Practical Foundations for Programming Languages SECOND EDITION

Robert Harper Carnegie Mellon University Copyright © 2016 by Robert Harper.

All Rights Reserved.

This is an abbreviated version of a book published by Cambridge University Press (http://www.cambridge.org). This draft is made available for the personal use of a single individual. The reader may make one copy for personal use. No unauthorized distribution of any kind is allowed. No alterations are permitted.

Preface to the Second Edition

Writing the second edition to a text book incurs the same risk as building the second version of a software system. It is difficult to make substantive improvements, while avoiding the temptation to overburden and undermine the foundation on which one is building. With the hope of avoiding the second system effect, I have sought to make corrections, revisions, expansions, and deletions that improve the coherence of the development, remove some topics that distract from the main themes, add new topics that were omitted from the first edition, and include exercises for almost every chapter.

The revision removes a number of typographical errors, corrects a few material errors (especially the formulation of the parallel abstract machine and of concurrency in Algol), and improves the writing throughout. Some chapters have been deleted (general pattern matching and polarization, restricted forms of polymorphism), some have been completely rewritten (the chapter on higher kinds), some have been substantially revised (general and parametric inductive definitions, concurrent and distributed Algol), several have been reorganized (to better distinguish partial from total type theories), and a new chapter has been added (on type refinements). Titular attributions on several chapters have been removed, not to diminish credit, but to avoid confusion between the present and the original formulations of several topics. A new system of (pronounceable!) language names has been introduced throughout. The exercises generally seek to expand on the ideas in the main text, and their solutions often involve significant technical ideas that merit study. Routine exercises of the kind one might include in a homework assignment are deliberately few.

My purpose in writing this book is to establish a comprehensive framework for formulating and analyzing a broad range of ideas in programming languages. If language design and programming methodology are to advance from a trade-craft to a rigorous discipline, it is essential that we first get the definitions right. Then, and only then, can there be meaningful analysis and consolidation of ideas. My hope is that I have helped to build such a foundation.

I am grateful to Stephanie Balzer, Stephen Brookes, Evan Cavallo, Karl Crary, Jon Sterling, and Todd Wilson for their help in critiquing drafts of this edition and for their suggestions for revision. I thank my department head, Frank Pfenning, for his support of my work on the completion of this edition. Thanks also to my editors, Ada Brunstein and Lauren Cowles, for their guidance and assistance. And thanks to Evan Cavallo and Andrew Shulaev for corrections to the draft.

Neither the author nor the publisher make any warranty, express or implied, that the definitions, theorems, and proofs contained in this volume are free of error, or are consistent with

any particular standard of merchantability, or that they will meet requirements for any particular application. They should not be relied on for solving a problem whose incorrect solution could result in injury to a person or loss of property. If you do use this material in such a manner, it is at your own risk. The author and publisher disclaim all liability for direct or consequential damage resulting from its use.

Pittsburgh July, 2015

Preface to the First Edition

Types are the central organizing principle of the theory of programming languages. Language features are manifestations of type structure. The syntax of a language is governed by the constructs that define its types, and its semantics is determined by the interactions among those constructs. The soundness of a language design—the absence of ill-defined programs—follows naturally.

The purpose of this book is to explain this remark. A variety of programming language features are analyzed in the unifying framework of type theory. A language feature is defined by its *statics*, the rules governing the use of the feature in a program, and its *dynamics*, the rules defining how programs using this feature are to be executed. The concept of *safety* emerges as the coherence of the statics and the dynamics of a language.

In this way we establish a foundation for the study of programming languages. But why these particular methods? The main justification is provided by the book itself. The methods we use are both *precise* and *intuitive*, providing a uniform framework for explaining programming language concepts. Importantly, these methods *scale* to a wide range of programming language concepts, supporting rigorous analysis of their properties. Although it would require another book in itself to justify this assertion, these methods are also *practical* in that they are *directly applicable* to implementation and *uniquely effective* as a basis for mechanized reasoning. No other framework offers as much.

Being a consolidation and distillation of decades of research, this book does not provide an exhaustive account of the history of the ideas that inform it. Suffice it to say that much of the development is not original, but rather is largely a reformulation of what has gone before. The notes at the end of each chapter signpost the major developments, but are not intended as a complete guide to the literature. For further information and alternative perspectives, the reader is referred to such excellent sources as Constable (1986), Constable (1998), Girard (1989), Martin-Löf (1984), Mitchell (1996), Pierce (2002, 2004), and Reynolds (1998).

The book is divided into parts that are, in the main, independent of one another. Parts I and II, however, provide the foundation for the rest of the book, and must therefore be considered prior to all other parts. On first reading it may be best to skim Part I, and begin in earnest with Part II, returning to Part I for clarification of the logical framework in which the rest of the book is cast.

Numerous people have read and commented on earlier editions of this book, and have suggested corrections and improvements to it. I am particularly grateful to Umut Acar, Jesper Louis Andersen, Carlo Angiuli, Andrew Appel, Stephanie Balzer, Eric Bergstrom, Guy E. Blelloch, Iliano Cervesato, Lin Chase, Karl Crary, Rowan Davies, Derek Dreyer, Dan Licata, Zhong Shao,

Rob Simmons, and Todd Wilson for their extensive efforts in reading and criticizing the book. I also thank the following people for their suggestions: Joseph Abrahamson, Arbob Ahmad, Zena Ariola, Eric Bergstrome, William Byrd, Alejandro Cabrera, Luis Caires, Luca Cardelli, Manuel Chakravarty, Richard C. Cobbe, James Cooper, Yi Dai, Daniel Dantas, Anupam Datta, Jake Donham, Bill Duff, Matthias Felleisen, Kathleen Fisher, Dan Friedman, Peter Gammie, Maia Ginsburg, Byron Hawkins, Kevin Hely, Kuen-Bang Hou (Favonia), Justin Hsu, Wojciech Jedynak, Cao Jing, Salil Joshi, Gabriele Keller, Scott Kilpatrick, Danielle Kramer, Dan Kreysa, Akiva Leffert, Ruy Ley-Wild, Karen Liu, Dave MacQueen, Chris Martens, Greg Morrisett, Stefan Muller, Tom Murphy, Aleksandar Nanevski, Georg Neis, David Neville, Adrian Trejo Nuñez, Cyrus Omar, Doug Perkins, Frank Pfenning, Jean Pichon, Benjamin Pierce, Andrew M. Pitts, Gordon Plotkin, David Renshaw, John Reynolds, Andreas Rossberg, Carter Schonwald, Dale Schumacher, Dana Scott, Shayak Sen, Pawel Sobocinski, Kristina Sojakova, Daniel Spoonhower, Paulo Tanimoto, Joe Tassarotti, Peter Thiemann, Bernardo Toninho, Michael Tschantz, Kami Vaniea, Carsten Varming, David Walker, Dan Wang, Jack Wileden, Sergei Winitzki, Roger Wolff, Omer Zach, Luke Zarko, and Yu Zhang. I am very grateful to the students of 15–312 and 15–814 at Carnegie Mellon who have provided the impetus for the preparation of this book and who have endured the many revisions to it over the last ten years.

I thank the Max Planck Institute for Software Systems for its hospitality and support. I also thank Espresso a Mano in Pittsburgh, CB2 Cafe in Cambridge, and Thonet Cafe in Saarbrücken for providing a steady supply of coffee and a conducive atmosphere for writing.

This material is, in part, based on work supported by the National Science Foundation under Grant Nos. 0702381 and 0716469. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

Robert Harper Pittsburgh March, 2012

Contents

Pr	eface	e to the Second Edition	iii
Pr	eface	e to the First Edition	v
Ι	Juc	dgments and Rules	1
1	Abs	stract Syntax	3
	1.1	Abstract Syntax Trees	
	1.2	Abstract Binding Trees	6
	1.3	Notes	10
2	Ind	uctive Definitions	13
-	2.1	Judgments	
	2.2	Inference Rules	
	2.3	Derivations	
	2.4	Rule Induction	
	2.5	Iterated and Simultaneous Inductive Definitions	18
	2.6		
	2.7		
3	Hyr	pothetical and General Judgments	23
	3.1	Hypothetical Judgments	
	0.1	3.1.1 Derivability	
		3.1.2 Admissibility	
	3.2	Hypothetical Inductive Definitions	
	3.3	General Judgments	
	3.4	Generic Inductive Definitions	
	3.5	Notes	

II	Statics and Dynamics	33
4	Statics 4.1 Syntax 4.2 Type System 4.3 Structural Properties 4.4 Notes	35 35 36 37 39
5	Dynamics5.1 Transition Systems5.2 Structural Dynamics5.3 Contextual Dynamics5.4 Equational Dynamics5.5 Notes	42 44 46
6	Type Safety 6.1 Preservation 6.2 Progress 6.3 Run-Time Errors 6.4 Notes	53
7	Evaluation Dynamics 7.1 Evaluation Dynamics	58 59 60
III	I Total Functions	63
8	Function Definitions and Values 8.1 First-Order Functions	67 69 70
9	System T of Higher-Order Recursion 9.1 Statics 9.2 Dynamics 9.3 Definability 9.4 Undefinability 9.5 Notes	73 73 74 76 77 79

CONTENTS	ix
----------	----

IV	Finite Data Types	81
10	Product Types	83
	10.1 Nullary and Binary Products	
	10.2 Finite Products	
	10.3 Primitive Mutual Recursion	
	10.4 Notes	
11	Sum Types	89
	11.1 Nullary and Binary Sums	. 89
	11.2 Finite Sums	. 91
	11.3 Applications of Sum Types	. 92
	11.3.1 Void and Unit	. 92
	11.3.2 Booleans	. 92
	11.3.3 Enumerations	
	11.3.4 Options	. 94
	11.4 Notes	. 95
V	Types and Propositions	97
12	Constructive Logic	99
	12.1 Constructive Semantics	
	12.2 Constructive Logic	
	12.2.1 Provability	
	12.2.2 Proof Terms	
	12.3 Proof Dynamics	
	12.4 Propositions as Types	
	12.5 Notes	. 105
12	Classical Logis	109
13	Classical Logic	
	13.1 Classical Logic	
	13.1.1 Provability and Refutability	
	13.1.2 Proofs and Refutations	
	13.3 Proof Dynamics	
	13.5 The Double-Negation Translation	
	13.6 Notes	. 119
VI	Infinite Data Types	121
• 1	Time Zuck Types	***
14	Generic Programming	123
	14.1 Introduction	. 123

	4.2 Polynomial Type Operators	126
	nductive and Coinductive Types 5.1 Motivating Examples 5.2 Statics 15.2.1 Types 15.2.2 Expressions 5.3 Dynamics 5.4 Solving Type Equations 5.5 Notes	132 132 133 134 135
VI	Variable Types	139
	ystem F of Polymorphic Types 6.1 Polymorphic Abstraction	145 145 146 147
	Abstract Types 7.1 Existential Types 17.1.1 Statics 17.1.2 Dynamics 17.1.3 Safety 7.2 Data Abstraction 7.3 Definability of Existential Types 7.4 Representation Independence 7.5 Notes	152 153 153 154 155
	Ligher Kinds 8.1 Constructors and Kinds 8.2 Constructor Equality 8.3 Expressions and Types 8.4 Notes	161 162
VI	Partiality and Recursive Types	165
	ystem PCF of Recursive Functions 9.1 Statics	167 169

CONTENTS	xi

	19.2 Dynamics 19.3 Definability 19.4 Finite and Infinite Data Structures 19.5 Totality and Partiality 19.6 Notes	171 173 173
	System FPC of Recursive Types 20.1 Solving Type Equations	179
IX	C Dynamic Types	185
21	The Untyped λ -Calculus 21.1 The λ -Calculus	188 190 192
22	Dynamic Typing22.1 Dynamically Typed PCF22.2 Variations and Extensions22.3 Critique of Dynamic Typing22.4 Notes	201
23	Hybrid Typing 23.1 A Hybrid Language	207208210
X	Subtyping	213
24	Structural Subtyping 24.1 Subsumption	216218

25	Behavioral Typing 25.1 Statics	. 234
XI	Dynamic Dispatch	239
26	Classes and Methods 26.1 The Dispatch Matrix 26.2 Class-Based Organization 26.3 Method-Based Organization 26.4 Self-Reference 26.5 Notes	244245246
27	Inheritance 27.1 Class and Method Extension	. 252. 254
ΧI	I Control Flow	257
	Control Flow Control Stacks 28.1 Machine Definition 28.2 Safety 28.3 Correctness of the K Machine 28.3.1 Completeness 28.3.2 Soundness 28.4 Notes	259 . 259 . 261 . 262 . 263 . 263
28	Control Stacks 28.1 Machine Definition	259 259 261 262 263 265 267 267 269

xiii
X

XI	III Symbolic Data	281
	Symbols 31.1 Symbol Declaration	284 285 286 286 287
	31.2.3 Safety	
	Fluid Binding 32.1 Statics 32.2 Dynamics 32.3 Type Safety 32.4 Some Subtleties 32.5 Fluid References 32.6 Notes	292293294295
	Dynamic Classification 33.1 Dynamic Classes 33.1.1 Statics 33.1.2 Dynamics 33.1.3 Safety 33.2 Class References 33.3 Definability of Dynamic Classes 33.4 Applications of Dynamic Classification 33.4.1 Classifying Secrets 33.4.2 Exception Values 33.5 Notes	299 300 301 301 302 303 303 304
XI	V Mutable State	307
	Modernized Algol 34.1 Basic Commands 34.1.1 Statics 34.1.2 Dynamics 34.1.3 Safety 34.2 Some Programming Idioms 34.3 Typed Commands and Typed Assignables 34.4 Notes	310 311 313 314 315

35.1 Capabilities 35.2 Scoped Assignables 35.3 Free Assignables 35.4 Safety 35.5 Benign Effects 35.6 Notes 36.1 PCF By-Need 36.2 Safety of PCF By-Need 36.3 FPC By-Need 36.4 Suspension Types 36.5 Notes	. 322 . 324 . 326 . 328 . 329 . 331 . 332 . 334 . 336
XV Parallelism	
XV Parallelism	341
37.1 Binary Fork-Join 37.2 Cost Dynamics 37.3 Multiple Fork-Join 37.4 Bounded Implementations 37.5 Scheduling 37.6 Notes 38 Futures and Speculations 38.1 Futures	. 346 . 349 . 351 . 354 . 356 . 360
38.1.1 Statics 38.1.2 Sequential Dynamics 38.2 Speculations 38.2.1 Statics 38.2.2 Sequential Dynamics 38.3 Parallel Dynamics 38.4 Pipelining With Futures 38.5 Notes	. 360. 361. 361. 362. 364
39 Process Calculus 39.1 Actions and Events 39.2 Interaction 39.3 Replication	. 371

CONTENTS	XV
CONTENTS	XV

	39.4	Allocating Channels	374
	39.5	Communication	376
	39.6	Channel Passing	379
	39.7	Universality	381
	39.8	Notes	382
40		current Algol	385
		Concurrent Algol	
		Broadcast Communication	
		Selective Communication	
		Free Assignables as Processes	
	40.5	Notes	394
41		ributed Algol	395
	41.1	Statics	395
		Dynamics	
		Safety	
	41.4	Notes	400
X	VII	Modularity	403
			403
-			403
		Iularity and Linking	405
	Mod 42.1	Iularity and Linking Simple Units and Linking	405 405
	Mod 42.1	lularity and Linking	405 405
	Mod 42.1 42.2	Iularity and Linking Simple Units and Linking	405 405 406
42	Mod 42.1 42.2 42.3	Simple Units and Linking	405 405 406 408
42	Mod 42.1 42.2 42.3 Sing	Iularity and Linking Simple Units and Linking Initialization and Effects Notes Ideton Kinds and Subkinding	405 405 406 408 409
42	Mod 42.1 42.2 42.3 Sing 43.1	Iularity and Linking Simple Units and Linking Initialization and Effects Notes Ideton Kinds and Subkinding Overview	405 405 406 408 409 410
42	Mod 42.1 42.2 42.3 Sing 43.1 43.2	Iularity and Linking Simple Units and Linking Initialization and Effects Notes Ideton Kinds and Subkinding Overview Singletons	405 406 408 409 410 410
42	Mod 42.1 42.2 42.3 Sing 43.1 43.2 43.3	Iularity and Linking Simple Units and Linking Initialization and Effects Notes Ileton Kinds and Subkinding Overview Singletons Dependent Kinds	405 406 408 409 410 410 412
42	Mod 42.1 42.2 42.3 Sing 43.1 43.2 43.3 43.4	Iularity and Linking Simple Units and Linking Initialization and Effects Notes Ideton Kinds and Subkinding Overview Singletons Dependent Kinds Higher Singletons	405 405 406 408 409 410 410 412 415
42	Mod 42.1 42.2 42.3 Sing 43.1 43.2 43.3 43.4	Iularity and Linking Simple Units and Linking Initialization and Effects Notes Ileton Kinds and Subkinding Overview Singletons Dependent Kinds	405 405 406 408 409 410 410 412 415
42	Mod 42.1 42.2 42.3 Sing 43.1 43.2 43.3 43.4 43.5	Iularity and Linking Simple Units and Linking Initialization and Effects Notes Ileton Kinds and Subkinding Overview Singletons Dependent Kinds Higher Singletons Notes	405 406 408 409 410 410 412 415 417
42	Mod 42.1 42.2 42.3 Sing 43.1 43.2 43.3 43.4 43.5	Iularity and Linking Simple Units and Linking Initialization and Effects Notes Ileton Kinds and Subkinding Overview Singletons Dependent Kinds Higher Singletons Notes Abstractions and Type Classes	405 406 408 409 410 410 412 415 417
42	Mod 42.1 42.2 42.3 Sing 43.1 43.2 43.3 43.4 43.5 Type 44.1	Simple Units and Linking Initialization and Effects Notes Cleton Kinds and Subkinding Overview Singletons Dependent Kinds Higher Singletons Notes Abstractions and Type Classes Type Abstraction	405 406 408 409 410 410 412 415 417 419 420
42	Mod 42.1 42.2 42.3 Sing 43.1 43.2 43.3 43.4 43.5 Type 44.1 44.2	Iularity and Linking Simple Units and Linking Initialization and Effects Notes Ileton Kinds and Subkinding Overview Singletons Dependent Kinds Higher Singletons Notes Abstractions and Type Classes Type Abstraction Type Classes	405 406 408 409 410 410 412 415 417 419 420 422
42	Mod 42.1 42.2 42.3 Sing 43.1 43.2 43.3 43.4 43.5 Type 44.1 44.2 44.3	Iularity and Linking Simple Units and Linking Initialization and Effects Notes Ileton Kinds and Subkinding Overview Singletons Dependent Kinds Higher Singletons Notes Abstractions and Type Classes Type Abstraction Type Classes A Module Language	405 406 408 409 410 410 412 415 417 419 420 422 424
42	Mod 42.1 42.2 42.3 Sing 43.1 43.2 43.3 43.4 43.5 Type 44.1 44.2 44.3 44.4	Iularity and Linking Simple Units and Linking Initialization and Effects Notes Ileton Kinds and Subkinding Overview Singletons Dependent Kinds Higher Singletons Notes Abstractions and Type Classes Type Abstraction Type Classes	405 406 408 409 410 410 412 415 417 419 420 422 424 428

45	Hierarchy and Parameterization	431
	45.1 Hierarchy	
	45.2 Abstraction	434
	45.3 Hierarchy and Abstraction	
	45.4 Applicative Functors	438
	45.5 Notes	
		7
XV	III Equational Reasoning	141
46	Equality for System T	443
	46.1 Observational Equivalence	
	46.2 Logical Equivalence	446
	46.3 Logical and Observational Equivalence Coincide	448
	46.4 Some Laws of Equality	
	46.4.1 General Laws	
	46.4.2 Equality Laws	
	46.4.3 Induction Law	451
	46.5 Notes	452
47	Equality for System PCF	453
	47.1 Observational Equivalence	453
	47.2 Logical Equivalence	
	47.3 Logical and Observational Equivalence Coincide	
	47.4 Compactness	457
	47.5 Lazy Natural Numbers	460
	47.6 Notes	461
48	Parametricity	463
	48.1 Overview	463
	48.2 Observational Equivalence	464
	48.3 Logical Equivalence	465
	48.4 Parametricity Properties	47 0
	48.5 Representation Independence, Revisited	473
	48.6 Notes	474
49	Process Equivalence	475
	•	475
	49.2 Strong Equivalence	477
	49.3 Weak Equivalence	
	49.4 Notes	481

	CONTENTS	xvii
	XIX Appendices	483
	A Answers to the Exercises	485
	B Background on Finite Sets	539
-		



Part I Judgments and Rules



Chapter 1

Abstract Syntax

Programming languages express computations in a form comprehensible to both people and machines. The syntax of a language specifies how various sorts of phrases (expressions, commands, declarations, and so forth) may be combined to form programs. But what are these phrases? What is a program made of?

The informal concept of syntax involves several distinct concepts. The *surface*, or *concrete*, *syntax* is concerned with how phrases are entered and displayed on a computer. The surface syntax is usually thought of as given by strings of characters from some alphabet (say, ASCII or Unicode). The *structural*, or *abstract*, *syntax* is concerned with the structure of phrases, specifically how they are composed from other phrases. At this level a phrase is a tree, called an *abstract syntax tree*, whose nodes are operators that combine several phrases to form another phrase. The *binding* structure of syntax is concerned with the introduction and use of identifiers: how they are declared, and how declared identifiers can be used. At this level phrases are *abstract binding trees*, which enrich abstract syntax trees with the concepts of binding and scope.

We will not concern ourselves in this book with concrete syntax, but will instead consider pieces of syntax to be finite trees augmented with a means of expressing the binding and scope of identifiers within a syntax tree. To prepare the ground for the rest of the book, we define in this chapter what is a "piece of syntax" in two stages. First, we define abstract syntax trees, or asts, which capture the hierarchical structure of a piece of syntax, while avoiding commitment to their concrete representation as a string. Second, we augment abstract syntax trees with the means of specifying the binding (declaration) and scope (range of significance) of an identifier. Such enriched forms of abstract syntax are called abstract binding trees, or abts for short.

Several functions and relations on abts are defined that give precise meaning to the informal ideas of binding and scope of identifiers. The concepts are infamously difficult to define properly, and are the mother lode of bugs for language implementors. Consequently, precise definitions are essential, but they are also fairly technical and take some getting used to. It is probably best to skim this chapter on first reading to get the main ideas, and return to it for clarification as necessary.

1.1 Abstract Syntax Trees

An abstract syntax tree, or ast for short, is an ordered tree whose leaves are variables, and whose interior nodes are operators whose arguments are its children. Asts are classified into a variety of sorts corresponding to different forms of syntax. A variable stands for an unspecified, or generic, piece of syntax of a specified sort. Asts can be combined by an operator, which has an arity specifying the sort of the operator and the number and sorts of its arguments. An operator of sort s and arity s_1, \ldots, s_n combines $n \ge 0$ asts of sort s_1, \ldots, s_n , respectively, into a compound ast of sort s.

The concept of a variable is central, and therefore deserves special emphasis. A variable is an *unknown* object drawn from some domain. The unknown can become known by *substitution* of a particular object for all occurrences of a variable in a formula, thereby specializing a general formula to a particular instance. For example, in school algebra variables range over real numbers, and we may form polynomials, such as $x^2 + 2x + 1$, that can be specialized by substitution of, say, 7 for x to obtain $7^2 + (2 \times 7) + 1$, which can be simplified according to the laws of arithmetic to obtain 64, which is $(7 + 1)^2$.

Abstract syntax trees are classified by *sorts* that divide asts into syntactic categories. For example, familiar programming languages often have a syntactic distinction between expressions and commands; these are two sorts of abstract syntax trees. Variables in abstract syntax trees range over sorts in the sense that only asts of the specified sort of the variable can be plugged in for that variable. Thus it would make no sense to replace an expression variable by a command, nor a command variable by an expression, the two being different sorts of things. But the core idea carries over from school mathematics, namely that *a variable is an unknown, or a place-holder, whose meaning is given by substitution*.

As an example, consider a language of arithmetic expressions built from numbers, addition, and multiplication. The abstract syntax of such a language consists of a single sort Exp generated by these operators:

- 1. An operator num[n] of sort Exp for each $n \in \mathbb{N}$;
- 2. Two operators, plus and times, of sort Exp, each with two arguments of sort Exp.

The expression $2 + (3 \times x)$, which involves a variable, x, would be represented by the ast

of sort Exp, under the assumption that x is also of this sort. Because, say, num[4], is an ast of sort Exp, we may plug it in for x in the above ast to obtain the ast

which is written informally as $2 + (3 \times 4)$. We may, of course, plug in more complex asts of sort Exp for x to obtain other asts as result.

The tree structure of asts provides a very useful principle of reasoning, called *structural induction*. Suppose that we wish to prove that some property $\mathcal{P}(a)$ holds of all asts a of a given sort. To show this it is enough to consider all the ways in which a can be generated, and show that the property holds in each case under the assumption that it holds for its constituent asts (if any). So, in the case of the sort Exp just described, we must show

- 1. The property holds for any variable x of sort Exp: prove that $\mathcal{P}(x)$.
- 2. The property holds for any number, num[n]: for every $n \in \mathbb{N}$, prove that $\mathcal{P}(num[n])$.
- 3. Assuming that the property holds for a_1 and a_2 , prove that it holds for plus($a_1; a_2$) and times($a_1; a_2$): if $\mathcal{P}(a_1)$ and $\mathcal{P}(a_2)$, then $\mathcal{P}(\text{plus}(a_1; a_2))$ and $\mathcal{P}(\text{times}(a_1; a_2))$.

Because these cases exhaust all possibilities for the formation of a, we are assured that $\mathcal{P}(a)$ holds for any ast a of sort Exp.

It is common to apply the principle of structural induction in a form that takes account of the interpretation of variables as place-holders for asts of the appropriate sort. Informally, it is often useful to prove a property of an ast involving variables in a form that is conditional on the property holding for the variables. Doing so anticipates that the variables will be replaced with asts that ought to have the property assumed for them, so that the result of the replacement will have the property as well. This amounts to applying the principle of structural induction to properties $\mathcal{P}(a)$ of the form "if a involves variables x_1, \ldots, x_k , and \mathcal{Q} holds of each x_i , then \mathcal{Q} holds of a", so that a proof of $\mathcal{P}(a)$ for all asts a by structural induction is just a proof that $\mathcal{Q}(a)$ holds for all asts a under the assumption that \mathcal{Q} holds for its variables. When there are no variables, there are no assumptions, and the proof of \mathcal{P} is a proof that \mathcal{Q} holds for all closed asts. On the other hand if x is a variable in a, and we replace it by an ast b for which \mathcal{Q} holds, then \mathcal{Q} will hold for the result of replacing x by b in a.

For the sake of precision, we now give precise definitions of these concepts. Let S be a finite set of sorts. For a given set S of sorts, an *arity* has the form $(s_1, \ldots, s_n)s$, which specifies the sort $s \in S$ of an operator taking $n \geq 0$ arguments, each of sort $s_i \in S$. Let $\mathcal{O} = \{\mathcal{O}_\alpha\}$ be an arity-indexed family of disjoint sets of *operators* \mathcal{O}_α of arity α . If σ is an operator of arity $(s_1, \ldots, s_n)s$, we say that σ has sort s and has s arguments of sorts s_1, \ldots, s_n .

Fix a set \mathcal{S} of sorts and an arity-indexed family \mathcal{O} of sets of operators of each arity. Let $\mathcal{X} = \{\mathcal{X}_s\}_{s \in \mathcal{S}}$ be a sort-indexed family of disjoint finite sets \mathcal{X}_s of variables x of sort s. When \mathcal{X} is clear from context, we say that a variable x is of sort s if $x \in \mathcal{X}_s$, and we say that x is fresh for \mathcal{X} , or just fresh when \mathcal{X} is understood, if $x \notin \mathcal{X}_s$ for any sort s. If x is fresh for \mathcal{X} and s is a sort, then \mathcal{X} , x is the family of sets of variables obtained by adding x to \mathcal{X}_s . The notation is ambiguous in that the sort s is not explicitly stated, but determined from context.

The family $A[\mathcal{X}] = \{A[\mathcal{X}]_s\}_{s \in \mathcal{S}}$ of abstract syntax trees, or asts, of sort s is the smallest family satisfying the following conditions:

- 1. A variable of sort s is an ast of sort s: if $x \in \mathcal{X}_s$, then $x \in \mathcal{A}[\mathcal{X}]_s$.
- 2. Operators combine asts: if o is an operator of arity $(s_1, ..., s_n)s$, and if $a_1 \in \mathcal{A}[\mathcal{X}]_{s_1}, ..., a_n \in \mathcal{A}[\mathcal{X}]_{s_n}$, then $o(a_1; ...; a_n) \in \mathcal{A}[\mathcal{X}]_s$.

It follows from this definition that the principle of *structural induction* can be used to prove that some property \mathcal{P} holds of every ast. To show $\mathcal{P}(a)$ holds for every $a \in \mathcal{A}[\mathcal{X}]$, it is enough to show:

- 1. If $x \in \mathcal{X}_s$, then $\mathcal{P}_s(x)$.
- 2. If o has arity $(s_1, \ldots, s_n)s$ and $\mathcal{P}_{s_1}(a_1)$ and \ldots and $\mathcal{P}_{s_n}(a_n)$, then $\mathcal{P}_{s}(o(a_1; \ldots; a_n))$.

For example, it is easy to prove by structural induction that $\mathcal{A}[\mathcal{X}] \subseteq \mathcal{A}[\mathcal{Y}]$ whenever $\mathcal{X} \subseteq \mathcal{Y}$.

Variables are given meaning by *substitution*. If $a \in \mathcal{A}[\mathcal{X}, x]_{s'}$, and $b \in \mathcal{A}[\mathcal{X}]_s$, then $[b/x]a \in \mathcal{A}[\mathcal{X}]_{s'}$ is the result of substituting b for every occurrence of x in a. The ast a is called the *target*, and x is called the *subject*, of the substitution. Substitution is defined by the following equations:

- 1. [b/x]x = b and [b/x]y = y if $x \neq y$.
- 2. $[b/x]o(a_1;...;a_n) = o([b/x]a_1;...;[b/x]a_n)$.

For example, we may check that

$$[\text{num}[2]/x]\text{plus}(x;\text{num}[3]) = \text{plus}(\text{num}[2];\text{num}[3]).$$

We may prove by structural induction that substitution on asts is well-defined.

Theorem 1.1. If $a \in A[X, x]$, then for every $b \in A[X]$ there exists a unique $c \in A[X]$ such that [b/x]a = c

Proof. By structural induction on a. If a = x, then c = b by definition, otherwise if $a = y \neq x$, then c = y, also by definition. Otherwise, $a = o(a_1; ...; a_n)$, and we have by induction unique $c_1, ..., c_n$ such that $[b/x]a_1 = c_1$ and ... $[b/x]a_n = c_n$, and so c is $c = o(c_1; ...; c_n)$, by definition of substitution.

1.2 Abstract Binding Trees

Abstract binding trees, or abts, enrich asts with the means to introduce new variables and symbols, called a binding, with a specified range of significance, called its scope. The scope of a binding is an abt within which the bound identifier can be used, either as a place-holder (in the case of a variable declaration) or as the index of some operator (in the case of a symbol declaration). Thus the set of active identifiers can be larger within a subtree of an abt than it is within the surrounding tree. Moreover, different subtrees may introduce identifiers with disjoint scopes. The crucial principle is that any use of an identifier should be understood as a reference, or abstract pointer, to its binding. One consequence is that the choice of identifiers is immaterial, so long as we can always associate a unique binding with each use of an identifier.

As a motivating example, consider the expression $let x be a_1 in a_2$, which introduces a variable x for use within the expression a_2 to stand for the expression a_1 . The variable x is bound by the let expression for use within a_2 ; any use of x within a_1 refers to a different variable that happens to have the same name. For example, in the expression let x be 7 in x + x occurrences of x in the addition refer to the variable introduced by the let. On the other hand in the expression let x be x * x in x + x, occurrences of x within the multiplication refer to a different variable than those occurring within the addition. The latter occurrences refer to the binding introduced by the let, whereas the former refer to some outer binding not displayed here.

The names of bound variables are immaterial insofar as they determine the same binding. So, for example, let x be x * x in x + x could just as well have been written let y be x * x in y + y, without changing its meaning. In the former case the variable x is bound within the addition, and

in the latter it is the variable y, but the "pointer structure" remains the same. On the other hand the expression let x be y * y in x + x has a different meaning to these two expressions, because now the variable y within the multiplication refers to a different surrounding variable. Renaming of bound variables is constrained to the extent that it must not alter the reference structure of the expression. For example, the expression

let
$$x$$
 be 2 in let y be 3 in $x + x$

has a different meaning than the expression

let
$$y$$
 be 2 in let y be 3 in $y + y$,

because the y in the expression y + y in the second case refers to the inner declaration, not the outer one as before.

The concept of an ast can be enriched to account for binding and scope of a variable. These enriched asts are called *abstract binding trees*, or *abts* for short. Abts generalize asts by allowing an operator to bind any finite number (possibly zero) of variables in each argument. An argument to an operator is called an *abstractor*, and has the form $x_1, \ldots, x_k.a$. The sequence of variables x_1, \ldots, x_k are bound within the abt a. (When k is zero, we elide the distinction between a and a itself.) Written in the form of an abt, the expression let x be a_1 in a_2 has the form let $(a_1; x.a_2)$, which more clearly specifies that the variable x is bound within a_2 , and not within a_1 . We often write \vec{x} to stand for a finite sequence x_1, \ldots, x_n of distinct variables, and write $\vec{x}.a$ to mean $x_1, \ldots, x_n.a$.

To account for binding, operators are assigned *generalized arities* of the form $(v_1, \ldots, v_n)s$, which specifies operators of sort s with n arguments of *valence* v_1, \ldots, v_n . In general a valence v has the form $s_1, \ldots, s_k.s$, which specifies the sort of an argument as well as the number and sorts of the variables bound within it. We say that a sequence \vec{x} of variables is of sort \vec{s} to mean that the two sequences have the same length k and that the variable x_i is of sort s_i for each $1 \le i \le k$.

Thus, to specify that the operator let has arity (Exp, Exp.Exp)Exp indicates that it is of sort Exp whose first argument is of sort Exp and binds no variables, and whose second argument is also of sort Exp, within which is bound one variable of sort Exp. The informal expression let x be 2+2 in $x\times x$ may then be written as the abt

in which the operator let has two arguments, the first of which is an expression, and the second of which is an abstractor that binds one expression variable.

Fix a set \mathcal{S} of sorts, and a family \mathcal{O} of disjoint sets of operators indexed by their generalized arities. For a given family of disjoint sets of variables \mathcal{X} , the family of *abstract binding trees*, or *abts* $\mathcal{B}[\mathcal{X}]$ is defined similarly to $\mathcal{A}[\mathcal{X}]$, except that \mathcal{X} is not fixed throughout the definition, but rather changes as we enter the scopes of abstractors.

This simple idea is surprisingly hard to make precise. A first attempt at the definition is as the least family of sets closed under the following conditions:

```
1. If x \in \mathcal{X}_s, then x \in \mathcal{B}[\mathcal{X}]_s.
```

2. For each operator o of arity $(\vec{s}_1.s_1,...,\vec{s}_n.s_n)s$, if $a_1 \in \mathcal{B}[\mathcal{X},\vec{x}_1]_{s_1},...$, and $a_n \in \mathcal{B}[\mathcal{X},\vec{x}_n]_{s_n}$, then $o(\vec{x}_1.a_1;...;\vec{x}_n.a_n) \in \mathcal{B}[\mathcal{X}]_s$.

The bound variables are adjoined to the set of active variables within each argument, with the sort of each variable determined by the valence of the operator.

This definition is *almost* correct, but fails to properly account for renaming of bound variables. An abt of the form $let(a_1; x.let(a_2; x.a_3))$ is ill-formed according to this definition, because the first binding adds x to \mathcal{X} , which implies that the second cannot also add x to \mathcal{X} , x, because it is not fresh for \mathcal{X} , x. The solution is to ensure that each of the arguments is well-formed regardless of the choice of bound variable names, which is achieved using *fresh renamings*, which are bijections between sequences of variables. Specifically, a fresh renaming (relative to \mathcal{X}) of a finite sequence of variables \vec{x} is a bijection $\rho: \vec{x} \leftrightarrow \vec{x}'$ between \vec{x} and \vec{x}' , where \vec{x}' is fresh for \mathcal{X} . We write $\widehat{\rho}(a)$ for the result of replacing each occurrence of x_i in a by $\rho(x_i)$, its fresh counterpart.

This is achieved by altering the second clause of the definition of abts using fresh renamings as follows:

```
For each operator o of arity (\vec{s}_1.s_1,...,\vec{s}_n.s_n)s, if for each 1 \le i \le n and each fresh renaming \rho_i : \vec{x}_i \leftrightarrow \vec{x}_i', we have \hat{\rho}_i(a_i) \in \mathcal{B}[\mathcal{X}, \vec{x}_i'], then o(\vec{x}_1.a_1;...;\vec{x}_n.a_n) \in \mathcal{B}[\mathcal{X}]_s.
```

The renaming, $\hat{\rho}_i(a_i)$, of each a_i ensures that collisions cannot occur, and that the abt is valid for almost all renamings of any bound variables that occur within it.

The principle of structural induction extends to abts, and is called *structural induction modulo fresh renaming*. It states that to show that $\mathcal{P}[\mathcal{X}](a)$ holds for every $a \in \mathcal{B}[\mathcal{X}]$, it is enough to show the following:

- 1. if $x \in \mathcal{X}_s$, then $\mathcal{P}[\mathcal{X}]_s(x)$.
- 2. For every o of arity $(\vec{s}_1.s_1, \ldots, \vec{s}_n.s_n)s$, if for each $1 \leq i \leq n$, $\mathcal{P}[\mathcal{X}, \vec{x}_i']_{s_i}(\widehat{\rho_i}(a_i))$ holds for every $\rho_i : \vec{x}_i \leftrightarrow \vec{x}_i'$ with $\vec{x}_i' \notin \mathcal{X}$, then $\mathcal{P}[\mathcal{X}]_s(o(\vec{x}_1.a_1; \ldots; \vec{x}_n.a_n))$.

The second condition ensures that the inductive hypothesis holds for *all* fresh choices of bound variable names, and not just the ones actually given in the abt.

As an example let us define the judgment $x \in a$, where $a \in \mathcal{B}[\mathcal{X}, x]$, to mean that x occurs free in a. Informally, this means that x is bound somewhere outside of a, rather than within a itself. If x is bound within a, then those occurrences of x are different from those occurring outside the binding. The following definition ensures that this is the case:

- 1. $x \in x$.
- 2. $x \in o(\vec{x}_1.a_1; ...; \vec{x}_n.a_n)$ if there exists $1 \le i \le n$ such that for every fresh renaming $\rho : \vec{x}_i \leftrightarrow \vec{z}_i$ we have $x \in \hat{\rho}(a_i)$.

The first condition states that x is free in x, but not free in y for any variable y other than x. The second condition states that if x is free in some argument, independently of the choice of bound variable names in that argument, then it is free in the overall abt.

The relation $a =_{\alpha} b$ of α -equivalence (so-called for historical reasons), means that a and b are identical up to the choice of bound variable names. The α -equivalence relation is the strongest congruence containing the following two conditions:

- 1. $x =_{\alpha} x$.
- 2. $o(\vec{x}_1.a_1; \dots; \vec{x}_n.a_n) =_{\alpha} o(\vec{x}_1'.a_1'; \dots; \vec{x}_n'.a_n')$ if for every $1 \leq i \leq n$, $\widehat{\rho_i}(a_i) =_{\alpha} \widehat{\rho_i'}(a_i')$ for all fresh renamings $\rho_i : \vec{x}_i \leftrightarrow \vec{z}_i$ and $\rho_i' : \vec{x}_i' \leftrightarrow \vec{z}_i$.

The idea is that we rename \vec{x}_i and \vec{x}_i' consistently, avoiding confusion, and check that a_i and a_i' are α -equivalent. If $a =_{\alpha} b$, then a and b are α -variants of each other.

Some care is required in the definition of *substitution* of an abt b of sort s for free occurrences of a variable x of sort s in some abt a of some sort, written $\lfloor b/x \rfloor a$. Substitution is partially defined by the following conditions:

- 1. [b/x]x = b, and [b/x]y = y if $x \neq y$.
- 2. $[b/x]o(\vec{x}_1.a_1;...;\vec{x}_n.a_n) = o(\vec{x}_1.a'_1;...;\vec{x}_n.a'_n)$, where, for each $1 \le i \le n$, we require that $\vec{x}_i \notin b$, and we set $a'_i = [b/x]a_i$ if $x \notin \vec{x}_i$, and $a'_i = a_i$ otherwise.

The definition of [b/x]a is quite delicate, and merits careful consideration.

One trouble spot for substitution is to notice that if x is bound by an abstractor within a, then x does not occur free within the abstractor, and hence is unchanged by substitution. For example, $[b/x] \text{let}(a_1; x.a_2) = \text{let}([b/x]a_1; x.a_2)$, there being no free occurrences of x in $x.a_2$. Another trouble spot is the *capture* of a free variable of b during substitution. For example, if $y \in b$, and $x \neq y$, then $[b/x] \text{let}(a_1; y.a_2)$ is undefined, rather than being $\text{let}([b/x]a_1; y.[b/x]a_2)$, as one might at first suspect. For example, provided that $x \neq y$, [y/x] let(num[0]; y.plus(x; y)) is undefined, not let(num[0]; y.plus(y; y)), which confuses two different variables named y.

Although capture avoidance is an essential characteristic of substitution, it is, in a sense, merely a technical nuisance. If the names of bound variables have no significance, then capture can always be avoided by first renaming the bound variables in a to avoid any free variables in b. In the foregoing example if we rename the bound variable y to y' to obtain $a' \triangleq \mathtt{let}(\mathtt{num}[0]; y'.\mathtt{plus}(x; y'))$, then [b/x]a' is defined, and is equal to $\mathtt{let}(\mathtt{num}[0]; y'.\mathtt{plus}(b; y'))$. The price for avoiding capture in this way is that substitution is only determined up to α -equivalence, and so we may no longer think of substitution as a function, but only as a proper relation.

To restore the functional character of substitution, it is sufficient to adopt the *identification convention*, which is stated as follows:

Abstract binding trees are always identified up to α -equivalence.

That is, α -equivalent abts are regarded as identical. Substitution can be extended to α -equivalence classes of abts to avoid capture by choosing representatives of the equivalence classes of b and a in such a way that substitution is defined, then forming the equivalence class of the result. Any two choices of representatives for which substitution is defined gives α -equivalent results, so that substitution becomes a well-defined total function. We will adopt the identification convention for abts throughout this book.

It will often be necessary to consider languages whose abstract syntax cannot be specified by a fixed set of operators, but rather requires that the available operators be sensitive to the context in which they occur. For our purposes it will suffice to consider a set of *symbolic parameters*, or *symbols*, that index families of operators so that as the set of symbols varies, so does the set of operators. An

10 1.3 Notes

indexed operator o is a family of operators indexed by symbols u, so that o[u] is an operator when u is an available symbol. If \mathcal{U} is a finite set of symbols, then $\mathcal{B}[\mathcal{U};\mathcal{X}]$ is the sort-indexed family of abts that are generated by operators and variables as before, admitting all indexed operator instances by symbols $u \in \mathcal{U}$. Whereas a variable is a place-holder that stands for an unknown abt of its sort, a symbol *does not stand for anything*, and is not, itself, an abt. The only significance of symbol is whether it is the same as or differs from another symbol; the operator instances o[u] and o[u'] are the same exactly when u is u', and are the same symbol.

The set of symbols is extended by introducing a *new*, or *fresh*, symbol within a scope using the abstractor *u.a*, which binds the symbol *u* within the abt *a*. An abstracted symbol is "new" in the same sense as for an abstracted variable: the name of the bound symbol can be varied at will provided that no conflicts arise. This renaming property ensures that an abstracted symbol is distinct from all others in scope. The only difference between symbols and variables is that the only operation on symbols is renaming; there is no notion of substitution for a symbol.

Finally, a word about notation: to help improve the readability we often "group" and "stage" the arguments to an operator, using round brackets and braces to show grouping, and generally regarding stages to progress from right to left. All arguments in a group are considered to occur at the same stage, though their order is significant, and successive groups are considered to occur in sequential stages. Staging and grouping is often a helpful mnemonic device, but has no fundamental significance. For example, the abt $o\{a_1; a_2\}(a_3; x.a_4)$ is the same as the abt $o\{a_1; a_2; a_3; x.a_4\}$, as would be any other order-preserving grouping or staging of its arguments.

1.3 Notes

The concept of abstract syntax has its origins in the pioneering work of Church, Turing, and Gödel, who first considered writing programs that act on representations of programs. Originally programs were represented by natural numbers, using encodings, now called *Gödel-numberings*, based on the prime factorization theorem. Any standard text on mathematical logic, such as Kleene (1952), has a thorough account of such representations. The Lisp language (McCarthy, 1965; Allen, 1978) introduced a much more practical and direct representation of syntax as *symbolic expressions*. These ideas were developed further in the language ML (Gordon et al., 1979), which featured a type system capable of expressing abstract syntax trees. The AUTOMATH project (Nederpelt et al., 1994) introduced the idea of using Church's λ notation (Church, 1941) to account for the binding and scope of variables. These ideas were developed further in LF (Harper et al., 1993).

The concept of abstract binding trees presented here was inspired by the system of notation developed in the NuPRL Project, which is described in Constable (1986) and from Martin-Löf's system of arities, which is described in Nordstrom et al. (1990). Their enrichment with symbol binders is influenced by Pitts and Stark (1993).

Exercises

1.1. Prove by structural induction on abstract syntax trees that if $\mathcal{X} \subseteq \mathcal{Y}$, then $\mathcal{A}[\mathcal{X}] \subseteq \mathcal{A}[\mathcal{Y}]$.

1.3 Notes 11

1.2. Prove by structural induction modulo renaming on abstract binding trees that if $\mathcal{X} \subseteq \mathcal{Y}$, then $\mathcal{B}[\mathcal{X}] \subseteq \mathcal{B}[\mathcal{Y}]$.

- **1.3.** Show that if $a =_{\alpha} a'$ and $b =_{\alpha} b'$ and both [b/x]a and [b'/x]a' are defined, then $[b/x]a =_{\alpha} [b'/x]a'$.
- **1.4.** Bound variables can be seen as the formal analogs of pronouns in natural languages. The binding occurrence of a variable at an abstractor fixes a "fresh" pronoun for use within its body that refers unambiguously to that variable (in contrast to English, in which the referent of a pronoun can often be ambiguous). This observation suggests an alternative representation of abts, called *abstract binding graphs*, or *abg's* for short, as directed graphs constructed as follows:
 - (a) Free variables are atomic nodes with no outgoing edges.
 - (b) Operators with *n* arguments are *n*-ary nodes, with one outgoing edge directed at each of their children.
 - (c) Abstractors are nodes with one edge directed to the scope of the abstracted variable.
 - (d) Bound variables are back edges directed at the abstractor that introduced it.

Notice that asts, thought of as abts with no abstractors, are *acyclic* directed graphs (more precisely, variadic trees), whereas general abts can be *cyclic*. Draw a few examples of abg's corresponding to the example abts given in this chapter. Give a precise definition of the sort-indexed family $\mathcal{G}[\mathcal{X}]$ of abstract binding graphs. What representation would you use for bound variables (back edges)?

1.3 Notes



Chapter 2

Inductive Definitions

Inductive definitions are an indispensable tool in the study of programming languages. In this chapter we will develop the basic framework of inductive definitions, and give some examples of their use. An inductive definition consists of a set of *rules* for deriving *judgments*, or *assertions*, of a variety of forms. Judgments are statements about one or more abstract binding trees of some sort. The rules specify necessary and sufficient conditions for the validity of a judgment, and hence fully determine its meaning.

2.1 Judgments

We start with the notion of a *judgment*, or *assertion*, about an abstract binding tree. We shall make use of many forms of judgment, including examples such as these:

```
n nat n is a natural number n_1 + n_2 = n n is the sum of n_1 and n_2 \tau type t is a type t expression t has type t expression t has value t
```

A judgment states that one or more abstract binding trees have a property or stand in some relation to one another. The property or relation itself is called a *judgment form*, and the judgment that an object or objects have that property or stand in that relation is said to be an *instance* of that judgment form. A judgment form is also called a *predicate*, and the objects constituting an instance are its *subjects*. We write a J or J a, for the judgment asserting that J holds of the abt a. Correspondingly, we sometimes notate the judgment form J by - J, or J -, using a dash to indicate the absence of an argument to J. When it is not important to stress the subject of the judgment, we write J to stand for an unspecified judgment, that is, an instance of some judgment form. For particular judgment forms, we freely use prefix, infix, or mix-fix notation, as illustrated by the above examples, in order to enhance readability.

14 2.2 Inference Rules

2.2 Inference Rules

An inductive definition of a judgment form consists of a collection of rules of the form

$$\frac{J_1 \quad \cdots \quad J_k}{J} \tag{2.1}$$

in which J and J_1, \ldots, J_k are all judgments of the form being defined. The judgments above the horizontal line are called the *premises* of the rule, and the judgment below the line is called its *conclusion*. If a rule has no premises (that is, when k is zero), the rule is called an *axiom*; otherwise it is called a *proper rule*.

An inference rule can be read as stating that the premises are *sufficient* for the conclusion: to show J, it is enough to show J_1, \ldots, J_k . When k is zero, a rule states that its conclusion holds unconditionally. Bear in mind that there may be, in general, many rules with the same conclusion, each specifying sufficient conditions for the conclusion. Consequently, if the conclusion of a rule holds, then it is not necessary that the premises hold, for it might have been derived by another rule.

For example, the following rules form an inductive definition of the judgment form - nat:

$$(2.2a)$$

$$\frac{a \text{ nat}}{\text{succ}(a) \text{ nat}}$$
 (2.2b)

These rules specify that a nat holds whenever either a is zero, or a is succ(b) where b nat for some b. Taking these rules to be exhaustive, it follows that a nat iff a is a natural number.

Similarly, the following rules constitute an inductive definition of the judgment form - tree:

$$\frac{a_1 \text{ tree} \quad a_2 \text{ tree}}{\text{node}(a_1; a_2) \text{ tree}}$$
 (2.3b)

These rules specify that a tree holds if either a is empty, or a is node($a_1;a_2$), where a_1 tree and a_2 tree. Taking these to be exhaustive, these rules state that a is a binary tree, which is to say it is either empty, or a node consisting of two children, each of which is also a binary tree.

The judgment form a is b expressing the equality of two abts a and b such that a nat and b nat is inductively defined by the following rules:

$$\frac{a \text{ is } b}{\text{succ}(a) \text{ is succ}(b)}$$
 (2.4b)

2.3 Derivations 15

In each of the preceding examples we have made use of a notational convention for specifying an infinite family of rules by a finite number of patterns, or *rule schemes*. For example, rule (2.2b) is a rule scheme that determines one rule, called an *instance* of the rule scheme, for each choice of object *a* in the rule. We will rely on context to determine whether a rule is stated for a *specific* object *a* or is instead intended as a rule scheme specifying a rule for *each choice* of objects in the rule.

A collection of rules is considered to define the *strongest* judgment form that is *closed under*, or *respects*, those rules. To be closed under the rules simply means that the rules are *sufficient* to show the validity of a judgment: *J* holds *if* there is a way to obtain it using the given rules. To be the *strongest* judgment form closed under the rules means that the rules are also *necessary*: *J* holds *only if* there is a way to obtain it by applying the rules. The sufficiency of the rules means that we may show that *J* holds by *deriving* it by composing rules. Their necessity means that we may reason about it using *rule induction*.

2.3 Derivations

To show that an inductively defined judgment holds, it is enough to exhibit a *derivation* of it. A derivation of a judgment is a finite composition of rules, starting with axioms and ending with that judgment. It can be thought of as a tree in which each node is a rule whose children are derivations of its premises. We sometimes say that a derivation of *J* is evidence for the validity of an inductively defined judgment *J*.

We usually depict derivations as trees with the conclusion at the bottom, and with the children of a node corresponding to a rule appearing above it as evidence for the premises of that rule. Thus, if

$$J_1 \quad \cdots \quad J_k$$

is an inference rule and $\nabla_1, \dots, \nabla_k$ are derivations of its premises, then

$$\frac{\nabla_1 \cdots \nabla_k}{I}$$

is a derivation of its conclusion. In particular, if k = 0, then the node has no children. For example, this is a derivation of succ(succ(succ(zero))) nat:

$$\frac{\overline{\text{zero nat}}}{\text{succ}(\text{zero}) \text{ nat}} \\
\underline{\frac{\text{succ}(\text{succ}(\text{zero})) \text{ nat}}{\text{succ}(\text{succ}(\text{zero}))) \text{ nat}}}.$$

16 2.4 Rule Induction

Similarly, here is a derivation of node (node (empty; empty); empty) tree:

To show that an inductively defined judgment is derivable we need only find a derivation for it. There are two main methods for finding derivations, called *forward chaining*, or *bottom-up construction*, and *backward chaining*, or *top-down construction*. Forward chaining starts with the axioms and works forward towards the desired conclusion, whereas backward chaining starts with the desired conclusion and works backwards towards the axioms.

More precisely, forward chaining search maintains a set of derivable judgments, and continually extends this set by adding to it the conclusion of any rule all of whose premises are in that set. Initially, the set is empty; the process terminates when the desired judgment occurs in the set. Assuming that all rules are considered at every stage, forward chaining will eventually find a derivation of any derivable judgment, but it is impossible (in general) to decide algorithmically when to stop extending the set and conclude that the desired judgment is not derivable. We may go on and on adding more judgments to the derivable set without ever achieving the intended goal. It is a matter of understanding the global properties of the rules to determine that a given judgment is not derivable.

Forward chaining is undirected in the sense that it does not take account of the end goal when deciding how to proceed at each step. In contrast, backward chaining is goal-directed. Backward chaining search maintains a queue of current goals, judgments whose derivations are to be sought. Initially, this set consists solely of the judgment we wish to derive. At each stage, we remove a judgment from the queue, and consider all rules whose conclusion is that judgment. For each such rule, we add the premises of that rule to the back of the queue, and continue. If there is more than one such rule, this process must be repeated, with the same starting queue, for each candidate rule. The process terminates whenever the queue is empty, all goals having been achieved; any pending consideration of candidate rules along the way can be discarded. As with forward chaining, backward chaining will eventually find a derivation of any derivable judgment, but there is, in general, no algorithmic method for determining in general whether the current goal is derivable. If it is not, we may futilely add more and more judgments to the goal set, never reaching a point at which all goals have been satisfied.

2.4 Rule Induction

Because an inductive definition specifies the *strongest* judgment form closed under a collection of rules, we may reason about them by *rule induction*. The principle of rule induction states that to show that a property a \mathcal{P} holds whenever a J is derivable, it is enough to show that \mathcal{P} is *closed under*, or *respects*, the rules defining the judgment form J. More precisely, the property \mathcal{P} respects the rule

$$\frac{a_1 \ \mathsf{J} \quad \dots \quad a_k \ \mathsf{J}}{a \ \mathsf{J}}$$

2.4 Rule Induction 17

if $\mathcal{P}(a)$ holds whenever $\mathcal{P}(a_1), \ldots, \mathcal{P}(a_k)$ do. The assumptions $\mathcal{P}(a_1), \ldots, \mathcal{P}(a_k)$ are called the *inductive hypotheses*, and $\mathcal{P}(a)$ is called the *inductive conclusion* of the inference.

The principle of rule induction is simply the expression of the definition of an inductively defined judgment form as the *strongest* judgment form closed under the rules comprising the definition. Thus, the judgment form defined by a set of rules is both (a) closed under those rules, and (b) sufficient for any other property also closed under those rules. The former means that a derivation is evidence for the validity of a judgment; the latter means that we may reason about an inductively defined judgment form by rule induction.

When specialized to rules (2.2), the principle of rule induction states that to show $\mathcal{P}(a)$ whenever a nat, it is enough to show:

```
1. \mathcal{P}(zero).
```

2. for every a, if $\mathcal{P}(a)$, then $\mathcal{P}(\operatorname{succ}(a))$.

The sufficiency of these conditions is the familiar principle of *mathematical induction*.

Similarly, rule induction for rules (2.3) states that to show $\mathcal{P}(a)$ whenever a tree, it is enough to show

```
1. \mathcal{P}(\text{empty}).
```

2. for every a_1 and a_2 , if $\mathcal{P}(a_1)$, and if $\mathcal{P}(a_2)$, then $\mathcal{P}(\text{node}(a_1;a_2))$.

The sufficiency of these conditions is called the principle of *tree induction*.

We may also show by rule induction that the predecessor of a natural number is also a natural number. Although this may seem self-evident, the point of the example is to show how to derive this from first principles.

```
Lemma 2.1. If succ(a) nat, then a nat.
```

Proof. Define $\mathcal{P}(a)$ to mean that if a = succ(b), then b nat. It suffices to show that \mathcal{P} is closed under rules (2.2).

Rule (2.2a) Vacuous, because zero is not of the form succ(-).

Rule (2.2b) The premise of the rule ensures that b nat when a = succ(b).

Using rule induction we may show that equality, as defined by rules (2.4) is reflexive.

Lemma 2.2. *If a* nat, then a is a.

Proof. By rule induction on rules (2.2):

Rule (2.2a) Applying rule (2.4a) we obtain zero is zero.

Rule (2.2b) Assume that a is a. It follows that succ(a) is succ(a) by an application of rule (2.4b).

Similarly, we may show that the successor operation is injective.

Lemma 2.3. *If* $succ(a_1)$ is $succ(a_2)$, then a_1 is a_2 .

Proof. Similar to the proof of Lemma 2.1.

2.5 Iterated and Simultaneous Inductive Definitions

Inductive definitions are often *iterated*, meaning that one inductive definition builds on top of another. In an iterated inductive definition the premises of a rule

$$\frac{J_1 \quad \dots \quad J_k}{J}$$

may be instances of either a previously defined judgment form, or the judgment form being defined. For example, the following rules define the judgment form - list, which states that a is a list of natural numbers:

$$\frac{a \text{ nat } b \text{ list}}{\cos(a;b) \text{ list}}$$
 (2.7b)

The first premise of rule (2.7b) is an instance of the judgment form a nat, which was defined previously, whereas the premise b list is an instance of the judgment form being defined by these rules.

Frequently two or more judgments are defined at once by a *simultaneous inductive definition*. A simultaneous inductive definition consists of a set of rules for deriving instances of several different judgment forms, any of which may appear as the premise of any rule. Because the rules defining each judgment form may involve any of the others, none of the judgment forms can be taken to be defined prior to the others. Instead we must understand that all of the judgment forms are being defined at once by the entire collection of rules. The judgment forms defined by these rules are, as before, the strongest judgment forms that are closed under the rules. Therefore the principle of proof by rule induction continues to apply, albeit in a form that requires us to prove a property of each of the defined judgment forms simultaneously.

For example, consider the following rules, which constitute a simultaneous inductive definition of the judgments *a* even, stating that *a* is an even natural number, and *a* odd, stating that *a* is an odd natural number:

$$\frac{b \text{ odd}}{\text{succ}(b) \text{ even}} \tag{2.8b}$$

$$\frac{a \text{ even}}{\text{succ}(a) \text{ odd}} \tag{2.8c}$$

The principle of rule induction for these rules states that to show simultaneously that $\mathcal{P}(a)$ whenever a even and $\mathcal{Q}(b)$ whenever b odd, it is enough to show the following:

- 1. $\mathcal{P}(\text{zero})$;
- 2. if Q(b), then $P(\operatorname{succ}(b))$;
- 3. if $\mathcal{P}(a)$, then $\mathcal{Q}(\operatorname{succ}(a))$.

As an example, we may use simultaneous rule induction to prove that (1) if a even, then either a is zero or a is succ(b) with b odd, and (2) if a odd, then a is succ(b) with b even. We define $\mathcal{P}(a)$ to hold iff a is zero or a is succ(b) for some b with b odd, and define $\mathcal{Q}(b)$ to hold iff b is succ(a) for some a with a even. The desired result follows by rule induction, because we can prove the following facts:

- 1. $\mathcal{P}(zero)$, which holds because zero is zero.
- 2. If Q(b), then succ(b) is succ(b') for some b' with Q(b'). Take b' to be b and apply the inductive assumption.
- 3. If $\mathcal{P}(a)$, then $\operatorname{succ}(a)$ is $\operatorname{succ}(a')$ for some a' with $\mathcal{P}(a')$. Take a' to be a and apply the inductive assumption.

2.6 Defining Functions by Rules

A common use of inductive definitions is to define a function by giving an inductive definition of its graph relating inputs to outputs, and then showing that the relation uniquely determines the outputs for given inputs. For example, we may define the addition function on natural numbers as the relation sum(a;b;c), with the intended meaning that c is the sum of a and b, as follows:

$$\frac{b \text{ nat}}{\text{sum}(\text{zero};b;b)} \tag{2.9a}$$

$$\frac{\operatorname{sum}(a;b;c)}{\operatorname{sum}(\operatorname{succ}(a);b;\operatorname{succ}(c))} \tag{2.9b}$$

The rules define a ternary (three-place) relation sum(a;b;c) among natural numbers a, b, and c. We may show that c is determined by a and b in this relation.

Theorem 2.4. For every a nat and b nat, there exists a unique c nat such that sum(a;b;c).

Proof. The proof decomposes into two parts:

- 1. (Existence) If a nat and b nat, then there exists c nat such that sum(a;b;c).
- 2. (Uniqueness) If sum(a;b;c), and sum(a;b;c'), then c is c'.

For existence, let $\mathcal{P}(a)$ be the proposition *if* b nat *then there exists* c nat *such that* sum(a;b;c). We prove that if a nat then $\mathcal{P}(a)$ by rule induction on rules (2.2). We have two cases to consider:

20 2.7 Notes

Rule (2.2a) We are to show $\mathcal{P}(\mathsf{zero})$. Assuming b nat and taking c to be b, we obtain $\mathsf{sum}(\mathsf{zero};b;c)$ by rule (2.9a).

Rule (2.2b) Assuming $\mathcal{P}(a)$, we are to show $\mathcal{P}(\operatorname{succ}(a))$. That is, we assume that if b nat then there exists c such that $\operatorname{sum}(a;b;c)$, and are to show that if b' nat, then there exists c' such that $\operatorname{sum}(\operatorname{succ}(a);b';c')$. To this end, suppose that b' nat. Then by induction there exists c such that $\operatorname{sum}(a;b';c)$. Taking c' to be $\operatorname{succ}(c)$, and applying rule (2.9b), we obtain $\operatorname{sum}(\operatorname{succ}(a);b';c')$, as required.

For uniqueness, we prove that *if* sum(a;b;c₁), *then if* sum(a;b;c₂), *then* c₁ is c₂ by rule induction based on rules (2.9).

Rule (2.9a) We have a is zero and c_1 is b. By an inner induction on the same rules, we may show that if $sum(zero;b;c_2)$, then c_2 is b. By Lemma 2.2 we obtain b is b.

Rule (2.9b) We have that a is succ(a') and c_1 is $succ(c'_1)$, where $sum(a';b;c'_1)$. By an inner induction on the same rules, we may show that if $sum(a;b;c_2)$, then c_2 is $succ(c'_2)$ where $sum(a';b;c'_2)$. By the outer inductive hypothesis c'_1 is c'_2 and so c_1 is c_2 .

2.7 Notes

Aczel (1977) provides a thorough account of the theory of inductive definitions on which the present account is based. A significant difference is that we consider inductive definitions of judgments over abts as defined in Chapter 1, rather than with natural numbers. The emphasis on judgments is inspired by Martin-Löf's logic of judgments (Martin-Löf, 1983, 1987).

Exercises

- **2.1**. Give an inductive definition of the judgment max(m;n;p), where m nat, n nat, and p nat, with the meaning that p is the larger of m and n. Prove that every m and n are related to a unique p by this judgment.
- **2.2.** Consider the following rules, which define the judgment hgt(t;n) stating that the binary tree t has *height n*.

$$\frac{}{\mathsf{hgt}(\mathsf{empty};\mathsf{zero})} \tag{2.10a}$$

$$\frac{\operatorname{hgt}(t_1;n_1) \quad \operatorname{hgt}(t_2;n_2) \quad \max(n_1;n_2;n)}{\operatorname{hgt}(\operatorname{node}(t_1;t_2);\operatorname{succ}(n))}$$
(2.10b)

Prove that the judgment hgt defines a function from trees to natural numbers.

2.7 Notes 21

2.3. Given an inductive definition of *ordered variadic trees* whose nodes have a finite, but variable, number of children with a specified left-to-right ordering among them. Your solution should consist of a simultaneous definition of two judgments, *t* tree, stating that *t* is a variadic tree, and *f* forest, stating that *f* is a "forest" (finite sequence) of variadic trees.

- **2.4.** Give an inductive definition of the height of a variadic tree of the kind defined in Exercise **2.3**. Your definition should make use of an auxiliary judgment defining the height of a forest of variadic trees, and will be defined simultaneously with the height of a variadic tree. Show that the two judgments so defined each define a function.
- **2.5**. Give an inductive definition of the *binary natural numbers*, which are either zero, twice a binary number, or one more than twice a binary number. The size of such a representation is logarithmic, rather than linear, in the natural number it represents.
- **2.6**. Give an inductive definition of addition of binary natural numbers as defined in Exercise **2.5**. *Hint*: Proceed by analyzing both arguments to the addition, and make use of an auxiliary function to compute the successor of a binary number. *Hint*: Alternatively, define both the sum and the sum-plus-one of two binary numbers mutually recursively.

22 2.7 Notes



Chapter 3

Hypothetical and General Judgments

A hypothetical judgment expresses an entailment between one or more hypotheses and a conclusion. We will consider two notions of entailment, called *derivability* and *admissibility*. Both express a form of entailment, but they differ in that derivability is stable under extension with new rules, admissibility is not. A *general judgment* expresses the universality, or genericity, of a judgment. There are two forms of general judgment, the *generic* and the *parametric*. The generic judgment expresses generality with respect to all substitution instances for variables in a judgment. The parametric judgment expresses generality with respect to renamings of symbols.

3.1 Hypothetical Judgments

The hypothetical judgment codifies the rules for expressing the validity of a conclusion conditional on the validity of one or more hypotheses. There are two forms of hypothetical judgment that differ according to the sense in which the conclusion is conditional on the hypotheses. One is stable under extension with more rules, and the other is not.

3.1.1 Derivability

For a given set \mathcal{R} of rules, we define the *derivability* judgment, written $J_1, \ldots, J_k \vdash_{\mathcal{R}} K$, where each J_i and K are basic judgments, to mean that we may derive K from the *expansion* $\mathcal{R} \cup \{J_1, \ldots, J_k\}$ of the rules \mathcal{R} with the axioms

$$\overline{J_1}$$
 \cdots $\overline{J_k}$

We treat the *hypotheses*, or *antecedents*, of the judgment, J_1, \ldots, J_k as "temporary axioms", and derive the *conclusion*, or *consequent*, by composing rules in \mathcal{R} . Thus, evidence for a hypothetical judgment consists of a derivation of the conclusion from the hypotheses using the rules in \mathcal{R} .

We use capital Greek letters, usually Γ or Δ , to stand for a finite set of basic judgments, and write $\mathcal{R} \cup \Gamma$ for the expansion of \mathcal{R} with an axiom corresponding to each judgment in Γ . The

judgment $\Gamma \vdash_{\mathcal{R}} K$ means that K is derivable from rules $\mathcal{R} \cup \Gamma$, and the judgment $\vdash_{\mathcal{R}} \Gamma$ means that $\vdash_{\mathcal{R}} J$ for each J in Γ . An equivalent way of defining $J_1, \ldots, J_n \vdash_{\mathcal{R}} J$ is to say that the rule

$$\frac{J_1 \quad \dots \quad J_n}{I} \tag{3.1}$$

is *derivable* from \mathcal{R} , which means that there is a derivation of J composed of the rules in \mathcal{R} augmented by treating J_1, \ldots, J_n as axioms.

For example, consider the derivability judgment

$$a \text{ nat } \vdash_{(2,2)} \operatorname{succ}(\operatorname{succ}(a)) \text{ nat}$$
 (3.2)

relative to rules (2.2). This judgment is valid for any choice of object a, as shown by the derivation

$$\frac{a \text{ nat}}{\operatorname{succ}(a) \text{ nat}}$$

$$\frac{succ(\operatorname{succ}(a)) \text{ nat}}{\operatorname{succ}(\operatorname{succ}(a)) \text{ nat}}$$
(3.3)

which composes rules (2.2), starting with a nat as an axiom, and ending with succ(succ(a)) nat. Equivalently, the validity of (3.2) may also be expressed by stating that the rule

$$\frac{a \text{ nat}}{\text{succ}(\text{succ}(a)) \text{ nat}}$$
 (3.4)

is derivable from rules (2,2).

It follows directly from the definition of derivability that it is stable under extension with new rules.

Theorem 3.1 (Stability). *If* $\Gamma \vdash_{\mathcal{R}} J$, then $\Gamma \vdash_{\mathcal{R} \cup \mathcal{R}'} J$.

Proof. Any derivation of J from $\mathcal{R} \cup \Gamma$ is also a derivation from $(\mathcal{R} \cup \mathcal{R}') \cup \Gamma$, because any rule in \mathcal{R} is also a rule in $\mathcal{R} \cup \mathcal{R}'$.

Derivability enjoys a number of *structural properties* that follow from its definition, independently of the rules R in question.

Reflexivity Every judgment is a consequence of itself: Γ , $J \vdash_{\mathcal{R}} J$. Each hypothesis justifies itself as conclusion.

Weakening If $\Gamma \vdash_{\mathcal{R}} J$, then $\Gamma, K \vdash_{\mathcal{R}} J$. Entailment is not influenced by un-exercised options.

Transitivity If Γ , $K \vdash_{\mathcal{R}} J$ and $\Gamma \vdash_{\mathcal{R}} K$, then $\Gamma \vdash_{\mathcal{R}} J$. If we replace an axiom by a derivation of it, the result is a derivation of its consequent without that hypothesis.

Reflexivity follows directly from the meaning of derivability. Weakening follows directly from the definition of derivability. Transitivity is proved by rule induction on the first premise.

3.1.2 Admissibility

Admissibility, written $\Gamma \models_{\mathcal{R}} J$, is a weaker form of hypothetical judgment stating that $\vdash_{\mathcal{R}} \Gamma$ implies $\vdash_{\mathcal{R}} J$. That is, the conclusion J is derivable from rules \mathcal{R} when the assumptions Γ are all derivable from rules \mathcal{R} . In particular if any of the hypotheses are *not* derivable relative to \mathcal{R} , then the judgment is *vacuously* true. An equivalent way to define the judgment $J_1, \ldots, J_n \models_{\mathcal{R}} J$ is to state that the rule

$$\frac{J_1 \quad \dots \quad J_n}{J} \tag{3.5}$$

is *admissible* relative to the rules in \mathcal{R} . Given any derivations of J_1, \ldots, J_n using the rules in \mathcal{R} , we may build a derivation of J using the rules in \mathcal{R} .

For example, the admissibility judgment

$$\operatorname{succ}(a) \operatorname{even} \models_{(2.8)} a \operatorname{odd}$$
 (3.6)

is valid, because any derivation of succ(a) even from rules (2.8) must contain a sub-derivation of a odd from the same rules, which justifies the conclusion. This fact can be proved by induction on rules (2.8). That judgment (3.6) is valid may also be expressed by saying that the rule

$$\frac{\operatorname{succ}(a) \operatorname{even}}{a \operatorname{odd}} \tag{3.7}$$

is *admissible* relative to rules (2.8).

In contrast to derivability the admissibility judgment is *not* stable under extension to the rules. For example, if we enrich rules (2.8) with the axiom

$$\frac{}{\text{succ}(\text{zero}) \text{ even}}$$
, (3.8)

then rule (3.6) is *inadmissible*, because there is no composition of rules deriving zero odd. Admissibility is as sensitive to which rules are *absent* from an inductive definition as it is to which rules are *present* in it.

The structural properties of derivability ensure that derivability is stronger than admissibility.

Theorem 3.2. *If*
$$\Gamma \vdash_{\mathcal{R}} J$$
, then $\Gamma \models_{\mathcal{R}} J$.

Proof. Repeated application of the transitivity of derivability shows that if $\Gamma \vdash_{\mathcal{R}} J$ and $\vdash_{\mathcal{R}} \Gamma$, then $\vdash_{\mathcal{R}} J$.

To see that the converse fails, note that

$$succ(zero)$$
 even $\forall_{(2.8)}$ zero odd,

because there is no derivation of the right-hand side when the left-hand side is added as an axiom to rules (2.8). Yet the corresponding admissibility judgment

$$succ(zero)$$
 even $\models_{(2.8)}$ zero odd

is valid, because the hypothesis is false: there is no derivation of succ(zero) even from rules (2.8). Even so, the derivability

$$succ(zero)$$
 even $\vdash_{(2.8)} succ(succ(zero))$ odd

is valid, because we may derive the right-hand side from the left-hand side by composing rules (2.8).

Evidence for admissibility can be thought of as a mathematical function transforming derivations $\nabla_1, \ldots, \nabla_n$ of the hypotheses into a derivation ∇ of the consequent. Therefore, the admissibility judgment enjoys the same structural properties as derivability, and hence is a form of hypothetical judgment:

Reflexivity If *J* is derivable from the original rules, then *J* is derivable from the original rules: $J \models_{\mathcal{R}} J$.

Weakening If *J* is derivable from the original rules assuming that each of the judgments in Γ are derivable from these rules, then *J* must also be derivable assuming that Γ and *K* are derivable from the original rules: if $\Gamma \models_{\mathcal{R}} J$, then $\Gamma, K \models_{\mathcal{R}} J$.

Transitivity If Γ , $K \models_{\mathcal{R}} J$ and $\Gamma \models_{\mathcal{R}} K$, then $\Gamma \models_{\mathcal{R}} J$. If the judgments in Γ are derivable, so is K, by assumption, and hence so are the judgments in Γ , K, and hence so is J.

Theorem 3.3. The admissibility judgment $\Gamma \models_{\mathcal{R}} J$ enjoys the structural properties of entailment.

Proof. Follows immediately from the definition of admissibility as stating that if the hypotheses are derivable relative to \mathcal{R} , then so is the conclusion.

If a rule r is admissible with respect to a rule set \mathcal{R} , then $\vdash_{\mathcal{R},r} J$ is equivalent to $\vdash_{\mathcal{R}} J$. For if $\vdash_{\mathcal{R}} J$, then obviously $\vdash_{\mathcal{R},r} J$, by simply disregarding r. Conversely, if $\vdash_{\mathcal{R},r} J$, then we may replace any use of r by its expansion in terms of the rules in \mathcal{R} . It follows by rule induction on \mathcal{R} , r that every derivation from the expanded set of rules \mathcal{R} , r can be transformed into a derivation from \mathcal{R} alone. Consequently, if we wish to prove a property of the judgments derivable from \mathcal{R} , r, when r is admissible with respect to \mathcal{R} , it suffices show that the property is closed under rules \mathcal{R} alone, because its admissibility states that the consequences of rule r are implicit in those of rules \mathcal{R} .

3.2 Hypothetical Inductive Definitions

It is useful to enrich the concept of an inductive definition to allow rules with derivability judgments as premises and conclusions. Doing so lets us introduce *local hypotheses* that apply only in the derivation of a particular premise, and also allows us to constrain inferences based on the *global hypotheses* in effect at the point where the rule is applied.

A hypothetical inductive definition consists of a set of hypothetical rules of the following form:

$$\frac{\Gamma\Gamma_1 \vdash J_1 \quad \dots \quad \Gamma\Gamma_n \vdash J_n}{\Gamma\vdash J} \quad . \tag{3.9}$$

The hypotheses Γ are the *global hypotheses* of the rule, and the hypotheses Γ_i are the *local hypotheses* of the *i*th premise of the rule. Informally, this rule states that J is a derivable consequence of Γ when each J_i is a derivable consequence of Γ , augmented with the hypotheses Γ_i . Thus, one way to show that J is derivable from Γ is to show, in turn, that each J_i is derivable from $\Gamma \Gamma_i$. The derivation of each premise involves a "context switch" in which we extend the global hypotheses with the local hypotheses of that premise, establishing a new set of global hypotheses for use within that derivation.

We require that all rules in a hypothetical inductive definition be *uniform* in the sense that they are applicable in *all* global contexts. Uniformity ensures that a rule can be presented in *implicit*, or *local form*,

$$\frac{\Gamma_1 \vdash J_1 \quad \dots \quad \Gamma_n \vdash J_n}{I} \quad , \tag{3.10}$$

in which the global context has been suppressed with the understanding that the rule applies for any choice of global hypotheses.

A hypothetical inductive definition is to be regarded as an ordinary inductive definition of a formal derivability judgment $\Gamma \vdash J$ consisting of a finite set of basic judgments Γ and a basic judgment J. A set of hypothetical rules \mathcal{R} defines the strongest formal derivability judgment that is *structural* and *closed* under uniform rules \mathcal{R} . Structurality means that the formal derivability judgment must be closed under the following rules:

$$\Gamma I \vdash I$$
 (3.11a)

$$\frac{\Gamma \vdash J}{\Gamma, K \vdash J} \tag{3.11b}$$

$$\frac{\Gamma \vdash K \quad \Gamma, K \vdash J}{\Gamma \vdash I} \tag{3.11c}$$

These rules ensure that formal derivability behaves like a hypothetical judgment. We write $\Gamma \vdash_{\mathcal{R}} J$ to mean that $\Gamma \vdash J$ is derivable from rules \mathcal{R} .

The principle of *hypothetical rule induction* is just the principle of rule induction applied to the formal hypothetical judgment. So to show that $\mathcal{P}(\Gamma \vdash J)$ when $\Gamma \vdash_{\mathcal{R}} J$, it is enough to show that \mathcal{P} is closed under the rules of \mathcal{R} and under the structural rules.¹ Thus, for each rule of the form (3.9), whether structural or in \mathcal{R} , we must show that

if
$$\mathcal{P}(\Gamma \Gamma_1 \vdash J_1)$$
 and ... and $\mathcal{P}(\Gamma \Gamma_n \vdash J_n)$, then $\mathcal{P}(\Gamma \vdash J)$.

But this is just a restatement of the principle of rule induction given in Chapter 2, specialized to the formal derivability judgment $\Gamma \vdash J$.

In practice we usually dispense with the structural rules by the method described in Section 3.1.2. By proving that the structural rules are admissible any proof by rule induction may restrict attention to the rules in \mathcal{R} alone. If all rules of a hypothetical inductive definition are uniform, the structural rules (3.11b) and (3.11c) are clearly admissible. Usually, rule (3.11a) must be postulated explicitly as a rule, rather than shown to be admissible on the basis of the other rules.

 $^{^1}$ Writing $\mathcal{P}(\Gamma \vdash J)$ is a mild abuse of notation in which the turnstile is used to separate the two arguments to \mathcal{P} for the sake of readability.

3.3 General Judgments

General judgments codify the rules for handling variables in a judgment. As in mathematics in general, a variable is treated as an *unknown* ranging over a specified set of objects. A *generic* judgment states that a judgment holds for any choice of objects replacing designated variables in the judgment. Another form of general judgment codifies the handling of symbolic parameters. A *parametric* judgment expresses generality over any choice of fresh renamings of designated symbols of a judgment. To keep track of the active variables and symbols in a derivation, we write $\Gamma \vdash_{\mathcal{R}}^{\mathcal{U};\mathcal{X}} J$ to say that J is derivable from Γ according to rules \mathcal{R} , with objects consisting of abts over symbols \mathcal{U} and variables \mathcal{X} .

The concept of uniformity of a rule must be extended to require that rules be *closed under renaming and substitution* for variables and *closed under renaming* for parameters. More precisely, if \mathcal{R} is a set of rules containing a free variable x of sort s then it must also contain all possible substitution instances of abts a of sort s for s, including those that contain other free variables. Similarly, if s contains rules with a parameter s, then it must contain all instances of that rule obtained by renaming s of a sort to any s of the same sort. Uniformity rules out stating a rule for a variable, without also stating it all instances of that variable. It also rules out stating a rule for a parameter without stating it for all possible renamings of that parameter.

Generic derivability judgment is defined by

$$\mathcal{Y} \mid \Gamma \vdash^{\mathcal{X}}_{\mathcal{R}} J$$
 iff $\Gamma \vdash^{\mathcal{X}\mathcal{Y}}_{\mathcal{R}} J$,

where $\mathcal{Y} \cap \mathcal{X} = \emptyset$. Evidence for generic derivability consists of a *generic derivation* ∇ involving the variables $\mathcal{X} \mathcal{Y}$. So long as the rules are uniform, the choice of \mathcal{Y} does not matter, in a sense to be explained shortly.

For example, the generic derivation ∇ ,

$$\frac{\overline{x \text{ nat}}}{\operatorname{succ}(x) \text{ nat}}$$
$$\operatorname{succ}(\operatorname{succ}(x)) \text{ nat}$$

is evidence for the judgment

$$x\mid x$$
 nat $dash_{(2.2)}^{\mathcal{X}}$ succ $($ succ (x) $)$ nat

provided $x \notin \mathcal{X}$. Any other choice of x would work just as well, as long as all rules are uniform.

The generic derivability judgment enjoys the following *structural properties* governing the behavior of variables, provided that \mathcal{R} is uniform.

Proliferation If
$$\mathcal{Y} \mid \Gamma \vdash^{\mathcal{X}}_{\mathcal{R}} J$$
, then $\mathcal{Y}, y \mid \Gamma \vdash^{\mathcal{X}}_{\mathcal{R}} J$.

Renaming If
$$\mathcal{Y}, y \mid \Gamma \vdash^{\mathcal{X}}_{\mathcal{R}} J$$
, then $\mathcal{Y}, y' \mid [y \leftrightarrow y']\Gamma \vdash^{\mathcal{X}}_{\mathcal{R}} [y \leftrightarrow y']J$ for any $y' \notin \mathcal{X} \mathcal{Y}$.

Substitution If
$$\mathcal{Y}, y \mid \Gamma \vdash_{\mathcal{R}}^{\mathcal{X}} J$$
 and $a \in \mathcal{B}[\mathcal{X} \mathcal{Y}]$, then $\mathcal{Y} \mid [a/y]\Gamma \vdash_{\mathcal{R}}^{\mathcal{X}} [a/y]J$.

Proliferation is guaranteed by the interpretation of rule schemes as ranging over all expansions of the universe. Renaming is built into the meaning of the generic judgment. It is left implicit in the principle of substitution that the substituting abt is of the same sort as the substituted variable.

Parametric derivability is defined analogously to generic derivability, albeit by generalizing over symbols, rather than variables. Parametric derivability is defined by

$$V \parallel \mathcal{Y} \mid \Gamma \vdash_{\mathcal{R}}^{\mathcal{U};\mathcal{X}} J \text{ iff } \mathcal{Y} \mid \Gamma \vdash_{\mathcal{R}}^{\mathcal{U}\mathcal{V};\mathcal{X}} J,$$

where $V \cap U = \emptyset$. Evidence for parametric derivability consists of a derivation ∇ involving the symbols V. Uniformity of \mathcal{R} ensures that any choice of parameter names is as good as any other; derivability is stable under renaming.

3.4 Generic Inductive Definitions

A *generic inductive definition* admits generic hypothetical judgments in the premises of rules, with the effect of augmenting the variables, as well as the rules, within those premises. A *generic rule* has the form

$$\frac{\mathcal{Y}\,\mathcal{Y}_1\mid\Gamma\,\Gamma_1\vdash J_1\quad\ldots\quad\mathcal{Y}\,\mathcal{Y}_n\mid\Gamma\,\Gamma_n\vdash J_n}{\mathcal{Y}\mid\Gamma\vdash J}\,.$$
(3.12)

The variables \mathcal{Y} are the *global variables* of the inference, and, for each $1 \le i \le n$, the variables \mathcal{Y}_i are the *local variables* of the *i*th premise. In most cases a rule is stated for *all* choices of global variables and global hypotheses. Such rules can be given in *implicit form*,

$$\frac{\mathcal{Y}_1 \mid \Gamma_1 \vdash J_1 \quad \dots \quad \mathcal{Y}_n \mid \Gamma_n \vdash J_n}{I} \quad . \tag{3.13}$$

A generic inductive definition is just an ordinary inductive definition of a family of *formal* generic judgments of the form $\mathcal{Y} \mid \Gamma \vdash J$. Formal generic judgments are identified up to renaming of variables, so that the latter judgment is treated as identical to the judgment $\mathcal{Y}' \mid \widehat{\rho}(\Gamma) \vdash \widehat{\rho}(J)$ for any renaming $\rho: \mathcal{Y} \leftrightarrow \mathcal{Y}'$. If \mathcal{R} is a collection of generic rules, we write $\mathcal{Y} \mid \Gamma \vdash_{\mathcal{R}} J$ to mean that the formal generic judgment $\mathcal{Y} \mid \Gamma \vdash_{\mathcal{I}} J$ is derivable from rules \mathcal{R} .

When specialized to a set of generic rules, the principle of rule induction states that to show $\mathcal{P}(\mathcal{Y} \mid \Gamma \vdash J)$ when $\mathcal{Y} \mid \Gamma \vdash_{\mathcal{R}} J$, it is enough to show that \mathcal{P} is closed under the rules \mathcal{R} . Specifically, for each rule in \mathcal{R} of the form (3.12), we must show that

if
$$\mathcal{P}(\mathcal{Y}\mathcal{Y}_1 \mid \Gamma\Gamma_1 \vdash I_1) \dots \mathcal{P}(\mathcal{Y}\mathcal{Y}_n \mid \Gamma\Gamma_n \vdash I_n)$$
 then $\mathcal{P}(\mathcal{Y} \mid \Gamma \vdash I)$.

By the identification convention (stated in Chapter 1) the property \mathcal{P} must respect renamings of the variables in a formal generic judgment.

To ensure that the formal generic judgment behaves like a generic judgment, we must always ensure that the following *structural rules* are admissible:

$$\frac{1}{\mathcal{Y} \mid \Gamma, I \vdash I} \tag{3.14a}$$

30 3.5 Notes

$$\frac{\mathcal{Y} \mid \Gamma \vdash J}{\mathcal{Y} \mid \Gamma, J' \vdash J} \tag{3.14b}$$

$$\frac{\mathcal{Y} \mid \Gamma \vdash J}{\mathcal{Y}, x \mid \Gamma \vdash J} \tag{3.14c}$$

$$\frac{\mathcal{Y}, x' \mid [x \leftrightarrow x']\Gamma \vdash [x \leftrightarrow x']J}{\mathcal{Y}, x \mid \Gamma \vdash J} \tag{3.14d}$$

$$\frac{\mathcal{Y} \mid \Gamma \vdash J \quad \mathcal{Y} \mid \Gamma, J \vdash J'}{\mathcal{Y} \mid \Gamma \vdash J'} \tag{3.14e}$$

$$\frac{\mathcal{Y}, x \mid \Gamma \vdash J \quad a \in \mathcal{B}[\mathcal{Y}]}{\mathcal{Y} \mid [a/x]\Gamma \vdash [a/x]J}$$
(3.14f)

The admissibility of rule (3.14a) is, in practice, ensured by explicitly including it. The admissibility of rules (3.14b) and (3.14c) is assured if each of the generic rules is uniform, because we may assimilate the added variable x to the global variables, and the added hypothesis J, to the global hypotheses. The admissibility of rule (3.14d) is ensured by the identification convention for the formal generic judgment. Rule (3.14f) must be verified explicitly for each inductive definition.

The concept of a generic inductive definition extends to parametric judgments as well. Briefly, rules are defined on formal parametric judgments of the form $\mathcal{V} \parallel \mathcal{Y} \mid \Gamma \vdash J$, with symbols \mathcal{V} , as well as variables, \mathcal{Y} . Such formal judgments are identified up to renaming of its variables and its symbols to ensure that the meaning is independent of the choice of variable and symbol names.

3.5 Notes

The concepts of entailment and generality are fundamental to logic and programming languages. The formulation given here builds on Martin-Löf (1983, 1987) and Avron (1991). Hypothetical and general reasoning are consolidated into a single concept in the AUTOMATH languages (Nederpelt et al., 1994) and in the LF Logical Framework (Harper et al., 1993). These systems allow arbitrarily nested combinations of hypothetical and general judgments, whereas the present account considers only general hypothetical judgments over basic judgment forms. On the other hand we consider here symbols, as well as variables, which are not present in these previous accounts. Parametric judgments are required for specifying languages that admit the dynamic creation of "new" objects (see Chapter 34).

Exercises

3.1. *Combinators* are inductively defined by the rule set C given as follows:

$$\frac{}{\text{s comb}}$$
 (3.15a)

3.5 Notes 31

$$\frac{}{\text{k comb}}$$
 (3.15b)

$$\frac{a_1 \text{ comb}}{\text{ap}(a_1; a_2) \text{ comb}}$$
 (3.15c)

Give an inductive definition of the *length* of a combinator defined as the number of occurrences of S and K within it.

3.2. The general judgment

$$x_1, \ldots, x_n \mid x_1 \text{ comb}, \ldots, x_n \text{ comb } \vdash_{\mathcal{C}} A \text{ comb}$$

3.3. *Conversion,* or *equivalence,* of combinators is expressed by the judgment $A \equiv B$ defined by the rule set \mathcal{E} extending \mathcal{C} as follows:²

$$\frac{a \text{ comb}}{a \equiv a} \tag{3.16a}$$

$$\frac{a_2 \equiv a_1}{a_1 \equiv a_2} \tag{3.16b}$$

$$\begin{array}{ccc}
a_1 \equiv a_2 & a_2 \equiv a_3 \\
a_1 \equiv a_3
\end{array}
\tag{3.16c}$$

$$\frac{a_1 \equiv a_1' \quad a_2 \equiv a_2'}{a_1 a_2 \equiv a_1' a_2'} \tag{3.16d}$$

$$\frac{a_1 \text{ comb} \quad a_2 \text{ comb}}{\text{k } a_1 a_2 \equiv a_1} \tag{3.16e}$$

$$\frac{a_1 \text{ comb}}{\text{s } a_1 a_2 a_3 \equiv (a_1 a_3) (a_2 a_3)}$$
(3.16f)

The no-doubt mysterious motivation for the last two equations will become clearer in a moment. For now, show that

$$x \mid x \text{ comb } \vdash_{\mathcal{C} \cup \mathcal{E}} \mathtt{skk} x \equiv x.$$

3.4. Show that if $x \mid x$ comb $\vdash_{\mathcal{C}} a$ comb, then there is a combinator a', written [x]a and called *bracket abstraction*, such that

$$x \mid x \text{ comb } \vdash_{\mathcal{C} \cup \mathcal{E}} a' x \equiv a.$$

Consequently, by Exercise 3.2, if a'' comb, then

$$([x]a)a'' \equiv [a''/x]a.$$

²The combinator $ap(a_1;a_2)$ is written $a_1 a_2$ for short, left-associatively when used in succession.

32 3.5 Notes

Hint: Inductively define the judgment

$$x \mid x \text{ comb} \vdash \text{abs}_x \ a \text{ is } a'$$
,

where $x \mid x \text{ comb } \vdash a \text{ comb}$. Then argue that it defines a' as a binary function of x and a. The motivation for the conversion axioms governing k and k should become clear while developing the proof of the desired equivalence.

3.5. Prove that bracket abstraction, as defined in Exercise **3.4**, is *non-compositional* by exhibiting *a* and *b* such that *a* comb and

$$xy \mid x \text{ comb } y \text{ comb } \vdash_{\mathcal{C}} b \text{ comb}$$

such that $[a/y]([x]b) \neq [x]([a/y]b)$. *Hint*: Consider the case that *b* is *y*.

Suggest a modification to the definition of bracket abstraction that is *compositional* by showing under the same conditions given above that

$$[a/y]([x]b) = [x]([a/y]b).$$

3.6. Consider the set $\mathcal{B}[\mathcal{X}]$ of abts generated by the operators ap, with arity $(\mathsf{Exp}, \mathsf{Exp})\mathsf{Exp}$, and λ , with arity $(\mathsf{Exp}, \mathsf{Exp})\mathsf{Exp}$, and possibly involving variables in \mathcal{X} , all of which are of sort Exp. Give an inductive definition of the judgment b closed, which specifies that b has no free occurrences of the variables in \mathcal{X} . *Hint*: it is essential to give an inductive definition of the hypothetical, general judgment

$$x_1, \ldots, x_n \mid x_1 \text{ closed}, \ldots, x_n \text{ closed} \vdash b \text{ closed}$$

in order to account for the binding of a variable by the λ operator. The hypothesis that a variable is closed seems self-contradictory in that a variable obviously occurs free in itself. Explain why this is not the case by examining carefully the meaning of the hypothetical and general judgments.

Part II Statics and Dynamics



Chapter 4

Statics

Most programming languages exhibit a *phase distinction* between the *static* and *dynamic* phases of processing. The static phase consists of parsing and type checking to ensure that the program is well-formed; the dynamic phase consists of execution of well-formed programs. A language is said to be *safe* exactly when well-formed programs are well-behaved when executed.

The static phase is specified by a *statics* comprising a set of rules for deriving *typing judgments* stating that an expression is well-formed of a certain type. Types mediate the interaction between the constituent parts of a program by "predicting" some aspects of the execution behavior of the parts so that we may ensure they fit together properly at run-time. Type safety tells us that these predictions are correct; if not, the statics is considered to be improperly defined, and the language is deemed *unsafe* for execution.

In this chapter we present the statics of a simple expression language, **E**, as an illustration of the method that we will employ throughout this book.

4.1 Syntax

When defining a language we shall be primarily concerned with its abstract syntax, specified by a collection of operators and their arities. The abstract syntax provides a systematic, unambiguous account of the hierarchical and binding structure of the language, and is considered the official presentation of the language. However, for the sake of clarity, it is also useful to specify minimal concrete syntax conventions, without going through the trouble to set up a fully precise grammar for it.

We will accomplish both of these purposes with a syntax chart, whose meaning is best illus-

36 4.2 Type System

trated by example. The following chart summarizes the abstract and concrete syntax of E.

Тур	τ	::=	num	num	numbers
			str	str	strings
Exp	е	::=	\boldsymbol{x}	\boldsymbol{x}	variable
			$\mathtt{num}[n]$	n	numeral
			$\mathtt{str}[s]$	"s"	literal
			$plus(e_1;e_2)$	$e_1 + e_2$	addition
			$times(e_1;e_2)$	$e_1 * e_2$	multiplication
			$\mathtt{cat}(\mathit{e}_1;\mathit{e}_2)$	$e_1 \hat{e}_2$	concatenation
			len(e)	e	length
			$let(e_1; x.e_2)$	let x be e_1 in e_2	definition

This chart defines two sorts, Typ, ranged over by τ , and Exp, ranged over by e. The chart defines a set of operators and their arities. For example, it specifies that the operator let has arity (Exp, Exp.Exp)Exp, which specifies that it has two arguments of sort Exp, and binds a variable of sort Exp in the second argument.

4.2 Type System

The role of a type system is to impose constraints on the formations of phrases that are sensitive to the context in which they occur. For example, whether the expression $\mathtt{plus}(x; \mathtt{num}[n])$ is sensible depends on whether the variable x is restricted to have type \mathtt{num} in the surrounding context of the expression. This example is, in fact, illustrative of the general case, in that the *only* information required about the context of an expression is the type of the variables within whose scope the expression lies. Consequently, the statics of \mathbf{E} consists of an inductive definition of generic hypothetical judgments of the form

$$\mathcal{X} \mid \Gamma \vdash e : \tau$$
,

where \mathcal{X} is a finite set of variables, and Γ is a *typing context* consisting of hypotheses of the form $x:\tau$, one for each $x\in\mathcal{X}$. We rely on typographical conventions to determine the set of variables, using the letters x and y to stand for them. We write $x\notin dom(\Gamma)$ to say that there is no assumption in Γ of the form $x:\tau$ for any type τ , in which case we say that the variable x is *fresh* for Γ .

The rules defining the statics of **E** are as follows:

$$\overline{\Gamma, x : \tau \vdash x : \tau} \tag{4.1a}$$

$$\overline{\Gamma \vdash \mathsf{str}[s] : \mathsf{str}} \tag{4.1b}$$

$$\overline{\Gamma \vdash \text{num}[n] : \text{num}}$$
 (4.1c)

$$\frac{\Gamma \vdash e_1 : \text{num} \quad \Gamma \vdash e_2 : \text{num}}{\Gamma \vdash \text{plus}(e_1; e_2) : \text{num}}$$
(4.1d)

$$\frac{\Gamma \vdash e_1 : \text{num} \quad \Gamma \vdash e_2 : \text{num}}{\Gamma \vdash \text{times}(e_1; e_2) : \text{num}}$$
(4.1e)

$$\frac{\Gamma \vdash e_1 : \operatorname{str} \quad \Gamma \vdash e_2 : \operatorname{str}}{\Gamma \vdash \operatorname{cat}(e_1; e_2) : \operatorname{str}}$$
(4.1f)

$$\frac{\Gamma \vdash e : \mathtt{str}}{\Gamma \vdash \mathtt{len}(e) : \mathtt{num}} \tag{4.1g}$$

$$\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma, x : \tau_1 \vdash e_2 : \tau_2}{\Gamma \vdash \mathsf{let}(e_1; x.e_2) : \tau_2} \tag{4.1h}$$

In rule (4.1h) we tacitly assume that the variable x is not already declared in Γ . This condition may always be met by choosing a suitable representative of the α -equivalence class of the let expression.

It is easy to check that every expression has at most one type by *induction on typing*, which is rule induction applied to rules (4.1).

Lemma 4.1 (Unicity of Typing). *For every typing context* Γ *and expression e, there exists at most one* τ *such that* $\Gamma \vdash e : \tau$.

Proof. By rule induction on rules (4.1), making use of the fact that variables have at most one type in any typing context.

The typing rules are *syntax-directed* in the sense that there is exactly one rule for each form of expression. Consequently it is easy to give necessary conditions for typing an expression that invert the sufficient conditions expressed by the corresponding typing rule.

Lemma 4.2 (Inversion for Typing). Suppose that $\Gamma \vdash e : \tau$. If $e = \text{plus}(e_1; e_2)$, then $\tau = \text{num}$, $\Gamma \vdash e_1 : \text{num}$, and $\Gamma \vdash e_2 : \text{num}$, and similarly for the other constructs of the language.

Proof. These may all be proved by induction on the derivation of the typing judgment $\Gamma \vdash e : \tau$. \Box

In richer languages such inversion principles are more difficult to state and to prove.

4.3 Structural Properties

The statics enjoys the structural properties of the generic hypothetical judgment.

Lemma 4.3 (Weakening). If $\Gamma \vdash e' : \tau'$, then $\Gamma, x : \tau \vdash e' : \tau'$ for any $x \notin dom(\Gamma)$ and any type τ .

Proof. By induction on the derivation of $\Gamma \vdash e' : \tau'$. We will give one case here, for rule (4.1h). We have that $e' = \mathtt{let}(e_1; z.e_2)$, where by the conventions on variables we may assume z is chosen such that $z \notin dom(\Gamma)$ and $z \neq x$. By induction we have

- 1. $\Gamma, x : \tau \vdash e_1 : \tau_1$,
- 2. $\Gamma, x : \tau, z : \tau_1 \vdash e_2 : \tau'$

from which the result follows by rule (4.1h).

Lemma 4.4 (Substitution). *If* Γ , $x : \tau \vdash e' : \tau'$ *and* $\Gamma \vdash e : \tau$, *then* $\Gamma \vdash [e/x]e' : \tau'$.

Proof. By induction on the derivation of Γ , $x:\tau\vdash e':\tau'$. We again consider only rule (4.1h). As in the preceding case, $e'=\mathtt{let}(e_1;z.e_2)$, where z is chosen so that $z\neq x$ and $z\notin dom(\Gamma)$. We have by induction and Lemma 4.3 that

- 1. $\Gamma \vdash [e/x]e_1 : \tau_1$,
- 2. $\Gamma, z : \tau_1 \vdash [e/x]e_2 : \tau'$.

By the choice of z we have

$$[e/x]$$
let $(e_1; z.e_2) =$ let $([e/x]e_1; z.[e/x]e_2).$

It follows by rule (4.1h) that $\Gamma \vdash [e/x] \texttt{let}(e_1; z.e_2) : \tau'$, as desired.

From a programming point of view, Lemma 4.3 allows us to use an expression in any context that binds its free variables: if e is well-typed in a context Γ , then we may "import" it into any context that includes the assumptions Γ . In other words introducing new variables beyond those required by an expression e does not invalidate e itself; it remains well-formed, with the same type. More importantly, Lemma 4.4 expresses the important concepts of *modularity* and *linking*. We may think of the expressions e and e' as two *components* of a larger system in which e' is a *client* of the *implementation* e. The client declares a variable specifying the type of the implementation, and is type checked knowing only this information. The implementation must be of the specified type to satisfy the assumptions of the client. If so, then we may link them to form the composite system [e/x]e'. This implementation may itself be the client of another component, represented by a variable e0 that is replaced by that component during linking. When all such variables have been implemented, the result is a *closed expression* that is ready for execution (evaluation).

The converse of Lemma 4.4 is called *decomposition*. It states that any (large) expression can be decomposed into a client and implementor by introducing a variable to mediate their interaction.

Lemma 4.5 (Decomposition). *If* $\Gamma \vdash [e/x]e' : \tau'$, then for every type τ such that $\Gamma \vdash e : \tau$, we have $\Gamma, x : \tau \vdash e' : \tau'$.

Proof. The typing of [e/x]e' depends only on the type of e wherever it occurs, if at all.

Lemma 4.5 tells us that any sub-expression can be isolated as a separate module of a larger system. This property is especially useful when the variable x occurs more than once in e', because then one copy of e suffices for all occurrences of x in e'.

The statics of **E** given by rules (4.1) exemplifies a recurrent pattern. The constructs of a language are classified into one of two forms, the *introduction* and the *elimination*. The introduction forms for a type determine the *values*, or *canonical forms*, of that type. The elimination forms determine how to manipulate the values of a type to form a computation of another (possibly the same) type.

¹This point may seem so obvious that it is not worthy of mention, but, surprisingly, there are useful type systems that lack this property. Because they do not necessarily validate the structural principle of weakening, they are called *substructural* type systems.

4.4 Notes 39

In the language **E** the introduction forms for the type num are the numerals, and those for the type str are the literals. The elimination forms for the type num are addition and multiplication, and those for the type str are concatenation and length.

The importance of this classification will become clear once we have defined the dynamics of the language in Chapter 5. Then we will see that the elimination forms are *inverse* to the introduction forms in that they "take apart" what the introduction forms have "put together." The coherence of the statics and dynamics of a language expresses the concept of *type safety*, the subject of Chapter 6.

4.4 Notes

The concept of the statics of a programming language was historically slow to develop, perhaps because the earliest languages had relatively few features and only very weak type systems. Statics in the sense considered here was introduced in the definition of the Standard ML programming language (Milner et al., 1997), building on earlier work by Church and others on the typed λ -calculus (Barendregt, 1992). The concept of introduction and elimination, and the associated inversion principle, was introduced by Gentzen in his pioneering work on natural deduction (Gentzen, 1969). These principles were applied to the structure of programming languages by Martin-Löf (1984, 1980).

Exercises

4.1. It is sometimes useful to give the typing judgment $\Gamma \vdash e : \tau$ an "operational" reading that specifies more precisely the flow of information required to derive a typing judgment (or determine that it is not derivable). The *analytic* mode corresponds to the context, expression, and type being given, with the goal to determine whether the typing judgment is derivable. The *synthetic* mode corresponds to the context and expression being given, with the goal to find the unique type τ , if any, possessed by the expression in that context. These two readings can be made explicit as judgments of the form $e \downarrow \tau$, corresponding to the analytic mode, and $e \uparrow \tau$, corresponding to the synthetic mode.

Give a simultaneous inductive definition of these two judgments according to the following guidelines:

- (a) Variables are introduced in synthetic form.
- (b) If we can synthesize a unique type for an expression, then we can analyze it with respect to a given type by checking type equality.
- (c) Definitions need care, because the type of the defined expression is not given, even when the type of the result is given.

There is room for variation; the point of the exercise is to explore the possibilities.

4.4 Notes

4.2. One way to limit the range of possibilities in the solution to Exercise **4.1** is to restrict and extend the syntax of the language so that every expression is either synthetic or analytic according to the following suggestions:

- (a) Variables are analytic.
- (b) Introduction forms are analytic, elimination forms are synthetic.
- (c) An analytic expression can be made synthetic by introducing a *type cast* of the form $cast\{\tau\}(e)$ specifying that e must check against the specified type τ , which is synthesized for the whole expression.
- (d) The defining expression of a definition must be synthetic, but the scope of the definition can be either synthetic or analytic.

Reformulate your solution to Exercise 4.1 to take account of these guidelines.

Chapter 5

Dynamics

The *dynamics* of a language describes how programs are executed. The most important way to define the dynamics of a language is by the method of *structural dynamics*, which defines a *transition system* that inductively specifies the step-by-step process of executing a program. Another method for presenting dynamics, called *contextual dynamics*, is a variation of structural dynamics in which the transition rules are specified in a slightly different way. An *equational dynamics* presents the dynamics of a language by a collection of rules defining when one program is *definitionally equivalent* to another.

5.1 Transition Systems

A transition system is specified by the following four forms of judgment:

- 1. *s* state, asserting that *s* is a *state* of the transition system.
- 2. *s* final, where *s* state, asserting that *s* is a *final* state.
- 3. *s* initial, where *s* state, asserting that *s* is an *initial* state.
- 4. $s \mapsto s'$, where s state and s' state, asserting that state s may transition to state s'.

In practice we always arrange things so that no transition is possible from a final state: if s final, then there is no s' state such that $s \mapsto s'$. A state from which no transition is possible is stuck. Whereas all final states are, by convention, stuck, there may be stuck states in a transition system that are not final. A transition system is deterministic iff for every state s there exists at most one state s' such that $s \mapsto s'$, otherwise it is non-deterministic.

A transition sequence is a sequence of states s_0, \ldots, s_n such that s_0 initial, and $s_i \mapsto s_{i+1}$ for every $0 \le i < n$. A transition sequence is *maximal* iff there is no s such that $s_n \mapsto s$, and it is complete iff it is maximal and s_n final. Thus every complete transition sequence is maximal, but maximal sequences are not necessarily complete. The judgment $s \downarrow$ means that there is a complete transition sequence starting from s, which is to say that there exists s' final such that $s \mapsto^* s'$.

The *iteration* of transition judgment $s \mapsto^* s'$ is inductively defined by the following rules:

$$\overline{s \longmapsto^* s}$$
 (5.1a)

$$\frac{s \longmapsto s' \quad s' \longmapsto^* s''}{s \longmapsto^* s''} \tag{5.1b}$$

When applied to the definition of iterated transition, the principle of rule induction states that to show that P(s,s') holds when $s \mapsto^* s'$, it is enough to show these two properties of P:

- 1. P(s,s).
- 2. if $s \mapsto s'$ and P(s', s''), then P(s, s'').

The first requirement is to show that *P* is reflexive. The second is to show that *P* is *closed under* head expansion, or closed under inverse evaluation. Using this principle, it is easy to prove that \mapsto^* is reflexive and transitive.

The *n*-times iterated transition judgment $s \mapsto^* ns'$, where $n \ge 0$, is inductively defined by the following rules.

$$\overline{s} \xrightarrow{*} 0s$$
 (5.2a)

$$\frac{s \longmapsto^* 0s}{s \longmapsto^* s' \longmapsto^* ns''} \\
\frac{s \longmapsto^* n + 1s''}{s \longmapsto^* n + 1s''}$$
(5.2a)

Theorem 5.1. For all states s and s', $s \mapsto^* s'$ iff $s \mapsto^* ks'$ for some $k \ge 0$.

Proof. From left to right, by induction on the definition of multi-step transition. From right to left, by mathematical induction on $k \geq 0$.

Structural Dynamics 5.2

A structural dynamics for the language E is given by a transition system whose states are closed expressions. All states are initial. The final states are the (closed) values, which represent the completed computations. The judgment e val, which states that e is a value, is inductively defined by the following rules:

$$\overline{\operatorname{num}[n] \text{ val}} \tag{5.3a}$$

$$\overline{\operatorname{str}[s]\operatorname{val}}$$
 (5.3b)

The transition judgment $e \mapsto e'$ between states is inductively defined by the following rules:

$$\frac{n_1 + n_2 = n}{\operatorname{plus}(\operatorname{num}[n_1]; \operatorname{num}[n_2]) \longmapsto \operatorname{num}[n]}$$
(5.4a)

$$\frac{e_1 \longmapsto e'_1}{\mathtt{plus}(e_1; e_2) \longmapsto \mathtt{plus}(e'_1; e_2)} \tag{5.4b}$$

$$\frac{e_1 \text{ val} \quad e_2 \longmapsto e_2'}{\text{plus}(e_1; e_2) \longmapsto \text{plus}(e_1; e_2')}$$
(5.4c)

$$\frac{s_1 \hat{s}_2 = s}{\operatorname{cat}(\operatorname{str}[s_1]; \operatorname{str}[s_2]) \longmapsto \operatorname{str}[s]}$$
(5.4d)

$$\frac{e_1 \longmapsto e'_1}{\operatorname{cat}(e_1; e_2) \longmapsto \operatorname{cat}(e'_1; e_2)} \tag{5.4e}$$

$$\frac{e_1 \text{ val} \quad e_2 \longmapsto e_2'}{\cot(e_1; e_2) \longmapsto \cot(e_1; e_2')} \tag{5.4f}$$

$$\left[\frac{e_1 \longmapsto e'_1}{\operatorname{let}(e_1; x.e_2) \longmapsto \operatorname{let}(e'_1; x.e_2)}\right] \tag{5.4g}$$

$$\frac{[e_1 \text{ val}]}{\text{let}(e_1; x.e_2) \longmapsto [e_1/x]e_2}$$
(5.4h)

We have omitted rules for multiplication and computing the length of a string, which follow a similar pattern. Rules (5.4a), (5.4d), and (5.4h) are *instruction transitions*, because they correspond to the primitive steps of evaluation. The remaining rules are *search transitions* that determine the order of execution of instructions.

The bracketed rule, rule (5.4g), and bracketed premise on rule (5.4h), are included for a *by-value* interpretation of let, and omitted for a *by-name* interpretation. The by-value interpretation evaluates an expression before binding it to the defined variable, whereas the by-name interpretation binds it in unevaluated form. The by-value interpretation saves work if the defined variable is used more than once, but wastes work if it is not used at all. Conversely, the by-name interpretation saves work if the defined variable is not used, and wastes work if it is used more than once.

A derivation sequence in a structural dynamics has a two-dimensional structure, with the number of steps in the sequence being its "width" and the derivation tree for each step being its "height." For example, consider the following evaluation sequence.

$$\begin{array}{l} \operatorname{let}(\operatorname{plus}(\operatorname{num}[1];\operatorname{num}[2]);x.\operatorname{plus}(\operatorname{plus}(x;\operatorname{num}[3]);\operatorname{num}[4])) \\ \longmapsto \operatorname{let}(\operatorname{num}[3];x.\operatorname{plus}(\operatorname{plus}(x;\operatorname{num}[3]);\operatorname{num}[4])) \\ \longmapsto \operatorname{plus}(\operatorname{plus}(\operatorname{num}[3];\operatorname{num}[3]);\operatorname{num}[4]) \\ \longmapsto \operatorname{plus}(\operatorname{num}[6];\operatorname{num}[4]) \\ \longmapsto \operatorname{num}[10] \end{array}$$

Each step in this sequence of transitions is justified by a derivation according to rules (5.4). For example, the third transition in the preceding example is justified by the following derivation:

$$\frac{\overline{\mathtt{plus}(\mathtt{num}[3];\mathtt{num}[3])} \longmapsto \mathtt{num}[6]}{\mathtt{plus}(\mathtt{plus}(\mathtt{num}[3];\mathtt{num}[4]) \longmapsto \mathtt{plus}(\mathtt{num}[6];\mathtt{num}[4])} \tag{5.4b}$$

The other steps are similarly justified by composing rules.

The principle of rule induction for the structural dynamics of **E** states that to show $\mathcal{P}(e \longmapsto e')$ when $e \longmapsto e'$, it is enough to show that \mathcal{P} is closed under rules (5.4). For example, we may show by rule induction that the structural dynamics of **E** is *determinate*, which means that an expression may transition to at most one other expression. The proof requires a simple lemma relating transition to values.

Lemma 5.2 (Finality of Values). For no expression e do we have both e val and $e \mapsto e'$ for some e'.

Proof. By rule induction on rules (5.3) and (5.4).

Lemma 5.3 (Determinacy). *If* $e \mapsto e'$ and $e \mapsto e''$, then e' and e'' are α -equivalent.

Proof. By rule induction on the premises $e \mapsto e'$ and $e \mapsto e''$, carried out either simultaneously or in either order. The primitive operators, such as addition, are assumed to have a unique value when applied to values.

Rules (5.4) exemplify the *inversion principle* of language design, which states that the elimination forms are *inverse* to the introduction forms of a language. The search rules determine the *principal arguments* of each elimination form, and the instruction rules specify how to evaluate an elimination form when all of its principal arguments are in introduction form. For example, rules (5.4) specify that both arguments of addition are principal, and specify how to evaluate an addition once its principal arguments are evaluated to numerals. The inversion principle is central to ensuring that a programming language is properly defined, the exact statement of which is given in Chapter 6.

5.3 Contextual Dynamics

A variant of structural dynamics, called *contextual dynamics*, is sometimes useful. There is no fundamental difference between contextual and structural dynamics, rather one of style. The main idea is to isolate instruction steps as a special form of judgment, called *instruction transition*, and to formalize the process of locating the next instruction using a device called an *evaluation context*. The judgment *e* val, defining whether an expression is a value, remains unchanged.

The instruction transition judgment $e_1 \longrightarrow e_2$ for **E** is defined by the following rules, together with similar rules for multiplication of numbers and the length of a string.

$$\frac{m+n=p}{\mathtt{plus}(\mathtt{num}[m];\mathtt{num}[n]) \longrightarrow \mathtt{num}[p]}$$
(5.5a)

$$\frac{s^{\hat{}}t = u}{\operatorname{cat}(\operatorname{str}[s]; \operatorname{str}[t]) \longrightarrow \operatorname{str}[u]}$$
(5.5b)

$$\overline{\operatorname{let}(e_1; x.e_2) \longrightarrow [e_1/x]e_2} \tag{5.5c}$$

The judgment \mathcal{E} ectx determines the location of the next instruction to execute in a larger expression. The position of the next instruction step is specified by a "hole", written \circ , into which

the next instruction is placed, as we shall detail shortly. (The rules for multiplication and length are omitted for concision, as they are handled similarly.)

$$\overline{\circ}$$
 ectx (5.6a)

$$\frac{\mathcal{E}_1 \text{ ectx}}{\text{plus}(\mathcal{E}_1; e_2) \text{ ectx}}$$
 (5.6b)

$$\frac{e_1 \text{ val } \mathcal{E}_2 \text{ ectx}}{\text{plus}(e_1; \mathcal{E}_2) \text{ ectx}}$$
(5.6c)

The first rule for evaluation contexts specifies that the next instruction may occur "here", at the occurrence of the hole. The remaining rules correspond one-for-one to the search rules of the structural dynamics. For example, rule (5.6c) states that in an expression plus (e_1 ; e_2), if the first argument, e_1 , is a value, then the next instruction step, if any, lies at or within the second argument,

An evaluation context is a template that is instantiated by replacing the hole with an instruction to be executed. The judgment $e' = \mathcal{E}\{e\}$ states that the expression e' is the result of filling the hole in the evaluation context \mathcal{E} with the expression e. It is inductively defined by the following rules:

$$\overline{e = \circ\{e\}} \tag{5.7a}$$

$$e = o\{e\}$$

$$\frac{e_1 = \mathcal{E}_1\{e\}}{\operatorname{plus}(e_1; e_2) = \operatorname{plus}(\mathcal{E}_1; e_2)\{e\}}$$
(5.7b)

$$\frac{e_1 \text{ val} \quad e_2 = \mathcal{E}_2\{e\}}{\text{plus}(e_1; e_2) = \text{plus}(e_1; \mathcal{E}_2)\{e\}}$$
(5.7c)

There is one rule for each form of evaluation context. Filling the hole with *e* results in *e*; otherwise we proceed inductively over the structure of the evaluation context.

Finally, the contextual dynamics for **E** is defined by a single rule:

$$\frac{e = \mathcal{E}\{e_0\} \quad e_0 \longrightarrow e'_0 \quad e' = \mathcal{E}\{e'_0\}}{e \longmapsto e'} \tag{5.8}$$

Thus, a transition from e to e' consists of (1) decomposing e into an evaluation context and an instruction, (2) execution of that instruction, and (3) replacing the instruction by the result of its execution in the same spot within e to obtain e'.

The structural and contextual dynamics define the same transition relation. For the sake of the proof, let us write $e \mapsto e'$ for the transition relation defined by the structural dynamics (rules (5.4)), and $e \mapsto e'$ for the transition relation defined by the contextual dynamics (rules (5.8)).

Theorem 5.4. $e \mapsto_{str} e'$ if, and only if, $e \mapsto_{ctx} e'$.

Proof. From left to right, proceed by rule induction on rules (5.4). It is enough in each case to exhibit an evaluation context \mathcal{E} such that $e = \mathcal{E}\{e_0\}$, $e' = \mathcal{E}\{e'_0\}$, and $e_0 \longrightarrow e'_0$. For example, for rule (5.4a), take $\mathcal{E} = \circ$, and note that $e \longrightarrow e'$. For rule (5.4b), we have by induction that there exists an evaluation context \mathcal{E}_1 such that $e_1 = \mathcal{E}_1\{e_0\}$, $e'_1 = \mathcal{E}_1\{e'_0\}$, and $e_0 \longrightarrow e'_0$. Take $\mathcal{E} = \text{plus}(\mathcal{E}_1; e_2)$, and note that $e = \text{plus}(\mathcal{E}_1; e_2)\{e_0\}$ and $e' = \text{plus}(\mathcal{E}_1; e_2)\{e'_0\}$ with $e_0 \longrightarrow e'_0$.

an evaluation context \mathcal{E}_1 such that $e_1 = \mathcal{E}_1\{e_0\}$, $e'_1 = \mathcal{E}_1\{e'_0\}$, and $e_0 \longrightarrow e'_0$. Take $\mathcal{E} = \mathtt{plus}(\mathcal{E}_1; e_2)$, and note that $e = \mathtt{plus}(\mathcal{E}_1; e_2)\{e_0\}$ and $e' = \mathtt{plus}(\mathcal{E}_1; e_2)\{e'_0\}$ with $e_0 \longrightarrow e'_0$. From right to left, note that if $e \longmapsto_{\mathsf{ctx}} e'$, then there exists an evaluation context \mathcal{E} such that $e = \mathcal{E}\{e_0\}$, $e' = \mathcal{E}\{e'_0\}$, and $e_0 \longrightarrow e'_0$. We prove by induction on rules (5.7) that $e \longmapsto_{\mathsf{str}} e'$. For example, for rule (5.7a), e_0 is e, e'_0 is e', and $e \longrightarrow e'$. Hence $e \longmapsto_{\mathsf{str}} e'$. For rule (5.7b), we have that $\mathcal{E} = \mathtt{plus}(\mathcal{E}_1; e_2)$, $e_1 = \mathcal{E}_1\{e_0\}$, $e'_1 = \mathcal{E}_1\{e'_0\}$, and $e_1 \longmapsto_{\mathsf{str}} e'_1$. Therefore e is $\mathtt{plus}(e_1; e_2)$, e' is $\mathtt{plus}(e'_1; e_2)$, and therefore by rule (5.4b), $e \mapsto_{\mathsf{str}} e'$.

Because the two transition judgments coincide, contextual dynamics can be considered an alternative presentation of a structural dynamics. It has two advantages over structural dynamics, one relatively superficial, one rather less so. The superficial advantage stems from writing rule (5.8) in the simpler form

$$\frac{e_0 \longrightarrow e'_0}{\mathcal{E}\{e_0\} \longmapsto \mathcal{E}\{e'_0\}} \ . \tag{5.9}$$

This formulation is superficially simpler in that it does not make explicit how an expression is decomposed into an evaluation context and a reducible expression. The deeper advantage of contextual dynamics is that all transitions are between complete programs. One need never consider a transition between expressions of any type other than the observable type, which simplifies certain arguments, such as the proof of Lemma 47.16.

5.4 Equational Dynamics

Another formulation of the dynamics of a language regards computation as a form of equational deduction, much in the style of elementary algebra. For example, in algebra we may show that the polynomials $x^2 + 2x + 1$ and $(x + 1)^2$ are equivalent by a simple process of calculation and re-organization using the familiar laws of addition and multiplication. The same laws are enough to determine the value of any polynomial, given the values of its variables. So, for example, we may plug in 2 for x in the polynomial $x^2 + 2x + 1$ and calculate that $2^2 + 2 \times 2 + 1 = 9$, which is indeed $(2 + 1)^2$. We thus obtain a model of computation in which the value of a polynomial for a given value of its variable is determined by substitution and simplification.

Very similar ideas give rise to the concept of *definitional*, or *computational*, *equivalence* of expressions in **E**, which we write as $\mathcal{X} \mid \Gamma \vdash e \equiv e' : \tau$, where Γ consists of one assumption of the form $x : \tau$ for each $x \in \mathcal{X}$. We only consider definitional equality of well-typed expressions, so that when considering the judgment $\Gamma \vdash e \equiv e' : \tau$, we tacitly assume that $\Gamma \vdash e : \tau$ and $\Gamma \vdash e' : \tau$. Here, as usual, we omit explicit mention of the variables \mathcal{X} when they can be determined from the forms of the assumptions Γ .

Definitional equality of expressions in **E** under the by-name interpretation of let is inductively defined by the following rules:

$$\frac{\Gamma \vdash e : \tau}{\Gamma \vdash e \equiv e : \tau} \tag{5.10a}$$

$$\frac{\Gamma \vdash e' \equiv e : \tau}{\Gamma \vdash e \equiv e' : \tau} \tag{5.10b}$$

$$\frac{\Gamma \vdash e \equiv e' : \tau \quad \Gamma \vdash e' \equiv e'' : \tau}{\Gamma \vdash e \equiv e'' : \tau}$$
(5.10c)

$$\frac{\Gamma \vdash e_1 \equiv e_1' : \text{num} \quad \Gamma \vdash e_2 \equiv e_2' : \text{num}}{\Gamma \vdash \text{plus}(e_1; e_2) \equiv \text{plus}(e_1'; e_2') : \text{num}}$$
(5.10d)

$$\frac{\Gamma \vdash e_1 \equiv e'_1 : \operatorname{str} \quad \Gamma \vdash e_2 \equiv e'_2 : \operatorname{str}}{\Gamma \vdash \operatorname{cat}(e_1; e_2) \equiv \operatorname{cat}(e'_1; e'_2) : \operatorname{str}}$$
(5.10e)

$$\frac{\Gamma \vdash e_1 \equiv e_1' : \tau_1 \quad \Gamma, x : \tau_1 \vdash e_2 \equiv e_2' : \tau_2}{\Gamma \vdash \mathsf{let}(e_1; x.e_2) \equiv \mathsf{let}(e_1'; x.e_2') : \tau_2}$$
(5.10f)

$$\frac{n_1 + n_2 = n}{\Gamma \vdash \mathtt{plus}(\mathtt{num}[n_1]; \mathtt{num}[n_2]) \equiv \mathtt{num}[n] : \mathtt{num}}$$
 (5.10g)

$$\frac{s_1 \hat{s}_2 = s}{\Gamma \vdash \mathsf{cat}(\mathsf{str}[s_1]; \mathsf{str}[s_2]) \equiv \mathsf{str}[s] : \mathsf{str}}$$
(5.10h)

$$\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma, x : \tau_1 \vdash e : \tau_2}{\Gamma \vdash \mathsf{let}(e_1; x.e_2) \equiv [e_1/x]e_2 : \tau_2}$$
(5.10i)

Rules (5.10a) through (5.10c) state that definitional equality is an *equivalence relation*. Rules (5.10d) through (5.10f) state that it is a *congruence relation*, which means that it is compatible with all expression-forming constructs in the language. Rules (5.10g) through (5.10i) specify the meanings of the primitive constructs of **E**. We say that rules (5.10) define the *strongest congruence* closed under rules (5.10g), (5.10h), and (5.10i).

Rules (5.10) suffice to calculate the value of an expression by a deduction similar to that used in high school algebra. For example, we may derive the equation

let
$$x$$
 be $1 + 2$ in $x + 3 + 4 \equiv 10$: num

by applying rules (5.10). Here, as in general, there may be many different ways to derive the same equation, but we need find only one derivation in order to carry out an evaluation.

Definitional equality is rather weak in that many equivalences that we might intuitively think are true are not derivable from rules (5.10). A prototypical example is the putative equivalence

$$x_1 : \text{num}, x_2 : \text{num} \vdash x_1 + x_2 \equiv x_2 + x_1 : \text{num},$$
 (5.11)

which, intuitively, expresses the commutativity of addition. Although we shall not prove this here, this equivalence is *not* derivable from rules (5.10). And yet we *may* derive all of its closed instances,

$$n_1 + n_2 \equiv n_2 + n_1 : \text{num},$$
 (5.12)

48 5.5 Notes

where n_1 nat and n_2 nat are particular numbers.

The "gap" between a general law, such as Equation (5.11), and all of its instances, given by Equation (5.12), may be filled by enriching the notion of equivalence to include a principle of proof by mathematical induction. Such a notion of equivalence is sometimes called *semantic equivalence*, because it expresses relationships that hold by virtue of the dynamics of the expressions involved. (Semantic equivalence is developed rigorously for a related language in Chapter 46.)

Theorem 5.5. For the expression language **E**, the relation $e \equiv e' : \tau$ holds iff there exists e_0 val such that $e \mapsto^* e_0$ and $e' \mapsto^* e_0$.

Proof. The proof from right to left is direct, because every transition step is a valid equation. The converse follows from the following, more general, proposition, which is proved by induction on rules (5.10): if $x_1 : \tau_1, \ldots, x_n : \tau_n \vdash e \equiv e' : \tau$, then when $e_1 : \tau_1, e'_1 : \tau_1, \ldots, e_n : \tau_n, e'_n : \tau_n$, if for each $1 \le i \le n$ the expressions e_i and e'_i evaluate to a common value v_i , then there exists e_0 val such that

$$[e_1,\ldots,e_n/x_1,\ldots,x_n]e \mapsto^* e_0$$

and

$$[e'_1,\ldots,e'_n/x_1,\ldots,x_n]e' \mapsto^* e_0$$

5.5 Notes

The use of transition systems to specify the behavior of programs goes back to the early work of Church and Turing on computability. Turing's approach emphasized the concept of an abstract machine consisting of a finite program together with unbounded memory. Computation proceeds by changing the memory in accordance with the instructions in the program. Much early work on the operational semantics of programming languages, such as the SECD machine (Landin, 1965), emphasized machine models. Church's approach emphasized the language for expressing computations, and defined execution in terms of the programs themselves, rather than in terms of auxiliary concepts such as memories or tapes. Plotkin's elegant formulation of structural operational semantics (Plotkin, 1981), which we use heavily throughout this book, was inspired by Church's and Landin's ideas (Plotkin, 2004). Contextual semantics, which was introduced by Felleisen and Hieb (1992), may be seen as an alternative formulation of structural semantics in which "search rules" are replaced by "context matching". Computation viewed as equational deduction goes back to the early work of Herbrand, Gödel, and Church.

Exercises

5.1. Prove that if $s \mapsto^* s'$ and $s' \mapsto^* s''$, then $s \mapsto^* s''$.

5.2. Complete the proof of Theorem 5.1 along the lines suggested there.

5.5 Notes 49

- 5.3. Complete the proof of Theorem 5.5 along the lines suggested there.
- **5.4.** Prove that if $\Gamma \vdash e \equiv e' : \tau$ according to Rules (5.10), then $\Gamma \vdash e : \tau$ and $\Gamma \vdash e' : \tau$ according to Rules (4.1).



5.5 Notes



Chapter 6

Type Safety

Most programming languages are *safe* (or, *type safe*, or *strongly typed*). Informally, this means that certain kinds of mismatches cannot arise during execution. For example, type safety for **E** states that it will never arise that a number is added to a string, or that two numbers are concatenated, neither of which is meaningful.

In general type safety expresses the coherence between the statics and the dynamics. The statics may be seen as predicting that the value of an expression will have a certain form so that the dynamics of that expression is well-defined. Consequently, evaluation cannot "get stuck" in a state for which no transition is possible, corresponding in implementation terms to the absence of "illegal instruction" errors at execution time. Safety is proved by showing that each step of transition preserves typability and by showing that typable states are well-defined. Consequently, evaluation can never "go off into the weeds," and hence can never encounter an illegal instruction.

Type safety for the language **E** is stated precisely as follows:

Theorem 6.1 (Type Safety).

- 1. If $e: \tau$ and $e \mapsto e'$, then $e': \tau$.
- 2. If $e:\tau$, then either e val, or there exists e' such that $e\longmapsto e'$.

The first part, called *preservation*, says that the steps of evaluation preserve typing; the second, called *progress*, ensures that well-typed expressions are either values or can be further evaluated. Safety is the conjunction of preservation and progress.

We say that an expression e is stuck iff it is not a value, yet there is no e' such that $e \mapsto e'$. It follows from the safety theorem that a stuck state is necessarily ill-typed. Or, putting it the other way around, that well-typed states do not get stuck.

6.1 Preservation

The preservation theorem for **E** defined in Chapters 4 and 5 is proved by rule induction on the transition system (rules (5.4)).

52 6.2 Progress

Theorem 6.2 (Preservation). *If* $e : \tau$ *and* $e \mapsto e'$, *then* $e' : \tau$.

Proof. We will give the proof in two cases, leaving the rest to the reader. Consider rule (5.4b),

$$\frac{e_1 \longmapsto e_1'}{\mathtt{plus}(e_1; e_2) \longmapsto \mathtt{plus}(e_1'; e_2)} .$$

Assume that $plus(e_1;e_2):\tau$. By inversion for typing, we have that $\tau=\text{num},e_1:\text{num},$ and $e_2:\text{num}.$ By induction we have that $e_1':\text{num},$ and hence $plus(e_1';e_2):\text{num}.$ The case for concatenation is handled similarly.

Now consider rule (5.4h),

$$\frac{1}{\text{let}(e_1; x.e_2) \longmapsto [e_1/x]e_2} .$$

Assume that $let(e_1; x.e_2) : \tau_2$. By the inversion lemma 4.2, $e_1 : \tau_1$ for some τ_1 such that $x : \tau_1 \vdash e_2 : \tau_2$. By the substitution lemma 4.4 $[e_1/x]e_2 : \tau_2$, as desired.

It is easy to check that the primitive operations are all type-preserving; for example, if a nat and b nat and a + b = c, then c nat.

The proof of preservation is naturally structured as an induction on the transition judgment, because the argument hinges on examining all possible transitions from a given expression. In some cases we may manage to carry out a proof by structural induction on e, or by an induction on typing, but experience shows that this often leads to awkward arguments, or, sometimes, cannot be made to work at all.

6.2 Progress

The progress theorem captures the idea that well-typed programs cannot "get stuck". The proof depends crucially on the following lemma, which characterizes the values of each type.

Lemma 6.3 (Canonical Forms). *If* e *val* and e : τ , then

- 1. If $\tau = \text{num}$, then e = num[n] for some number n.
- 2. If $\tau = \operatorname{str}$, then $e = \operatorname{str}[s]$ for some string s.

Proof. By induction on rules (4.1) and (5.3).

Progress is proved by rule induction on rules (4.1) defining the statics of the language.

Theorem 6.4 (Progress). *If* $e : \tau$, then either e val, or there exists e' such that $e \mapsto e'$.

Proof. The proof proceeds by induction on the typing derivation. We will consider only one case, for rule (4.1d),

$$\frac{e_1: \mathtt{num} \quad e_2: \mathtt{num}}{\mathtt{plus}(e_1; e_2): \mathtt{num}} ,$$

6.3 Run-Time Errors 53

where the context is empty because we are considering only closed terms.

By induction we have that either e_1 val, or there exists e_1' such that $e_1 \mapsto e_1'$. In the latter case it follows that $\mathtt{plus}(e_1;e_2) \mapsto \mathtt{plus}(e_1';e_2)$, as required. In the former we also have by induction that either e_2 val, or there exists e_2' such that $e_2 \mapsto e_2'$. In the latter case we have that $\mathtt{plus}(e_1;e_2) \mapsto \mathtt{plus}(e_1;e_2')$, as required. In the former, we have, by the Canonical Forms Lemma 6.3, $e_1 = \mathtt{num}[n_1]$ and $e_2 = \mathtt{num}[n_2]$, and hence

$$plus(num[n_1]; num[n_2]) \longrightarrow num[n_1 + n_2].$$

Because the typing rules for expressions are syntax-directed, the progress theorem could equally well be proved by induction on the structure of *e*, appealing to the inversion theorem at each step to characterize the types of the parts of *e*. But this approach breaks down when the typing rules are not syntax-directed, that is, when there is more than one rule for a given expression form. Such rules present no difficulites, so long as the proof proceeds by induction on the typing rules, and not on the structure of the expression.

Summing up, the combination of preservation and progress together constitute the proof of safety. The progress theorem ensures that well-typed expressions do not "get stuck" in an ill-defined state, and the preservation theorem ensures that if a step is taken, the result remains well-typed (with the same type). Thus the two parts work together to ensure that the statics and dynamics are coherent, and that no ill-defined states can ever be encountered while evaluating a well-typed expression.

6.3 Run-Time Errors

Suppose that we wish to extend **E** with, say, a quotient operation that is undefined for a zero divisor. The natural typing rule for quotients is given by the following rule:

$$\frac{e_1: \mathtt{num} \quad e_2: \mathtt{num}}{\mathtt{div}(e_1; e_2): \mathtt{num}}$$

But the expression div(num[3]; num[0]) is well-typed, yet stuck! We have two options to correct this situation:

- 1. Enhance the type system, so that no well-typed program may divide by zero.
- 2. Add dynamic checks, so that division by zero signals an error as the outcome of evaluation.

Either option is, in principle, practical, but the most common approach is the second. The first requires that the type checker prove that an expression be non-zero before permitting it to be used in the denominator of a quotient. It is difficult to do this without ruling out too many programs as ill-formed. We cannot predict statically whether an expression will be non-zero when evaluated, so the second approach is most often used in practice.

The overall idea is to distinguish *checked* from *unchecked* errors. An unchecked error is one that is ruled out by the type system. No run-time checking is performed to ensure that such an

54 6.4 Notes

error does not occur, because the type system rules out the possibility of it arising. For example, the dynamics need not check, when performing an addition, that its two arguments are, in fact, numbers, as opposed to strings, because the type system ensures that this is the case. On the other hand the dynamics for quotient *must* check for a zero divisor, because the type system does not rule out the possibility.

One approach to modeling checked errors is to give an inductive definition of the judgment e err stating that the expression e incurs a checked run-time error, such as division by zero. Here are some representative rules that would be present in a full inductive definition of this judgment:

$$\frac{e_1 \text{ val}}{\text{div}(e_1; \text{num}[0]) \text{ err}}$$
 (6.1a)

$$\frac{e_1 \text{ err}}{\text{div}(e_1; e_2) \text{ err}} \tag{6.1b}$$

$$\frac{e_1 \text{ val} \quad e_2 \text{ err}}{\text{div}(e_1; e_2) \text{ err}} \tag{6.1c}$$

Rule (6.1a) signals an error condition for division by zero. The other rules propagate this error upwards: if an evaluated sub-expression is a checked error, then so is the overall expression.

Once the error judgment is available, we may also consider an expression, error, which forcibly induces an error, with the following statics and dynamics:

$$\frac{}{\Gamma \vdash \mathsf{error} : \tau} \tag{6.2a}$$

The preservation theorem is not affected by checked errors. However, the statement (and proof) of progress is modified to account for checked errors.

Theorem 6.5 (Progress With Error). *If* $e : \tau$, then either e err, or e val, or there exists e' such that $e \mapsto e'$.

Proof. The proof is by induction on typing, and proceeds similarly to the proof given earlier, except that there are now three cases to consider at each point in the proof. \Box

6.4 Notes

The concept of type safety was first formulated by Milner (1978), who invented the slogan "well-typed programs do not go wrong." Whereas Milner relied on an explicit notion of "going wrong" to express the concept of a type error, Wright and Felleisen (1994) observed that we can instead show that ill-defined states cannot arise in a well-typed program, giving rise to the slogan "well-typed programs do not get stuck." However, their formulation relied on an analysis showing that no stuck state is well-typed. The progress theorem, which relies on the characterization of canonical forms in the style of Martin-Löf (1980), eliminates this analysis.

6.4 Notes 55

Exercises

- **6.1**. Complete the proof of Theorem 6.2 in full detail.
- **6.2**. Complete the proof of Theorem 6.4 in full detail.
- **6.3**. Give several cases of the proof of Theorem 6.5 to illustrate how checked errors are handled in type safety proofs.



56 6.4 Notes



Chapter 7

Evaluation Dynamics

In Chapter 5 we defined evaluation of expressions in **E** using a structural dynamics. Structural dynamics is very useful for proving safety, but for some purposes, such as writing a user manual, another formulation, called *evaluation dynamics* is preferable. An evaluation dynamics is a relation between a phrase and its value that is defined without detailing the step-by-step process of evaluation. A *cost dynamics* enriches an evaluation dynamics with a *cost measure* specifying the resource usage of evaluation. A prime example is time, measured as the number of transition steps required to evaluate an expression according to its structural dynamics.

7.1 Evaluation Dynamics

An *evaluation dynamics*, consists of an inductive definition of the evaluation judgment $e \Downarrow v$ stating that the closed expression e evaluates to the value v. The evaluation dynamics of \mathbf{E} is defined by the following rules:

$$\overline{\operatorname{num}[n] \Downarrow \operatorname{num}[n]} \tag{7.1a}$$

$$\overline{\operatorname{str}[s] \Downarrow \operatorname{str}[s]} \tag{7.1b}$$

$$\frac{e_1 \Downarrow \operatorname{num}[n_1] \quad e_2 \Downarrow \operatorname{num}[n_2] \quad n_1 + n_2 = n}{\operatorname{plus}(e_1; e_2) \Downarrow \operatorname{num}[n]}$$
(7.1c)

$$\frac{e_1 \Downarrow \operatorname{str}[s_1] \quad e_2 \Downarrow \operatorname{str}[s_2] \quad s_1 \hat{\ } s_2 = s}{\operatorname{cat}(e_1; e_2) \Downarrow \operatorname{str}[s]}$$
(7.1d)

$$\frac{e \Downarrow \text{str}[s] \quad |s| = n}{\text{len}(e) \Downarrow \text{num}[n]}$$
 (7.1e)

$$\frac{[e_1/x]e_2 \Downarrow v_2}{\operatorname{let}(e_1; x.e_2) \Downarrow v_2} \tag{7.1f}$$

The value of a let expression is determined by substitution of the binding into the body. The rules are not syntax-directed, because the premise of rule (7.1f) is not a sub-expression of the expression in the conclusion of that rule.

Rule (7.1f) specifies a by-name interpretation of definitions. For a by-value interpretation the following rule should be used instead:

$$\frac{e_1 \Downarrow v_1 \quad [v_1/x]e_2 \Downarrow v_2}{\operatorname{let}(e_1; x.e_2) \Downarrow v_2} \tag{7.2}$$

Because the evaluation judgment is inductively defined, we prove properties of it by rule induction. Specifically, to show that the property $\mathcal{P}(e \downarrow v)$ holds, it is enough to show that \mathcal{P} is closed under rules (7.1):

- 1. Show that $\mathcal{P}(\text{num}[n] \Downarrow \text{num}[n])$.
- 2. Show that $\mathcal{P}(\mathsf{str}[s] \Downarrow \mathsf{str}[s])$.
- 3. Show that $\mathcal{P}(\mathtt{plus}(e_1; e_2) \Downarrow \mathtt{num}[n])$, if $\mathcal{P}(e_1 \Downarrow \mathtt{num}[n_1])$, $\mathcal{P}(e_2 \Downarrow \mathtt{num}[n_2])$, and $n_1 + n_2 = n$.
- 4. Show that $\mathcal{P}(\mathsf{cat}(e_1; e_2) \Downarrow \mathsf{str}[s])$, if $\mathcal{P}(e_1 \Downarrow \mathsf{str}[s_1])$, $\mathcal{P}(e_2 \Downarrow \mathsf{str}[s_2])$, and $s_1 \hat{s}_2 = s$.
- 5. Show that $\mathcal{P}(\text{let}(e_1; x.e_2) \Downarrow v_2)$, if $\mathcal{P}([e_1/x]e_2 \Downarrow v_2)$.

This induction principle is *not* the same as structural induction on *e* itself, because the evaluation rules are not syntax-directed.

Lemma 7.1. *If* $e \Downarrow v$, then v val.

Proof. By induction on rules (7.1). All cases except rule (7.1f) are immediate. For the latter case, the result follows directly by an appeal to the inductive hypothesis for the premise of the evaluation rule. \Box

7.2 Relating Structural and Evaluation Dynamics

We have given two different forms of dynamics for **E**. It is natural to ask whether they are equivalent, but to do so first requires that we consider carefully what we mean by equivalence. The structural dynamics describes a step-by-step process of execution, whereas the evaluation dynamics suppresses the intermediate states, focusing attention on the initial and final states alone. This remark suggests that the right correspondence is between *complete* execution sequences in the structural dynamics and the evaluation judgment in the evaluation dynamics.

Theorem 7.2. For all closed expressions e and values v, $e \mapsto^* v$ iff $e \Downarrow v$.

How might we prove such a theorem? We will consider each direction separately. We consider the easier case first.

Lemma 7.3. *If* $e \Downarrow v$, then $e \longmapsto^* v$.

Proof. By induction on the definition of the evaluation judgment. For example, suppose that $plus(e_1; e_2) \Downarrow num[n]$ by the rule for evaluating additions. By induction we know that $e_1 \mapsto^* num[n_1]$ and $e_2 \mapsto^* num[n_2]$. We reason as follows:

$$\begin{array}{ccc} \mathtt{plus}(\,e_1;e_2\,) & \longmapsto^* & \mathtt{plus}(\,\mathtt{num}[\,n_1\,];e_2\,) \\ & \longmapsto^* & \mathtt{plus}(\,\mathtt{num}[\,n_1\,];\mathtt{num}[\,n_2\,]\,) \\ & \longmapsto & \mathtt{num}[\,n_1+n_2\,] \end{array}$$

Therefore plus $(e_1; e_2) \mapsto^* \text{num}[n_1 + n_2]$, as required. The other cases are handled similarly. \square

For the converse, recall from Chapter 5 the definitions of multi-step evaluation and complete evaluation. Because $v \downarrow v$ when v val, it suffices to show that evaluation is closed under converse evaluation:¹

Lemma 7.4. *If* $e \mapsto e'$ *and* $e' \Downarrow v$, *then* $e \Downarrow v$.

Proof. By induction on the definition of the transition judgment. For example, suppose that $plus(e_1;e_2) \mapsto plus(e_1';e_2)$, where $e_1 \mapsto e_1'$. Suppose further that $plus(e_1';e_2) \Downarrow v$, so that $e_1' \Downarrow num[n_1]$, and $e_2 \Downarrow num[n_2]$, and $e_1 + e_2 = n$, and $e_2 \Downarrow num[n_2]$, and $e_3 \Downarrow num[n_3]$, and hence $plus(e_1;e_2) \Downarrow num[n]$, as required.

7.3 Type Safety, Revisited

Type safety is defined in Chapter 6 as preservation and progress (Theorem 6.1). These concepts are meaningful when applied to a dynamics given by a transition system, as we shall do throughout this book. But what if we had instead given the dynamics as an evaluation relation? How is type safety proved in that case?

The answer, unfortunately, is that we cannot. Although there is an analog of the preservation property for an evaluation dynamics, there is no clear analog of the progress property. Preservation may be stated as saying that if $e \Downarrow v$ and $e : \tau$, then $v : \tau$. It can be readily proved by induction on the evaluation rules. But what is the analog of progress? We might be tempted to phrase progress as saying that if $e : \tau$, then $e \Downarrow v$ for some v. Although this property is true for \mathbf{E} , it demands much more than just progress — it requires that every expression evaluate to a value! If \mathbf{E} were extended to admit operations that may result in an error (as discussed in Section 6.3), or to admit non-terminating expressions, then this property would fail, even though progress would remain valid.

One possible attitude towards this situation is to conclude that type safety cannot be properly discussed in the context of an evaluation dynamics, but only by reference to a structural dynamics. Another point of view is to instrument the dynamics with explicit checks for dynamic type errors, and to show that any expression with a dynamic type fault must be statically ill-typed. Re-stated in the contrapositive, this means that a statically well-typed program cannot incur a dynamic type error. A difficulty with this point of view is that we must explicitly account for a form of error

¹Converse evaluation is also known as *head expansion*.

60 7.4 Cost Dynamics

solely to prove that it cannot arise! Nevertheless, a semblance of type safety can be established using evaluation dynamics.

We define a judgment *e* ?? stating that the expression *e goes wrong* when executed. The exact definition of "going wrong" is given by a set of rules, but the intention is that it should cover all situations that correspond to type errors. The following rules are representative of the general case:

$$\overline{\mathsf{plus}(\mathsf{str}[s]; e_2)??} \tag{7.3a}$$

$$\frac{e_1 \text{ val}}{\text{plus}(e_1; \text{str}[s])??} \tag{7.3b}$$

These rules explicitly check for the misapplication of addition to a string; similar rules govern each of the primitive constructs of the language.

Theorem 7.5. *If* e ??, then there is no τ such that e : τ .

Proof. By rule induction on rules (7.3). For example, for rule (7.3a), we note that str[s] : str, and hence $plus(str[s]; e_2)$ is ill-typed.

Corollary 7.6. *If* $e : \tau$, then $\neg (e ??)$.

Apart from the inconvenience of having to define the judgment e?? only to show that it is irrelevant for well-typed programs, this approach suffers a very significant methodological weakness. If we should omit one or more rules defining the judgment e??, the proof of Theorem 7.5 remains valid; there is nothing to ensure that we have included sufficiently many checks for run-time type errors. We can prove that the ones we define cannot arise in a well-typed program, but we cannot prove that we have covered all possible cases. By contrast the structural dynamics does not specify any behavior for ill-typed expressions. Consequently, any ill-typed expression will "get stuck" without our explicit intervention, and the progress theorem rules out all such cases. Moreover, the transition system corresponds more closely to implementation—a compiler need not make any provisions for checking for run-time type errors. Instead, it relies on the statics to ensure that these cannot arise, and assigns no meaning to any ill-typed program. Therefore, execution is more efficient, and the language definition is simpler.

7.4 Cost Dynamics

A structural dynamics provides a natural notion of *time complexity* for programs, namely the number of steps required to reach a final state. An evaluation dynamics, however, does not provide such a direct notion of time. Because the individual steps required to complete an evaluation are suppressed, we cannot directly read off the number of steps required to evaluate to a value. Instead we must augment the evaluation relation with a cost measure, resulting in a *cost dynamics*.

Evaluation judgments have the form $e \downarrow^k v$, with the meaning that e evaluates to v in k steps.

$$\frac{1}{\operatorname{num}[n] \downarrow^{0} \operatorname{num}[n]} \tag{7.4a}$$

7.5 Notes 61

$$\frac{e_1 \downarrow^{k_1} \text{num}[n_1] \quad e_2 \downarrow^{k_1} \text{num}[n_2]}{\text{plus}(e_1; e_2) \downarrow^{k_1+k_2+1} \text{num}[n_1+n_2]}$$
(7.4b)

$$\frac{1}{\operatorname{str}[s] \Downarrow^{0} \operatorname{str}[s]} \tag{7.4c}$$

$$\frac{e_1 \downarrow^{k_1} s_1 \quad e_2 \downarrow^{k_2} s_2}{\operatorname{cat}(e_1; e_2) \downarrow^{k_1 + k_2 + 1} \operatorname{str}[s_1 \hat{s}_2]}$$
(7.4d)

$$\frac{[e_1/x]e_2 \downarrow^{k_2} v_2}{\text{let}(e_1; x.e_2) \downarrow^{k_2+1} v_2}$$
 (7.4e)

For a by-value interpretation of let, rule (7.4e) is replaced by the following rule:

$$\frac{e_1 \Downarrow^{k_1} v_1 \quad [v_1/x]e_2 \Downarrow^{k_2} v_2}{\text{let}(e_1; x.e_2) \Downarrow^{k_1+k_2+1} v_2}$$
(7.5)

Theorem 7.7. For any closed expression e and closed value v of the same type, $e \downarrow k$ v iff $e \longmapsto^* kv$.

Proof. From left to right proceed by rule induction on the definition of the cost dynamics. From right to left proceed by induction on k, with an inner rule induction on the definition of the structural dynamics.

7.5 Notes

The structural similarity between evaluation dynamics and typing rules was first developed in *The Definition of Standard ML* (Milner et al., 1997). The advantage of evaluation dynamics is its directness; its disadvantage is that it is not well-suited to proving properties such as type safety. Robin Milner introduced the apt phrase "going wrong" as a description of a type error. Cost dynamics was introduced by Blelloch and Greiner (1996) in a study of parallel computation (see Chapter 37).

Exercises

- **7.1**. Show that evaluation is deterministic: if $e \Downarrow v_1$ and $e \Downarrow v_2$, then $v_1 = v_2$.
- **7.2**. Complete the proof of Lemma 7.3.
- **7.3**. Complete the proof of Lemma 7.4. Then show that if $e \mapsto^* e'$ with e' val, then $e \downarrow e'$.
- **7.4.** Augment the evaluation dynamics with checked errors, along the lines sketched in Chapter 5, using *e* ?? to say that *e* incurs a checked (or an unchecked) error. What remains unsatisfactory about the type safety proof? Can you think of a better alternative?

62 7.5 Notes

7.5. Consider generic hypothetical judgments of the form

$$x_1 \Downarrow v_1, \ldots, x_n \Downarrow v_n \vdash e \Downarrow v$$

where v_1 val,..., v_n val, and v val. The hypotheses, written Δ , are called the *environment* of the evaluation; they provide the values of the free variables in e. The hypothetical judgment $\Delta \vdash e \Downarrow v$ is called an *environmental evaluation dynamics*.

Give a hypothetical inductive definition of the environmental evaluation dynamics *without* making any use of substitution. In particular, you should include the rule

$$\overline{\Delta, x \Downarrow v \vdash x \Downarrow v}$$

defining the evaluation of a free variable.

Show that $x_1 \Downarrow v_1, \ldots, x_n \Downarrow v_n \vdash e \Downarrow v$ iff $[v_1, \ldots, v_n/x_1, \ldots, x_n]e \Downarrow v$ (using the by-value form of evaluation).

Part III Total Functions





Chapter 9

System T of Higher-Order Recursion

System **T**, well-known as *Gödel's T*, is the combination of function types with the type of natural numbers. In contrast to **E**, which equips the naturals with some arbitrarily chosen arithmetic operations, the language **T** provides a general mechanism, called *primitive recursion*, from which these primitives may be defined. Primitive recursion captures the essential inductive character of the natural numbers, and hence may be seen as an intrinsic termination proof for each program in the language. Consequently, we may only define *total* functions in the language, those that always return a value for each argument. In essence every program in **T** "comes equipped" with a proof of its termination. Although this may seem like a shield against infinite loops, it is also a weapon that can be used to show that some programs cannot be written in **T**. To do so would demand a master termination proof for every possible program in the language, something that we shall prove does not exist.

9.1 Statics

The syntax of **T** is given by the following grammar:

We write \overline{n} for the expression $s(\dots s(z))$, in which the successor is applied $n \ge 0$ times to zero. The expression $rec\{e_0; x.y.e_1\}(e)$ is called the *recursor*. It represents the *e*-fold iteration of the

74 9.2 Dynamics

transformation $x.y.e_1$ starting from e_0 . The bound variable x represents the predecessor and the bound variable y represents the result of the x-fold iteration. The "with" clause in the concrete syntax for the recursor binds the variable y to the result of the recursive call, as will become clear shortly.

Sometimes the *iterator*, $iter\{e_0; y.e_1\}(e)$, is considered as an alternative to the recursor. It has essentially the same meaning as the recursor, except that only the result of the recursive call is bound to y in e_1 , and no binding is made for the predecessor. Clearly the iterator is a special case of the recursor, because we can always ignore the predecessor binding. Conversely, the recursor is definable from the iterator, provided that we have product types (Chapter 10) at our disposal. To define the recursor from the iterator, we simultaneously compute the predecessor while iterating the specified computation.

The statics of **T** is given by the following typing rules:

$$\frac{\Gamma, x : \tau \vdash x : \tau}{\Gamma, x : \tau} \tag{9.1a}$$

$$\Gamma \vdash \sigma \cdot \mathsf{nat}$$
 (9.1b)

$$\frac{\Gamma \vdash e : \text{nat}}{\Gamma \vdash s(e) : \text{nat}} \tag{9.1c}$$

$$\frac{\Gamma \vdash e : \mathtt{nat} \quad \Gamma \vdash e_0 : \tau \quad \Gamma, x : \mathtt{nat}, y : \tau \vdash e_1 : \tau}{\Gamma \vdash \mathtt{rec}\{e_0; x.y.e_1\}(e) : \tau} \tag{9.1d}$$

$$\frac{\Gamma, x : \tau_1 \vdash e : \tau_2}{\Gamma \vdash \lambda \{\tau_1\}(x.e) : \rightarrow (\tau_1; \tau_2)}$$

$$(9.1e)$$

$$\frac{\Gamma \vdash e_1 : \rightarrow (\tau_2; \tau) \quad \Gamma \vdash e_2 : \tau_2}{\Gamma \vdash \operatorname{ap}(e_1; e_2) : \tau}$$
(9.1f)

As usual, admissibility of the structural rule of substitution is crucially important.

Lemma 9.1. *If* $\Gamma \vdash e : \tau$ *and* $\Gamma, x : \tau \vdash e' : \tau'$, *then* $\Gamma \vdash [e/x]e' : \tau'$.

9.2 Dynamics

The closed values of **T** are defined by the following rules:

$$\frac{}{z \text{ val}}$$
 (9.2a)

$$\frac{[e \text{ val}]}{s(e) \text{ val}} \tag{9.2b}$$

$$\frac{1}{\lambda\{\tau\}(x.e) \text{ val}} \tag{9.2c}$$

9.2 Dynamics 75

The premise of rule (9.2b) is included for an *eager* interpretation of successor, and excluded for a *lazy* interpretation.

The transition rules for the dynamics of **T** are as follows:

$$\left[\frac{e \longmapsto e'}{s(e) \longmapsto s(e')}\right] \tag{9.3a}$$

$$\frac{e_1 \longmapsto e'_1}{\operatorname{ap}(e_1; e_2) \longmapsto \operatorname{ap}(e'_1; e_2)} \tag{9.3b}$$

$$\left[\frac{e_1 \text{ val} \quad e_2 \longmapsto e_2'}{\operatorname{ap}(e_1; e_2) \longmapsto \operatorname{ap}(e_1; e_2')} \right]$$
(9.3c)

$$\frac{[e_2 \text{ val}]}{\operatorname{ap}(\lambda\{\tau\}(x.e); e_2) \longmapsto [e_2/x]e} \tag{9.3d}$$

$$\frac{e \mapsto e'}{\operatorname{rec}\{e_0; x.y.e_1\}(e') \mapsto \operatorname{rec}\{e_0; x.y.e_1\}(e')}$$
(9.3e)

$$\frac{}{\operatorname{rec}\{e_0; x.y.e_1\}(\mathbf{z}) \longmapsto e_0} \tag{9.3f}$$

$$\frac{\mathbf{s}(e) \text{ val}}{\operatorname{rec}\{e_0; x.y.e_1\}(\mathbf{s}(e)) \longmapsto [e, \operatorname{rec}\{e_0; x.y.e_1\}(e)/x, y]e_1}$$
(9.3g)

The bracketed rules and premises are included for an eager successor and call-by-value application, and omitted for a lazy successor and call-by-name application. Rules (9.3f) and (9.3g) specify the behavior of the recursor on z and s(e). In the former case the recursor reduces to e_0 , and in the latter case the variable x is bound to the predecessor e and y is bound to the (unevaluated) recursion on e. If the value of y is not required in the rest of the computation, the recursive call is not evaluated.

Lemma 9.2 (Canonical Forms). *If* $e : \tau$ *and* e *val, then*

- 1. If $\tau = \text{nat}$, then either e = z or e = s(e') for some e'.
- 2. If $\tau = \tau_1 \rightarrow \tau_2$, then $e = \lambda (x : \tau_1) e_2$ for some e_2 .

Theorem 9.3 (Safety). 1. If $e : \tau$ and $e \mapsto e'$, then $e' : \tau$.

2. If $e: \tau$, then either e val or $e \mapsto e'$ for some e'.

76 9.3 Definability

9.3 Definability

A mathematical function $f : \mathbb{N} \to \mathbb{N}$ on the natural numbers is *definable* in **T** iff there exists an expression e_f of type nat \to nat such that for every $n \in \mathbb{N}$,

$$e_f(\overline{n}) \equiv \overline{f(n)}$$
: nat. (9.4)

That is, the numeric function $f: \mathbb{N} \to \mathbb{N}$ is definable iff there is an expression e_f of type $\mathtt{nat} \to \mathtt{nat}$ such that, when applied to the numeral representing the argument $n \in \mathbb{N}$, the application is definitionally equal to the numeral corresponding to $f(n) \in \mathbb{N}$.

Definitional equality for **T**, written $\Gamma \vdash e \equiv e' : \tau$, is the strongest congruence containing these axioms:

$$\frac{\Gamma, x : \tau_1 \vdash e_2 : \tau_2 \quad \Gamma \vdash e_1 : \tau_1}{\Gamma \vdash \operatorname{ap}(\lambda \{\tau_1\}(x.e_2); e_1) \equiv [e_1/x]e_2 : \tau_2}$$

$$(9.5a)$$

$$\frac{\Gamma \vdash e_0 : \tau \quad \Gamma, x : \tau \vdash e_1 : \tau}{\Gamma \vdash \operatorname{rec}\{e_0; x.y.e_1\}(z) \equiv e_0 : \tau}$$

$$(9.5b)$$

$$\frac{\Gamma \vdash e_0 : \tau \quad \Gamma, x : \tau \vdash e_1 : \tau}{\Gamma \vdash \operatorname{rec}\{e_0; x.y.e_1\}(s(e)) \equiv [e, \operatorname{rec}\{e_0; x.y.e_1\}(e)/x, y]e_1 : \tau}$$
(9.5c)

For example, the doubling function, $d(n) = 2 \times n$, is definable in **T** by the expression e_d : nat \rightarrow nat given by

$$\lambda(x: \mathtt{nat}) \mathtt{rec} x \{z \hookrightarrow z \mid \mathtt{s}(u) \mathtt{with} v \hookrightarrow \mathtt{s}(\mathtt{s}(v))\}.$$

To check that this defines the doubling function, we proceed by induction on $n \in \mathbb{N}$. For the basis, it is easy to check that

$$e_d(\overline{0}) \equiv \overline{0}$$
: nat.

For the induction, assume that

$$e_d(\overline{n}) \equiv \overline{d(n)}$$
: nat.

Then calculate using the rules of definitional equality:

$$\begin{split} e_d(\,\overline{n+1}\,) &\equiv \mathtt{s}(\,\mathtt{s}(\,e_d(\,\overline{n}\,)\,)\,) \\ &\equiv \mathtt{s}(\,\mathtt{s}(\,\overline{2\times n}\,)\,) \\ &= \overline{2\times (n+1)} \\ &= \overline{d(n+1)}. \end{split}$$

As another example, consider the following function, called *Ackermann's function*, defined by the following equations:

$$A(0,n) = n+1$$

$$A(m+1,0) = A(m,1)$$

$$A(m+1,n+1) = A(m,A(m+1,n)).$$

9.4 Undefinability 77

The Ackermann function grows very quickly. For example, $A(4,2) \approx 2^{65,536}$, which is often cited as being larger than the number of atoms in the universe! Yet we can show that the Ackermann function is total by a lexicographic induction on the pair of arguments (m,n). On each recursive call, either m decreases, or else m remains the same, and n decreases, so inductively the recursive calls are well-defined, and hence so is A(m,n).

A first-order primitive recursive function is a function of type $\mathtt{nat} \to \mathtt{nat}$ that is defined using the recursor, but without using any higher order functions. Ackermann's function is defined so that it is not first-order primitive recursive, but is higher-order primitive recursive. The key to showing that it is definable in **T** is to note that A(m+1,n) iterates n times the function A(m,-), starting with A(m,1). As an auxiliary, let us define the higher-order function

$$\mathtt{it}: (\mathtt{nat} \to \mathtt{nat}) \to \mathtt{nat} \to \mathtt{nat} \to \mathtt{nat}$$

to be the λ -abstraction

$$\lambda (f: \mathtt{nat} \to \mathtt{nat}) \lambda (n: \mathtt{nat}) \mathtt{rec} n \{z \hookrightarrow \mathtt{id} \mid \mathtt{s}(_) \mathtt{with} g \hookrightarrow f \circ g\},$$

where $id = \lambda(x:nat)x$ is the identity, and $f \circ g = \lambda(x:nat)f(g(x))$ is the composition of f and g. It is easy to check that

$$it(f)(\overline{n})(\overline{m}) \equiv f^{(n)}(\overline{m}) : nat,$$

where the latter expression is the n-fold composition of f starting with \overline{m} . We may then define the Ackermann function

$$e_a: \mathtt{nat} \to \mathtt{nat} \to \mathtt{nat}$$

to be the expression

$$\lambda \, (\, m \, : \, \mathtt{nat} \,) \, \mathtt{rec} \, m \, \big\{ \mathtt{z} \, \hookrightarrow \, \mathtt{s} \, \big| \, \mathtt{s}(\, _) \, \mathtt{with} \, f \, \hookrightarrow \, \lambda \, \big(\, n \, : \, \mathtt{nat} \, \big) \, \mathtt{it}(\, f \, \big)(\, n \, \big)(\, f(\, \overline{1} \,) \, \big) \big\}.$$

It is instructive to check that the following equivalences are valid:

$$e_{a}(\overline{0})(\overline{n}) \equiv s(\overline{n}) \tag{9.6}$$

$$e_a(\overline{m+1})(\overline{0}) \equiv e_a(\overline{m})(\overline{1}) \tag{9.7}$$

$$e_a(\overline{m+1})(\overline{n+1}) \equiv e_a(\overline{m})(e_a(s(\overline{m}))(\overline{n})). \tag{9.8}$$

That is, the Ackermann function is definable in **T**.

9.4 Undefinability

It is impossible to define an infinite loop in **T**.

Theorem 9.4. *If* $e : \tau$, then there exists v val such that $e \equiv v : \tau$.

78 9.4 Undefinability

Consequently, values of function type in **T** behave like mathematical functions: if $e: \tau_1 \to \tau_2$ and $e_1: \tau_1$, then $e(e_1)$ evaluates to a value of type τ_2 . Moreover, if e: nat, then there exists a natural number n such that $e \equiv \overline{n}:$ nat.

Using this, we can show, using a technique called *diagonalization*, that there are functions on the natural numbers that are not definable in \mathbf{T} . We make use of a technique, called *Gödel-numbering*, that assigns a unique natural number to each closed expression of \mathbf{T} . By assigning a unique number to each expression, we may manipulate expressions as data values in \mathbf{T} so that \mathbf{T} is able to compute with its own programs.¹

The essence of Gödel-numbering is captured by the following simple construction on abstract syntax trees. (The generalization to abstract binding trees is slightly more difficult, the main complication being to ensure that all α -equivalent expressions are assigned the same Gödel number.) Recall that a general ast a has the form $o(a_1, \ldots, a_k)$, where o is an operator of arity k. Enumerate the operators so that every operator has an index $i \in \mathbb{N}$, and let m be the index of o in this enumeration. Define the $G\"{o}del$ number $^{\Gamma}a^{\Gamma}$ of a to be the number

$$2^m 3^{n_1} 5^{n_2} \dots p_k^{n_k}$$

where p_k is the kth prime number (so that $p_0 = 2$, $p_1 = 3$, and so on), and n_1, \ldots, n_k are the Gödel numbers of a_1, \ldots, a_k , respectively. This procedure assigns a natural number to each ast. Conversely, given a natural number, n, we may apply the prime factorization theorem to "parse" n as a unique abstract syntax tree. (If the factorization is not of the right form, which can only be because the arity of the operator does not match the number of factors, then n does not code any ast.)

Now, using this representation, we may define a (mathematical) function $f_{univ}: \mathbb{N} \to \mathbb{N} \to \mathbb{N}$ such that, for any $e: \mathtt{nat} \to \mathtt{nat}$, $f_{univ}(\lceil e \rceil)(m) = n$ iff $e(\overline{m}) \equiv \overline{n}: \mathtt{nat}.^2$ The determinacy of the dynamics, together with Theorem 9.4, ensure that f_{univ} is a well-defined function. It is called the *universal function* for \mathbf{T} because it specifies the behavior of any expression e of type $\mathtt{nat} \to \mathtt{nat}$. Using the universal function, let us define an auxiliary mathematical function, called the *diagonal function* $\delta: \mathbb{N} \to \mathbb{N}$, by the equation $\delta(m) = f_{univ}(m)(m)$. The δ function is chosen so that $\delta(\lceil e \rceil) = n$ iff $e(\lceil e \rceil) \equiv \overline{n}: \mathtt{nat}$. (The motivation for its definition will become clear in a moment.)

The function f_{univ} is not definable in **T**. Suppose that it were definable by the expression e_{univ} , then the diagonal function δ would be definable by the expression

$$e_{\delta} = \lambda \left(m : \mathtt{nat} \right) e_{univ}(m)(m).$$

But in that case we would have the equations

$$e_{\delta}(\lceil e \rceil) \equiv e_{univ}(\lceil e \rceil)(\lceil e \rceil)$$
$$\equiv e(\lceil e \rceil).$$

Now let e_{Δ} be the function expression

$$\lambda(x: \mathtt{nat}) s(e_{\delta}(x)),$$

¹The same technique lies at the heart of the proof of Gödel's celebrated incompleteness theorem. The undefinability of certain functions on the natural numbers within \mathbf{T} may be seen as a form of incompleteness like that considered by Gödel. ²The value of $f_{univ}(k)(m)$ may be chosen arbitrarily to be zero when k is not the code of any expression e.

9.5 Notes 79

so that we may deduce

$$e_{\Delta}(\lceil \overline{e_{\Delta}} \rceil) \equiv s(e_{\delta}(\lceil \overline{e_{\Delta}} \rceil))$$
$$\equiv s(e_{\Delta}(\lceil \overline{e_{\Delta}} \rceil)).$$

But the termination theorem implies that there exists n such that $e_{\Delta}(\lceil \overline{e_{\Delta}} \rceil) \equiv \overline{n}$, and hence we have $\overline{n} \equiv \mathfrak{s}(\overline{n})$, which is impossible.

We say that a language \mathcal{L} is *universal* if it is possible to write an interpreter for \mathcal{L} in \mathcal{L} itself. It is intuitively clear that f_{univ} is computable in the sense that we can define it in *some* sufficiently powerful programming language. But the preceding argument shows that \mathbf{T} is not up to the task; it is not a universal programming language. Examination of the foregoing proof reveals an inescapable tradeoff: by insisting that all expressions terminate, we necessarily lose universality—there are computable functions that are not definable in the language.

9.5 Notes

System **T** was introduced by Gödel in his study of the consistency of arithmetic (Gödel, 1980). He showed how to "compile" proofs in arithmetic into well-typed terms of system **T**, and to reduce the consistency problem for arithmetic to the termination of programs in **T**. It was perhaps the first programming language whose design was directly influenced by the verification (of termination) of its programs.

Exercises

- **9.1**. Prove Lemma 9.2.
- 9.2. Prove Theorem 9.3.
- **9.3.** Attempt to prove that if e: nat is closed, then there exists n such that $e \mapsto^* \overline{n}$ under the eager dynamics. Where does the proof break down?
- **9.4.** Attempt to prove termination for all well-typed closed terms: if $e:\tau$, then there exists e' val such that $e \mapsto^* e'$. You are free to appeal to Lemma 9.2 and Theorem 9.3 as necessary. Where does the attempt break down? Can you think of a stronger inductive hypothesis that might evade the difficulty?
- **9.5.** Define a closed term e of type τ in **T** to be *hereditarily terminating at type* τ by induction on the structure of τ as follows:
 - (a) If $\tau = \mathtt{nat}$, then e is hereditarily terminating at type τ iff e is terminating (that is, iff $e \mapsto^* \overline{n}$ for some n.)
 - (b) If $\tau = \tau_1 \to \tau_2$, then e is hereditarily terminating iff when e_1 is hereditarily terminating at type τ_1 , then $e(e_1)$ is hereditarily terminating at type τ_2 .

9.5 Notes

Attempt to prove hereditary termination for well-typed terms: if $e:\tau$, then e is hereditarily terminating at type τ . The stronger inductive hypothesis bypasses the difficulty that arose in Exercise 9.4, but introduces another obstacle. What is the complication? Can you think of an even stronger inductive hypothesis that would suffice for the proof?

- **9.6.** Show that if e is hereditarily terminating at type τ , e': τ , and $e' \mapsto e$, then e' is also hereditarily terminating at type τ . (The need for this result will become clear in the solution to Exercise **9.5**.)
- 9.7. Define an open, well-typed term

$$x_1:\tau_1,\ldots,x_n:\tau_n\vdash e:\tau$$

to be open hereditarily terminating iff every substitution instance

$$[e_1,\ldots,e_n/x_1,\ldots,x_n]e$$

is closed hereditarily terminating at type τ when each e_i is closed hereditarily terminating at type τ_i for each $1 \le i \le n$. Derive Exercise 9.3 from this result.

Part IV Finite Data Types



Chapter 10

Product Types

The *binary product* of two types consists of *ordered pairs* of values, one from each type in the order specified. The associated elimination forms are *projections*, which select the first and second component of a pair. The *nullary product*, or *unit*, type consists solely of the unique "null tuple" of no values, and has no associated elimination form. The product type admits both a *lazy* and an *eager* dynamics. According to the lazy dynamics, a pair is a value without regard to whether its components are values; they are not evaluated until (if ever) they are accessed and used in another computation. According to the eager dynamics, a pair is a value only if its components are values; they are evaluated when the pair is created.

More generally, we may consider the *finite product*, $\langle \tau_i \rangle_{i \in I}$, indexed by a finite set of *indices I*. The elements of the finite product type are *I-indexed tuples* whose *i*th component is an element of the type τ_i , for each $i \in I$. The components are accessed by *I-indexed projection* operations, generalizing the binary case. Special cases of the finite product include *n-tuples*, indexed by sets of the form $I = \{0, \ldots, n-1\}$, and *labeled tuples*, or *records*, indexed by finite sets of symbols. Similarly to binary products, finite products admit both an eager and a lazy interpretation.

10.1 Nullary and Binary Products

The abstract syntax of products is given by the following grammar:

There is no elimination form for the unit type, there being nothing to extract from the null tuple.

The statics of product types is given by the following rules.

$$\frac{}{\Gamma \vdash \langle \rangle : \mathtt{unit}} \tag{10.1a}$$

$$\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2}{\Gamma \vdash \langle e_1, e_2 \rangle : \tau_1 \times \tau_2}$$
(10.1b)

$$\frac{\Gamma \vdash e : \tau_1 \times \tau_2}{\Gamma \vdash e \cdot 1 : \tau_1} \tag{10.1c}$$

$$\frac{\Gamma \vdash e : \tau_1 \times \tau_2}{\Gamma \vdash e \cdot \mathbf{r} : \tau_2} \tag{10.1d}$$

The dynamics of product types is defined by the following rules:

$$\langle \rangle$$
 val (10.2a)

$$\frac{[e_1 \text{ val}] \quad [e_2 \text{ val}]}{\langle e_1, e_2 \rangle \text{ val}}$$
 (10.2b)

$$\left[\frac{e_1 \longmapsto e_1'}{\langle e_1, e_2 \rangle \longmapsto \langle e_1', e_2 \rangle} \right]$$
(10.2c)

$$\left[\frac{e_1 \text{ val} \quad e_2 \longmapsto e_2'}{\langle e_1, e_2 \rangle \longmapsto \langle e_1, e_2' \rangle} \right]$$
(10.2d)

$$\frac{e \longmapsto e'}{e \cdot 1 \longmapsto e' \cdot 1} \tag{10.2e}$$

$$\frac{e \longmapsto e'}{e \cdot \mathbf{r} \longmapsto e' \cdot \mathbf{r}} \tag{10.2f}$$

$$\frac{[e_1 \text{ val}] \quad [e_2 \text{ val}]}{\langle e_1, e_2 \rangle \cdot 1 \longmapsto e_1}$$
(10.2g)

$$\frac{[e_1 \text{ val}] \quad [e_2 \text{ val}]}{\langle e_1, e_2 \rangle \cdot \mathbf{r} \longmapsto e_2} \tag{10.2h}$$

The bracketed rules and premises are omitted for a lazy dynamics, and included for an eager dynamics of pairing.

The safety theorem applies to both the eager and the lazy dynamics, with the proof proceeding along similar lines in each case.

Theorem 10.1 (Safety). 1. If $e : \tau$ and $e \mapsto e'$, then $e' : \tau$.

2. If $e: \tau$ then either e val or there exists e' such that $e \mapsto e'$.

Proof. Preservation is proved by induction on transition defined by rules (10.2). Progress is proved by induction on typing defined by rules (10.1). \Box

10.2 Finite Products 85

10.2 Finite Products

The syntax of finite product types is given by the following grammar:

$$\begin{array}{llll} \mathsf{Typ} & \tau & ::= & \mathsf{prod}(\ \{i \hookrightarrow \tau_i\}_{i \in I}) & \langle \tau_i \rangle_{i \in I} & \mathsf{product} \\ \mathsf{Exp} & e & ::= & \mathsf{tpl}(\ \{i \hookrightarrow e_i\}_{i \in I}) & \langle e_i \rangle_{i \in I} & \mathsf{tuple} \\ & & \mathsf{pr}[\ i\](\ e) & e \cdot i & \mathsf{projection} \end{array}$$

The variable I stands for a finite index set over which products are formed. The type $\operatorname{prod}(\{i \hookrightarrow \tau_i\}_{i \in I})$, or $\prod_{i \in I} \tau_i$ for short, is the type of I-tuples of expressions e_i of type τ_i , one for each $i \in I$. An I-tuple has the form $\operatorname{tpl}(\{i \hookrightarrow e_i\}_{i \in I})$, or $\langle e_i \rangle_{i \in I}$ for short, and for each $i \in I$ the ith projection from an I-tuple e is written $\operatorname{pr}[i](e)$, or $e \cdot i$ for short.

When $I = \{i_1, \dots, i_n\}$, the *I*-tuple type may be written in the form

$$\langle i_1 \hookrightarrow \tau_1, \ldots, i_n \hookrightarrow \tau_n \rangle$$

where we make explicit the association of a type to each index $i \in I$. Similarly, we may write

$$\langle i_1 \hookrightarrow e_1, \ldots, i_n \hookrightarrow e_n \rangle$$

for the *I*-tuple whose *i*th component is e_i .

Finite products generalize empty and binary products by choosing I to be empty or the twoelement set $\{1,r\}$, respectively. In practice I is often chosen to be a finite set of symbols that serve as labels for the components of the tuple to enhance readability.

The statics of finite products is given by the following rules:

$$\frac{\Gamma \vdash e_1 : \tau_1 \quad \dots \quad \Gamma \vdash e_n : \tau_n}{\Gamma \vdash \langle i_1 \hookrightarrow e_1, \dots, i_n \hookrightarrow e_n \rangle : \langle i_1 \hookrightarrow \tau_1, \dots, i_n \hookrightarrow \tau_n \rangle}$$
(10.3a)

$$\frac{\Gamma \vdash e : \langle i_1 \hookrightarrow \tau_1, \dots, i_n \hookrightarrow \tau_n \rangle \quad (1 \le k \le n)}{\Gamma \vdash e \cdot i_k : \tau_k}$$
 (10.3b)

In rule (10.3b) the index $i_k \in I$ is a *particular* element of the index set I, whereas in rule (10.3a), the indices i_1, \ldots, i_n range over the entire index set I.

The dynamics of finite products is given by the following rules:

$$\frac{[e_1 \text{ val } \dots e_n \text{ val}]}{\langle i_1 \hookrightarrow e_1, \dots, i_n \hookrightarrow e_n \rangle \text{ val}}$$
 (10.4a)

$$\begin{bmatrix}
e_1 \text{ val } \dots & e_{j-1} \text{ val } e'_1 = e_1 & \dots & e'_{j-1} = e_{j-1} \\
e_j \longmapsto e'_j & e'_{j+1} = e_{j+1} & \dots & e'_n = e_n \\
\hline
\langle i_1 \hookrightarrow e_1, \dots, i_n \hookrightarrow e_n \rangle \longmapsto \langle i_1 \hookrightarrow e'_1, \dots, i_n \hookrightarrow e'_n \rangle
\end{bmatrix}$$
(10.4b)

$$\frac{e \longmapsto e'}{e \cdot i \longmapsto e' \cdot i} \tag{10.4c}$$

$$\frac{\left[\left\langle i_{1} \hookrightarrow e_{1}, \dots, i_{n} \hookrightarrow e_{n} \right\rangle \text{ val}\right]}{\left\langle i_{1} \hookrightarrow e_{1}, \dots, i_{n} \hookrightarrow e_{n} \right\rangle \cdot i_{k} \longmapsto e_{k}}$$

$$(10.4d)$$

As formulated, rule (10.4b) specifies that the components of a tuple are evaluated in *some* sequential order, without specifying the order in which the components are considered. It is not hard, but a bit technically complicated, to impose an evaluation order by imposing a total ordering on the index set and evaluating components according to this ordering.

Theorem 10.2 (Safety). *If* $e : \tau$, then either e val or there exists e' such that $e' : \tau$ and $e \mapsto e'$.

Proof. The safety theorem is decomposed into progress and preservation lemmas, which are proved as in Section 10.1. \Box

10.3 Primitive Mutual Recursion

Using products we may simplify the primitive recursion construct of **T** so that only the recursive result on the predecessor, and not the predecessor itself, is passed to the successor branch. Writing this as $iter\{e_0; x.e_1\}(e)$, we may define $rec\{e_0; x.y.e_1\}(e)$ to be $e' \cdot r$, where e' is the expression

$$iter\{\langle z, e_0 \rangle; x'.\langle s(x'\cdot 1), [x'\cdot 1, x'\cdot r/x, y]e_1 \rangle\}(e).$$

The idea is to compute inductively both the number n and the result of the recursive call on n, from which we can compute both n + 1 and the result of another recursion using e_1 . The base case is computed directly as the pair of zero and e_0 . It is easy to check that the statics and dynamics of the recursor are preserved by this definition.

We may also use product types to implement *mutual primitive recursion*, in which we define two functions simultaneously by primitive recursion. For example, consider the following recursion equations defining two mathematical functions on the natural numbers:

$$e(0) = 1$$

$$o(0) = 0$$

$$e(n+1) = o(n)$$

$$o(n+1) = e(n)$$

Intuitively, e(n) is non-zero if and only if n is even, and o(n) is non-zero if and only if n is odd.

To define these functions in **T** enriched with products, we first define an auxiliary function e_{eo} of type

$$\mathtt{nat} \rightarrow (\mathtt{nat} \times \mathtt{nat})$$

that computes both results simultaneously by swapping back and forth on recursive calls:

$$\lambda$$
 (n :nat) iter n { $z \hookrightarrow \langle 1, 0 \rangle \mid s(b) \hookrightarrow \langle b \cdot r, b \cdot 1 \rangle$ }.

We may then define e_{ev} and e_{od} as follows:

$$e_{\text{ev}} \triangleq \lambda (n: \text{nat}) e_{\text{eo}}(n) \cdot 1$$

 $e_{\text{od}} \triangleq \lambda (n: \text{nat}) e_{\text{eo}}(n) \cdot r.$

10.4 Notes 87

10.4 Notes

Product types are the most basic form of structured data. All languages have some form of product type, but often in a form that is combined with other, separable, concepts. Common manifestations of products include: (1) functions with "multiple arguments" or "multiple results"; (2) "objects" represented as tuples of mutually recursive functions; (3) "structures," which are tuples with mutable components. There are many papers on finite product types, which include record types as a special case. Pierce (2002) provides a thorough account of record types, and their subtyping properties (for which, see Chapter 24). Allen et al. (2006) analyzes many of the key ideas in the framework of dependent type theory.

Exercises

- **10.1.** A *database schema* may be thought of as a finite product type $\prod_{i \in I} \tau$, in which the *columns*, or *attributes*. are labeled by the indices I whose values are restricted to *atomic* types, such as nat and str. A value of a schema type is called a *tuple*, or *instance*, of that schema. A *database* may be thought of as a finite sequence of such tuples, called the *rows* of the database. Give a representation of a database using function, product, and natural numbers types, and define the *project* operation that sends a database with columns I to a database with columns $I' \subseteq I$ by restricting each row to the specified columns.
- **10.2.** Rather than choose between a lazy and an eager dynamics for products, we can instead distinguish two forms of product type, called the *positive* and the *negative*. The statics of the negative product is given by rules (10.1), and the dynamics is lazy. The statics of the positive product, written $\tau_1 \otimes \tau_2$, is given by the following rules:

$$\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2}{\Gamma \vdash \otimes (e_1; e_2) : \tau_1 \otimes \tau_2}$$
(10.5a)

$$\frac{\Gamma \vdash e_0 : \tau_1 \otimes \tau_2 \quad \Gamma x_1 : \tau_1 x_2 : \tau_2 \vdash e : \tau}{\Gamma \vdash \mathsf{split}(e_0; x_1, x_2.e) : \tau}$$
(10.5b)

The dynamics of fuse, the introduction form for the positive pair, is eager, essentially because the elimination form, split, extracts both components simultaneously.

Show that the negative product is definable in terms of the positive product using the unit and function types to express the lazy dynamics of negative pairing. Show that the positive product is definable in terms of the negative product, provided that we have at our disposal a let expression with a by-value dynamics so that we may enforce eager evaluation of positive pairs.

10.3. Specializing Exercise **10.2** to nullary products, we obtain a positive and a negative unit type. The negative unit type is given by rules (10.1), with no elimination forms and one introduction form. Give the statics and dynamics for a positive unit type, and show that the positive and negative unit types are inter-definable without any further assumptions.

88 10.4 Notes



Chapter 11

Sum Types

Most data structures involve alternatives such as the distinction between a leaf and an interior node in a tree, or a choice in the outermost form of a piece of abstract syntax. Importantly, the choice determines the structure of the value. For example, nodes have children, but leaves do not, and so forth. These concepts are expressed by *sum types*, specifically the *binary sum*, which offers a choice of two things, and the *nullary sum*, which offers a choice of no things. *Finite sums* generalize nullary and binary sums to allow an arbitrary number of cases indexed by a finite index set. As with products, sums come in both eager and lazy variants, differing in how values of sum type are defined.

11.1 Nullary and Binary Sums

The abstract syntax of sums is given by the following grammar:

```
\begin{array}{lllll} \text{Typ} & \tau & ::= & \text{void} & \text{void} & \text{nullary sum} \\ & & +(\tau_1;\tau_2) & \tau_1+\tau_2 & \text{binary sum} \\ \text{Exp} & e & ::= & \text{case}\{\tau\}(e) & \text{case}\,e\,\{\,\} & \text{null case} \\ & & \text{in}[1]\{\tau_1;\tau_2\}(e) & 1\cdot e & \text{left injection} \\ & & & \text{in}[\tau]\{\tau_1;\tau_2\}(e) & \text{r}\cdot e & \text{right injection} \\ & & & \text{case}\{x_1.e_1;x_2.e_2\}(e) & \text{case}\,e\,\{1\cdot x_1\hookrightarrow e_1\,|\, \textbf{r}\cdot x_2\hookrightarrow e_2\} & \text{case analysis} \end{array}
```

The nullary sum represents a choice of zero alternatives, and hence admits no introduction form. The elimination form, case e $\{$ $\}$, expresses the contradiction that e is a value of type void. The elements of the binary sum type are labeled to show whether they are drawn from the left or the right summand, either $1 \cdot e$ or $r \cdot e$. A value of the sum type is eliminated by case analysis.

The statics of sum types is given by the following rules.

$$\frac{\Gamma \vdash e : \mathtt{void}}{\Gamma \vdash \mathsf{case}\, e\, \{\,\} : \tau} \tag{11.1a}$$

$$\frac{\Gamma \vdash e : \tau_1}{\Gamma \vdash 1 \cdot e : \tau_1 + \tau_2} \tag{11.1b}$$

$$\frac{\Gamma \vdash e : \tau_2}{\Gamma \vdash \mathbf{r} \cdot e : \tau_1 + \tau_2} \tag{11.1c}$$

$$\frac{\Gamma \vdash e : \tau_1 + \tau_2 \quad \Gamma, x_1 : \tau_1 \vdash e_1 : \tau \quad \Gamma, x_2 : \tau_2 \vdash e_2 : \tau}{\Gamma \vdash \mathsf{case}\, e\, \{1 \cdot x_1 \hookrightarrow e_1 \mid \mathbf{r} \cdot x_2 \hookrightarrow e_2\} : \tau} \tag{11.1d}$$

For the sake of readability, in rules (11.1b) and (11.1c) we have written $1 \cdot e$ and $\mathbf{r} \cdot e$ in place of the abstract syntax $\inf[1]\{\tau_1;\tau_2\}(e)$ and $\inf[\mathbf{r}]\{\tau_1;\tau_2\}(e)$, which includes the types τ_1 and τ_2 explicitly. In rule (11.1d) both branches of the case analysis must have the same type. Because a type expresses a static "prediction" on the form of the value of an expression, and because an expression of sum type could evaluate to either form at run-time, we must insist that both branches yield the same type.

The dynamics of sums is given by the following rules:

$$\frac{e \longmapsto e'}{\operatorname{case} e \{\} \longmapsto \operatorname{case} e' \{\}}$$
 (11.2a)

$$\frac{[e \text{ val}]}{1 \cdot e \text{ val}} \tag{11.2b}$$

$$\frac{[e \text{ val}]}{r \cdot e \text{ val}} \tag{11.2c}$$

$$\left[\frac{e \mapsto e'}{1 \cdot e \mapsto 1 \cdot e'} \right]$$
(11.2d)

$$\left[\frac{e \longmapsto e'}{\mathbf{r} \cdot e \longmapsto \mathbf{r} \cdot e'}\right] \tag{11.2e}$$

$$\frac{e \longmapsto e'}{\operatorname{case} e \left\{ 1 \cdot x_1 \hookrightarrow e_1 \mid \mathbf{r} \cdot x_2 \hookrightarrow e_2 \right\} \longmapsto \operatorname{case} e' \left\{ 1 \cdot x_1 \hookrightarrow e_1 \mid \mathbf{r} \cdot x_2 \hookrightarrow e_2 \right\}}$$
 (11.2f)

$$\frac{[e \text{ val}]}{\text{case } 1 \cdot e \{1 \cdot x_1 \hookrightarrow e_1 \mid \mathbf{r} \cdot x_2 \hookrightarrow e_2\} \longmapsto [e/x_1]e_1}$$
 (11.2g)

$$\frac{[e \text{ val}]}{\operatorname{caser} \cdot e \left\{ 1 \cdot x_1 \hookrightarrow e_1 \mid r \cdot x_2 \hookrightarrow e_2 \right\} \longmapsto [e/x_2]e_2} \tag{11.2h}$$

The bracketed premises and rules are included for an eager dynamics, and excluded for a lazy dynamics.

The coherence of the statics and dynamics is stated and proved as usual.

Theorem 11.1 (Safety). 1. If $e : \tau$ and $e \mapsto e'$, then $e' : \tau$.

2. If $e: \tau$, then either e val or $e \mapsto e'$ for some e'.

Proof. The proof proceeds by induction on rules (11.2) for preservation, and by induction on rules (11.1) for progress.

91 11.2 Finite Sums

Finite Sums 11.2

Just as we may generalize nullary and binary products to finite products, so may we also generalize nullary and binary sums to finite sums. The syntax for finite sums is given by the following grammar:

$$\begin{array}{lll} \mathsf{Typ} & \tau & ::= & \mathsf{sum}(\{i \hookrightarrow \tau_i\}_{i \in I}) & [\tau_i]_{i \in I} & \mathsf{sum} \\ \mathsf{Exp} & e & ::= & \mathsf{in}[i]\{\vec{\tau}\}(e) & i \cdot e & \mathsf{injection} \\ & & \mathsf{case}(e;\{i \hookrightarrow x_i.e_i\}_{i \in I}) & \mathsf{case}\,e\,\{i \cdot x_i \hookrightarrow e_i\}_{i \in I} & \mathsf{case}\,\,\mathsf{analysis} \end{array}$$

The variable *I* stands for a finite index set over which sums are formed. The notation $\vec{\tau}$ stands for a finite function $\{i \hookrightarrow \tau_i\}_{i \in I}$ for some index set I. The type $sum(\{i \hookrightarrow \tau_i\}_{i \in I})$, or $\sum_{i \in I} \tau_i$ for short, is the type of *I*-classified values of the form $\inf[i]\{I\}(e_i)$, or $i \cdot e_i$ for short, where $i \in I$ and e_i is an expression of type τ_i . An *I*-classified value is analyzed by an *I*-way case analysis of the form case(e; { $i \hookrightarrow x_i.e_i$ } $_{i \in I}$).

When $I = \{i_1, \dots, i_n\}$, the type of *I*-classified values may be written

$$[i_1 \hookrightarrow \tau_1, \ldots, i_n \hookrightarrow \tau_n]$$

specifying the type associated with each class $l_i \in I$. Correspondingly, the *I*-way case analysis has the form

case
$$e\{i_1 \cdot x_1 \hookrightarrow e_1 \mid \ldots \mid i_n \cdot x_n \hookrightarrow e_n\}$$
.

Finite sums generalize empty and binary sums by choosing *I* to be empty or the two-element set $\{1,r\}$, respectively. In practice I is often chosen to be a finite set of symbols that serve as names for the classes so as to enhance readability.

The statics of finite sums is defined by the following rules:

$$\frac{\Gamma \vdash e : \tau_k \quad (1 \le k \le n)}{\Gamma \vdash i_k \cdot e : [i_1 \hookrightarrow \tau_1, \dots, i_n \hookrightarrow \tau_n]}$$
(11.3a)

$$\frac{\Gamma \vdash e : [i_1 \hookrightarrow \tau_1, \dots, i_n \hookrightarrow \tau_n] \quad \Gamma, x_1 : \tau_1 \vdash e_1 : \tau \quad \dots \quad \Gamma, x_n : \tau_n \vdash e_n : \tau}{\Gamma \vdash \mathsf{case} \, e \, \{i_1 \cdot x_1 \hookrightarrow e_1 \mid \dots \mid i_n \cdot x_n \hookrightarrow e_n\} : \tau}$$
(11.3b)

These rules generalize the statics for nullary and binary sums given in Section 11.1.

The dynamics of finite sums is defined by the following rules:

$$\frac{[e \text{ val}]}{i \cdot e \text{ val}} \tag{11.4a}$$

$$\left[\frac{e \longmapsto e'}{i \cdot e \longmapsto i \cdot e'}\right] \tag{11.4b}$$

$$\frac{\left[\frac{e \longmapsto e}{i \cdot e \longmapsto i \cdot e'}\right]}{e \longmapsto e'}$$

$$\frac{e \longmapsto e'}{\text{case } e \left\{i \cdot x_i \hookrightarrow e_i\right\}_{i \in I} \longmapsto \text{case } e' \left\{i \cdot x_i \hookrightarrow e_i\right\}_{i \in I}}$$
(11.4c)

$$\frac{i \cdot e \text{ val}}{\operatorname{case} i \cdot e \left\{ i \cdot x_i \hookrightarrow e_i \right\}_{i \in I} \longmapsto [e/x_i]e_i} \tag{11.4d}$$

These again generalize the dynamics of binary sums given in Section 11.1.

Theorem 11.2 (Safety). *If* $e : \tau$, then either e val or there exists $e' : \tau$ such that $e \mapsto e'$.

Proof. The proof is like that for the binary case, as described in Section 11.1.

11.3 Applications of Sum Types

Sum types have many uses, several of which we outline here. More interesting examples arise once we also have inductive and recursive types, which are introduced in Parts VI and Part VIII.

11.3.1 Void and Unit

It is instructive to compare the types unit and void, which are often confused with one another. The type unit has exactly one element, $\langle \rangle$, whereas the type void has no elements at all. Consequently, if e: unit, then if e evaluates to a value, that value is $\langle \rangle$ — in other words, e has no interesting value. On the other hand, if e: void, then e must not yield a value; if it were to have a value, it would have to be a value of type void, of which there are none. Thus what is called the void type in many languages is really the type unit because it indicates that an expression has no interesting value, not that it has no value at all!

11.3.2 Booleans

Perhaps the simplest example of a sum type is the familiar type of Booleans, whose syntax is given by the following grammar:

Typ $ au$::=	bool	bool	booleans
Exp e	∷≢	true	true	truth
		false	false	falsity
		$if(e;e_1;e_2)$	if e then e_1 else e_2	conditional

The expression if $(e; e_1; e_2)$ branches on the value of e: bool.

The statics of Booleans is given by the following typing rules:

$$\frac{}{\Gamma \vdash \mathsf{true} : \mathsf{bool}} \tag{11.5a}$$

$$\frac{}{\Gamma \vdash \mathsf{false:bool}} \tag{11.5b}$$

$$\frac{\Gamma \vdash e : \text{bool} \quad \Gamma \vdash e_1 : \tau \quad \Gamma \vdash e_2 : \tau}{\Gamma \vdash \text{if } e \text{ then } e_1 \text{ else } e_2 : \tau}$$
(11.5c)

The dynamics is given by the following value and transition rules:

$$\frac{}{\text{if true then } e_1 \text{ else } e_2 \longmapsto e_1} \tag{11.6c}$$

$$\frac{}{\text{if false then } e_1 \text{ else } e_2 \longmapsto e_2} \tag{11.6d}$$

$$\frac{e \longmapsto e'}{\text{if } e \text{ then } e_1 \text{ else } e_2 \longmapsto \text{if } e' \text{ then } e_1 \text{ else } e_2} \tag{11.6e}$$

The type bool is definable in terms of binary sums and nullary products:

$$bool = unit + unit (11.7a)$$

$$true = 1 \cdot \langle \rangle \tag{11.7b}$$

$$false = r \cdot \langle \rangle \tag{11.7c}$$

$$false = \mathbf{r} \cdot \langle \rangle$$
 (11.7c)
$$if e then e_1 else e_2 = case e \{1 \cdot x_1 \hookrightarrow e_1 \mid \mathbf{r} \cdot x_2 \hookrightarrow e_2\}$$
 (11.7d)

In the last equation above the variables x_1 and x_2 are chosen arbitrarily such that $x_1 \notin e_1$ and $x_2 \notin e_2$. It is a simple matter to check that the readily-defined statics and dynamics of the type bool are engendered by these definitions.

11.3.3 **Enumerations**

More generally, sum types can be used to define finite enumeration types, those whose values are one of an explicitly given finite set, and whose elimination form is a case analysis on the elements of that set. For example, the type suit, whose elements are \clubsuit , \diamondsuit , \heartsuit , and \spadesuit , has as elimination form the case analysis

case
$$e \{ A \hookrightarrow e_0 \mid \Diamond \hookrightarrow e_1 \mid \heartsuit \hookrightarrow e_2 \mid A \hookrightarrow e_3 \}$$
,

which distinguishes among the four suits. Such finite enumerations are easily representable as sums. For example, we may define suit $= [unit]_{\in I}$, where $I = \{ \clubsuit, \lozenge, \heartsuit, \spadesuit \}$ and the type family is constant over this set. The case analysis form for a labeled sum is almost literally the desired case analysis for the given enumeration, the only difference being the binding for the uninteresting value associated with each summand, which we may ignore.

Other examples of enumeration types abound. For example, most languages have a type char of characters, which is a large enumeration type containing all possible Unicode (or other such standard classification) characters. Each character is assigned a code (such as UTF-8) used for interchange among programs. The type char is equipped with operations such as chcode(n) that yield the char associated to the code n, and codech(c) that yield the code of character c. Using the linear ordering on codes we may define a total ordering of characters, called the collating sequence determined by that code.

Options 11.3.4

Another use of sums is to define the *option* types, which have the following syntax:

```
Typ \tau ::= opt(\tau)
                                                                option
                                                  \tau opt
Exp e ::= null
                                                  null
                                                                nothing
                                                \mathtt{just}(\mathit{e}\,)
                  just(e)
                                                               something
                  ifnull\{\tau\}\{e_1; x.e_2\}(e) ifnulle\{null\hookrightarrow e_1| just(x)\hookrightarrow e_2\}
                                                               null test
```

The type opt(τ) represents the type of "optional" values of type τ . The introduction forms are null, corresponding to "no value", and just (e), corresponding to a specified value of type τ . The elimination form discriminates between the two possibilities.

The option type is definable from sums and nullary products according to the following equations:1

$$\tau \, \text{opt} = \text{unit} + \tau \tag{11.8a}$$

$$null = l \cdot \langle \rangle \tag{11.8b}$$

$$null = 1 \cdot \langle \rangle$$

$$just(e) = r \cdot e$$
(11.8b)
(11.8c)

$$ifnull e \{null \hookrightarrow e_1 \mid just(x_2) \hookrightarrow e_2\} = case e \{1 \cdot _ \hookrightarrow e_1 \mid r \cdot x_2 \hookrightarrow e_2\}$$
 (11.8d)

We leave it to the reader to check the statics and dynamics implied by these definitions.

The option type is the key to understanding a common misconception, the *null pointer fallacy*. This fallacy arises from two related errors. The first error is to deem values of certain types to be mysterious entities called *pointers*. This terminology arises from suppositions about how these values might be represented at run-time, rather than on their semantic role in the language. The second error compounds the first. A particular value of a pointer type is distinguished as the null pointer, which, unlike the other elements of that type, does not stand for a value of that type at all, but rather rejects all attempts to use it.

To help avoid such failures, such languages usually include a function, say null: $\tau \to bool$, that yields true if its argument is null, and false otherwise. Such a test allows the programmer to take steps to avoid using null as a value of the type it purports to inhabit. Consequently, programs are riddled with conditionals of the form

$$if null(e) then ... error ... else ... proceed$$
 (11.9)

Despite this, "null pointer" exceptions at run-time are rampant, in part because it is quite easy to overlook the need for such a test, and in part because detection of a null pointer leaves little recourse other than abortion of the program.

The underlying problem is the failure to distinguish the type τ from the type τ opt. Rather than think of the elements of type τ as pointers, and thereby have to worry about the null pointer, we instead distinguish between a *genuine* value of type τ and an *optional* value of type τ . An optional

¹We often write an underscore in place of a bound variable that is not used within its scope.

11.4 Notes 95

value of type τ may or may not be present, but, if it is, the underlying value is truly a value of type τ (and cannot be null). The elimination form for the option type,

if
$$\operatorname{null} e\left\{\operatorname{null} \hookrightarrow e_{error} \mid \operatorname{just}(x) \hookrightarrow e_{ok}\right\},$$
 (11.10)

propagates the information that e is present into the non-null branch by binding a genuine value of type τ to the variable x. The case analysis effects a change of type from "optional value of type τ " to "genuine value of type τ ", so that within the non-null branch no further null checks, explicit or implicit, are necessary. Note that such a change of type is not achieved by the simple Boolean-valued test exemplified by expression (11.9); the advantage of option types is precisely that they do so.

11.4 Notes

Heterogeneous data structures are ubiquitous. Sums codify heterogeneity, yet few languages support them in the form given here. The best approximation in commercial languages is the concept of a class in object-oriented programming. A class is an injection into a sum type, and dispatch is case analysis on the class of the data object. (See Chapter 26 for more on this correspondence.) The absence of sums is the origin of C.A.R. Hoare's self-described "billion dollar mistake," the null pointer (Hoare, 2009). Bad language designs put the burden of managing "null" values entirely at run-time, instead of making the possibility or the impossibility of "null" apparent at compile time.

Exercises

- **11.1**. Complete the definition of a finite enumeration type sketched in Section 11.3.3. Derive enumeration types from finite sum types.
- **11.2.** The essence of Hoare's mistake is the misidentification of the type τ opt with the type bool \times τ . Values of the latter type are pairs consisting of a boolean "flag" and a value of type τ . The idea is that the flag indicates whether the associated value is "present". When the flag is true, the second component is present, and, when the flag is false, the second component is absent.

Analyze Hoare's mistake by attempting to define τ opt to be the type bool $\times \tau$ by filling in the following chart:

$$\label{eq:null} \verb"null"\,\triangleq\, \verb"?"$$

$$\verb"just"(\,e\,) \triangleq \verb"?"$$

$$\verb"ifnull"\,e\, \{\verb"null"\,\hookrightarrow\, e_1 \mid \verb"just"(\,x\,) \hookrightarrow e_2\} \triangleq \verb"?"$$

Argue that *even if* we adopt Hoare's convention of admitting a "null" value of every type, the chart cannot be properly filled.

96 11.4 Notes

11.3. Databases have a version of the "null pointer" problem that arises when not every tuple provides a value for every attribute (such as a person's middle name). More generally, many commercial databases are limited to a single atomic type for each attribute, presenting problems when the value of that attribute may have several types (for example, one may have different sorts of postal codes depending on the country). Consider how to address these problems using the methods discussed in Exercise 10.1. Suggest how to handle null values and heterogeneous values that avoids some of the complications that arise in traditional formulations of databases.

11.4. A combinational circuit is an open expression of type

$$x_1 : bool, \ldots, x_n : bool \vdash e : bool,$$

which computes a boolean value from n boolean inputs. Define a NOR and a NAND gate as boolean circuits with two inputs and one output. There is no reason to restrict to a single output. For example, define an HALF-ADDER that takes two boolean inputs, but produces two boolean outputs, the sum and the carry outputs of the HALF-ADDER. Then define a FULL-ADDER that takes three inputs, the addends and an incoming carry, and produces two outputs, the sum and the outgoing carry. Define the type NYBBLE to be the product bool \times bool \times bool \times bool. Define the combinational circuit NYBBLE-ADDER that takes two nybbles as input and produces a nybble and a carry-out bit as output.

11.5. A *signal* is a time-varying sequence of booleans, representing the status of the signal at each time instant. An RS latch is a fundamental digital circuit with two input signals and two output signals. Define the type signal of signals to be the function type nat → bool of infinite sequences of booleans. Define an RS latch as a function of type

$$(signal \times signal) \rightarrow (signal \times signal).$$

Part V Types and Propositions



Part VI Infinite Data Types





Part VII Variable Types



Chapter 16

System F of Polymorphic Types

The languages we have considered so far are all *monomorphic* in that every expression has a unique type, given the types of its free variables, if it has a type at all. Yet it is often the case that essentially the same behavior is required, albeit at several different types. For example, in **T** there is a *distinct* identity function for each type τ , namely λ (x: τ) x, even though the behavior is the same for each choice of τ . Similarly, there is a distinct composition operator for each triple of types, namely

$$\circ_{\tau_1,\tau_2,\tau_3} = \lambda \left(f : \tau_2 \to \tau_3 \right) \lambda \left(g : \tau_1 \to \tau_2 \right) \lambda \left(x : \tau_1 \right) f(g(x)).$$

Each choice of the three types requires a *different* program, even though they all have the same behavior when executed.

Obviously it would be useful to capture the pattern once and for all, and to instantiate this pattern each time we need it. The expression patterns codify generic (type-independent) behaviors that are shared by all instances of the pattern. Such generic expressions are *polymorphic*. In this chapter we will study the language \mathbf{F} , which was introduced by Girard under the name *System F* and by Reynolds under the name *polymorphic typed* λ -calculus. Although motivated by a simple practical problem (how to avoid writing redundant code), the concept of polymorphism is central to an impressive variety of seemingly disparate concepts, including the concept of data abstraction (the subject of Chapter 17), and the definability of product, sum, inductive, and coinductive types considered in the preceding chapters. (Only general recursive types extend the expressive power of the language.)

16.1 Polymorphic Abstraction

The language **F** is a variant of **T** in which we eliminate the type of natural numbers, but add, in compensation, polymorphic types:¹

A type abstraction $\Lambda(t.e)$ defines a generic, or polymorphic, function with type variable t standing for an unspecified type within e. A type application, or instantiation $App\{\tau\}(e)$ applies a polymorphic function to a specified type, which is plugged in for the type variable to obtain the result. The universal type, $\forall (t.\tau)$, classifies polymorphic functions.

The statics of **F** consists of two judgment forms, the *type formation* judgment,

$$\Delta \vdash \tau$$
 type,

and the typing judgment,

$$\Lambda \Gamma \vdash \rho \cdot \tau$$

The hypotheses Δ have the form t type, where t is a variable of sort Typ, and the hypotheses Γ have the form $x : \tau$, where x is a variable of sort Exp.

The rules defining the type formation judgment are as follows:

$$\overline{\Delta}, t \text{ type} \vdash t \text{ type}$$
 (16.1a)

$$\frac{\Delta \vdash \tau_1 \text{ type } \Delta \vdash \tau_2 \text{ type}}{\Delta \vdash \rightarrow (\tau_1; \tau_2) \text{ type}}$$
(16.1b)

$$\frac{\Delta, t \text{ type} \vdash \tau \text{ type}}{\Delta \vdash \forall (t.\tau) \text{ type}}$$
 (16.1c)

The rules defining the typing judgment are as follows:

$$\overline{\Delta \Gamma, x : \tau \vdash x : \tau} \tag{16.2a}$$

$$\frac{\Delta \vdash \tau_1 \text{ type } \Delta \Gamma, x : \tau_1 \vdash e : \tau_2}{\Delta \Gamma \vdash \lambda \{\tau_1\}(x.e) : \rightarrow (\tau_1; \tau_2)}$$
(16.2b)

$$\frac{\Delta \Gamma \vdash e_1 : \rightarrow (\tau_2; \tau) \quad \Delta \Gamma \vdash e_2 : \tau_2}{\Delta \Gamma \vdash \operatorname{ap}(e_1; e_2) : \tau}$$
(16.2c)

¹Girard's original version of System F included the natural numbers as a basic type.

$$\frac{\Delta, t \text{ type } \Gamma \vdash e : \tau}{\Delta \Gamma \vdash \Lambda(t.e) : \forall (t.\tau)}$$
 (16.2d)

$$\frac{\Delta \Gamma \vdash e : \forall (t.\tau') \quad \Delta \vdash \tau \text{ type}}{\Delta \Gamma \vdash \mathsf{App}\{\tau\}(e) : [\tau/t]\tau'}$$
(16.2e)

Lemma 16.1 (Regularity). *If* $\Delta \Gamma \vdash e : \tau$, and if $\Delta \vdash \tau_i$ type for each assumption $x_i : \tau_i$ in Γ , then $\Delta \vdash \tau$ type.

Proof. By induction on rules (16.2).

The statics admits the structural rules for a general hypothetical judgment. In particular, we have the following critical substitution property for type formation and expression typing.

Lemma 16.2 (Substitution). 1. If Δ , t type $\vdash \tau'$ type and $\Delta \vdash \tau$ type, then $\Delta \vdash [\tau/t]\tau'$ type.

- 2. If Δ , t type $\Gamma \vdash e' : \tau'$ and $\Delta \vdash \tau$ type, then $\Delta [\tau/t]\Gamma \vdash [\tau/t]e' : [\tau/t]\tau'$.
- 3. If $\Delta \Gamma, x : \tau \vdash e' : \tau'$ and $\Delta \Gamma \vdash e : \tau$, then $\Delta \Gamma \vdash [e/x]e' : \tau'$.

The second part of the lemma requires substitution into the context Γ as well as into the term and its type, because the type variable t may occur freely in any of these positions.

Returning to the motivating examples from the introduction, the polymorphic identity function, *I*, is written

$$\Lambda(t)\lambda(x:t)x;$$

it has the polymorphic type

$$\forall (t.t \rightarrow t)$$

Instances of the polymorphic identity are written $I[\tau]$, where τ is some type, and have the type $\tau \to \tau$.

Similarly, the polymorphic composition function, *C*, is written

$$\Lambda(t_1)\Lambda(t_2)\Lambda(t_3)\lambda(f:t_2\to t_3)\lambda(g:t_1\to t_2)\lambda(x:t_1)f(g(x)).$$

The function *C* has the polymorphic type

$$\forall (t_1.\forall (t_2.\forall (t_3.(t_2 \rightarrow t_3) \rightarrow (t_1 \rightarrow t_2) \rightarrow (t_1 \rightarrow t_3)))).$$

Instances of *C* are obtained by applying it to a triple of types, written $C[\tau_1][\tau_2][\tau_3]$. Each such instance has the type

$$(\tau_2 \rightarrow \tau_3) \rightarrow (\tau_1 \rightarrow \tau_2) \rightarrow (\tau_1 \rightarrow \tau_3).$$

Dynamics

The dynamics of **F** is given as follows:

$$\frac{1}{\lambda\{\tau\}(x.e) \text{ val}} \tag{16.3a}$$

$$\Lambda(t.e) \text{ val}$$
 (16.3b)

$$\frac{[e_2 \text{ val}]}{\operatorname{ap}(\lambda\{\tau_1\}(x.e); e_2) \longmapsto [e_2/x]e}$$
(16.3c)

$$\frac{e_1 \longmapsto e_1'}{\operatorname{ap}(e_1; e_2) \longmapsto \operatorname{ap}(e_1'; e_2)} \tag{16.3d}$$

$$\left[\frac{e_1 \text{ val} \quad e_2 \longmapsto e_2'}{\operatorname{ap}(e_1; e_2) \longmapsto \operatorname{ap}(e_1; e_2')} \right]$$
(16.3e)

$$\frac{1}{\text{App}\{\tau\}(\Lambda(t,e)) \longmapsto [\tau/t]e} \tag{16.3f}$$

$$\frac{\operatorname{App}\{\tau\}(\Lambda(t.e)) \longmapsto [\tau/t]e}{\operatorname{App}\{\tau\}(e) \longmapsto \operatorname{App}\{\tau\}(e')}$$
(16.3g)

The bracketed premises and rule are included for a call-by-value interpretation, and omitted for a call-by-name interpretation of **F**.

It is a simple matter to prove safety for **F**, using familiar methods.

Lemma 16.3 (Canonical Forms). *Suppose that e* : τ *and e val, then*

1. If
$$\tau = \rightarrow (\tau_1; \tau_2)$$
, then $e = \lambda \{\tau_1\}(x.e_2)$ with $x : \tau_1 \vdash e_2 : \tau_2$.

2. If
$$\tau = \forall (t.\tau')$$
, then $e = \Lambda(t.e')$ with t type $\vdash e' : \tau'$.

Proof. By rule induction on the statics.

Theorem 16.4 (Preservation). *If* $e : \tau$ *and* $e \mapsto e'$, *then* $e' : \tau$.

Proof. By rule induction on the dynamics.

Theorem 16.5 (Progress). *If* $e : \tau$, then either e val or there exists e' such that $e \mapsto e'$.

Proof. By rule induction on the statics.

16.2 Polymorphic Definability

The language **F** is astonishingly expressive. Not only are all (lazy) finite products and sums definable in the language, but so are all (lazy) inductive and coinductive types. Their definability is most naturally expressed using definitional equality, which is the least congruence containing the following two axioms:

$$\frac{\Delta \Gamma, x : \tau_1 \vdash e_2 : \tau_2 \quad \Delta \Gamma \vdash e_1 : \tau_1}{\Delta \Gamma \vdash (\lambda (x : \tau_1) e_2)(e_1) \equiv [e_1/x]e_2 : \tau_2}$$
(16.4a)

$$\frac{\Delta, t \text{ type } \Gamma \vdash e : \tau \quad \Delta \vdash \rho \text{ type}}{\Delta \Gamma \vdash \Lambda(t) e[\rho] \equiv [\rho/t]e : [\rho/t]\tau}$$
(16.4b)

In addition there are rules omitted here specifying that definitional equality is a congruence relation (that is, an equivalence relation respected by all expression-forming operations).

16.2.1 Products and Sums

The nullary product, or unit, type is definable in **F** as follows:

unit
$$\triangleq \forall (r.r \rightarrow r)$$

 $\langle \rangle \triangleq \Lambda(r) \lambda(x:r) x$

The identity function plays the role of the null tuple, because it is the only closed value of this type.

Binary products are definable in **F** by using encoding tricks similar to those described in Chapter 21 for the untyped λ -calculus:

$$\tau_{1} \times \tau_{2} \triangleq \forall (r.(\tau_{1} \to \tau_{2} \to r) \to r)$$

$$\langle e_{1}, e_{2} \rangle \triangleq \Lambda(r) \lambda (x : \tau_{1} \to \tau_{2} \to r) x(e_{1})(e_{2})$$

$$e \cdot 1 \triangleq e[\tau_{1}] (\lambda (x : \tau_{1}) \lambda (y : \tau_{2}) x)$$

$$e \cdot \mathbf{r} \triangleq e[\tau_{2}] (\lambda (x : \tau_{1}) \lambda (y : \tau_{2}) y)$$

The statics given in Chapter 10 is derivable according to these definitions. Moreover, the following definitional equalities are derivable in **F** from these definitions:

$$\langle e_1, e_2 \rangle \cdot 1 \equiv e_1 : \tau_1$$

and

$$\langle e_1, e_2 \rangle \cdot \mathbf{r} \equiv e_2 : \tau_2.$$

The nullary sum, or void, type is definable in F:

$$\mathtt{void} \triangleq \forall (\mathit{r.r})$$
 $\mathtt{case}\,e\,\{\,\} \triangleq e[\,\rho\,]$

Binary sums are also definable in **F**:

$$\begin{split} \tau_1 + \tau_2 &\triangleq \forall (\mathit{r}.(\tau_1 \to \mathit{r}\,) \to (\tau_2 \to \mathit{r}\,) \to \mathit{r}\,) \\ 1 \cdot e &\triangleq \Lambda(\mathit{r}\,) \, \lambda \, (\mathit{x} : \tau_1 \to \mathit{r}\,) \, \lambda \, (\mathit{y} : \tau_2 \to \mathit{r}\,) \, x(\mathit{e}\,) \\ \mathbf{r} \cdot e &\triangleq \Lambda(\mathit{r}\,) \, \lambda \, (\mathit{x} : \tau_1 \to \mathit{r}\,) \, \lambda \, (\mathit{y} : \tau_2 \to \mathit{r}\,) \, y(\mathit{e}\,) \\ &\qquad \qquad \mathsf{case} \, e \, \{1 \cdot \mathit{x}_1 \hookrightarrow e_1 \, | \, \mathbf{r} \cdot \mathit{x}_2 \hookrightarrow e_2\} \triangleq \\ e[\,\rho\,] (\,\lambda \, (\mathit{x}_1 : \tau_1\,) \, e_1\,) (\,\lambda \, (\mathit{x}_2 : \tau_2\,) \, e_2\,) \end{split}$$

provided that the types make sense. It is easy to check that the following equivalences are derivable in \mathbf{F} :

$$case 1 \cdot d_1 \{1 \cdot x_1 \hookrightarrow e_1 \mid r \cdot x_2 \hookrightarrow e_2\} \equiv [d_1/x_1]e_1 : \rho$$

and

$$\operatorname{case} \mathbf{r} \cdot d_2 \left\{ 1 \cdot x_1 \hookrightarrow e_1 \mid \mathbf{r} \cdot x_2 \hookrightarrow e_2 \right\} \equiv [d_2/x_2]e_2 : \rho.$$

Thus the dynamic behavior specified in Chapter 11 is correctly implemented by these definitions.

16.2.2 Natural Numbers

As we remarked above, the natural numbers (under a lazy interpretation) are also definable in **F**. The key is the iterator, whose typing rule we recall here for reference:

$$\frac{e_0: \mathtt{nat} \quad e_1: \tau \quad x: \tau \vdash e_2: \tau}{\mathtt{iter}\{e_1; x.e_2\}(e_0): \tau}$$

Because the result type τ is arbitrary, this means that if we have an iterator, then we can use it to define a function of type

$$\mathtt{nat} \to \forall (t.t \to (t \to t) \to t).$$

This function, when applied to an argument n, yields a polymorphic function that, for any result type, t, given the initial result for z and a transformation from the result for x into the result for s(x), yields the result of iterating the transformation n times, starting with the initial result.

Because the *only* operation we can perform on a natural number is to iterate up to it, we may simply *identify* a natural number, n, with the polymorphic iterate-up-to-n function just described. Thus we may define the type of natural numbers in \mathbf{F} by the following equations:

$$\begin{split} \operatorname{nat} &\triangleq \forall (t.t \to (t \to t) \to t) \\ \mathbf{z} &\triangleq \Lambda(t) \, \lambda \, (z \colon \! t) \, \lambda \, (s \colon \! t \to t) \, z \\ \mathbf{s}(e) &\triangleq \Lambda(t) \, \lambda \, (z \colon \! t) \, \lambda \, (s \colon \! t \to t) \, s(e[t](z)(s)) \end{split} \\ & \\ \operatorname{iter}\{e_1; x.e_2\}(e_0) &\triangleq e_0[\tau](e_1) \, (\lambda \, (x \colon \! \tau) \, e_2) \end{split}$$

It is easy to check that the statics and dynamics of the natural numbers type given in Chapter 9 are derivable in **F** under these definitions. The representations of the numerals in **F** are called the *polymorphic Church numerals*.

The encodability of the natural numbers shows that **F** is *at least as expressive* as **T**. But is it *more* expressive? Yes! It is possible to show that the evaluation function for **T** is definable in **F**, even though it is not definable in **T** itself. However, the same diagonal argument given in Chapter 9 applies here, showing that the evaluation function for **F** is not definable in **F**. We may enrich **F** a bit more to define the evaluator for **F**, but as long as all programs in the enriched language terminate, we will once again have an undefinable function, the evaluation function for that extension.

16.3 Parametricity Overview

A remarkable property of **F** is that polymorphic types severely constrain the behavior of their elements. We may prove useful theorems about an expression knowing *only* its type—that is, without ever looking at the code. For example, if i is any expression of type $\forall (t.t \to t)$, then it is the identity function. Informally, when i is applied to a type, τ , and an argument of type τ , it returns a value of type τ . But because τ is not specified until i is called, the function has no choice but to return its argument, which is to say that it is essentially the identity function. Similarly, if b is any expression of type $\forall (t.t \to t \to t)$, then b is equivalent to either $\Lambda(t) \lambda(x:t) \lambda(y:t) x$ or $\Lambda(t) \lambda(x:t) \lambda(y:t) y$. Intuitively, when b is applied to two arguments of a given type, the only value it can return is one of the givens.

Properties of a program in **F** that can be proved knowing only its type are called *parametricity properties*. The facts about the functions *i* and *b* stated above are examples of parametricity properties. Such properties are sometimes called "free theorems," because they come from typing "for free", without any knowledge of the code itself. It bears repeating that in **F** we prove non-trivial behavioral properties of programs without ever examining the program text. The key to this incredible fact is that we are able to prove a deep property, called *parametricity*, about the language **F**, that then applies to every program written in **F**. One may say that the type system "pre-verifies" programs with respect to a broad range of useful properties, eliminating the need to prove those properties about every program separately. The parametricity theorem for **F** explains the remarkable experience that if a piece of code type checks, then it "just works." Parametricity narrows the space of well-typed programs sufficiently that the opportunities for programmer error are reduced to almost nothing.

So how does the parametricity theorem work? Without getting into too many technical details (but see Chapter 48 for a full treatment), we can give a brief summary of the main idea. Any function $i : \forall (t.t \to t)$ in **F** enjoys the following property:

For any type τ and any property \mathcal{P} of the type τ , then if \mathcal{P} holds of $x : \tau$, then \mathcal{P} holds of $i[\tau](x)$.

To show that for any type τ , and any x of type τ , the expression $i[\tau](x)$ is equivalent to x, it suffices to fix $x_0 : \tau$, and consider the property \mathcal{P}_{x_0} that holds of $y : \tau$ iff y is equivalent to x_0 . Obviously \mathcal{P} holds of x_0 itself, and hence by the above-displayed property of i, it sends any argument satisfying \mathcal{P}_{x_0} to a result satisfying \mathcal{P}_{x_0} , which is to say that it sends x_0 to x_0 . Because x_0 is an arbitrary element of τ , it follows that $i[\tau]$ is the identity function, $\lambda(x:\tau)x$, on the type τ , and because τ is itself arbitrary, i is the polymorphic identity function, $\lambda(x:\tau)x$.

148 16.4 Notes

A similar argument suffices to show that the function b, defined above, is either $\Lambda(t) \lambda(x:t) \lambda(y:t) x$ or $\Lambda(t) \lambda(x:t) \lambda(y:t) y$. By virtue of its type the function b enjoys the parametricity property

```
For any type \tau and any property \mathcal{P} of \tau, if \mathcal{P} holds of x : \tau and of y : \tau, then \mathcal{P} holds of b[\tau](x)(y).
```

Choose an arbitrary type τ and two arbitrary elements x_0 and y_0 of type τ . Define \mathcal{Q}_{x_0,y_0} to hold of $z:\tau$ iff either z is equivalent to x_0 or z is equivalent to y_0 . Clearly \mathcal{Q}_{x_0,y_0} holds of both x_0 and y_0 themselves, so by the quoted parametricity property of b, it follows that \mathcal{Q}_{x_0,y_0} holds of $b[\tau](x_0)(y_0)$, which is to say that it is equivalent to either x_0 or y_0 . Since τ , x_0 , and y_0 are arbitrary, it follows that b is equivalent to either $A(t)\lambda(y:t)\lambda(y:t)x$ or $A(t)\lambda(x:t)\lambda(y:t)y$.

The parametricity theorem for **F** implies even stronger properties of functions such as i and b considered above. For example, the function i of type $\forall (t.t \rightarrow t)$ also satisfies the following condition:

```
If \tau and \tau' are any two types, and \mathcal{R} is a binary relation between \tau and \tau', then for any x : \tau and x' : \tau', if \mathcal{R} relates x to x', then \mathcal{R} relates i[\tau](x) to i[\tau'](x').
```

Using this property we may again prove that i is equivalent to the polymorphic identity function. Specifically, if τ is any type and $g:\tau\to\tau$ is any function on that type, then it follows from the type of i alone that $i[\tau](g(x))$ is equivalent to $g(i[\tau](x))$ for any $x:\tau$. To prove this, simply choose $\mathcal R$ to the be graph of the function g, the relation $\mathcal R_g$ that holds of x and x' iff x' is equivalent to g(x). The parametricity property of i, when specialized to $\mathcal R_g$, states that if x' is equivalent to g(x), then $i[\tau](x')$ is equivalent to $g(i[\tau](x))$, which is to say that $i[\tau](g(x))$ is equivalent to $g(i[\tau](x))$. To show that i is equivalent to the identity function, choose $x_0:\tau$ arbitrarily, and consider the constant function g_0 on τ that always returns x_0 . Because x_0 is equivalent to $g_0(x_0)$, it follows that $i[\tau](x_0)$ is equivalent to x_0 , which is to say that i behaves like the polymorphic identity function.

16.4 Notes

System F was introduced by Girard (1972) in the context of proof theory and by Reynolds (1974) in the context of programming languages. The concept of parametricity was originally isolated by Strachey, but was not fully developed until the work of Reynolds (1983). The phrase "free theorems" for parametricity theorems was introduced by Wadler (1989).

Exercises

- **16.1**. Give polymorphic definitions and types to the s and k combinators defined in Exercise 3.1.
- **16.2**. Define in **F** the type bool of *Church booleans*. Define the type bool, and define true and false of this type, and the conditional if e then e_0 else e_1 , where e is of this type.

16.4 Notes 149

16.3. Define in **F** the inductive type of lists of natural numbers as defined in Chapter 15. *Hint*: Define the representation in terms of the recursor (elimination form) for lists, following the pattern for defining the type of natural numbers.

- **16.4**. Define in **F** an arbitrary inductive type, $\mu(t.\tau)$. *Hint*: generalize your answer to Exercise **16.3**.
- **16.5**. Define the type t list as in Exercise **16.3**, with the element type, t, unspecified. Define the finite set of *elements* of a list l to be those x given by the head of some number of tails of l. Now suppose that $f: \forall (t.t \texttt{list} \to t \texttt{list})$ is an arbitrary function of the stated type. Show that the elements of $f[\tau](l)$ are a subset of those of l. Thus f may only permute, replicate, or drop elements from its input list to obtain its output list.

150 16.4 Notes



Chapter 17

Abstract Types

Data abstraction is perhaps the most important technique for structuring programs. The main idea is to introduce an *interface* that serves as a contract between the *client* and the *implementor* of an abstract type. The interface specifies what the client may rely on for its own work, and, simultaneously, what the implementor must provide to satisfy the contract. The interface serves to isolate the client from the implementor so that each may be developed in isolation from the other. In particular one implementation can be replaced by another without affecting the behavior of the client, provided that the two implementations meet the same interface and that each simulates the other with respect to the operations of the interface. This property is called *representation independence* for an abstract type.

Data abstraction is formalized by extending the language **F** with *existential types*. Interfaces are existential types that provide a collection of operations acting on an unspecified, or abstract, type. Implementations are packages, the introduction form for existential types, and clients are uses of the corresponding elimination form. It is remarkable that the programming concept of data abstraction is captured so naturally and directly by the logical concept of existential type quantification. Existential types are closely connected with universal types, and hence are often treated together. The superficial reason is that both are forms of type quantification, and hence both require the machinery of type variables. The deeper reason is that existential types are *definable* from universals — surprisingly, data abstraction is actually just a form of polymorphism! Consequently, representation independence is an application of the parametricity properties of polymorphic functions discussed in Chapter 16.

17.1 Existential Types

The syntax of **FE** extends **F** with the following constructs:

The introduction form $\exists (t,\tau)$ is a *package* of the form pack ρ with e as $\exists (t,\tau)$, where ρ is a type and e is an expression of type $[\rho/t]\tau$. The type ρ is the representation type of the package, and the expression e is the *implementation* of the package. The elimination form is the expression open e_1 as t with x: τ in e_2 , which opens the package e_1 for use within the client e_2 by binding its representation type to t and its implementation to x for use within e_2 . Crucially, the typing rules ensure that the client is type-correct independently of the actual representation type used by the implementor, so that it can be varied without affecting the type correctness of the client.

The abstract syntax of the open construct specifies that the type variable t and the expression variable x are bound within the client. They may be renamed at will by α -equivalence without affecting the meaning of the construct, provided, of course, that the names do not conflict with any others in scope. In other words the type t is a "new" type, one that is distinct from all other types, when it is introduced. This principle is sometimes called generativity of abstract types: the use of an abstract type by a client "generates" a "new" type within that client. This behavior relies on the identification convention stated in Chapter 1.

17.1.1 Statics

The statics of **FE** is given by these rules:

$$\frac{\Delta, t \text{ type} \vdash \tau \text{ type}}{\Delta \vdash \exists (t.\tau) \text{ type}}$$
 (17.1a)

$$\frac{\Delta \vdash \rho \text{ type } \Delta, t \text{ type} \vdash \tau \text{ type } \Delta \Gamma \vdash e : [\rho/t]\tau}{\Delta \Gamma \vdash \text{pack}\{t, \tau\}\{\rho\}(e) : \exists (t, \tau)}$$
(17.1b)

$$\frac{\Delta \vdash \rho \text{ type } \Delta, t \text{ type} \vdash \tau \text{ type } \Delta \Gamma \vdash e : [\rho/t]\tau}{\Delta \Gamma \vdash \text{pack}\{t.\tau\}\{\rho\}(e) : \exists (t.\tau)}$$

$$\frac{\Delta \Gamma \vdash e_1 : \exists (t.\tau) \quad \Delta, t \text{ type } \Gamma, x : \tau \vdash e_2 : \tau_2 \quad \Delta \vdash \tau_2 \text{ type}}{\Delta \Gamma \vdash \text{open}\{t.\tau\}\{\tau_2\}(e_1; t, x.e_2) : \tau_2}$$
(17.1c)

Rule (17.1c) is complex, so study it carefully! There are two important things to notice:

- 1. The type of the client, τ_2 , must not involve the abstract type t. This restriction prevents the client from attempting to export a value of the abstract type outside of the scope of its definition.
- 2. The body of the client, e_2 , is type checked without knowledge of the representation type, t. The client is, in effect, polymorphic in the type variable t.

Lemma 17.1 (Regularity). Suppose that $\Delta \Gamma \vdash e : \tau$. If $\Delta \vdash \tau_i$ type for each $x_i : \tau_i$ in Γ , then $\Delta \vdash \tau$ type. *Proof.* By induction on rules (17.1), using substitution for expressions and types.

17.1.2 Dynamics

The dynamics of FE is defined by the following rules (including the bracketed material for an eager interpretation, and omitting it for a lazy interpretation):

$$\frac{[e \text{ val}]}{\text{pack}\{t.\tau\}\{\rho\}(e) \text{ val}}$$
 (17.2a)

17.2 Data Abstraction 153

$$\left[\frac{e \longmapsto e'}{\operatorname{pack}\{t.\tau\}\{\rho\}(e') \longmapsto \operatorname{pack}\{t.\tau\}\{\rho\}(e')}\right] \tag{17.2b}$$

$$\frac{e_1 \longmapsto e_1'}{\operatorname{open}\{t.\tau\}\{\tau_2\}(e_1;t,x.e_2) \longmapsto \operatorname{open}\{t.\tau\}\{\tau_2\}(e_1';t,x.e_2)}$$
(17.2c)

$$\frac{[e \text{ val}]}{\text{open}\{t.\tau\}\{\tau_2\}(\operatorname{pack}\{t.\tau\}\{\rho\}(e);t,x.e_2) \longmapsto [\rho,e/t,x]e_2}$$
(17.2d)

It is important to see that, according to these rules, *there are no abstract types at run time!* The representation type is propagated to the client by substitution when the package is opened, thereby eliminating the abstraction boundary between the client and the implementor. Thus, data abstraction is a *compile-time discipline* that leaves no traces of its presence at execution time.

17.1.3 Safety

Safety of FE is stated and proved by decomposing it into progress and preservation.

Theorem 17.2 (Preservation). *If* $e : \tau$ *and* $e \mapsto e'$, *then* $e' : \tau$.

Proof. By rule induction on $e \mapsto e'$, using substitution for both expression- and type variables. \Box

Lemma 17.3 (Canonical Forms). *If* $e : \exists (t.\tau)$ *and* e *val, then* $e = \text{pack}\{t.\tau\}\{\rho\}(e')$ *for some type* ρ *and some* e' *such that* $e' : [\rho/t]\tau$.

Proof. By rule induction on the statics, using the definition of closed values. \Box

Theorem 17.4 (Progress). If $e:\tau$ then either e val or there exists e' such that $e\longmapsto e'$.

Proof. By rule induction on $e:\tau$, using the canonical forms lemma.

17.2 Data Abstraction

To illustrate the use of **FE**, we consider an abstract type of queues of natural numbers supporting three operations:

- 1. Forming the empty queue.
- 2. Inserting an element at the tail of the queue.
- 3. Removing the head of the queue, which is assumed non-empty.

This is clearly a bare-bones interface, but suffices to illustrate the main ideas of data abstraction. Queue elements are natural numbers, but nothing depends on this choice.

The crucial property of this description is that nowhere do we specify what queues actually *are*, only what we can *do* with them. The behavior of a queue is expressed by the existential type $\exists (t.\tau)$ which serves as the interface of the queue abstraction:

$$\exists (t. \langle \texttt{emp} \hookrightarrow t, \texttt{ins} \hookrightarrow \texttt{nat} \times t \to t, \texttt{rem} \hookrightarrow t \to (\texttt{nat} \times t) \texttt{opt} \rangle).$$

The representation type t of queues is abstract — all that is known about it is that it supports the operations emp, ins, and rem, with the given types.

An implementation of queues consists of a package specifying the representation type, together with the implementation of the associated operations in terms of that representation. Internally to the implementation, the representation of queues is known and relied upon by the operations. Here is a very simple implementation e_l in which queues are represented as lists:

The elided body of e_i conses the first component of x, the element, onto the second component of x, the queue, and the elided body of e_r reverses its argument, and returns the head element paired with the reversal of the tail. Both of these operations "know" that queues are represented as values of type natlist, and are programmed accordingly.

It is also possible to give another implementation e_p of the same interface $\exists (t.\tau)$, but in which queues are represented as pairs of lists, consisting of the "back half" of the queue paired with the reversal of the "front half". This two-part representation avoids the need for reversals on each call, and, as a result, achieves amortized constant-time behavior:

```
 \begin{array}{l} \operatorname{pack} \operatorname{natlist} \times \operatorname{natlist} \operatorname{with} \langle \operatorname{emp} \hookrightarrow \langle \operatorname{nil}, \operatorname{nil} \rangle, \operatorname{ins} \hookrightarrow e_i, \operatorname{rem} \hookrightarrow e_r \rangle \operatorname{as} \exists (\ t.\tau \ ). \\ \\ \operatorname{In} \operatorname{this} \operatorname{case} e_i \operatorname{has} \operatorname{type} \\ \\ \operatorname{nat} \times (\operatorname{natlist} \times \operatorname{natlist}) \to (\operatorname{natlist} \times \operatorname{natlist}), \\ \\ \operatorname{and} e_r \operatorname{has} \operatorname{type} \\ \\ (\operatorname{natlist} \times \operatorname{natlist}) \to \operatorname{nat} \times (\operatorname{natlist} \times \operatorname{natlist}). \\ \end{array}
```

These operations "know" that queues are represented as values of type $natlist \times natlist$, and are implemented accordingly.

The important point is that the *same* client type checks regardless of which implementation of queues we choose, because the representation type is hidden, or *held abstract*, from the client during type checking. Consequently, it cannot rely on whether it is natlist or natlist × natlist or some other type. That is, the client is *independent* of the representation of the abstract type.

17.3 Definability of Existential Types

The language **FE** is not a proper extension of **F**, because existential types (under a lazy dynamics) are definable in terms of universal types. Why should this be possible? Note that the client of an abstract type is *polymorphic* in the representation type. The typing rule for

```
open e_1 as t with x:\tau in e_2:\tau_2,
```

where $e_1: \exists (t.\tau)$, specifies that $e_2: \tau_2$ under the assumptions t type and $x: \tau$. In essence, the client is a polymorphic function of type

$$\forall (t.\tau \rightarrow \tau_2),$$

where t may occur in τ (the type of the operations), but not in τ_2 (the type of the result). This suggests the following encoding of existential types:

$$\exists (t.\tau) \triangleq \forall (u.\forall (t.\tau \to u) \to u)$$

$$\operatorname{pack} \rho \operatorname{with} e \operatorname{as} \exists (t.\tau) \triangleq \Lambda(u) \lambda (x : \forall (t.\tau \to u)) x [\rho](e)$$

$$\operatorname{open} e_1 \operatorname{as} t \operatorname{with} x : \tau \operatorname{in} e_2 \triangleq e_1[\tau_2] (\Lambda(t) \lambda (x : \tau) e_2)$$

An existential is encoded as a polymorphic function taking the overall result type u as argument, followed by a polymorphic function representing the client with result type u, and yielding a value of type u as overall result. Consequently, the open construct simply packages the client as such a polymorphic function, instantiates the existential at the result type, τ_2 , and applies it to the polymorphic client. (The translation therefore depends on knowing the overall result type τ_2 of the open construct.) Finally, a package consisting of a representation type ρ and an implementation e is a polymorphic function that, when given the result type u and the client x instantiates x with ρ and passes to it the implementation e.

17.4 Representation Independence

An important consequence of parametricity is that it ensures that clients are insensitive to the representations of abstract types. More precisely, there is a criterion, *bisimilarity*, for relating two implementations of an abstract type such that the behavior of a client is unaffected by swapping one implementation by another that is bisimilar to it. This principle leads to a simple method for proving the correctness of a *candidate* implementation of an abstract type, which is to show that it is bisimilar to an obviously correct *reference* implementation of it. Because the candidate and the reference implementations are bisimilar, no client may distinguish them from one another, and hence if the client behaves properly with the reference implementation, then it must also behave properly with the candidate.

To derive the definition of bisimilarity of implementations, it is helpful to examine the definition of existential types in terms of universals given in Section 17.3. It is immediately clear that the client of an abstract type is polymorphic in the representation of the abstract type. A client c of an abstract type $\exists (t.\tau)$ has type $\forall (t.\tau \to \tau_2)$, where t does not occur free in τ_2 (but may, of course, occur in τ). Applying the parametricity property described informally in Chapter 16 (and developed rigorously in Chapter 48), this says that if R is a bisimulation relation between any two implementations of the abstract type, then the client behaves identically on them. The fact that t does not occur in the result type ensures that the behavior of the client is independent of the choice of relation between the implementations, provided that this relation is preserved by the operations that implement it.

Explaining what is a bisimulation is best done by example. Consider the existential type $\exists (t.\tau)$, where τ is the labeled tuple type

$$\langle \texttt{emp} \hookrightarrow t, \texttt{ins} \hookrightarrow \texttt{nat} \times t \rightarrow t, \texttt{rem} \hookrightarrow t \rightarrow (\texttt{nat} \times t) \texttt{opt} \rangle.$$

This specifies an abstract type of queues. The operations emp, ins, and rem specify, respectively, the empty queue, an insert operation, and a remove operation. For the sake of simplicity the element type is the type of natural numbers. The result of removal is an optional pair, according to whether the queue is empty or not.

Theorem 48.12 ensures that if ρ and ρ' are any two closed types, and if R is a relation between expressions of these two types, then if the implementations $e: [\rho/x]\tau$ and $e': [\rho'/x]\tau$ respect R, then $c[\rho]e$ behaves the same as $c[\rho']e'$. It remains to define when two implementations respect the relation R. Let

$$e \triangleq \langle \texttt{emp} \,{\hookrightarrow}\, e_{\texttt{m}}, \texttt{ins} \,{\hookrightarrow}\, e_{\texttt{i}}, \texttt{rem} \,{\hookrightarrow}\, e_{\texttt{r}} \rangle$$

and

$$e' \triangleq \langle emp \hookrightarrow e'_{\mathsf{m}}, ins \hookrightarrow e'_{\mathsf{i}}, rem \hookrightarrow e'_{\mathsf{r}} \rangle.$$

For these implementations to respect *R* means that the following three conditions hold:

- 1. The empty queues are related: $R(e_{m}, e'_{m})$.
- 2. Inserting the same element on each of two related queues yields related queues: if $d : \tau$ and R(q, q'), then $R(e_i(d)(q), e_i'(d)(q'))$.
- 3. If two queues are related, then either they are both empty, or their front elements are the same and their back elements are related: if R(q, q'), then either

(a)
$$e_r(q) \cong \text{null} \cong e'_r(q')$$
, or
(b) $e_r(q) \cong \text{just}(\langle d, r \rangle)$ and $e'_r(q') \cong \text{just}(\langle d', r' \rangle)$, with $d \cong d'$ and $R(r, r')$.

If such a relation R exists, then the implementations e and e' are bisimilar. The terminology stems from the requirement that the operations of the abstract type preserve the relation: if it holds before an operation is performed, then it must also hold afterwards, and the relation must hold for the initial state of the queue. Thus each implementation simulates the other up to the relationship specified by R.

To see how this works in practice, let us consider informally two implementations of the abstract type of queues defined earlier. For the reference implementation we choose ρ to be the type natlist, and define the empty queue to be the empty list, define insert to add the given element to the head of the list, and define remove to remove the last element of the list. The code is as follows:

```
t \triangleq \mathtt{natlist} \mathtt{emp} \triangleq \mathtt{nil} \mathtt{ins} \triangleq \lambda \, (\, x : \mathtt{nat} \,) \, \lambda \, (\, q : t \,) \, \mathtt{cons} \, (\, x ; q \,) \mathtt{rem} \triangleq \lambda \, (\, q : t \,) \, \mathtt{case} \, \mathtt{rev} \, (\, q \,) \, \{ \mathtt{nil} \hookrightarrow \mathtt{null} \, | \, \mathtt{cons} \, (\, f ; q r \,) \hookrightarrow \mathtt{just} \, (\, \langle f, \mathtt{rev} \, (\, q r \,) \rangle \,) \}.
```

17.5 Notes 157

Removing an element takes time linear in the length of the list, because of the reversal.

For the candidate implementation we choose ρ' to be the type natlist \times natlist of pairs of lists $\langle b, f \rangle$ in which b is the "back half" of the queue, and f is the reversal of the "front half" of the queue. For this representation we define the empty queue to be a pair of empty lists, define insert to extend the back with that element at the head, and define remove based on whether the front is empty or not. If it is non-empty, the head element is removed from it, and returned along with the pair consisting of the back and the tail of the front. If it is empty, and the back is not, then we reverse the back, remove the head element, and return the pair consisting of the empty list and the tail of the now-reversed back. The code is as follows:

```
\begin{split} t &\triangleq \mathtt{natlist} \times \mathtt{natlist} \\ & \texttt{emp} \triangleq \langle \mathtt{nil}, \mathtt{nil} \rangle \\ & \texttt{ins} \triangleq \lambda \left( x : \mathtt{nat} \right) \lambda \left( \langle bs, fs \rangle : t \right) \langle \mathtt{cons}(x; bs), fs \rangle \\ & \texttt{rem} \triangleq \lambda \left( \langle bs, fs \rangle : t \right) \mathtt{case} \, fs \, \{ \mathtt{nil} \hookrightarrow e \, | \, \mathtt{cons}(f; fs') \hookrightarrow \mathtt{just}(\langle f, \langle bs, fs' \rangle \rangle \, ) \}, \, \text{where} \\ & e \triangleq \mathtt{case} \, \mathtt{rev}(bs) \, \{ \mathtt{nil} \hookrightarrow \mathtt{null} \, | \, \mathtt{cons}(b; bs') \hookrightarrow \mathtt{just}(\langle b, \langle \mathtt{nil}, bs' \rangle \rangle \, ) \}. \end{split}
```

The cost of the occasional reversal is amortized across the sequence of inserts and removes to show that each operation in a sequence costs unit time overall.

To show that the candidate implementation is correct, we show that it is bisimilar to the reference implementation. To do so, we specify a relation R between the types natlist and natlist \times natlist such that the two implementations satisfy the three simulation conditions given earlier. The required relation states that $R(l,\langle b,f\rangle)$ iff the list l is the list app $(b)(\operatorname{rev}(f))$, where app is the evident append function on lists. That is, thinking of l as the reference representation of the queue, the candidate must ensure that the elements of b followed by the elements of b in reverse order form precisely the list b. It is easy to check that the implementations just described preserve this relation. Having done so, we are assured that the client b behaves the same regardless of whether we use the reference or the candidate. Because the reference implementation is obviously correct (albeit inefficient), the candidate must also be correct in that the behavior of any client is not affected by using it instead of the reference.

17.5 Notes

The connection between abstract types in programming languages and existential types in logic was made by Mitchell and Plotkin (1988). Closely related ideas were already present in Reynolds (1974), but the connection with existential types was not explicitly drawn there. The present formulation of representation independence follows closely Mitchell (1986).

Exercises

17.1. Show that the statics and dynamics of existential types are correctly simulated using the interpretation given in Section 17.3.

158 17.5 Notes

17.2. Define in **FE** of the coinductive type of streams of natural numbers as defined in Chapter 15. *Hint*: Define the representation in terms of the generator (introduction form) for streams.

- **17.3**. Define in **FE** an arbitrary coinductive type $v(t.\tau)$. *Hint*: generalize your answer to Exercise **17.2**.
- **17.4.** Representation independence for abstract types is a corollary of the parametricity theorem for polymorphic types, using the interpretation of **FE** in **F** given in Section 17.3. Recast the proof of equivalence of the two implementations of queues given in Section 17.4 as an instance of parametricity as defined informally in Chapter 16.

164 17.5 Notes



Part VIII Partiality and Recursive Types



Chapter 19

System PCF of Recursive Functions

We introduced the language **T** as a basis for discussing total computations, those for which the type system guarantees termination. The language **M** generalizes **T** to admit inductive and coinductive types, while preserving totality. In this chapter we introduce **PCF** as a basis for discussing partial computations, those that may not terminate when evaluated, even when they are well-typed. At first blush this may seem like a disadvantage, but as we shall see in Chapter 20 it admits greater expressive power than is possible in **T**.

The source of partiality in **PCF** is the concept of *general recursion*, which permits the solution of equations between expressions. The price for admitting solutions to all such equations is that computations may not terminate—the solution to some equations might be undefined (divergent). In **PCF** the programmer must make sure that a computation terminates; the type system does not guarantee it. The advantage is that the termination proof need not be embedded into the code itself, resulting in shorter programs.

For example, consider the equations

$$f(0) \triangleq 1$$
$$f(n+1) \triangleq (n+1) \times f(n).$$

Intuitively, these equations define the factorial function. They form a system of simultaneous equations in the unknown f which ranges over functions on the natural numbers. The function we seek is a *solution* to these equations—a specific function $f: \mathbb{N} \to \mathbb{N}$ such that the above conditions are satisfied.

A solution to such a system of equations is a fixed point of an associated functional (higher-order function). To see this, let us re-write these equations in another form:

$$f(n) \triangleq \begin{cases} 1 & \text{if } n = 0 \\ n \times f(n') & \text{if } n = n' + 1. \end{cases}$$

Re-writing yet again, we seek *f* given by

$$n \mapsto \begin{cases} 1 & \text{if } n = 0 \\ n \times f(n') & \text{if } n = n' + 1. \end{cases}$$

Now define the *functional* F by the equation F(f) = f', where f' is given by

$$n \mapsto \begin{cases} 1 & \text{if } n = 0 \\ n \times f(n') & \text{if } n = n' + 1. \end{cases}$$

Note well that the condition on f' is expressed in terms of f, the argument to the functional F, and not in terms of f' itself! The function f we seek is a *fixed point* of F, a function $f: \mathbb{N} \to \mathbb{N}$ such that f = F(f). In other words f is defined to be fix(F), where fix is a higher-order operator on functionals F that computes a fixed point for it.

Why should an operator such as F have a fixed point? The key is that functions in **PCF** are *partial*, which means that they may diverge on some (or even all) inputs. Consequently, a fixed point of a functional F is the limit of a series of approximations of the desired solution obtained by iterating F. Let us say that a partial function ϕ on the natural numbers, is an *approximation* to a total function f if $\phi(m) = n$ implies that f(m) = n. Let $\bot : \mathbb{N} \to \mathbb{N}$ be the totally undefined partial function— $\bot(n)$ is undefined for every $n \in \mathbb{N}$. This is the "worst" approximation to the desired solution f of the recursion equations given above. Given any approximation ϕ of f, we may "improve" it to $\phi' = F(\phi)$. The partial function ϕ' is defined on 0 and on m+1 for every $m \ge 0$ on which ϕ is defined. Continuing, $\phi'' = F(\phi') = F(F(\phi))$ is an improvement on ϕ' , and hence a further improvement on ϕ . If we start with \bot as the initial approximation to f, then pass to the limit

$$\lim_{i>0}F^{(i)}(\bot),$$

we will obtain the least approximation to f that is defined for every $m \in \mathbb{N}$, and hence is the function f itself. Turning this around, if the limit exists, it is the solution we seek.

Because this construction works for any functional *F*, we conclude that *all* such operators have fixed points, and hence that *all* systems of equations such as the one given above have solutions. The solution is given by general recursion, but there is no guarantee that it is a total function (defined on all elements of its domain). For the above example it happens to be true, because we can prove by induction that this is so, but in general the solution is a partial function that may diverge on some inputs. It is our task as programmers to ensure that the functions defined by general recursion are total, or at least that we have a grasp of those inputs for which it is well-defined.

19.1 Statics 169

19.1 Statics

The syntax of **PCF** is given by the following grammar:

The expression $fix\{\tau\}(x.e)$ is *general recursion*; it is discussed in more detail below. The expression $ifz\{e_0; x.e_1\}(e)$ branches according to whether e evaluates to z or not, binding the predecessor to x in the case that it is not.

The statics of **PCF** is inductively defined by the following rules:

$$\overline{\Gamma, x : \tau \vdash x : \tau} \tag{19.1a}$$

$$\frac{}{\Gamma \vdash z : nat}$$
 (19.1b)

$$\frac{\Gamma \vdash e : \mathtt{nat}}{\Gamma \vdash \mathtt{s}(e) : \mathtt{nat}} \tag{19.1c}$$

$$\frac{\Gamma \vdash e : \mathtt{nat} \quad \Gamma \vdash e_0 : \tau \quad \Gamma, x : \mathtt{nat} \vdash e_1 : \tau}{\Gamma \vdash \mathtt{ifz}\{e_0; x.e_1\}(e) : \tau}$$
(19.1d)

$$\frac{\Gamma, x : \tau_1 \vdash e : \tau_2}{\Gamma \vdash \lambda \{\tau_1\}(x.e) : \rightharpoonup (\tau_1; \tau_2)}$$
(19.1e)

$$\frac{\Gamma \vdash e_1 : \rightharpoonup (\tau_2; \tau) \quad \Gamma \vdash e_2 : \tau_2}{\Gamma \vdash \operatorname{ap}(e_1; e_2) : \tau}$$
(19.1f)

$$\frac{\Gamma, x : \tau \vdash e : \tau}{\Gamma \vdash \text{fix}\{\tau\}(x.e) : \tau}$$
 (19.1g)

Rule (19.1g) reflects the self-referential nature of general recursion. To show that $fix\{\tau\}$ (x.e) has type τ , we *assume* that it is the case by assigning that type to the variable x, which stands for the recursive expression itself, and checking that the body, e, has type τ under this very assumption.

The structural rules, including in particular substitution, are admissible for the statics.

Lemma 19.1. *If* Γ , $x : \tau \vdash e' : \tau'$, $\Gamma \vdash e : \tau$, then $\Gamma \vdash [e/x]e' : \tau'$.

170 19.2 Dynamics

19.2 Dynamics

The dynamics of **PCF** is defined by the judgments e val, specifying the closed values, and $e \mapsto e'$, specifying the steps of evaluation.

The judgment e val is defined by the following rules:

$$\frac{}{z \text{ val}}$$
 (19.2a)

$$\frac{[e \text{ val}]}{\mathsf{s}(e) \text{ val}} \tag{19.2b}$$

$$\frac{1}{\lambda\{\tau\}(x.e) \text{ val}} \tag{19.2c}$$

The bracketed premise on rule (19.2b) is included for the *eager* interpretation of the successor operation, and omitted for the *lazy* interpretation. (See Chapter 36 for a further discussion of laziness.) The transition judgment $e \mapsto e^t$ is defined by the following rules:

$$\frac{e \longmapsto e'}{\mathbf{s}(e) \longmapsto \mathbf{s}(e')}$$
 (19.3a)

$$\frac{e \longmapsto e'}{\text{ifz}\{e_0; x.e_1\}(e) \longmapsto \text{ifz}\{e_0; x.e_1\}(e')}$$
(19.3b)

$$\frac{}{\text{ifz}\{e_0; x.e_1\}(z) \longmapsto e_0} \tag{19.3c}$$

$$\frac{\mathbf{s}(e) \text{ val}}{\mathbf{ifz}\{e_0; x.e_1\}(\mathbf{s}(e)) \longmapsto [e/x]e_1}$$
(19.3d)

$$\frac{e_1 \longmapsto e_1'}{\operatorname{ap}(e_1; e_2) \longmapsto \operatorname{ap}(e_1'; e_2)}$$
(19.3e)

$$\left[\frac{e_1 \text{ val} \quad e_2 \longmapsto e_2'}{\operatorname{ap}(e_1; e_2) \longmapsto \operatorname{ap}(e_1; e_2')} \right]$$
(19.3f)

$$\frac{[e_2 \text{ val}]}{\operatorname{ap}(\lambda\{\tau\}(x.e); e_2) \longmapsto [e_2/x]e}$$
 (19.3g)

$$\frac{1}{\operatorname{fix}\{\tau\}(x.e) \longmapsto [\operatorname{fix}\{\tau\}(x.e)/x]e} \tag{19.3h}$$

The bracketed rule (19.3a) is included for an eager interpretation of the successor, and omitted otherwise. Bracketed rule (19.3f) and the bracketed premise on rule (19.3g) are included for a call-by-value interpretation, and omitted for a call-by-name interpretation, of function application. Rule (19.3h) implements self-reference by substituting the recursive expression itself for the variable x in its body; this is called unwinding the recursion.

19.3 Definability 171

Theorem 19.2 (Safety).

- 1. If $e : \tau$ and $e \longmapsto e'$, then $e' : \tau$.
- 2. If $e: \tau$, then either e val or there exists e' such that $e \mapsto e'$.

Proof. The proof of preservation is by induction on the derivation of the transition judgment. Consider rule (19.3h). Suppose that $fix\{\tau\}(x.e):\tau$. By inversion and substitution we have $[fix\{\tau\}(x.e)/x]e:\tau$, as required. The proof of progress proceeds by induction on the derivation of the typing judgment. For example, for rule (19.1g) the result follows because we may make progress by unwinding the recursion.

It is easy to check that if e val, then e is irreducible in that there is no e' such that $e \mapsto e'$. The safety theorem implies the converse, that an irreducible expression is a value, provided that it is closed and well-typed.

Definitional equality for the call-by-name variant of **PCF**, written $\Gamma \vdash e_1 \equiv e_2 : \tau$, is the strongest congruence containing the following axioms:

$$\frac{\Gamma \vdash \text{ifz}\{e_0; x.e_1\}(\mathbf{z}) \equiv e_0 : \tau}{\Gamma \vdash \text{ifz}\{e_0; x.e_1\}(\mathbf{z}) \equiv e_0 : \tau}$$
(19.4a)

$$\frac{\Gamma \vdash \text{ifz}\{e_0; x.e_1\}(s(e)) \equiv [e/x|e_1:\tau]}{\Gamma \vdash \text{ifz}\{e_0; x.e_1\}(s(e)) \equiv [e/x|e_1:\tau]}$$
(19.4b)

$$\frac{1}{\Gamma \vdash \text{fix}\{\tau\}(x.e) \equiv [\text{fix}\{\tau\}(x.e)/x]e : \tau}$$
(19.4c)

$$\Gamma \vdash \operatorname{ap}(\lambda\{\tau_1\}(x.e_2); e_1) \equiv [e_1/x]e_2 : \tau$$
(19.4d)

These rules suffice to calculate the value of any closed expression of type nat: if e: nat, then $e \equiv \overline{n}$: nat iff $e \mapsto^* \overline{n}$.

19.3 Definability

Let us write $\operatorname{fun} x(y:\tau_1):\tau_2$ is e for a recursive function within whose body, $e:\tau_2$, are bound two variables, $y:\tau_1$ standing for the argument and $x:\tau_1 \rightharpoonup \tau_2$ standing for the function itself. The dynamics of this construct is given by the axiom

$$\overline{(\operatorname{fun} x(y:\tau_1):\tau_2 \operatorname{is} e)(e_1) \longmapsto [\operatorname{fun} x(y:\tau_1):\tau_2 \operatorname{is} e, e_1/x, y]e}$$

That is, to apply a recursive function, we substitute the recursive function itself for x and the argument for y in its body.

172 19.3 Definability

Recursive functions are defined in **PCF** using fixed points, writing

$$fix x : \tau_1 \rightharpoonup \tau_2 is \lambda (y : \tau_1) e$$

for $\operatorname{fun} x(y:\tau_1):\tau_2$ is e. We may easily check that the static and dynamics of recursive functions are derivable from this definition.

The primitive recursion construct of **T** is defined in **PCF** using recursive functions by taking the expression

$$\operatorname{rec} e \left\{ \mathbf{z} \hookrightarrow e_0 \mid \mathbf{s}(x) \text{ with } y \hookrightarrow e_1 \right\}$$

to stand for the application e'(e), where e' is the general recursive function

$$fun f(u:nat):\tau is ifz u \{z \hookrightarrow e_0 \mid s(x) \hookrightarrow [f(x)/y]e_1\}.$$

The static and dynamics of primitive recursion are derivable in **PCF** using this expansion.

In general, functions definable in **PCF** are partial in that they may be undefined for some arguments. A partial (mathematical) function, $\phi: \mathbb{N} \to \mathbb{N}$, is *definable* in **PCF** iff there is an expression $e_{\phi}: \mathtt{nat} \to \mathtt{nat}$ such that $\phi(m) = n$ iff $e_{\phi}(\overline{m}) \equiv \overline{n}: \mathtt{nat}$. So, for example, if ϕ is the totally undefined function, then e_{ϕ} is any function that loops without returning when it is applied.

It is informative to classify those partial functions ϕ that are definable in **PCF**. The *partial recursive functions* are defined to be the primitive recursive functions extended with the *minimization* operation: given $\phi(m,n)$, define $\psi(n)$ to be the least $m \geq 0$ such that (1) for m' < m, $\phi(m',n)$ is defined and non-zero, and (2) $\phi(m,n) = 0$. If no such m exists, then $\psi(n)$ is undefined.

Theorem 19.3. A partial function ϕ on the natural numbers is definable in **PCF** iff it is partial recursive.

Proof sketch. Minimization is definable in **PCF**, so it is at least as powerful as the set of partial recursive functions. Conversely, we may, with some tedium, define an evaluator for expressions of **PCF** as a partial recursive function, using Gödel-numbering to represent expressions as numbers. Therefore **PCF** does not exceed the power of the set of partial recursive functions. □

Church's Law states that the partial recursive functions coincide with the set of effectively computable functions on the natural numbers—those that can be carried out by a program written in any programming language that is or will ever be defined.¹ Therefore **PCF** is as powerful as any other programming language with respect to the set of definable functions on the natural numbers.

The universal function ϕ_{univ} for **PCF** is the partial function on the natural numbers defined by

$$\phi_{univ}(\lceil e \rceil)(m) = n \text{ iff } e(\overline{m}) \equiv \overline{n} : \text{nat.}$$

In contrast to \mathbf{T} , the universal function ϕ_{univ} for **PCF** is partial (might be undefined for some inputs). It is, in essence, an interpreter that, given the code $\overline{}e^{-}$ of a closed expression of type nat \rightarrow nat, simulates the dynamics to calculate the result, if any, of applying it to the \overline{m} , obtaining \overline{n} . Because this process may fail to terminate, the universal function is not defined for all inputs.

¹See Chapter 21 for further discussion of Church's Law.

By Church's Law the universal function is definable in **PCF**. In contrast, we proved in Chapter 9 that the analogous function is *not* definable in **T** using the technique of diagonalization. It is instructive to examine why that argument does not apply in the present setting. As in Section 9.4, we may derive the equivalence

$$e_{\Delta}(\overline{\lceil e_{\Delta} \rceil}) \equiv \mathtt{s}(\,e_{\Delta}(\overline{\lceil e_{\Delta} \rceil}\,)\,)$$

for **PCF**. But now, instead of concluding that the universal function, e_{univ} , does not exist as we did for **T**, we instead conclude for **PCF** that e_{univ} diverges on the code for e_{Λ} applied to its own code.

19.4 Finite and Infinite Data Structures

Finite data types (products and sums), including their use in pattern matching and generic programming, carry over verbatim to **PCF**. However, the distinction between the eager and lazy dynamics for these constructs becomes more important. Rather than being a matter of preference, the decision to use an eager or lazy dynamics affects the meaning of a program: the "same" types mean different things in a lazy dynamics than in an eager dynamics. For example, the elements of a product type in an eager language are pairs of values of the component types. In a lazy language they are instead pairs of unevaluated, possibly divergent, computations of the component types, a very different thing indeed. And similarly for sums.

The situation grows more acute for infinite types such as the type nat of "natural numbers." The scare quotes are warranted, because the "same" type has a very different meaning under an eager dynamics than under a lazy dynamics. In the former case the type nat is, indeed, the authentic type of natural numbers—the least type containing zero and closed under successor. The principle of mathematical induction is valid for reasoning about the type nat in an eager dynamics. It corresponds to the inductive type nat defined in Chapter 15.

On the other hand, under a lazy dynamics the type nat is no longer the type of natural numbers at all. For example, it includes the value

$$\omega \triangleq \texttt{fix}\,x:\texttt{natiss}(x),$$

which has itself as predecessor! It is, intuitively, an "infinite stack of successors", growing without end. It is clearly not a natural number (it is larger than all of them), so the principle of mathematical induction does not apply. In a lazy setting nat could be renamed lnat to remind us of the distinction; it corresponds to the type conat defined in Chapter 15.

19.5 Totality and Partiality

The advantage of a total programming language such as **T** is that it ensures, by type checking, that every program terminates, and that every function is total. There is no way to have a well-typed program that goes into an infinite loop. This prohibition may seem appealing, until one considers that the upper bound on the time to termination may be large, so large that it might as well diverge for all practical purposes. But let us grant for the moment that it is a virtue of **T**

that it precludes divergence. Why, then, bother with a language such as **PCF** that does not rule out divergence? After all, infinite loops are invariably bugs, so why not rule them out by type checking? The notion seems appealing until one tries to write a program in a language such as **T**.

Consider the computation of the greatest common divisor (gcd) of two natural numbers. It can be programmed in **PCF** by solving the following equations using general recursion:

```
gcd(m,0) = m
gcd(0,n) = n
gcd(m,n) = gcd(m-n,n) \text{ if } m > n
gcd(m,n) = gcd(m,n-m) \text{ if } m < n
```

The type of gcd defined this way is $(nat \times nat) \rightarrow nat$, which suggests that it may not terminate for some inputs. But we may prove by induction on the sum of the pair of arguments that it is, in fact, a total function.

Now consider programming this function in **T**. It is, in fact, programmable using only primitive recursion, but the code to do it is rather painful (try it!). One way to see the problem is that in **T** the only form of looping is one that reduces a natural number by one on each recursive call; it is not (directly) possible to make a recursive call on a smaller number other than the immediate predecessor. In fact one may code up more general patterns of terminating recursion using only primitive recursion as a primitive, but if you check the details, you will see that doing so comes at a price in performance and program complexity. Program complexity can be mitigated by building libraries that codify standard patterns of reasoning whose cost of development should be amortized over all programs, not just one in particular. But there is still the problem of performance. Indeed, the encoding of more general forms of recursion into primitive recursion means that, deep within the encoding, there must be a "timer" that goes down by ones to ensure that the program terminates. The result will be that programs written with such libraries will be slower than necessary.

But, one may argue, **T** is simply not a serious language. A more serious total programming language would admit sophisticated patterns of control without performance penalty. Indeed, one could easily envision representing the natural numbers in binary, rather than unary, and allowing recursive calls by halving to get logarithmic complexity. Such a formulation is possible, as would be quite a number of analogous ideas that avoid the awkwardness of programming in **T**. Could we not then have a practical language that rules out divergence?

We can, but at a cost. We have already seen one limitation of total programming languages: they are not universal. You cannot write an interpreter for T within T, and this limitation extends to any total language whatever. If this does not seem important, then consider the *Blum Size Theorem (BST)*, which places another limitation on total languages. Fix *any* total language \mathcal{L} that permits writing functions on the natural numbers. Pick any blowup factor, say 2^{2^n} . The BST states that there is a total function on the natural numbers that is programmable in \mathcal{L} , but whose shortest program in \mathcal{L} is larger by the given blowup factor than its shortest program in **PCF**!

The underlying idea of the proof is that in a total language the proof of termination of a program must be baked into the code itself, whereas in a partial language the termination proof is an external verification condition left to the programmer. There are, and always will be, programs whose termination proof is rather complicated to express, if you fix in advance the means of proving it total. (In **T** it was

19.6 Notes 175

primitive recursion, but one can be more ambitious, yet still get caught by the BST.) But if you leave room for ingenuity, then programs can be short, because they do not have to embed the proof of their termination in their own running code.

19.6 Notes

The solution to recursion equations described here is based on Kleene's fixed point theorem for complete partial orders, specialized to the approximation ordering of partial functions. The language **PCF** is derived from Plotkin (1977) as a laboratory for the study of semantics of programming languages. Many authors have used PCF as the subject of study of many problems in semantics. It has thereby become the *E. coli* of programming languages.

Exercises

- **19.1.** Consider the problem considered in Section 10.3 of how to define the mutually recursive "even" and "odd" functions. There we gave a solution in terms of primitive recursion. You are, instead, to give a solution in terms of general recursion. *Hint*: consider that a pair of mutually recursive functions is a recursive pair of functions.
- **19.2**. Show that minimization, as explained before the statement of Theorem 19.3, is definable in **PCF**.
- **19.3**. Consider the partial function ϕ_{halts} such that if e: nat \rightarrow nat, then $\phi_{halts}(\lceil e \rceil)$ evaluates to zero iff $e(\lceil e \rceil)$ converges, and evaluates to one otherwise. Prove that ϕ_{halts} is not definable in **PCF**.
- **19.4**. Suppose that we changed the specification of minimization given prior to Theorem 19.3 so that $\psi(n)$ is the least m such that $\phi(m,n)=0$, and is undefined if no such m exists. Is this "simplified" form of minimization definable in **PCF**?
- **19.5**. Suppose that we wished to define, in the lazy variant of **PCF**, a version of the *parallel or* function specified a function of two arguments that returns z if either of its arguments is z, and s(z) otherwise. That is, we wish to find an expression *e* satisfying the following properties:

$$e(e_1)(e_2) \longrightarrow^* \mathbf{z} \text{ if } e_1 \longrightarrow^* \mathbf{z}$$

 $e(e_1)(e_2) \longrightarrow^* \mathbf{z} \text{ if } e_2 \longrightarrow^* \mathbf{z}$
 $e(e_1)(e_2) \longrightarrow^* \mathbf{s}(\mathbf{z}) \text{ otherwise}$

Thus, *e* defines a total function of its two arguments, *even if* one of the arguments diverges. Clearly such a function cannot be defined in the call-by-value variant of **PCF**, but can it be defined in the call-by-name variant? If so, show how; if not, prove that it cannot be, and suggest an extension of **PCF** that would allow it to be defined.

176 19.6 Notes

19.6. We appealed to Church's Law to argue that the universal function for **PCF** is definable in **PCF**. See what is behind this claim by considering two aspects of the problem: (1) Gödelnumbering, the representation of abstract syntax by a number; (2) evaluation, the process of interpreting a function on its inputs. Part (1) is a technical issue arising from the limited data structures available in **PCF**. Part (2) is the heart of the matter; explore its implementation in terms of a solution to Part (1).



Chapter 20

System FPC of Recursive Types

In this chapter we study **FPC**, a language with products, sums, partial functions, and *recursive types*. Recursive types are solutions to type equations $t \cong \tau$ where there is no restriction on where t may occur in τ . Equivalently, a recursive type is a *fixed point* up to isomorphism of the associated unrestricted type operator $t.\tau$. By removing the restrictions on the type operator we may consider the solution of a type equation such as $t \cong t \rightharpoonup t$, which describes a type that is isomorphic to the type of partial functions defined on itself. If types were sets, such an equation could not be solved, because there are more partial functions on a set than there are elements of that set. But *types are not sets*: they classify *computable* functions, not *arbitrary* functions. With types we may solve such "dubious" type equations, even though we cannot expect to do so with sets. The penalty is that we must admit non-termination. For one thing, type equations involving functions have solutions only if the functions involved are partial.

A benefit of working in the setting of partial functions is that type equations have *unique* solutions (up to isomorphism). Therefore it makes sense, as we shall do in this chapter, to speak of *the* solution to a type equation. But what about the *distinct* solutions to a type equation given in Chapter 15? These turn out to coincide for any fixed dynamics, but give rise to different solutions according to whether the dynamics is eager or lazy (as illustrated in Section 19.4 for the special case of the natural numbers). Under a lazy dynamics (where all constructs are evaluated lazily), recursive types have a coinductive flavor, and the inductive analogs are inaccessible. Under an eager dynamics (where all constructs are evaluated eagerly), recursive types have an inductive flavor. But the coinductive analogs are accessible as well, using function types to selectively impose laziness. It follows that the eager dynamics is *more expressive* than the lazy dynamics, because it is impossible to go the other way around (one cannot define inductive types in a lazy language).

20.1 Solving Type Equations

The language **FPC** has products, sums, and partial functions inherited from the preceding development, extended with the new concept of recursive types. The syntax of recursive types is

defined as follows:

Recursive types have the same general form as the inductive and coinductive types discussed in Chapter 15, but without restriction on the type operator involved. Recursive type are formed according to the rule:

$$\frac{\Delta, t \text{ type} \vdash \tau \text{ type}}{\Delta \vdash \text{rec}(t.\tau) \text{ type}}$$
 (20.1)

The statics of folding and unfolding is given by the following rules:

$$\frac{\Gamma \vdash e : [\operatorname{rec}(t.\tau)/t]\tau}{\Gamma \vdash \operatorname{fold}\{t.\tau\}(e) : \operatorname{rec}(t.\tau)}$$
 (20.2a)

$$\frac{\Gamma \vdash e : \operatorname{rec}(t.\tau)}{\Gamma \vdash \operatorname{unfold}(e) : [\operatorname{rec}(t.\tau)/t]\tau}$$
 (20.2b)

The dynamics of folding and unfolding is given by these rules:

$$\frac{[e \text{ val}]}{\text{fold}\{t.\tau\}(e) \text{ val}}$$
 (20.3a)

$$\left[\frac{e \longmapsto e'}{\operatorname{fold}\{t.\tau\}(e) \longmapsto \operatorname{fold}\{t.\tau\}(e')}\right] \tag{20.3b}$$

$$\frac{e \longmapsto e'}{\operatorname{unfold}(e) \longmapsto \operatorname{unfold}(e')}$$
 (20.3c)

$$\frac{\operatorname{fold}\{t.\tau\}(e)\operatorname{val}}{\operatorname{unfold}(\operatorname{fold}\{t.\tau\}(e))\longmapsto e} \tag{20.3d}$$

The bracketed premise and rule are included for an *eager* interpretation of the introduction form, and omitted for a *lazy* interpretation. As mentioned in the introduction, the choice of eager or lazy dynamics affects the meaning of recursive types.

Theorem 20.1 (Safety). 1. If $e : \tau$ and $e \mapsto e'$, then $e' : \tau$.

2. If $e: \tau$, then either e val, or there exists e' such that $e \mapsto e'$.

20.2 Inductive and Coinductive Types

Recursive types may be used to represent inductive types such as the natural numbers. Using an *eager* dynamics for **FPC**, the recursive type

$$\rho = \mathtt{rec}\,t\,\mathtt{is}\,[\,\mathtt{z}\,{\hookrightarrow}\,\mathtt{unit},\mathtt{s}\,{\hookrightarrow}\,t\,]$$

satisfies the type equation

$$\rho \cong [z \hookrightarrow \mathtt{unit}, s \hookrightarrow \rho],$$

and is isomorphic to the type of eager natural numbers. The introduction and elimination forms are defined on ρ by the following equations:¹

$$\begin{split} \mathbf{z} &\triangleq \mathtt{fold}(\,\mathbf{z} \cdot \langle \rangle \,) \\ \mathbf{s}(\,e\,) &\triangleq \mathtt{fold}(\,\mathbf{s} \cdot e\,) \\ \mathtt{ifz}\,e\, \{\mathbf{z} \,\hookrightarrow \, e_0 \mid \mathbf{s}(\,x\,) \,\hookrightarrow \, e_1\} &\triangleq \, \mathtt{case}\, \mathtt{unfold}(\,e\,)\, \{\mathbf{z} \,\cdot\, _\, \hookrightarrow \, e_0 \mid \, \mathbf{s} \,\cdot\, x \,\hookrightarrow \, e_1\}. \end{split}$$

It is a good exercise to check that the eager dynamics of natural numbers in **PCF** is correctly simulated by these definitions.

On the other hand, under a lazy dynamics for **FPC**, the same recursive type ρ' ,

$$rec t is [z \hookrightarrow unit, s \hookrightarrow t],$$

satisfies the same type equation,

$$ho'\cong [\, \mathtt{z}\hookrightarrow \mathtt{unit}, \mathtt{s}\hookrightarrow
ho'\,],$$

but is not the type of natural numbers! Rather, it is the type lnat of lazy natural numbers introduced in Section 19.4. As discussed there, the type ρ' contains the "infinite number" ω , which is of course not a natural number.

Similarly, using an eager dynamics for **FPC**, the type natlist of lists of natural numbers is defined by the recursive type

$$rec t is [n \hookrightarrow unit, c \hookrightarrow nat \times t],$$

which satisfies the type equation

$$\mathtt{natlist} \cong [\mathtt{n} \hookrightarrow \mathtt{unit}, \mathtt{c} \hookrightarrow \mathtt{nat} \times \mathtt{natlist}].$$

The list introduction operations are given by the following equations:

$$\begin{split} \text{nil} &\triangleq \texttt{fold}(\, \texttt{n} \cdot \left\langle \right\rangle \,) \\ &\texttt{cons}(\, e_1; e_2 \,) \triangleq \texttt{fold}(\, \texttt{c} \cdot \left\langle e_1, e_2 \right\rangle \,). \end{split}$$

¹The "underscore" stands for a variable that does not occur free in e_0 .

180 20.3 Self-Reference

A conditional list elimination form is given by the following equation:

$$\mathsf{case}\,e\,\{\mathsf{nil}\,\hookrightarrow\,e_0\,|\,\mathsf{cons}(\,x;y\,)\,\hookrightarrow\,e_1\}\,\triangleq\,\,\mathsf{case}\,\mathsf{unfold}(\,e\,)\,\{\mathsf{n}\,\cdot\,_\,\hookrightarrow\,e_0\,|\,\,\mathsf{c}\,\cdot\,\langle x,y\rangle\,\hookrightarrow\,e_1\},$$

where we have used pattern-matching syntax to bind the components of a pair for the sake of clarity.

Now consider the *same* recursive type, but in the context of a lazy dynamics for **FPC**. What type is it? If all constructs are lazy, then a value of the recursive type

$$\mathtt{rec}\,t\,\mathtt{is}\,[\,\mathtt{n}\,\hookrightarrow\,\mathtt{unit},\mathtt{c}\,\hookrightarrow\,\mathtt{nat}\, imes\,t\,],$$

has the form fold(e), where e is an unevaluated computation of the sum type, whose values are injections of unevaluated computations of either the unit type or of the product type $nat \times t$. And the latter consists of pairs of an unevaluated computation of a (lazy!) natural number, and an unevaluated computation of another value of this type. In particular, this type contains infinite lists whose tails go on without end, as well as finite lists that eventually reach an end. The type is, in fact, a version of the type of infinite streams defined in Chapter 15, rather than a type of finite lists as is the case under an eager dynamics.

It is common in textbooks to depict data structures using "box-and-pointer" diagrams. These work well in the eager setting, provided that no functions are involved. For example, an eager list of eager natural numbers may be depicted using this notation. We may think of fold as an abstract pointer to a tagged cell consisting of either (a) the tag n with no associated data, or (b) the tag c attached to a pair consisting of an authentic natural number and another list, which is an abstract pointer of the same type. But this notation does not scale well to types involving functions, or to languages with a lazy dynamics. For example, the recursive type of "lists" in lazy **FPC** cannot be depicted using boxes and pointers, because of the unevaluated computations occurring in values of this type. It is a mistake to limit one's conception of data structures to those that can be drawn on the blackboard using boxes and pointers or similar informal notations. There is no substitute for a programming language to express data structures fully and accurately.

It is deceiving that the "same" recursive type can have two different meanings according to whether the underlying dynamics is eager or lazy. For example, it is common for lazy languages to use the name "list" for the recursive type of streams, or the name "nat" for the type of lazy natural numbers. This terminology is misleading, considering that such languages do not (and can not) have a proper type of finite lists or a type of natural numbers. *Caveat emptor!*

20.3 Self-Reference

In the general recursive expression $fix\{\tau\}(x.e)$ the variable x stands for the expression itself. Self-reference is effected by the unrolling transition

$$fix\{\tau\}(x.e) \longmapsto [fix\{\tau\}(x.e)/x]e$$
,

which substitutes the expression itself for x in its body during execution. It is useful to think of x as an *implicit argument* to e that is instantiated to itself when the expression is used. In many well-known languages this implicit argument has a special name, such as this or self, to emphasize its self-referential interpretation.

20.3 Self-Reference

Using this intuition as a guide, we may derive general recursion from recursive types. This derivation shows that general recursion may, like other language features, be seen as a manifestation of type structure, instead of as an *ad hoc* language feature. The derivation isolates a type of self-referential expressions given by the following grammar:

The statics of these constructs is given by the following rules:

$$\frac{\Gamma, x : self(\tau) \vdash e : \tau}{\Gamma \vdash self(\tau) (x.e) : self(\tau)}$$
 (20.4a)

$$\frac{\Gamma \vdash e : \mathtt{self}(\tau)}{\Gamma \vdash \mathtt{unroll}(e) : \tau} \tag{20.4b}$$

The dynamics is given by the following rule for unrolling the self-reference:

$$\frac{}{\operatorname{self}\{\tau\}(x.e)\operatorname{val}} \tag{20.5a}$$

$$\frac{e \longmapsto e'}{\operatorname{unroll}(e) \longmapsto \operatorname{unroll}(e')} \tag{20.5b}$$

$$\frac{1}{\text{unroll}(\text{self}\{\tau\}(x.e)) \longmapsto [\text{self}\{\tau\}(x.e)/x]e}$$
 (20.5c)

The main difference, compared to general recursion, is that we distinguish a type of self-referential expressions, instead of having self-reference at every type. However, as we shall see, the self-referential type suffices to implement general recursion, so the difference is a matter of taste.

The type $self(\tau)$ is definable from recursive types. As suggested earlier, the key is to consider a self-referential expression of type τ to depend on the expression itself. That is, we seek to define the type $self(\tau)$ so that it satisfies the isomorphism

$$\operatorname{self}(\tau) \cong \operatorname{self}(\tau) \rightharpoonup \tau.$$

We seek a fixed point of the type operator $t.t \rightarrow \tau$, where $t \notin \tau$ is a type variable standing for the type in question. The required fixed point is just the recursive type

$$rec(t.t \rightarrow \tau)$$
,

which we take as the definition of $self(\tau)$.

The self-referential expression $self\{\tau\}(x.e)$ is the expression

$$fold(\lambda(x:self(\tau))e).$$

We may check that rule (20.4a) is derivable according to this definition. The expression unroll(e) is correspondingly the expression

$$unfold(e)(e)$$
.

It is easy to check that rule (20.4b) is derivable from this definition. Moreover, we may check that

unroll(self{
$$\tau$$
}($y.e$)) \longrightarrow^* [self{ τ }($y.e$)/ y] e .

This completes the derivation of the type $self(\tau)$ of self-referential expressions of type τ .

The self-referential type $self(\tau)$ can be used to define general recursion for *any* type. We may define $fix\{\tau\}(x.e)$ to stand for the expression

$$unroll(self\{\tau\}(y.[unroll(y)/x]e))$$

where the recursion at each occurrence of x is unrolled within e. It is easy to check that this verifies the statics of general recursion given in Chapter 19. Moreover, it also validates the dynamics, as shown by the following derivation:

$$\begin{split} \operatorname{fix} \{\tau\}(\, x.e \,) &= \operatorname{unroll}(\, \operatorname{self} \{\tau\}(\, y.[\operatorname{unroll}(\, y \,)/x]e \,) \,) \\ &\longmapsto^* [\operatorname{unroll}(\, \operatorname{self} \{\tau\}(\, y.[\operatorname{unroll}(\, y \,)/x]e \,) \,)/x]e \\ &= [\operatorname{fix} \{\tau\}(\, x.e \,)/x]e. \end{split}$$

It follows that recursive types can be used to define a non-terminating expression of every type, $fix\{\tau\}(x.x)$.

20.4 The Origin of State

The concept of *state* in a computation—which will be discussed in Part XIV—has its origins in the concept of recursion, or self-reference, which, as we have just seen, arises from the concept of recursive types. For example, the concept of a *flip-flop* or a *latch* is a circuit built from combinational logic elements (typically, nor or nand gates) that have the characteristic that they maintain an alterable state over time. An RS latch, for example, maintains its output at the logical level of zero or one in response to a signal on the R or S inputs, respectively, after a brief settling delay. This behavior is achieved using *feedback*, which is just a form of self-reference, or recursion: the output of the gate feeds back into its input so as to convey the current state of the gate to the logic that determines its next state.

We can implement an RS latch using recursive types. The idea is to use self-reference to model the passage of time, with the current output being computed from its input and its previous outputs. Specifically, an RS latch is a value of type τ_{rsl} given by

$$\operatorname{rec} t \text{ is } \langle \mathtt{X} \hookrightarrow \mathtt{bool}, \mathtt{Q} \hookrightarrow \mathtt{bool}, \mathtt{N} \hookrightarrow t \rangle.$$

The X and Q components of the latch represent its current outputs (of which Q represents the current state of the latch), and the N component represents the next state of the latch. If e is of type τ_{rsl} ,

20.5 Notes 183

then we define e @ X to mean unfold(e) · X, and define e @ Q and e @ N similarly. The expressions e @ X and e @ Q evaluate to the boolean outputs of the latch e, and e @ N evaluates to another latch representing its evolution over time based on these inputs.

For given values r and s, a new latch is computed from an old latch by the recursive function rsl defined as follows:²

$$fix rsl is \lambda(l:\tau_{rsl})e_{rsl}$$

where e_{rsl} is the expression

$$\texttt{fix}\,\textit{this}\, \texttt{is}\, \texttt{fold}(\, \langle \texttt{X} \hookrightarrow e_{nor}(\, \langle s, l @ \, \mathsf{Q} \rangle \,), \mathsf{Q} \hookrightarrow e_{nor}(\, \langle r, l \, @ \, \mathsf{X} \rangle \,), \mathsf{N} \hookrightarrow \textit{rsl}(\, \textit{this} \,) \rangle \,),$$

where e_{nor} is the obvious binary function on booleans. The outputs of the latch are computed in terms of the r and s inputs and the outputs of the previous state of the latch. To get the construction started, we define an initial state of the latch in which the outputs are arbitrarily set to false, and whose next state is determined by applying the recursive function rsl to that state:

$$fix this is fold(\langle X \hookrightarrow false, Q \hookrightarrow false, N \hookrightarrow rsl(this)\rangle).$$

Selection of the N component causes the outputs to be recalculated based on their current values. Notice the role of self-reference in maintaining the state of the latch.

20.5 Notes

The systematic study of recursive types in programming was initiated by Scott (1976, 1982) to give a mathematical model of the untyped λ -calculus. The derivation of recursion from recursive types is an application of Scott's theory. The category-theoretic view of recursive types was developed by Wand (1979) and Smyth and Plotkin (1982). Implementing state using self-reference is fundamental to digital logic (Ward and Halstead, 1990). The example given in Section 20.4 is inspired by Cook (2009) and Abadi and Cardelli (1996). The account of signals as streams (explored in the exercises) is inspired by the pioneering work of Kahn (MacQueen, 2009). The language name **FPC** is taken from Gunter (1992).

Exercises

20.1. Show that the recursive type $D \triangleq \text{rec } t \text{ is } t \rightarrow t \text{ is non-trivial by interpreting the sk-combinators defined in Exercise$ **3.1**into it. Specifically, define elements <math>k:D and s:D and a (left-associative) "application" function

$$x : Dy : D \vdash x \cdot y : D$$

such that

²For convenience we assume that fold is evaluated lazily.

184 20.5 Notes

(a)
$$k \cdot x \cdot y \longmapsto^* x$$
;
(b) $s \cdot x \cdot y \cdot z \longmapsto^* (x \cdot z) \cdot (y \cdot z)$.

20.2. Recursive types admit the structure of both inductive and coinductive types. Consider the recursive type $\tau \triangleq \mathtt{rec}\,t$ is τ' and the associated inductive and coinductive types $\mu(t.\tau')$ and $\nu(t.\tau')$. Complete the following chart consistently with the statics of inductive and coinductive types on the left-hand side and with the statics of recursive types on the right:

$$ext{fold}\{t.t \operatorname{opt}\}(e) \triangleq ext{fold}(e)$$
 $ext{rec}\{t.t \operatorname{opt}\}(x.e';e) \triangleq ?$
 $ext{unfold}\{t.t \operatorname{opt}\}(e) \triangleq ext{unfold}(e)$
 $ext{gen}\{t.t \operatorname{opt}\}(x.e';e) \triangleq ?$

Check that the statics is derivable under these definitions. *Hint*: you will need to use general recursion on the right to fill in the missing cases. You may also find it useful to use generic programming.

Now consider the dynamics of these definitions, under both an eager and a lazy interpretation. What happens in each case?

20.3. Define the type signal of *signals* to be the coinductive type of infinite streams of booleans (bits). Define a *signal transducer* to be a function of type signal — signal. Combinational logic gates, such as the NOR gate, can be defined as signal transducers. Give a coinductive definition of the type signal, and define NOR as a signal transducer. Be sure to take account of the underlying dynamics of **PCF**.

The passage from combinational to digital logic (circuit elements that maintain state) relies on self-reference. For example, an RS latch can be built from NOR two nor gates in this way. Define an RS latch using general recursion and two of the NOR gates just defined.

20.4. The type τ_{rsl} given in Section 20.4 above is the type of streams of pairs of booleans. Give another formulation of an RS latch as a value of type τ_{rsl} , but this time using the coinductive interpretation of the recursive type proposed in Exercise **20.2** (using the lazy dynamics for **FPC**). Expand and simplify this definition using your solution to Exercise **20.2**, and compare it with the formulation given in Section 20.4. *Hint*: the internal state of the stream is a pair of booleans corresponding to the X and Q outputs of the latch.

Part IX Dynamic Types



Chapter 21

The Untyped λ -Calculus

In this chapter we study the premier example of a uni-typed programming language, the (untyped) λ -calculus. This formalism was introduced by Church in the 1930's as a universal language of computable functions. It is distinctive for its austere elegance. The λ -calculus has but one "feature", the higher-order function. Everything is a function, hence every expression may be applied to an argument, which must itself be a function, with the result also being a function. To borrow a turn of phrase, in the λ -calculus it's functions all the way down.

21.1 The λ -Calculus

The abstract syntax of the untyped λ -calculus, called Λ , is given by the following grammar:

Exp
$$u ::= x$$
 x variable $\lambda(x.u)$ $\lambda(x)u$ λ -abstraction $ap(u_1;u_2)$ $u_1(u_2)$ application

The statics of Λ is defined by general hypothetical judgments of the form x_1 ok,..., x_n ok \vdash u ok, stating that u is a well-formed expression involving the variables $x_1,...,x_n$. (As usual, we omit explicit mention of the variables when they can be determined from the form of the hypotheses.) This relation is inductively defined by the following rules:

$$\frac{}{\Gamma, x \text{ ok} \vdash x \text{ ok}} \tag{21.1a}$$

$$\frac{\Gamma \vdash u_1 \text{ ok} \quad \Gamma \vdash u_2 \text{ ok}}{\Gamma \vdash u_1(u_2) \text{ ok}}$$
 (21.1b)

$$\frac{\Gamma, x \text{ ok} \vdash u \text{ ok}}{\Gamma \vdash \lambda (x) u \text{ ok}}$$
 (21.1c)

188 21.2 Definability

The dynamics of Λ is given equationally, rather than via a transition system. Definitional equality for Λ is a judgment of the form $\Gamma \vdash u \equiv u'$, where $\Gamma \vdash u$ ok and $\Gamma \vdash u'$ ok. It is inductively defined by the following rules:

$$\frac{}{\Gamma \cdot u \text{ ok} \vdash u \equiv u} \tag{21.2a}$$

$$\frac{\Gamma \vdash u \equiv u'}{\Gamma \vdash u' \equiv u} \tag{21.2b}$$

$$\frac{\Gamma \vdash u \equiv u' \quad \Gamma \vdash u' \equiv u''}{\Gamma \vdash u \equiv u''}$$
 (21.2c)

$$\frac{\Gamma \vdash u_1 \equiv u_1' \quad \Gamma \vdash u_2 \equiv u_2'}{\Gamma \vdash u_1(u_2) \equiv u_1'(u_2')}$$
(21.2d)

$$\frac{\Gamma, x \text{ ok} \vdash u \equiv u'}{\Gamma \vdash \lambda(x) u \equiv \lambda(x) u'}$$
(21.2e)

$$\frac{\Gamma, x \text{ ok} \vdash u_2 \text{ ok} \quad \Gamma \vdash u_1 \text{ ok}}{\Gamma \vdash (\lambda(x)u_2)(u_1) \equiv [u_1/x]u_2}$$
(21.2f)

We often write just $u \equiv u'$ when the variables involved need not be emphasized or are clear from context.

21.2 Definability

Interest in the untyped λ -calculus stems from its surprising expressiveness. It is a *Turing-complete* language in the sense that it has the same capability to express computations on the natural numbers as does any other known programming language. Church's Law states that any conceivable notion of computable function on the natural numbers is equivalent to the λ -calculus. This assertion is true for all *known* means of defining computable functions on the natural numbers. The force of Church's Law is that it postulates that all future notions of computation will be equivalent in expressive power (measured by definability of functions on the natural numbers) to the λ -calculus. Church's Law is therefore a *scientific law* in the same sense as, say, Newton's Law of Universal Gravitation, which predicts the outcome of all future measurements of the acceleration in a gravitational field.

We will sketch a proof that the untyped λ -calculus is as powerful as the language PCF described in Chapter 19. The main idea is to show that the PCF primitives for manipulating the natural numbers are definable in the untyped λ -calculus. In particular, we must show that the natural numbers are definable as λ -terms in such a way that case analysis, which discriminates between zero and non-zero numbers, is definable. The principal difficulty is with computing the predecessor of a number, which requires a bit of cleverness. Finally, we show how to represent general recursion, completing the proof.

¹It is debatable whether there are any scientific laws in Computer Science. In the opinion of the author, Church's Law, which is usually called *Church's Thesis*, is a strong candidate for being a scientific law.

189 21.2 Definability

The first task is to represent the natural numbers as certain λ -terms, called the *Church numerals*.

$$\overline{0} \triangleq \lambda(b)\lambda(s)b \tag{21.3a}$$

$$\overline{n+1} \triangleq \lambda(b)\lambda(s)s(\overline{n}(b)(s))$$
 (21.3b)

It follows that

$$\overline{n}(u_1)(u_2) \equiv u_2(\ldots(u_2(u_1))),$$

the *n*-fold application of u_2 to u_1 . That is, \overline{n} iterates its second argument (the induction step) ntimes, starting with its first argument (the basis).

Using this definition it is not difficult to define the basic functions of arithmetic. For example, successor, addition, and multiplication are defined by the following untyped λ -terms:

$$succ \triangleq \lambda(x)\lambda(b)\lambda(s)s(x(b)(s))$$
 (21.4)

$$plus \triangleq \lambda(x)\lambda(y)y(x)(succ)$$
 (21.5)

$$times \triangleq \lambda(x)\lambda(y)y(\overline{0})(plus(x))$$
 (21.6)

It is easy to check that $succ(\bar{n}) \equiv \overline{n+1}$, and that similar correctness conditions hold for the representations of addition and multiplication.

To define $ifz\{u_0; x.u_1\}(u)$ requires a bit of ingenuity. The key is to define the "cut-off predecessor", pred, such that

$$\operatorname{pred}(\bar{0}) \equiv \bar{0} \tag{21.7}$$

$$\operatorname{pred}(\overline{0}) \equiv \overline{0}$$
 (21.7)
$$\operatorname{pred}(\overline{n+1}) \equiv \overline{n}.$$
 (21.8)

To compute the predecessor using Church numerals, we must show how to compute the result for $\overline{n+1}$ in terms of its value for \overline{n} . At first glance this seems simple—just take the successor until we consider the base case, in which we define the predecessor of $\overline{0}$ to be $\overline{0}$. This formulation invalidates the obvious strategy of taking successors at inductive steps, and necessitates some other approach.

What to do? A useful intuition is to think of the computation in terms of a pair of "shift registers" satisfying the invariant that on the nth iteration the registers contain the predecessor of n and n itself, respectively. Given the result for n, namely the pair (n-1,n), we pass to the result for n+1 by shifting left and incrementing to obtain (n, n+1). For the base case, we initialize the registers with (0,0), reflecting the stipulation that the predecessor of zero be zero. To compute the predecessor of n we compute the pair (n-1,n) by this method, and return the first component.

To make this precise, we must first define a Church-style representation of ordered pairs.

$$\langle u_1, u_2 \rangle \triangleq \lambda(f) f(u_1)(u_2) \tag{21.9}$$

$$u \cdot 1 \triangleq u(\lambda(x)\lambda(y)x) \tag{21.10}$$

$$u \cdot \mathbf{r} \triangleq u(\lambda(x)\lambda(y)y) \tag{21.11}$$

190 21.3 Scott's Theorem

It is easy to check that under this encoding $\langle u_1, u_2 \rangle \cdot 1 \equiv u_1$, and that a similar equivalence holds for the second projection. We may now define the required representation, u_p , of the predecessor function:

$$u_p' \triangleq \lambda(x) x(\langle \overline{0}, \overline{0} \rangle) (\lambda(y) \langle y \cdot \mathbf{r}, \text{succ}(y \cdot \mathbf{r}) \rangle)$$
 (21.12)

$$u_p \triangleq \lambda(x) u_p'(x) \cdot 1 \tag{21.13}$$

It is easy to check that this gives us the required behavior. Finally, define $ifz\{u_0; x.u_1\}(u)$ to be the untyped term

$$u(u_0)(\lambda(_-)[u_p(u)/x]u_1).$$

This definition gives us all the apparatus of PCF, apart from general recursion. But general recursion is also definable in Λ using a *fixed point combinator*. There are many choices of fixed point combinator, of which the best known is the Y *combinator*:

$$\mathbf{Y} \triangleq \lambda(F)(\lambda(f)F(f(f)))(\lambda(f)F(f(f))).$$

It is easy to check that

$$Y(F) \equiv F(Y(F)).$$

Using the Y combinator, we may define general recursion by writing $Y(\lambda(x)u)$, where x stands for the recursive expression itself.

Although it is clear that Y as just defined computes a fixed point of its argument, it is probably less clear why it works or how we might have invented it in the first place. The main idea is quite simple. If a function is recursive, it is given an extra first argument, which is arranged at call sites to be the function itself. Whenever we wish to call a self-referential function with an argument, we apply the function first to itself and then to its argument; this protocol is imposed on both the "external" calls to the function and on the "internal" calls that the function may make to itself. For this reason the first argument is often called this or self, to remind you that it will be, by convention, bound to the function itself.

With this in mind, it is easy to see how to derive the definition of Y. If F is the function whose fixed point we seek, then the function $F' = \lambda(f)F(f(f))$ is a variant of F in which the self-application convention has been imposed internally by substituting for each occurrence of f in F(f) the self-application f(f). Now check that $F'(F') \equiv F(F'(F'))$, so that F'(F') is the desired fixed point of F. Expanding the definition of F', we have derived that the desired fixed point of F is

$$\lambda(f)F(f(f))(\lambda(f)F(f(f)).$$

To finish the derivation, we need only note that nothing depends on the particular choice of F, which means that we can compute a fixed point for F uniformly in F. That is, we may define a *single* function, the term Y as defined above, that computes the fixed point of *any* F.

21.3 Scott's Theorem

Scott's Theorem states that definitional equality for the untyped λ -calculus is undecidable: there is no algorithm to determine whether or not two untyped terms are definitionally equal. The proof

21.3 Scott's Theorem 191

uses the concept of *inseparability*. Any two properties, A_0 and A_1 , of λ -terms are *inseparable* if there is no decidable property, \mathcal{B} , such that A_0 u implies that \mathcal{B} u holds, and A_1 u implies that \mathcal{B} u does not hold. We say that a property, \mathcal{A} , of untyped terms is behavioral iff whenever $u \equiv u'$, then \mathcal{A} u iff \mathcal{A} u'.

The proof of Scott's Theorem decomposes into two parts:

- 1. For any untyped λ -term u, we may find an untyped term v such that $u(\lceil v \rceil) \equiv v$, where $\lceil v \rceil$ is the Gödel number of v, and $\lceil v \rceil$ is its representation as a Church numeral. (See Chapter 9 for a discussion of Gödel-numbering.)
- 2. Any two non-trivial² behavioral properties A_0 and A_1 of untyped terms are inseparable.

Lemma 21.1. For any u there exists v such that $u(\overline{\ }\overline{v}) \equiv v$.

Proof Sketch. The proof relies on the definability of the following two operations in the untyped λ -calculus:

1.
$$\operatorname{ap}(\overline{\lceil u_1 \rceil})(\overline{\lceil u_2 \rceil}) \equiv \overline{\lceil u_1(u_2) \rceil}$$
.

2.
$$nm(\overline{n}) \equiv \overline{\overline{n}}$$
.

Intuitively, the first takes the representations of two untyped terms, and builds the representation of the application of one to the other. The second takes a numeral for n, and yields the representation of the Church numeral \overline{n} . Given these, we may find the required term v by defining $v \triangleq w(\lceil \overline{w} \rceil)$, where $w \triangleq \lambda(x) u(\operatorname{ap}(x)(\operatorname{nm}(x)))$. We have

$$v = w(\lceil \overline{w} \rceil)$$

$$\equiv u(\operatorname{ap}(\lceil \overline{w} \rceil)(\operatorname{nm}(\lceil \overline{w} \rceil)))$$

$$\equiv u(\lceil \overline{w} (\lceil \overline{w} \rceil) \rceil)$$

$$\equiv u(\lceil \overline{v} \rceil).$$

The definition is very similar to that of Y(u), except that u takes as input the representation of a term, and we find a v such that, when applied to the representation of v, the term u yields v itself.

Lemma 21.2. Suppose that A_0 and A_1 are two non-trivial behavioral properties of untyped terms. Then there is no untyped term w such that

- 1. For every u either $w(\overline{\lceil u \rceil}) \equiv \overline{0}$ or $w(\overline{\lceil u \rceil}) \equiv \overline{1}$.
- 2. If A_0 u, then $w(\overline{\lceil u \rceil}) \equiv \overline{0}$.
- 3. If A_1 u, then $w(\overline{\lceil u \rceil}) \equiv \overline{1}$.

²A property of untyped terms is *trivial* if it either holds for all untyped terms or never holds for any untyped term.

Proof. Suppose there is such an untyped term w. Let v be the untyped term

$$\lambda(x)$$
ifz $\{u_1; ...u_0\}(w(x))$,

where u_0 and u_1 are chosen such that \mathcal{A}_0 u_0 and \mathcal{A}_1 u_1 . (Such a choice must exist by non-triviality of the properties.) By Lemma 21.1 there is an untyped term t such that $v(\overline{t}) \equiv t$. If $w(\overline{t}) \equiv 0$, then $t \equiv v(\overline{t}) \equiv u_1$, and so \mathcal{A}_1 t, because \mathcal{A}_1 is behavioral and \mathcal{A}_1 u_1 . But then $w(\overline{t}) \equiv 0$ by the defining properties of w, which is a contradiction. Similarly, if $w(\overline{t}) \equiv 0$, again a contradiction.

Corollary 21.3. There is no algorithm to decide whether $u \equiv u'$

Proof. For fixed u, the property \mathcal{E}_u u' defined by $u' \equiv u$ is a non-trivial behavioral property of untyped terms. So it is inseparable from its negation, and hence is undecidable.

21.4 Untyped Means Uni-Typed

The untyped λ -calculus can be faithfully embedded in a typed language with recursive types. Thus every untyped λ -term has a representation as a typed expression in such a way that execution of the representation of a λ -term corresponds to execution of the term itself. This embedding is *not* a matter of writing an interpreter for the λ -calculus in **FPC**, but rather a direct representation of untyped λ -terms as typed expressions in a language with recursive types.

The key observation is that the *untyped* λ -calculus is really the *uni-typed* λ -calculus. It is not the *absence* of types that gives it its power, but rather that it has *only one* type, the recursive type

$$D \triangleq \operatorname{rec} t \operatorname{is} t \rightharpoonup t.$$

A value of type D is of the form fold(e) where e is a value of type $D \rightarrow D$ — a function whose domain and range are both D. Any such function can be regarded as a value of type D by "folding", and any value of type D can be turned into a function by "unfolding". As usual, a recursive type is a solution to a type equation, which in the present case is the equation

$$D \cong D \rightharpoonup D$$
.

This isomorphism specifies that *D* is a type that is isomorphic to the space of partial functions on *D* itself, which is impossible if types are just sets.

This isomorphism leads to the following translation, of Λ into **FPC**:

$$x^{\dagger} \triangleq x \tag{21.14a}$$

$$\lambda(x)u^{\dagger} \triangleq \text{fold}(\lambda(x:D)u^{\dagger}) \tag{21.14b}$$

$$u_1(u_2)^{\dagger} \triangleq \operatorname{unfold}(u_1^{\dagger})(u_2^{\dagger})$$
 (21.14c)

21.5 Notes 193

Note that the embedding of a λ -abstraction is a value, and that the embedding of an application exposes the function being applied by unfolding the recursive type. And so we have

$$\lambda(x) u_1(u_2)^{\dagger} = \operatorname{unfold}(\operatorname{fold}(\lambda(x:D) u_1^{\dagger}))(u_2^{\dagger})$$

$$\equiv \lambda(x:D) u_1^{\dagger}(u_2^{\dagger})$$

$$\equiv [u_2^{\dagger}/x] u_1^{\dagger}$$

$$= ([u_2/x]u_1)^{\dagger}.$$

The last step, stating that the embedding commutes with substitution, is proved by induction on the structure of u_1 . Thus β -reduction is implemented by evaluation of the embedded terms.

Thus we see that the canonical untyped language, Λ , which by dint of terminology stands in opposition to typed languages, turns out to be but a typed language after all. Rather than eliminating types, an untyped language consolidates an infinite collection of types into a single recursive type. Doing so renders static type checking trivial, at the cost of incurring dynamic overhead to coerce values to and from the recursive type. In Chapter 22 we will take this a step further by admitting many different types of data values (not just functions), each of which is a component of a "master" recursive type. This generalization shows that so-called $dynamically\ typed$ languages are, in fact, $statically\ typed$. Thus this traditional distinction cannot be considered an opposition, because dynamic languages are but particular forms of static languages in which undue emphasis is placed on a single recursive type.

21.5 Notes

The untyped λ -calculus was introduced by Church (1941) as a formalization of the informal concept of a computable function. Unlike the well-known machine models, such as the Turing machine or the random access machine, the λ -calculus codifies mathematical and programming practice. Barendregt (1984) is the definitive reference for all aspects of the untyped λ -calculus; the proof of Scott's theorem is adapted from Barendregt's account. Scott (1980a) gave the first model of the untyped λ -calculus in terms of an elegant theory of recursive types. This construction underlies Scott's apt description of the λ -calculus as "uni-typed", rather than "untyped." The idea to characterize Church's Law as such was communicated to the author, independently of each other, by Robert L. Constable and Mark Lillibridge.

Exercises

- **21.1**. Define an encoding of finite products as defined in Chapter 10 in Λ .
- **21.2**. Define the factorial function in Λ two ways, one without using Y, and one using Y. In both cases show that your solution, u, has the property that $u(\overline{n}) \equiv \overline{n!}$.
- **21.3**. Define the "Church booleans" in Λ by defining terms true and false such that

194 21.5 Notes

- (a) true(u_1)(u_2) $\equiv u_1$.
- (b) false(u_1)(u_2) $\equiv u_2$.

What is the encoding of if u then u_1 else u_2 ?

- **21.4**. Define an encoding of finite sums as defined in Chapter 11 in Λ .
- **21.5**. Define an encoding of finite lists of natural numbers as defined in Chapter 15 in Λ .
- **21.6**. Define an encoding of the infinite streams of natural numbers as defined in Chapter 15 in Λ .
- **21.7**. Show that Λ can be "compiled" to sk-combinators using bracket abstraction (see Exercises 3.4 and 3.5. Define a translation u^* from Λ into sk combinators such that

if
$$u_1 \equiv u_2$$
, then $u_1^* \equiv u_2^*$.

Hint: Define u^* by induction on the structure of u, using the compositional form of bracket abstraction considered in Exercise 3.5. Show that the translation is itself compositional in that it commutes with substitution:

$$([u_2/x]u_1)^* = [u_2^*/x]u^*.$$

Then proceed by rule induction on rules (21.2) to show the required correctness condition.

Chapter 22

Dynamic Typing

We saw in Chapter 21 that an untyped language is a uni-typed language in which "untyped" terms are just terms of single recursive type. Because all expressions of Λ are well-typed, type safety ensures that no misinterpretation of a value is possible. When spelled out for Λ , type safety follows from there being exactly one class of values, that of functions on values. No application can get stuck, because every value is a function that may be applied to an argument.

This safety property breaks down once more than one class of value is admitted. For example, if the natural numbers are added as a primitive to Λ , then it is possible to incur a run-time error by attempting to apply a number to an argument. One way to manage this is to embrace the possibility, treating class mismatches as checked errors, and weakening the progress theorem as outlined in Chapter 6. Such languages are called *dynamic languages* because an error such as the one described is postponed to run time, rather than precluded at compile time by type checking. Languages of the latter sort are called *static languages*.

Dynamic languages are often considered in opposition to static languages, but the opposition is illusory. Just as the untyped λ -calculus is uni-typed, so dynamic languages are but special cases of static languages in which there is only one recursive type (albeit with multiple classes of value).

22.1 Dynamically Typed PCF

To illustrate dynamic typing we formulate a dynamically typed version of **PCF**, called **DPCF**. The abstract syntax of **DPCF** is given by the following grammar:

```
Exp \quad d ::= x
                                                       variable
                                         \overline{n}
                                                       numeral
                 num[n]
                 zero
                                                       zero
                                         succ(d)
                                                       successor
                 succ(d)
                 ifz\{d_0; x.d_1\}(d) ifz d \{zero \hookrightarrow d_0 \mid succ(x) \hookrightarrow d_1\}
                                                       zero test
                 fun(x.d)
                                         \lambda(x)d
                                                       abstraction
                 ap(d_1; d_2)
                                         d_1(d_2)
                                                       application
                 fix(x.d)
                                         fix x is d
                                                        recursion
```

There are two classes of values in **DPCF**, the *numbers*, which have the form num[n], and the *functions*, which have the form fun(x.d). The expressions zero and succ(d) are not themselves values, but rather are *constructors* that evaluate to values. General recursion is definable using a fixed point combinator, but is taken as primitive here to simplify the analysis of the dynamics in Section 22.3.

As usual, the abstract syntax of **DPCF** is what matters, but we use the concrete syntax to improve readability. However, notational conveniences can obscure important details, such as the tagging of values with their class and the checking of these tags at run-time. For example, the concrete syntax for a number, \overline{n} , suggests a "bare" representation, the abstract syntax reveals that the number is labeled with the class num to distinguish it from a function. Correspondingly, the concrete syntax for a function is λ (x) d, but its abstract syntax, fun(x.d), shows that it also sports a class label. The class labels are required to ensure safety by run-time checking, and must not be overlooked when comparing static with dynamic languages.

The statics of **DPCF** is like that of Λ ; it merely checks that there are no free variables in the expression. The judgment

$$x_1$$
 ok,... x_n ok $\vdash d$ ok

states that d is a well-formed expression with free variables among those in the hypotheses. If the assumptions are empty, then we write just d ok to mean that d is a closed expression of **DPCF**.

The dynamics of **DPCF** must check for errors that would never arise in a language such as **PCF**. For example, evaluation of a function application must ensure that the value being applied is indeed a function, signaling an error if it is not. Similarly the conditional branch must ensure that its principal argument is a number, signaling an error if it is not. To account for these possibilities,

the dynamics is given by several judgment forms, as summarized in the following chart:

d val	d is a (closed) value
$d \longmapsto d'$	d evaluates in one step to d'
d err	<i>d</i> incurs a run-time error
d is_num n	d is of class num with value n
d isnt_num	d is not of class num
$d \text{ is_fun } x.d$	d is of class fun with body $x.d$
d isnt_fun	d is not of class fun

The last four judgment forms implement dynamic class checking. They are only relevant when *d* is already a value. The affirmative class-checking judgments have a second argument that represents the underlying structure of a value; this argument is *not* itself an expression of **DPCF**.

The value judgment *d* val states that *d* is a evaluated (closed) expression:

$$\overline{\operatorname{num}[n] \operatorname{val}} \tag{22.1a}$$

$$\overline{\text{fun}(x.d)\text{ val}}$$
 (22.1b)

The affirmative class-checking judgments are defined by the following rules:

$$\overline{\operatorname{num}[n] \text{ is_num } n} \tag{22.2a}$$

$$fun(x.d)$$
 is $fun(x.d)$ (22.2b)

The negative class-checking judgments are correspondingly defined by these rules:

$$\overline{\text{num}[n] \text{ isnt_fun}}$$
 (22.3a)

$$\overline{\text{fun}(x.d) \text{ isnt_num}}$$
 (22.3b)

The transition judgment $d \mapsto d'$ and the error judgment d err are defined simultaneously by the following rules:

$$\underline{\text{zero} \longmapsto \text{num}[z]}$$
(22.4a)

$$\frac{d \longmapsto d'}{\operatorname{succ}(d) \longmapsto \operatorname{succ}(d')} \tag{22.4b}$$

$$\frac{d \operatorname{err}}{\operatorname{succ}(d) \operatorname{err}} \tag{22.4c}$$

$$\frac{d \text{ is_num } n}{\text{succ}(d) \longmapsto \text{num}[s(n)]}$$
 (22.4d)

$$\frac{d \operatorname{isnt_num}}{\operatorname{succ}(d) \operatorname{err}} \tag{22.4e}$$

$$\frac{d \longmapsto d'}{\operatorname{ifz}\{d_0; x.d_1\}(d) \longmapsto \operatorname{ifz}\{d_0; x.d_1\}(d')}$$
 (22.4f)

$$\frac{d \operatorname{err}}{\operatorname{ifz}\{d_0; x.d_1\}(d) \operatorname{err}}$$
 (22.4g)

$$\frac{d \text{ is_num z}}{\text{ifz}\{d_0; x.d_1\}(d) \longmapsto d_0}$$
 (22.4h)

$$\frac{d \text{ is_num s}(n)}{\text{ifz}\{d_0; x.d_1\}(d) \longmapsto [\text{num}[n]/x]d_1}$$
(22.4i)

$$\frac{d \operatorname{isnt_num}}{\operatorname{ifz}\{d_0; x.d_1\}(d) \operatorname{err}}$$
 (22.4j)

$$\frac{d_1 \longmapsto d'_1}{\operatorname{ap}(d_1; d_2) \longmapsto \operatorname{ap}(d'_1; d_2)} \tag{22.4k}$$

$$\left[\frac{d_1 \text{ val} \quad d_2 \longmapsto d_2'}{\operatorname{ap}(d_1; d_2) \longmapsto \operatorname{ap}(d_1; d_2')} \right]$$
(22.41)

$$\frac{d_1 \text{ err}}{\operatorname{ap}(d_1; d_2) \text{ err}} \tag{22.4m}$$

$$\frac{d_1 \text{ is_fun } x.d \quad [d_2 \text{ val}]}{\operatorname{ap}(d_1; d_2) \longmapsto [d_2/x]d}$$
(22.4n)

$$\frac{d_1 \operatorname{isnt_fun}}{\operatorname{ap}(d_1; d_2) \operatorname{err}} \tag{22.40}$$

$$fix(x.d) \mapsto [fix(x.d)/x]d$$
 (22.4p)

Rule (22.4i) labels the predecessor with the class num to maintain the invariant that variables are bound to expressions of **DPCF**.

Lemma 22.1 (Class Checking). If d val, then

- 1. either d is_num n for some n, or d isnt_num;
- 2. either d is_fun x.d' for some x and d', or d isnt_fun.

Proof. By inspection of the rules defining the class-checking judgments.

Theorem 22.2 (Progress). If d ok, then either d val, or d err, or there exists d' such that $d \mapsto d'$.

Proof. By induction on the structure of d. For example, if $d = \verb+succ+(d')$, then we have by induction either d' val, or d' err, or $d' \longmapsto d''$ for some d''. In the last case we have by rule (22.4b) that $\verb+succ+(d') \longmapsto \verb+succ+(d'')$, and in the second-to-last case we have by rule (22.4c) that $\verb+succ+(d') \mapsto \verb+succ+(d') \mapsto \verb+su$

Lemma 22.3 (Exclusivity). For any d in **DPCF**, exactly one of the following holds: d val, or d err, or $d \mapsto d'$ for some d'.

Proof. By induction on the structure of d, making reference to rules (22.4).

22.2 Variations and Extensions

The dynamic language **DPCF** defined in Section 22.1 parallels the static language **PCF** defined in Chapter 19. One discrepancy, however, is in the treatment of natural numbers. Whereas in **PCF** the zero and successor operations are introduction forms for the type nat, in **DPCF** they are elimination forms that act on specially defined numerals. The present formulation uses only a single class of numbers.

One could instead treat zero and succ(d) as values of separate classes, and introduce the obvious class checking judgments for them. When written in this style, the dynamics of the conditional branch is given as follows:

$$\frac{d \longmapsto d'}{\text{ifz}\{d_0; x.d_1\}(d) \longmapsto \text{ifz}\{d_0; x.d_1\}(d')}$$
(22.5a)

$$\frac{d \text{ is_zero}}{\text{ifz}\{d_0; x.d_1\}(d) \longmapsto d_0}$$
 (22.5b)

$$\frac{d \text{ is_succ } d'}{\text{ifz}\{d_0; x.d_1\}(d) \longmapsto [d'/x]d_1}$$
 (22.5c)

$$\frac{d \text{ isnt_zero} \quad d \text{ isnt_succ}}{\text{ifz}\{d_0; x.d_1\}(d) \text{ err}}$$
(22.5d)

Notice that the predecessor of a value of the successor class need not be a number, whereas in the previous formulation this possibility does not arise.

DPCF can be extended with structured data similarly. A classic example is to consider a class nil, consisting of a "null" value, and a class cons, consisting of pairs of values.

The expression $ifnil(d; d_0; x, y.d_1)$ distinguishes the null value from a pair, and signals an error on any other class of value.

Lists (finite sequences) can be encoded using null and pairing. For example, the list consisting of three zeros can be represented by the value

But what to make of the following value?

$$cons(zero; cons(zero; cons(zero; \lambda(x)x)))$$

It is not a list, because it does not end with nil, but it is a permissible value in the enriched language.

mentation of Lisp.

A difficulty with encoding lists using null and pair emerges when defining functions on them. For example, here is a definition of the function append that concatenates two lists:

$$fix a is \lambda(x) \lambda(y) ifnil(x; y; x_1, x_2.cons(x_1; a(x_2)(y)))$$

Nothing prevents us from applying this function to any two values, regardless of whether they are lists. If the first argument is not a list, then execution aborts with an error. But because the function does not traverse its second argument, it can be any value at all. For example, we may apply append with a list and a function to obtain the "list" that ends with a λ given above.

It might be argued that the conditional branch that distinguishes null from a pair is inappropriate in **DPCF**, because there are more than just these two classes in the language. One approach that avoids this criticism is to abandon pattern matching on the class of data, replacing it by a general conditional branch that distinguishes null from all other values, and adding to the language *predicates*¹ that test the class of a value and *destructors* that invert the constructors of each class.

We could instead reformulate null and and pairing as follows:

The conditional $\operatorname{cond}(d;d_0;d_1)$ distinguishes d between nil and all other values. If d is not nil , the conditional evaluates to d_0 , and otherwise evaluates to d_1 . In other words the value nil represents boolean falsehood, and all other values represent boolean truth. The predicates $\operatorname{nil}?(d)$ and $\operatorname{cons}?(d)$ test the class of their argument, yielding nil if the argument is not of the specified class, and yielding some non-nil if so. The destructors $\operatorname{car}(d)$ and $\operatorname{cdr}(d)$ decompose $\operatorname{cons}(d_1;d_2)$ into d_1 and d_2 , respectively.²

Written in this form, the function append is given by the expression

fix
$$a$$
 is $\lambda(x)\lambda(y)$ cond(x ; cons(car(x); a (cdr(x))(y)); y).

The behavior of this formulation of append is no different from the earlier one; the only difference is that instead of dispatching on whether a value is either null or a pair, we instead allow discrimination on any predicate of the value, which includes such checks as special cases.

An alternative, which is not widely used, is to enhance, and not restrict, the conditional branch so that it includes cases for each possible class of value in the language. So in a language with numbers, functions, null, and pairing, the conditional would have four branches. The fourth branch, for pairing, would deconstruct the pair into its constituent parts. The difficulty with this approach is that in realistic languages there are many classes of data, and such a conditional would be rather unwieldy. Moreover, even once we have dispatched on the class of a value, it is nevertheless necessary for the primitive operations associated with that class to admit run-time checks. For example,

¹Predicates evaluate to the null value to mean that a condition is false, and some non-null value to mean that it is true.

²The terminology for the projections is archaic, but well-established. It is said that car originally stood for "contents of the address register" and that cdr stood for "contents of the data register", referring to the details of the original imple-

we may determine that a value d is of the numeric class, but there is no way to propagate this information into the branch of the conditional that then adds d to some other number. The addition operation must still check the class of d, recover the underlying number, and create a new value of numeric class. It is an inherent limitation of dynamic languages that they do not allow values other than classified values.

22.3 Critique of Dynamic Typing

The safety theorem for **DPCF** is an advantage of dynamic over static typing. Unlike static languages, which rule out some candidate programs as ill-typed, every piece of abstract syntax in **DPCF** is well-formed, and hence, by Theorem 22.2, has a well-defined dynamics (albeit one with checked errors). But this convenience is also a disadvantage, because errors that could be ruled out at compile time by type checking are not signaled until run time.

Consider, for example, the addition function in **DPCF**, whose specification is that, when passed two values of class num, returns their sum, which is also of class num:³

$$fun(x.fix(p.fun(y.ifz\{x;y'.succ(p(y'))\}(y)))).$$

The addition function may, deceptively, be written in concrete syntax as follows:

$$\lambda(x)$$
 fix p is $\lambda(y)$ ifz y {zero $\hookrightarrow x$ | succ (y') \hookrightarrow succ $(p(y'))$ }.

It is deceptive, because it obscures the class tags on values, and the operations that check the validity of those tags. Let us now examine the costs of these operations in a bit more detail.

First, note that the body of the fixed point expression is labeled with class fun. The dynamics of the fixed point construct binds p to this function. Consequently, the dynamic class check incurred by the application of p in the recursive call is guaranteed to succeed. But **DPCF** offers no means of suppressing the redundant check, because it cannot express the invariant that p is always bound to a value of class fun.

Second, note that the result of applying the inner λ -abstraction is either x, the argument of the outer λ -abstraction, or the successor of a recursive call to the function itself. The successor operation checks that its argument is of class num, even though this condition is guaranteed to hold for all but the base case, which returns the given x, which can be of any class at all. In principle we can check that x is of class num once, and note that it is otherwise a loop invariant that the result of applying the inner function is of this class. However, **DPCF** gives us no way to express this invariant; the repeated, redundant tag checks imposed by the successor operation cannot be avoided.

Third, the argument *y* to the inner function is either the original argument to the addition function, or is the predecessor of some earlier recursive call. But as long as the original call is to a value of class num, then the dynamics of the conditional will ensure that all recursive calls have this class. And again there is no way to express this invariant in **DPCF**, and hence there is no way to avoid the class check imposed by the conditional branch.

³This specification imposes no restrictions on the behavior of addition on arguments that are not classified as numbers, but we could make the further demand that the function abort when applied to arguments that are not classified by num.

202 22.4 Notes

Classification is not free—storage is required for the class label, and it takes time to detach the class from a value each time it is used and to attach a class to a value when it is created. Although the overhead of classification is not asymptotically significant (it slows down the program only by a constant factor), it is nevertheless non-negligible, and should be eliminated when possible. But this is impossible within **DPCF**, because it cannot enforce the restrictions required to express the required invariants. For that we need a static type system.

22.4 Notes

The earliest dynamically typed language is Lisp (McCarthy, 1965), which continues to influence language design a half century after its invention. Dynamic PCF is the core of Lisp, but with a proper treatment of variable binding, correcting what McCarthy himself has described as an error in the original design. Informal discussions of dynamic languages are often complicated by the elision of the dynamic checks that are made explicit here. Although the surface syntax of dynamic PCF is almost the same as that for PCF, minus the type annotations, the underlying dynamics is different. It is for this reason that static PCF cannot be seen as a restriction of dynamic PCF by the imposition of a type system.

Exercises

- **22.1.** Surface syntax can be deceiving; even simple arithmetic expressions do not have the same meaning in **DPCF** that they do in **PCF**. To see why, define the addition function, plus, in **DPCF**, and examine the dynamics of evaluating expressions such as $plus(\overline{5})(\overline{7})$. Even though this expression might be written as "5 + 7" in both static and dynamic languages, they have different meanings.
- **22.2.** Give a precise dynamics to the data structuring primitives described informally in Section 22.2. What class restrictions should cons impose on its arguments? Check the dynamics of the append function when called with two lists as arguments.
- **22.3**. To avoid the difficulties with the representation of lists using cons and nil, introduce a class of lists that are constructed using revised versions of nil and cons that operate on the class of lists. Revisit the dynamics of the append function when redefined using the class of lists.
- **22.4**. Allowing multiple arguments to, and multiple results from, functions is a notorious source of trouble in dynamic languages. The restriction to a single type makes it impossible even to distinguish *n* things from *m* things, let alone express more subtle properties of a program. Numerous workarounds have been proposed. Explore the problem yourself by enriching **DPCF** with multi-argument and multi-result functions. Be sure to consider these questions:
 - (a) If a function is defined with *n* parameters, what should happen if it is called with more or fewer than *n* arguments?

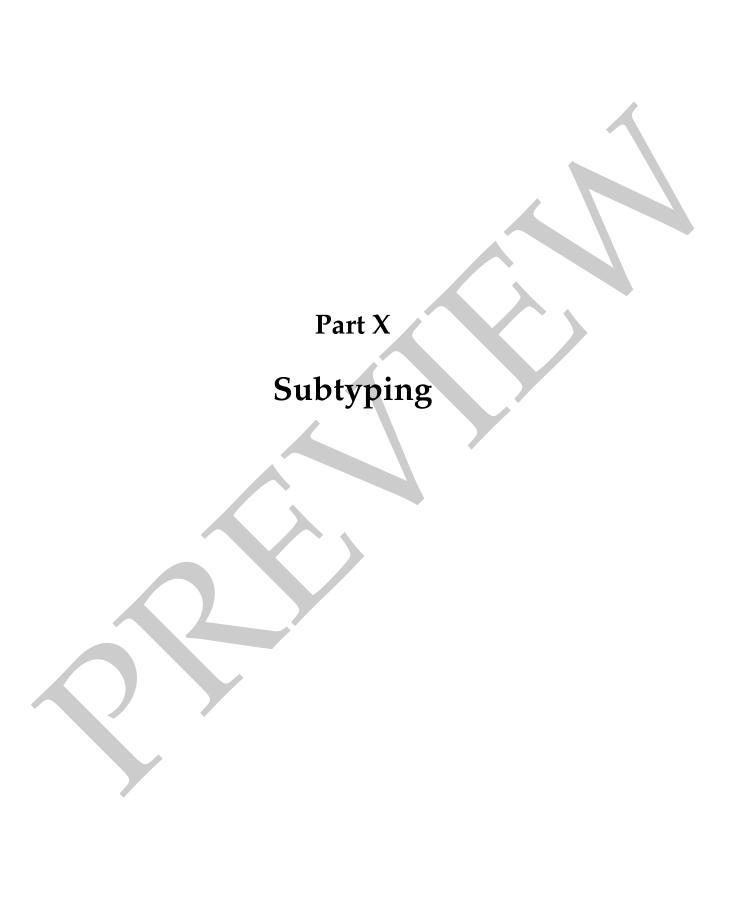
22.4 Notes 203

(b) What happens if one were to admit functions with a *varying* number of arguments? How would you refer to these arguments within the body of such a function? How does this relate to pattern matching?

- (c) What if one wished to admit keyword parameter passing by giving names to the arguments, and allowing them to be passed in any order by associating them with their names?
- (d) What notation would you suggest for functions returning multiple results? For example, a division function might return the quotient and the remainder. How might one notate this in the function body? How would a caller access the results individually or collectively?
- (e) How would one define the composition of two functions when either or both can take multiple arguments or return multiple results?

212 22.4 Notes







Chapter 24

Structural Subtyping

A *subtype* relation is a pre-order (reflexive and transitive relation) on types that validates the *sub-sumption principle*:

if τ' is a subtype of τ , then a value of type τ' may be provided when a value of type τ is required.

The subsumption principle relaxes the strictures of a type system to allow values of one type to be treated as values of another.

Experience shows that the subsumption principle, although useful as a general guide, can be tricky to apply correctly in practice. The key to getting it right is the principle of introduction and elimination. To see whether a candidate subtyping relationship is sensible, it suffices to consider whether every *introduction* form of the subtype can be safely manipulated by every *elimination* form of the supertype. A subtyping principle makes sense only if it passes this test; the proof of the type safety theorem for a given subtyping relation ensures that this is the case.

A good way to get a subtyping principle wrong is to think of a type merely as a set of values (generated by introduction forms), and to consider whether every value of the subtype can also be considered to be a value of the supertype. The intuition behind this approach is to think of subtyping as akin to the subset relation in ordinary mathematics. But, as we shall see, this can lead to serious errors, because it fails to take account of the elimination forms that are applicable to the supertype. It is not enough to think only of the introduction forms; subtyping is a matter of behavior, and not containment.

24.1 Subsumption

A *subtyping judgment* has the form $\tau' <: \tau$, and states that τ' is a subtype of τ . At the least we demand that the following *structural rules* of subtyping be admissible:

$$\overline{\tau <: \tau}$$
 (24.1a)

$$\frac{\tau'' <: \tau' \quad \tau' <: \tau}{\tau'' <: \tau} \tag{24.1b}$$

In practice we either tacitly include these rules as primitive, or prove that they are admissible for a given set of subtyping rules.

The point of a subtyping relation is to enlarge the set of well-typed programs, which is accomplished by the *subsumption rule*:

$$\frac{\Gamma \vdash e : \tau' \quad \tau' <: \tau}{\Gamma \vdash e : \tau} \tag{24.2}$$

In contrast to most other typing rules, the rule of subsumption is *not* syntax-directed, because it does not constrain the form of e. That is, the subsumption rule can be applied to *any* form of expression. In particular, to show that $e:\tau$, we have two choices: either apply the rule appropriate to the particular form of e, or apply the subsumption rule, checking that $e:\tau'$ and $\tau' <: \tau$.

24.2 Varieties of Subtyping

In this section we will informally explore several different forms of subtyping in the context of extensions of the language **FPC** introduced in Chapter 20.

Numeric Types

We begin with an informal discussion of numeric types such as are common in many programming languages. Our mathematical experience suggests subtyping relationships among numeric types. For example, in a language with types int, rat, and real, representing the integers, the rationals, and the reals, it is tempting to postulate the subtyping relationships

by analogy with the set containments

$$\mathbb{Z} \subseteq \mathbb{Q} \subseteq \mathbb{R}$$
.

But are these subtyping relationships sensible? The answer depends on the representations and interpretations of these types. Even in mathematics, the containments just mentioned are usually not true—or are true only in a rough sense. For example, the set of rational numbers can be considered to consist of ordered pairs (m,n), with $n \neq 0$ and $\gcd(m,n) = 1$, representing the ratio m/n. The set $\mathbb Z$ of integers can be isomorphically embedded within $\mathbb Q$ by identifying $n \in \mathbb Z$ with the ratio n/1. Similarly, the real numbers are often represented as convergent sequences of rationals, so that strictly speaking the rationals are not a subset of the reals, but rather can be embedded in them by choosing a canonical representative (a particular convergent sequence) of each rational.

For mathematical purposes it is entirely reasonable to overlook fine distinctions such as that between \mathbb{Z} and its embedding within \mathbb{Q} . Ignoring the difference is justified because the operations on rationals restrict to the embedding in the expected way: if we add two integers thought of as rationals in the canonical way, then the result is the rational associated with their sum. And similarly for the other operations, provided that we take some care in defining them to ensure that it all works out properly. For the purposes of computing, however, we must also take account of algorithmic efficiency and the finiteness of machine representations. For example, what are often

called "real numbers" in a programming language are, of course, floating point numbers, a finite subset of the rational numbers. Not every rational can be exactly represented as a floating point number, nor does floating point arithmetic restrict to rational arithmetic, even when its arguments are exactly represented as floating point numbers.

Product Types

Product types give rise to a form of subtyping based on the subsumption principle. The only elimination form applicable to a value of product type is a projection. Under mild assumptions about the dynamics of projections, we may consider one product type to be a subtype of another by considering whether the projections applicable to the supertype can be validly applied to values of the subtype.

Consider a context in which a value of type $\tau = \langle \tau_j \rangle_{j \in J}$ is required. The statics of finite products (rules (10.3)) ensures that the only operation we may perform on a value of type τ , other than to bind it to a variable, is to take the jth projection from it for some $j \in J$ to obtain a value of type τ_j . Now suppose that e is of type τ' . For the projection $e \cdot j$ to be well-formed, then τ' is a finite product type $\langle \tau'_i \rangle_{i \in I}$ such that $j \in I$. Moreover, for the projection to be of type τ_j , it is enough to require that $\tau'_j = \tau_j$. Because $j \in J$ is arbitrary, we arrive at the following subtyping rule for finite product types:

$$\frac{J \subseteq I}{\prod_{i \in I} \tau_i <: \prod_{j \in I} \tau_j}$$
 (24.3)

This rule sufices for the required subtyping, but not necessary; we will consider a more liberal form of this rule in Section 24.3. The justification for rule (24.3) is that we may evaluate $e \cdot i$ regardless of the actual form of e, provided only that it has a field indexed by $i \in I$.

Sum Types

By an argument dual to the one given for finite product types we may derive a related subtyping rule for finite sum types. If a value of type $\sum_{j\in J} \tau_j$ is required, the statics of sums (rules (11.3)) ensures that the only non-trivial operation that we may perform on that value is a J-indexed case analysis. If we provide a value of type $\sum_{i\in J} \tau_i'$ instead, no difficulty will arise so long as $I\subseteq J$ and each τ_i' is equal to τ_i . If the containment is strict, some cases cannot arise, but this does not disrupt safety.

$$\frac{I \subseteq J}{\sum_{i \in I} \tau_i <: \sum_{j \in J} \tau_j}$$
 (24.4)

Note well the reversal of the containment as compared to rule (24.3).

Dynamic Types

A popular form of subtyping is associated with the type dyn introduced in Chapter 23. The type dyn provides no information about the class of a value of this type. One might argue that it is whole the point of dynamic typing to suppress this information statically, making it available only

218 24.3 Variance

dynamically. On the other hand, it is not much trouble to introduce subtypes of dyn that specify the class of a value, relying on subsumption to "forget" the class when it cannot be determined statically.

Working in the context of Chapter 23 this amounts to introduce two new types, dyn[num] and dyn[fun], governed by the following two subtyping axioms:

$$\frac{1}{\operatorname{dyn}[\operatorname{num}] <: \operatorname{dyn}} \tag{24.5a}$$

$$\frac{1}{\operatorname{dyn}[fun]} <: dyn} \tag{24.5b}$$

Of course, in a richer language with more classes of dynamic values one would correspondingly introduce more such subtypes of dyn, one for each additional class. As a matter of notation, the type dyn is frequently spelled object, and its class-specific subtypes dyn[num] and dyn[fun], are often written as num and fun, respectively. But doing so invites confusion between the separate concepts of *class* and *type*, as discussed in detail in Chapters 22 and 23.

The class-specific subtypes of dyn come into play by reformulating the typing rules for introducing values of type dyn to note the class of the created value:

$$\frac{\Gamma \vdash e : \mathtt{nat}}{\Gamma \vdash \mathtt{new}[\mathtt{num}](e) : \mathtt{dyn}[\mathtt{num}]} \tag{24.6a}$$

$$\frac{\Gamma \vdash e : dyn \rightharpoonup dyn}{\Gamma \vdash new[fun](e) : dyn[fun]}$$
(24.6b)

Thus, in this formulation, classified values "start life" with class-specific types, because in those cases it is statically apparent what is the class of the introduced value. Subsumption is used to weaken the type to dyn in those cases where no static prediction can be made—for example, when the branches of a conditional evaluate to dynamic values of different classes it is necessary to weaken the type of the branches to dyn.

The advantage of such a subtyping mechanism is that we can express more precise types, such as the type $dyn[num] \rightarrow dyn[num]$ of functions mapping a value of type dyn with class num to another such value. This typing is more precise than, say, $dyn \rightarrow dyn$, which merely classifies functions that act on dynamically typed values. In this way weak invariants can be expressed and enforced, but only insofar as it is possible to track the classes of the values involved in a computation. Subtyping is not nearly a powerful enough mechanism for practical situations, rendering the additional specificity not worth the effort of including it. (A more powerful approach is developed in Chapter 25.)

24.3 Variance

In addition to basic subtyping principles such as those considered in Section 24.2, it is also important to consider the effect of subtyping on type constructors. A type constructor *covariant* in an

24.3 Variance 219

argument if subtyping in that argument is preserved by the constructor. It is *contravariant* if subtyping in that argument is reversed by the constructor. It is *invariant* in an argument if subtyping for the constructed type is not affected by subtyping in that argument.

Product and Sum Types

Finite product types are *covariant* in each field. For if e is of type $\prod_{i \in I} \tau'_i$, and the projection $e \cdot j$ is to be of type τ_j , then it suffices to require that $j \in I$ and $\tau'_i <: \tau_j$.

$$\frac{(\forall i \in I) \ \tau_i' <: \tau_i}{\prod_{i \in I} \tau_i' <: \prod_{i \in I} \tau_i}$$
(24.7)

It is implicit in this rule that the dynamics of projection cannot be sensitive to the precise type of any of the fields of a value of finite product type.

Finite sum types are also covariant, because each branch of a case analysis on a value of the supertype expects a value of the corresponding summand, for which it suffices to provide a value of the corresponding subtype summand:

$$\frac{(\forall i \in I) \ \tau_i' <: \tau_i}{\sum_{i \in I} \tau_i' <: \sum_{i \in I} \tau_i}$$
(24.8)

Partial Function Types

The variance of the function type constructors is a bit more subtle. Let us consider first the variance of the function type in its range. Suppose that $e: \tau_1 \rightharpoonup \tau_2'$. Then if $e_1: \tau_1$, then $e(e_1): \tau_2'$, and if $\tau_2' <: \tau_2$, then $e(e_1): \tau_2$ as well.

$$\frac{\tau_2' <: \tau_2}{\tau_1 \rightharpoonup \tau_2' <: \tau_1 \rightharpoonup \tau_2} \tag{24.9}$$

Every function that delivers a value of type τ_2' also delivers a value of type τ_2 , provided that $\tau_2' <: \tau_2$. Thus the function type constructor is covariant in its range.

Now let us consider the variance of the function type in its domain. Suppose again that $e: \tau_1 \rightarrow \tau_2$. Then e can be applied to any value of type τ_1 to obtain a value of type τ_2 . Hence, by the subsumption principle, it can be applied to any value of a subtype τ'_1 of τ_1 , and it will still deliver a value of type τ_2 . Consequently, we may just as well think of e as having type $\tau'_1 \rightarrow \tau_2$.

$$\frac{\tau_1' <: \tau_1}{\tau_1 \rightharpoonup \tau_2 <: \tau_1' \rightharpoonup \tau_2} \tag{24.10}$$

The function type is contravariant in its domain position. Note well the reversal of the subtyping relation in the premise as compared to the conclusion of the rule!

Combining these rules we obtain the following general principle of contra- and covariance for function types:

$$\frac{\tau_1' <: \tau_1 \quad \tau_2' <: \tau_2}{\tau_1 \rightharpoonup \tau_2' <: \tau_1' \rightharpoonup \tau_2} \tag{24.11}$$

220 24.3 Variance

Beware of the reversal of the ordering in the domain!

Recursive Types

The language **FPC** has a partial function types, which behave the same under subtyping as total function types, sums and products, which behave as described above, and recursive types, which introduce some subtleties that have been the source of error in language design. To gain some intuition, consider the type of labeled binary trees with natural numbers at each node,

$$\operatorname{rec} t \text{ is } [\operatorname{empty} \hookrightarrow \operatorname{unit}, \operatorname{binode} \hookrightarrow \langle \operatorname{data} \hookrightarrow \operatorname{nat}, \operatorname{lft} \hookrightarrow t, \operatorname{rht} \hookrightarrow t \rangle],$$

and the type of "bare" binary trees, without data attached to the nodes,

$$\mathtt{rec}\,t\,\mathtt{is}\,[\,\mathtt{empty}\,\hookrightarrow\mathtt{unit},\mathtt{binode}\,\hookrightarrow\langle\mathtt{lft}\,\hookrightarrow t,\mathtt{rht}\,\hookrightarrow t\rangle\,].$$

Is either a subtype of the other? Intuitively, we might expect the type of labeled binary trees to be a *subtype* of the type of bare binary trees, because any use of a bare binary tree can simply ignore the presence of the label.

Now consider the type of bare "two-three" trees with two sorts of nodes, those with two children, and those with three:

$$\operatorname{rec} t$$
 is $[\operatorname{empty} \hookrightarrow \operatorname{unit}, \operatorname{binode} \hookrightarrow \tau_2, \operatorname{trinode} \hookrightarrow \tau_3],$

where

$$au_2 \triangleq \langle \mathtt{lft} \hookrightarrow t, \mathtt{rht} \hookrightarrow t \rangle, \mathtt{and}$$

$$au_3 \triangleq \langle \mathtt{lft} \hookrightarrow t, \mathtt{mid} \hookrightarrow t, \mathtt{rht} \hookrightarrow t \rangle.$$

What subtype relationships should hold between this type and the preceding two tree types? Intuitively the type of bare two-three trees should be a *supertype* of the type of bare binary trees, because any use of a two-three tree proceeds by three-way case analysis, which covers both forms of binary tree.

To capture the pattern illustrated by these examples, we need a subtyping rule for recursive types. It is tempting to consider the following rule:

$$\frac{t \text{ type} \vdash \tau' <: \tau}{\text{rec} t \text{ is } \tau' <: \text{rec} t \text{ is } \tau} ?? \tag{24.12}$$

That is, to check whether one recursive type is a subtype of the other, we simply compare their bodies, with the bound variable treated as an argument. Notice that by reflexivity of subtyping, we have t <: t, and hence we may use this fact in the derivation of $\tau' <: \tau$.

Rule (24.12) validates the intuitively plausible subtyping between labeled binary tree and bare binary trees just described. Deriving this requires checking that the subtyping relationship

$$\langle \mathtt{data} \hookrightarrow \mathtt{nat}, \mathtt{lft} \hookrightarrow t, \mathtt{rht} \hookrightarrow t \rangle <: \langle \mathtt{lft} \hookrightarrow t, \mathtt{rht} \hookrightarrow t \rangle,$$

24.3 Variance 221

holds generically in *t*, which is evidently the case.

Unfortunately, Rule (24.12) also underwrites *incorrect* subtyping relationships, as well as some correct ones. As an example of what goes wrong, consider the recursive types

$$\tau' = \operatorname{rec} t \text{ is } \langle \mathtt{a} \hookrightarrow t \rightharpoonup \mathtt{nat}, \mathtt{b} \hookrightarrow t \rightharpoonup \mathtt{int} \rangle$$

and

$$\tau = \operatorname{rec} t \text{ is } \langle \mathtt{a} \hookrightarrow t \rightharpoonup \operatorname{int}, \mathtt{b} \hookrightarrow t \rightharpoonup \operatorname{int} \rangle.$$

We assume for the sake of the example that nat <: int, so that by using rule (24.12) we may derive $\tau' <: \tau$, which is incorrect. Let $e: \tau'$ be the expression

fold(
$$\langle a \hookrightarrow \lambda (x : \tau') 4, b \hookrightarrow \lambda (x : \tau') q((unfold(x) \cdot a)(x)) \rangle)$$
,

where q: nat \rightharpoonup nat is the discrete square root function. Because $\tau'<:\tau$, it follows that $e:\tau$ as well, and hence

$$unfold(e): \langle a \hookrightarrow \tau \rightharpoonup int, b \hookrightarrow \tau \rightharpoonup int \rangle.$$

Now let e': τ be the expression

fold(
$$\langle a \hookrightarrow \lambda (x : \tau) - 4, b \hookrightarrow \lambda (x : \tau) 0 \rangle$$
).

(The important point about e' is that the a method returns a negative number; the b method is of no significance.) To finish the proof, observe that

$$(\operatorname{unfold}(e) \cdot b)(e') \longmapsto^* q(-4),$$

which is a stuck state. We have derived a well-typed program that "gets stuck", refuting type safety!

Rule (24.12) is therefore incorrect. But what has gone wrong? The error lies in the choice of a single variable to stand for both recursive types, which does not correctly model self-reference. In effect we are treating two distinct recursive types as if they were equal while checking their bodies for a subtyping relationship. But this is clearly wrong! It fails to take account of the self-referential nature of recursive types. On the left side the bound variable stands for the subtype, whereas on the right the bound variable stands for the super-type. Confusing them leads to the unsoundness just illustrated.

As is often the case with self-reference, the solution is to *assume* what we are trying to prove, and check that this assumption can be maintained by examining the bodies of the recursive types. To do so we use hypothetical judgments of the form $\Delta \vdash \tau' <: \tau$, where Δ consists of hypotheses t type and $t <: \tau$ that declares a fresh type variable t that is not otherwise declared in Δ . Using such hypothetical judgments we may state the correct rule for subtyping recursive types as follows:

$$\frac{\Delta, t \text{ type, } t' \text{ type, } t' <: t \vdash \tau' <: \tau \quad \Delta, t' \text{ type} \vdash \tau' \text{ type} \quad \Delta, t \text{ type} \vdash \tau \text{ type}}{\Delta \vdash \text{rec } t' \text{ is } \tau' <: \text{rec } t \text{ is } \tau} . \tag{24.13}$$

That is, to check whether rec t' is $\tau' <: rec t$ is τ , we assume that t' <: t, because t' and t stand for the corresponding recursive types, and check that $\tau' <: \tau$ under this assumption. It is instructive to check that the unsound subtyping example given above is not derivable using this rule: the subtyping assumption is at odds with the contravariance of the function type in its domain.

222 24.3 Variance

Quantified Types

Consider extending **FPC** with the universal and existential quantified types discussed in Chapters 16 and 17. The variance principles for the quantifiers state that they are uniformly covariant in the quantified types:

$$\frac{\Delta, t \text{ type} \vdash \tau' <: \tau}{\Delta \vdash \forall (t.\tau') <: \forall (t.\tau)}$$
(24.14a)

$$\frac{\Delta, t \text{ type} \vdash \tau' <: \tau}{\Delta \vdash \exists (t.\tau') <: \exists (t.\tau)}$$
 (24.14b)

Consequently, we may derive the principle of substitution:

Lemma 24.1. *If* Δ , t type $\vdash \tau_1 <: \tau_2$, and $\Delta \vdash \tau$ type, then $\Delta \vdash [\tau/t]\tau_1 <: [\tau/t]\tau_2$.

Proof. By induction on the subtyping derivation.

It is easy to check that the above variance principles for the quantifiers are consistent with the principle of subsumption. For example, a package of the subtype $\exists (t.\tau')$ consists of a representation type ρ and an implementation e of type $[\rho/t]\tau'$. But if t type $\vdash \tau' <: \tau$, we have by substitution that $[\rho/t]\tau' <: [\rho/t]\tau$, and hence e is also an implementation of type $[\rho/t]\tau$. So the package is also of the supertype.

It is natural to extend subtyping to the quantifiers by allowing quantification over all subtypes of a specified type; this is called *bounded quantification*.

$$\Delta, t \text{ type, } t <: \tau \vdash t <: \tau$$
 (24.15a)

$$\frac{\Delta \vdash \tau :: T}{\Delta \vdash \tau <: \tau} \tag{24.15b}$$

$$\frac{\Delta \vdash \tau'' <: \tau' \quad \Delta \vdash \tau' <: \tau}{\Delta \vdash \tau'' <: \tau}$$
 (24.15c)

$$\frac{\Delta \vdash \tau_1' <: \tau_1 \quad \Delta, t \text{ type, } t <: \tau_1' \vdash \tau_2 <: \tau_2'}{\Delta \vdash \forall t <: \tau_1.\tau_2 <: \forall t <: \tau_1'.\tau_2'}$$

$$(24.15d)$$

$$\frac{\Delta \vdash \tau_1 <: \tau_1' \quad \Delta, t \text{ type, } t <: \tau_1 \vdash \tau_2 <: \tau_2'}{\Delta \vdash \exists \ t <: \tau_1.\tau_2 <: \exists \ t <: \tau_1'.\tau_2'}$$
 (24.15e)

Rule (24.15d) states that the universal quantifier is contravariant in its bound, whereas rule (24.15e) states that the existential quantifier is covariant in its bound.

24.4 Dynamics and Safety

There is a subtle assumption in the definition of product subtyping in Section 24.2, namely that the *same* projection operation from an *I*-tuple applies also to a *J*-tuple, provided $J \supseteq I$. But this need not be the case. One could represent *I*-tuples differently from *J*-tuples at will, so that the meaning of the projection at position $i \in I \subseteq J$ is different in the two cases. Nothing rules out this possibility, yet product subtyping relies on it not being the case. From this point of view product subtyping is not well-justified, but one may instead consider that subtyping limits possible implementations to ensure that it makes sense.

Similar considerations apply to sum types. An J-way case analysis need not be applicable to an I-way value of sum type, even when $I \subseteq J$ and all the types in common agree. For example, one might represent values of a sum type with a "small" index set in a way that is not applicable for a "large" index set. In that case the "large" case analysis would not make sense on a value of "small" sum type. Here again we may consider either that subtyping is not justified, or that it imposes limitations on the implementation that are not otherwise forced.

These considerations merit careful consideration of the safety of languages with subtyping. As an illustrative case we consider the safety of **FPC** enriched with product subtyping. The main concern is that the subsumption rule obscures the "true" type of a value, complicating the canonical forms lemma. Moreover, we assume that the same projection makes sense for a wider tuple than a narrower one, provided that it is within range.

Lemma 24.2 (Structurality).

- 1. The tuple subtyping relation is reflexive and transitive.
- 2. The typing judgment $\Gamma \vdash e : \tau$ is closed under weakening and substitution.

Proof.

- 1. Reflexivity is proved by induction on the structure of types. Transitivity is proved by induction on the derivations of the judgments $\tau'' <: \tau'$ and $\tau' <: \tau$ to obtain a derivation of $\tau'' <: \tau$.
- 2. By induction on rules (10.3), augmented by rule (24.2).

Lemma 24.3 (Inversion).

- 1. If $e \cdot j : \tau$, then $e : \prod_{i \in I} \tau_i$, $j \in I$, and $\tau_i <: \tau$.
- 2. If $\langle e_i \rangle_{i \in I} : \tau$, then $\prod_{i \in I} \tau'_i <: \tau$ where $e_i : \tau'_i$ for each $i \in I$.
- 3. If $\tau' <: \prod_{i \in I} \tau_i$, then $\tau' = \prod_{i \in I} \tau'_i$ for some I and some types τ'_i for $i \in I$.
- 4. If $\prod_{i \in I} \tau'_i <: \prod_{j \in J} \tau_j$, then $J \subseteq I$ and $\tau'_i <: \tau_j$ for each $j \in J$.

Proof. By induction on the subtyping and typing rules, paying special attention to rule (24.2).

224 24.5 Notes

Theorem 24.4 (Preservation). *If* $e : \tau$ *and* $e \mapsto e'$, *then* $e' : \tau$.

Proof. By induction on rules (10.4). For example, consider rule (10.4d), so that $e = \langle e_i \rangle_{i \in I} \cdot k$ and $e' = e_k$. By Lemma 24.3 we have $\langle e_i \rangle_{i \in I} : \prod_{j \in J} \tau_j$, with $k \in J$ and $\tau_k <: \tau$. By another application of Lemma 24.3 for each $i \in I$ there exists τ'_i such that $e_i : \tau'_i$ and $\prod_{i \in I} \tau'_i <: \prod_{j \in J} \tau_j$. By Lemma 24.3 again, we have $J \subseteq I$ and $\tau'_j <: \tau_j$ for each $j \in J$. But then $e_k : \tau_k$, as desired. The remaining cases are similar.

Lemma 24.5 (Canonical Forms). *If* e *val* and e : $\prod_{j \in J} \tau_j$, then e is of the form $\langle e_i \rangle_{i \in I}$, where $J \subseteq I$, and $e_j : \tau_j$ for each $j \in J$.

Proof. By induction on rules (10.3) augmented by rule (24.2). \Box

Theorem 24.6 (Progress). *If* $e:\tau$, then either e val or there exists e' such that $e\longmapsto e'$.

Proof. By induction on rules (10.3) augmented by rule (24.2). The rule of subsumption is handled by appeal to the inductive hypothesis on the premise of the rule. rule (10.4d) follows from Lemma 24.5. \Box

24.5 Notes

Subtyping is perhaps the most widely misunderstood concept in programming languages. Subtyping is principally a convenience, akin to type inference, that makes some programs simpler to write. But the subsumption rule cuts both ways. Inasmuch as it allows the implicit passage from τ' to τ when τ' is a subtype of τ , it also weakens the meaning of a type assertion $e:\tau$ to mean that e has some type contained in the type τ . Subsumption precludes expressing the requirement that e has exactly the type τ , or that two expressions jointly have the same type. And it is just this weakness that creates so many of the difficulties with subtyping.

Much has been written about subtyping, often in relation to object-oriented programming. Standard ML (Milner et al., 1997) is one of the first languages to make use of subtyping, in two forms, called *enrichment* and *realization*. The former corresponds to product subtyping, and the latter to the "forgetful" subtyping associated with type definitions (see Chapter 43). The first systematic studies of subtyping include those by Mitchell (1984); Reynolds (1980), and Cardelli (1988). Pierce (2002) give a thorough account of subtyping, especially of recursive and polymorphic types, and proves that subtyping for bounded impredicative universal quantification is undecidable.

Exercises

24.1. Check the variance of the type

(unit
$$\rightharpoonup \tau$$
) \times ($\tau \rightharpoonup$ unit).

When viewed as a constructor with argument τ , is it covariant or contravariant? Give a precise proof or counterexample in each case.

24.5 Notes 225

24.2. Consider the two recursive types,

$$\rho_1 \stackrel{\triangle}{=} \operatorname{rec} t \operatorname{is} \langle \operatorname{eq} \hookrightarrow (t \rightharpoonup \operatorname{bool}) \rangle,$$

and

$$\rho_2 \stackrel{\triangle}{=} \operatorname{rec} t \text{ is } \langle \operatorname{eq} \hookrightarrow (t \rightharpoonup \operatorname{bool}), \operatorname{f} \hookrightarrow \operatorname{bool} \rangle.$$

It is clear that ρ_1 could not be a subtype of ρ_2 , because, viewed as a product after unrolling, a value of the former type lacks a component that a value of the latter has. But is ρ_2 a subtype of ρ_1 ? If so, prove it by exhibiting a derivation of this fact using the rules given in Section 24.3. If not, give a counterexample showing that the suggested subtyping would violate type safety.

- **24.3**. Another approach to the dynamics of subtyping that ensures safety, but gives subsumption dynamic significance, associates a witness, called a *coercion*, to each subtyping relation, and inserts a coercion wherever subsumption is used. More precisely,
 - (a) Assign to each valid subtyping $\tau <: \tau'$ a coercion function $\chi : \tau \rightharpoonup \tau'$ that transforms a value of type τ into a value of type τ' .
 - (b) Interpret the subsumption rule as implicit coercion. Specifically, when $\tau <: \tau'$ is witnessed by $\chi : \tau \rightharpoonup \tau'$, applying subsumption to $e : \tau$ inserts an application of χ to obtain $\chi(e) : \tau'$.

Formulate this idea precisely for the case of a subtype relation generated by "width" subtyping for products, and the variance principles for product, sum and function types. Your solution should make clear that it evades the tacit projection assumption mentioned above.

But there may be more than one coercion $\chi: \tau \rightharpoonup \tau'$ corresponding to the subtyping $\tau <: \tau'$. The meaning of a program would then depend on which coercion is chosen when subsumption is used. If there is exactly one coercion for each subtyping relation, it is said to be *coherent*. Is your coercion interpretation of product subtyping coherent? (A proper treatment of coherence requires expression equivalence, which is discussed in Chapter 47.)

Part XI Dynamic Dispatch





Chapter 27

Inheritance

In this chapter we build on Chapter 26 and consider the process of defining the dispatch matrix that determines the behavior of each method on each class. A common strategy is to build the dispatch matrix incrementally by adding new classes or methods to an existing dispatch matrix. To add a class requires that we define the behavior of each method on objects of that class, and to define a method requires that we define the behavior of that method on objects of the classes. The definition of these behaviors can be given by any means available in the language. However, it is often suggested that a useful means of defining a new class is to *inherit* the behavior of another class on some methods, and to *override* its behavior on others, resulting in an amalgam of the old and new behaviors. The new class is often called a *subclass* of the old class, which is then called the *superclass*. Similarly, a new method can be defined by inheriting the behavior of another method on some classes, and overriding the behavior on others. By analogy we may call the new method a *sub-method* of a given *super-method*. For the sake of clarity we restrict attention to the non-self-referential case in the following development.

27.1 Class and Method Extension

We begin by extending a given dispatch matrix, e_{dm} , of type

$$\prod_{c \in C} \prod_{d \in D} \left(\tau^c \to \rho_d \right)$$

with a new class $c^* \notin C$ and a new method $d^* \notin D$ to obtain a new dispatch matrix $e^*_{\sf dm}$ of type

$$\prod_{c \in C^*} \prod_{d \in D^*} (\tau^c \to \rho_d),$$

where $C^* = C \cup \{c^*\}$ and $D^* = D \cup \{d^*\}$.

To add a new class c^* to the dispatch matrix, we must specify the following information:¹

¹The extension with a new method will be considered separately for the sake of clarity.

- 1. The instance type τ^{c^*} of the new class c^* .
- 2. The behavior $e_d^{c^*}$ of each method $d \in D$ on an object of the new class c^* , a function of type $\tau^{c^*} \to \rho_d$.

This data determines a new dispatch matrix e_{dm}^* such that the following conditions are satisfied:

- 1. For each $c \in C$ and $d \in D$, the behavior $e_{dm}^* \cdot c \cdot d$ is the same as the behavior $e_{dm} \cdot c \cdot d$.
- 2. For each $d \in D$, the behavior $e_{dm}^* \cdot c^* \cdot d$ is given by $e_d^{c^*}$.

To define c^* as a subclass of some class $c \in C$ means to define the behavior $e_d^{c^*}$ to be e_d^c for some (perhaps many) $d \in D$. It is sensible to inherit a method d in this way only if the subtype relationship

$$\tau^c \to \rho_d <: \tau^{c^*} \to \rho_d$$

is valid, which will be the case if $\tau^{c^*} <: \tau^c$. This subtyping condition ensures that the inherited behavior can be invoked on the instance data of the new class.

Similarly, to add a new method d^* to the dispatch matrix, we must specify the following information:

- 1. The result type ρ_{d^*} of the new method d^* .
- 2. The behavior $e^c_{d^*}$ of the new method d^* on an object of each class $c \in C$, a function of type $\tau^c \to \rho_{d^*}$.

This data determines a new dispatch matrix e_{dm}^* such that the following conditions are satisfied:

- 1. For each $c \in C$ and $d \in D$, the behavior $e_{\mathsf{dm}}^* \cdot c \cdot d$ is the same as $e_{\mathsf{dm}} \cdot c \cdot d$.
- 2. The behavior $e_{\sf dm}^* \cdot c \cdot d^*$ is given by $e_{d^*}^c$.

To define d^* as a sub-method of some $d \in D$ means to define the behavior $e^c_{d^*}$ to be e^c_d for some (perhaps many) classes $c \in C$. This definition is only sensible if the subtype relationship

$$\tau^c
ightarrow
ho_d <: \tau^c
ightarrow
ho_{d^*}$$

holds, which is the case if $\rho_d <: \rho_{d^*}$. This subtyping relationship ensures that the result of the old behavior suffices for the new behavior.

We will now consider how inheritance relates to the method- and class-based organizations of dynamic dispatch considered in Chapter 26.

27.2 Class-Based Inheritance

Recall that the class-based organization given in Chapter 26 consists of a class vector e_{cv} of type

$$au_{\mathsf{cv}} riangleq \prod_{c \in C} (au^c o
ho),$$

where the object type ρ is the finite product type $\prod_{d \in D} \rho_d$. The class vector consists of a tuple of constructors that specialize the methods to a given object of each class.

Let us consider the effect of adding a new class c^* as described in Section 27.1. The new class vector e_{cv}^* has type

$$au_{\mathsf{cv}}^* riangleq \prod_{c \in C^*} (au^c o
ho).$$

There is an isomorphism, written ()[†], between τ_{cv}^* and the type

$$\tau_{\sf cv} \times (\tau^{c^*} \to \rho),$$

which can be used to define the new class vector e_{cv}^* as follows:

$$\langle e_{\mathsf{cv}}, \lambda \left(u : \tau^{c^*} \right) \langle d \hookrightarrow e_d^{c^*} (u) \rangle_{d \in D} \rangle^{\dagger}.$$

This definition makes clear that the old class vector e_{cv} is reused intact in the new class vector, which extends the old class vector with a new constructor.

Although the object type ρ is the same both before and after the extension with the new class, the behavior of an object of class c^* may differ arbitrarily from that of any other object, even that of the superclass from which it inherits its behavior. So, knowing that c^* inherits from c tells us nothing about the behavior of its objects, but only about the means by which the class is defined. Inheritance carries no semantic significance, but is only a record of the history of how a class is defined.

Now let us consider the effect of adding a new method d^* as described in Section 27.1. The new class vector e_{cv}^* has type

$$au_{\mathsf{cv}}^* \triangleq \prod_{c \in C} (au^c o
ho^*),$$

where ρ^* is the product type $\prod_{d \in D^*} \rho_d$. There is an isomorphism, written ()[‡], between ρ^* and the type $\rho \times \rho_{d^*}$, where ρ is the old object type. Using this the new class vector e_{cv}^* is defined by

$$\langle c \hookrightarrow \lambda (u : \tau^c) \langle \langle d \hookrightarrow ((e_{cv} \cdot c)(u)) \cdot d \rangle_{d \in D}, e_{d^*}^c(u) \rangle^{\ddagger} \rangle_{c \in C}.$$

Observe that each constructor must be re-defined to account for the new method, but the definition makes use of the old class vector for the definitions of the old methods.

By this construction the new object type ρ^* is a subtype of the old object type ρ . Consequently, any objects with the new method can be used in situations expecting an object without the new method, as might be expected. To avoid redefining old classes when a new method is introduced, we may restrict inheritance so that new methods are only added to new subclasses. Subclasses may then have more methods than super-classes, and objects of the subclass can be provided when an object of the superclass is required.

27.3 Method-Based Inheritance

The method-based organization is dual to that of the class-based organization. Recall that the method-based organization given in Chapter 26 consists of a method vector e_{mv} of type

$$au_{\mathsf{mv}} riangleq \prod_{d \in D} au o
ho_d$$
,

where the instance type τ is the sum type $\sum_{c \in C} \tau^c$. The method vector consists of a tuple of functions that dispatch on the class of the object to determine their behavior.

Let us consider the effect of adding a new method d^* as described in Section 27.1. The new method vector e_{mv}^* has type

$$\tau_{\mathsf{mv}}^* \triangleq \prod_{d \in D^*} \tau \to \rho_d.$$

There is an isomorphism, written ()[‡], between τ_{mv}^* and the type

$$\tau_{\mathsf{mv}} \times (\tau \to \rho_{d^*}).$$

Using this isomorphism, the new method vector e_{mv}^* is defined as

$$\langle e_{\mathsf{mv}}, \lambda \, (\mathit{this}: \tau) \, \mathsf{case} \, \mathit{this} \, \{c \cdot u \hookrightarrow e^{c}_{d^*}(u)\}_{c \in C} \rangle^{\ddagger}.$$

The old method vector is re-used intact, extended with a dispatch function for the new method.

The object type does not change under the extension with a new method, but because $\rho^* <: \rho$, there is no difficulty using a new object in a context expecting an old object—the added method is ignored.

Finally, let us consider the effect of adding a new class c^* as described in Section 27.1. The new method vector, e_{mv}^* , has the type

$$au_{\mathsf{mv}}^* \triangleq \prod_{d \in D} au^* o
ho_d,$$

where τ^* is the new object type $\sum_{c \in C^*} \tau^c$, which is a super-type of the old object type τ . There is an isomorphism, written $()^{\dagger}$, between τ^* and the sum type $\tau + \tau^{c^*}$, which we may use to define the new method vector e_{mv}^* as follows:

$$\langle d \hookrightarrow \chi \, (\, this : \tau^* \,) \, \mathsf{case} \, this^\dagger \, \{ \mathbf{1} \cdot u \hookrightarrow (\, e_\mathsf{mv} \cdot d \,) (\, u \,) \, | \, \mathbf{r} \cdot u \hookrightarrow e_d^{c^*} (\, u \,) \} \rangle_{d \in D}.$$

Every method must be redefined to account for the new class, but the old method vector is reused.

27.4 Notes

Abadi and Cardelli (1996) and Pierce (2002) provide thorough accounts of the interaction of inheritance and subtyping. Liskov and Wing (1994) discuss it from a behavioral perspective. They propose to require that subclasses respect the behavior of the superclass when inheritance is used.

27.4 Notes 255

Exercises

27.1. Consider the case of extending a dispatch matrix with self-reference by a new class c^* in which a method d is inherited from an existing class c. What requirements ensure that such an inheritance is properly defined? What happens if we extend a self-referential dispatch matrix with a new method, d^* that inherits its behavior on class c from another method d?

- **27.2.** Consider the example of two mutually recursive methods given in Exercise **26.3**. Suppose that num* is a new class with instance type $\tau^{\text{num}} <: \tau^{\text{num}}$ that inherits the ev method from num, but defines its own version of the od method. What happens when message ev is sent to an instance of num*? Will the revised od method ever be invoked?
- **27.3.** *Method specialization* consists of defining a new class by inheriting methods from another class or classes, while redefining some of the methods that the inherited methods might invoke. The behavior of the inherited methods on instances of the new class is altered to the extent that they invoke a method that is specialized to the new class. Reconsider Exercise **26.3** in light of Exercise **27.2**, seeking to ensure that the specialization of od is invoked when the inherited method ev is invoked on instances of the new class.
 - (a) Redefine the class num along the following lines. The instance data of num is an object admitting methods ev and od. The class num admits these methods, and simply hands them off to the instance object.
 - (b) The classes zero or of succ admit both the ev and od methods, and are defined using message send to effect mutual recursion as necessary.
 - (c) Define a subclass succ* of succ that overrides the od method. Show that ev on an instance of succ* correctly invokes the overridden od method.

256 27.4 Notes



Part XII Control Flow



Chapter 28

Control Stacks

Structural dynamics is convenient for proving properties of languages, such as a type safety theorem, but is less convenient as a guide for implementation. A structural dynamics defines a transition relation using rules that determine where to apply the next instruction without spelling out how to find where the instruction lies within an expression. To make this process explicit we introduce a mechanism, called a *control stack*, that records the work that remains to be done after an instruction is executed. Using a stack eliminates the need for premises on the transition rules so that the transition system defines an *abstract machine* whose steps are determined by information explicit in its state, much as a concrete computer does.

In this chapter we develop an abstract machine **K** for evaluating expressions in **PCF**. The machine makes explicit the context in which primitive instruction steps are executed, and the process by which the results are propagated to determine the next step of execution. We prove that **K** and **PCF** are equivalent in the sense that both achieve the same outcomes for the same expressions.

28.1 Machine Definition

A state s of the stack machine **K** for **PCF** consists of a *control stack* k and a closed expression e. States take one of two forms:

- 1. An *evaluation* state of the form $k \triangleright e$ corresponds to the evaluation of a closed expression e on a control stack k.
- 2. A *return* state of the form $k \triangleleft e$, where e val, corresponds to the evaluation of a stack k on a closed value e.

As an aid to memory, note that the separator "points to" the focal entity of the state, the expression in an evaluation state and the stack in a return state.

The control stack represents the context of evaluation. It records the "current location" of evaluation, the context into which the value of the current expression is returned. Formally, a control

260 28.1 Machine Definition

stack is a list of *frames*:

$$\overline{\epsilon}$$
 stack (28.1a)

$$\frac{f \text{ frame } k \text{ stack}}{k; f \text{ stack}}$$
 (28.1b)

The frames of the **K** machine are inductively defined by the following rules:

$$\frac{}{s(-) \text{ frame}}$$
 (28.2a)

$$\frac{1}{\text{ifz}\{e_0; x.e_1\}(-) \text{ frame}}$$
 (28.2b)

$$\frac{}{\mathsf{ap}(-;e_2)\mathsf{ frame}} \tag{28.2c}$$

The frames correspond to search rules in the dynamics of **PCF**. Thus, instead of relying on the structure of the transition derivation to keep a record of pending computations, we make an explicit record of them in the form of a frame on the control stack.

The transition judgment between states of the **PCF** machine is inductively defined by a set of inference rules. We begin with the rules for natural numbers, using an eager dynamics for the successor.

$$(28.3a)$$

$$\frac{1}{k \triangleright \mathsf{s}(e) \longmapsto k; \mathsf{s}(-) \triangleright e} \tag{28.3b}$$

$$\frac{1}{k; \mathbf{s}(-) \triangleleft e \longmapsto k \triangleleft \mathbf{s}(e)} \tag{28.3c}$$

To evaluate z we simply return it. To evaluate s(e), we push a frame on the stack to record the pending successor, and evaluate e; when that returns with e', we return s(e') to the stack.

Next, we consider the rules for case analysis.

$$\overline{k} \triangleright \operatorname{ifz}\{e_0; x.e_1\}(e) \longmapsto k; \operatorname{ifz}\{e_0; x.e_1\}(-) \triangleright e \tag{28.4a}$$

$$\frac{1}{k; \text{ifz}\{e_0; x.e_1\}(-) \triangleleft z \longmapsto k \triangleright e_0}$$
 (28.4b)

$$\frac{1}{k; ifz\{e_0; x.e_1\}(-) \triangleleft s(e) \longmapsto k \triangleright [e/x]e_1}$$
(28.4c)

The test expression is evaluated, recording the pending case analysis on the stack. Once the value of the test expression is determined, the zero or non-zero branch of the condition is evaluated, substituting the predecessor in the latter case.

28.2 Safety 261

Finally, we give the rules for functions, which are evaluated by-name, and the rule for general recursion.

$$\frac{1}{k \triangleright \lambda \{\tau\}(x.e) \longmapsto k \triangleleft \lambda \{\tau\}(x.e)} \tag{28.5a}$$

$$\frac{1}{k \triangleright \operatorname{ap}(e_1; e_2) \longmapsto k; \operatorname{ap}(-; e_2) \triangleright e_1}$$
 (28.5b)

$$\frac{1}{k; \operatorname{ap}(-; e_2) \triangleleft \lambda\{\tau\}(x.e) \longmapsto k \triangleright [e_2/x]e}$$
 (28.5c)

$$\frac{1}{k \triangleright \text{fix}\{\tau\}(x.e) \longmapsto k \triangleright [\text{fix}\{\tau\}(x.e)/x]e}$$
 (28.5d)

It is important that evaluation of a general recursion requires no stack space.

The initial and final states of the **K** machine are defined by the following rules:

$$\frac{}{\epsilon \triangleright e \text{ initial}} \tag{28.6a}$$

$$\frac{e \text{ val}}{\epsilon \triangleleft e \text{ final}} \tag{28.6b}$$

Safety 28.2

To define and prove safety for the PCF machine requires that we introduce a new typing judgment, $k \div \tau$, which states that the stack k expects a value of type τ . This judgment is inductively defined by the following rules:

$$(28.7a)$$

$$\frac{k \div \tau' \quad f : \tau \leadsto \tau'}{k : f \div \tau} \tag{28.7a}$$

This definition makes use of an auxiliary judgment, $f: \tau \leadsto \tau'$, stating that a frame f transforms a value of type τ to a value of type τ' .

$$\frac{}{\mathsf{s}(\,-\,):\mathsf{nat}\rightsquigarrow\mathsf{nat}}\tag{28.8a}$$

$$\frac{e_0:\tau\quad x:\mathtt{nat}\vdash e_1:\tau}{\mathtt{ifz}\{e_0;x.e_1\}(-):\mathtt{nat}\leadsto\tau} \tag{28.8b}$$

$$\frac{e_2 : \tau_2}{\operatorname{ap}(-; e_2) : \rightharpoonup (\tau_2; \tau) \leadsto \tau} \tag{28.8c}$$

The states of the **PCF** machine are well-formed if their stack and expression components match:

$$\frac{k \div \tau \quad e : \tau}{k \triangleright e \text{ ok}} \tag{28.9a}$$

$$\frac{k \div \tau \quad e : \tau \quad e \text{ val}}{k \triangleleft e \text{ ok}} \tag{28.9b}$$

We leave the proof of safety of the **PCF** machine as an exercise.

Theorem 28.1 (Safety). 1. If s ok and $s \mapsto s'$, then s' ok.

2. If s ok, then either s final or there exists s' such that $s \mapsto s'$.

28.3 Correctness of the K Machine

Does evaluation of an expression *e* using the **K** machine yield the same result as does the structural dynamics of **PCF**? The answer to this question can be derived from the following facts.

Completeness If $e \mapsto^* e'$, where e' val, then $\epsilon \triangleright e \mapsto^* \epsilon \triangleleft e'$.

Soundness If $\epsilon \triangleright e \mapsto^* \epsilon \triangleleft e'$, then $e \mapsto^* e'$ with e' val.

To prove completeness a plausible first step is to consider a proof by induction on the definition of multi-step transition, which reduces the theorem to the following two lemmas:

- 1. If e val, then $\epsilon \triangleright e \mapsto^* \epsilon \triangleleft e$.
- 2. If $e \mapsto e'$, then, for every v val, if $\epsilon \triangleright e' \mapsto^* \epsilon \triangleleft v$, then $\epsilon \triangleright e \mapsto^* \epsilon \triangleleft v$.

The first can be proved easily by induction on the structure of e. The second requires an inductive analysis of the derivation of $e \mapsto e'$ that gives rise to two complications. The first complication is that we cannot restrict attention to the empty stack, for if e is, say, $ap(e_1; e_2)$, then the first step of the **K** machine is

$$\epsilon \triangleright \operatorname{ap}(e_1; e_2) \longmapsto \epsilon \operatorname{;ap}(-; e_2) \triangleright e_1.$$

To handle such situations we consider the evaluation of e_1 on any stack, not just the empty stack.

Specifically, we prove that if $e \mapsto e'$ and $k \triangleright e' \mapsto^* k \triangleleft v$, then $k \triangleright e \mapsto^* k \triangleleft v$. Reconsider the case $e = \operatorname{ap}(e_1; e_2)$, $e' = \operatorname{ap}(e'_1; e_2)$, with $e_1 \mapsto^* e'_1$. We are given that $k \triangleright \operatorname{ap}(e'_1; e_2) \mapsto^* k \triangleleft v$, and we are to show that $k \triangleright \operatorname{ap}(e_1; e_2) \mapsto^* k \triangleleft v$. It is easy to show that the first step of the former derivation is

$$k \rhd \operatorname{ap}(e_1'; e_2) \longmapsto k ; \operatorname{ap}(-; e_2) \rhd e_1'.$$

We would like to apply induction to the derivation of $e_1 \mapsto e'_1$, but to do so we need a value v_1 such that $e'_1 \mapsto^* v_1$, which is not at hand.

We therefore consider the value of each sub-expression of an expression. This information is given by the evaluation dynamics described in Chapter 7, which has the property that $e \Downarrow e'$ iff $e \mapsto^* e'$ and e' val.

Lemma 28.2. *If* $e \Downarrow v$, then for every k stack, $k \triangleright e \longmapsto^* k \triangleleft v$.

The desired result follows by the analog of Theorem 7.2 for **PCF**, which states that $e \Downarrow v$ iff $e \longmapsto^* v$.

To prove soundness, we note that it is awkward to reason inductively about a multi-step transition from $\epsilon \triangleright e \longmapsto^* \epsilon \triangleleft v$. The intermediate steps could involve alternations of evaluation and return states. Instead we consider a **K** machine state to encode an expression, and show that the machine transitions are simulated by the transitions of the structural dynamics.

To do so we define a judgment, $s \hookrightarrow e$, stating that state s "unravels to" expression e. It will turn out that for initial states, $s = e \triangleright e$, and final states, $s = e \triangleleft e$, we have $s \hookrightarrow e$. Then we show that if $s \longmapsto^* s'$, where s' final, $s \hookrightarrow e$, and $s' \hookrightarrow e'$, then e' val and $e \longmapsto^* e'$. For this it is enough to show the following two facts:

- 1. If $s \hookrightarrow e$ and s final, then e val.
- 2. If $s \mapsto s'$, $s \mapsto e$, $s' \mapsto e'$, and $e' \mapsto^* v$, where v val, then $e \mapsto^* v$.

The first is quite simple, we need only note that the unraveling of a final state is a value. For the second, it is enough to prove the following lemma.

Lemma 28.3. If $s \mapsto s'$, $s \mapsto e$, and $s' \mapsto e'$, then $e \mapsto^* e'$.

Corollary 28.4. $e \mapsto^* \overline{n} \text{ iff } e \triangleright e \mapsto^* e \triangleleft \overline{n}.$

28.3.1 Completeness

Proof of Lemma 28.2. The proof is by induction on an evaluation dynamics for PCF.

Consider the evaluation rule

$$\frac{e_1 \Downarrow \lambda\{\tau_2\}(x.e) \quad [e_2/x]e \Downarrow v}{\operatorname{ap}(e_1; e_2) \Downarrow v}$$
(28.10)

For an arbitrary control stack k we are to show that $k \triangleright ap(e_1; e_2) \longmapsto^* k \triangleleft v$. Applying both of the inductive hypotheses in succession, interleaved with steps of the **K** machine, we obtain

$$k \triangleright \operatorname{ap}(e_1; e_2) \longmapsto k ; \operatorname{ap}(-; e_2) \triangleright e_1$$

 $\longmapsto^* k ; \operatorname{ap}(-; e_2) \triangleleft \lambda \{\tau_2\}(x.e)$
 $\longmapsto k \triangleright [e_2/x]e$
 $\longmapsto^* k \triangleleft \tau.$

The other cases of the proof are handled similarly.

28.3.2 Soundness

The judgment $s \hookrightarrow e'$, where s is either $k \triangleright e$ or $k \triangleleft e$, is defined in terms of the auxiliary judgment $k \bowtie e = e'$ by the following rules:

$$\frac{k \bowtie e = e'}{k \bowtie e \hookrightarrow e'} \tag{28.11a}$$

$$\frac{k \bowtie e = e'}{k \triangleleft e \looparrowright e'} \tag{28.11b}$$

In words, to unravel a state we wrap the stack around the expression to form a complete program. The unraveling relation is inductively defined by the following rules:

$$\overline{\epsilon \bowtie e = e} \tag{28.12a}$$

$$\frac{k \bowtie s(e) = e'}{k; s(-) \bowtie e = e'}$$
 (28.12b)

$$\frac{k \bowtie ifz\{e_0; x.e_1\}(e) = e'}{k; ifz\{e_0; x.e_1\}(-) \bowtie e = e'}$$
(28.12c)

$$\frac{k \bowtie \operatorname{ap}(e_1; e_2) = e}{k; \operatorname{ap}(-; e_2) \bowtie e_1 = e}$$
(28.12d)

These judgments both define total functions.

Lemma 28.5. The judgment $s \hookrightarrow e$ relates every state s to a unique expression e, and the judgment $k \bowtie e = e'$ relates every stack k and expression e to a unique expression e'.

We are therefore justified in writing $k \bowtie e$ for the unique e' such that $k \bowtie e = e'$.

The following lemma is crucial. It states that unraveling preserves the transition relation.

Lemma 28.6. If
$$e \mapsto e'$$
, $k \bowtie e = d$, $k \bowtie e' = d'$, then $d \mapsto d'$.

Proof. The proof is by rule induction on the transition $e \mapsto e'$. The inductive cases, where the transition rule has a premise, follow easily by induction. The base cases, where the transition is an axiom, are proved by an inductive analysis of the stack k.

For an example of an inductive case, suppose that $e = \operatorname{ap}(e_1; e_2)$, $e' = \operatorname{ap}(e_1'; e_2)$, and $e_1 \longmapsto e_1'$. We have $k \bowtie e = d$ and $k \bowtie e' = d'$. It follows from rules (28.12) that k; $\operatorname{ap}(-; e_2) \bowtie e_1 = d$ and k; $\operatorname{ap}(-; e_2) \bowtie e_1' = d'$. So by induction $d \longmapsto d'$, as desired.

For an example of a base case, suppose that $e = \operatorname{ap}(\lambda\{\tau_2\}(x.e);e_2)$ and $e' = [e_2/x]e$ with $e \mapsto e'$ directly. Assume that $k \bowtie e = d$ and $k \bowtie e' = d'$; we are to show that $d \mapsto d'$. We proceed by an inner induction on the structure of k. If k = e, the result follows immediately. Consider, say, the stack k = k'; $\operatorname{ap}(-;c_2)$. It follows from rules (28.12) that $k' \bowtie \operatorname{ap}(e;c_2) = d$ and $k' \bowtie \operatorname{ap}(e';c_2) = d'$. But by the structural dynamics $\operatorname{ap}(e;c_2) \mapsto \operatorname{ap}(e';c_2)$, so by the inner inductive hypothesis we have $d \mapsto d'$, as desired.

We may now complete the proof of Lemma 28.3.

Proof of Lemma 28.3. The proof is by case analysis on the transitions of the **K** machine. In each case, after unraveling, the transition will correspond to zero or one transitions of the **PCF** structural dynamics.

Suppose that $s = k \triangleright s(e)$ and s' = k; $s(-) \triangleright e$. Note that $k \bowtie s(e) = e'$ iff k; $s(-) \bowtie e = e'$, from which the result follows immediately.

28.4 Notes 265

Suppose that s=k; ap($\lambda\{\tau\}(x.e_1)$; $-) \triangleleft e_2$ and $s'=k \triangleright [e_2/x]e_1$. Let e' be such that k; ap($\lambda\{\tau\}(x.e_1)$; $-) \bowtie e_2=e'$ and let e'' be such that $k\bowtie [e_2/x]e_1=e''$. Observe that $k\bowtie ap(\lambda\{\tau\}(x.e_1);e_2)=e'$. The result follows from Lemma 28.6.

28.4 Notes

The abstract machine considered here is typical of a wide class of machines that make control flow explicit in the state. The prototype is the SECD machine (Landin, 1965), which is a linearization of a structural operational semantics (Plotkin, 1981). The advantage of a machine model is that the explicit treatment of control is needed for languages that allow the control state to be manipulated (see Chapter 30 for a prime example). The disadvantage is that the control state of the computation must be made explicit, necessitating rules for manipulating it that are left implicit in a structural dynamics.

Exercises

- **28.1**. Give the proof of Theorem 28.1 for conditional expressions.
- **28.2**. Formulate a call-by-value variant of the **PCF** machine.
- **28.3.** Analyze the worst-case asymptotic complexity of executing each instruction of the **K** machine.
- **28.4**. Refine the proof of Lemma 28.2 by bounding the number of machine steps taken for each step of the **PCF** dynamics.

28.4 Notes



Chapter 29

Exceptions

Exceptions effect a non-local transfer of control from the point at which the exception is *raised* to an enclosing *handler* for that exception. This transfer interrupts the normal flow of control in a program in response to unusual conditions. For example, exceptions can be used to signal an error condition, or to signal the need for special handling in unusual circumstances. We could use conditionals to check for and process errors or unusual conditions, but using exceptions is often more convenient, particularly because the transfer to the handler is conceptually direct and immediate, rather than indirect via explicit checks.

In this chapter we will consider two extensions of **PCF** with exceptions. The first, **FPCF**, enriches **PCF** with the simplest form of exception, called a *failure*, with no associated data. A failure can be intercepted, and turned into a success (or another failure!) by transferring control to another expression. The second, **XPCF**, enriches **PCF** with *exceptions*, with associated data that is passed to an exception handler that intercepts it. The handler may analyze the associated data to determine how to recover from the exceptional condition. A key choice is to decide on the type of the data associated to an exception.

29.1 Failures

The syntax of **FPCF** is defined by the following extension of the grammar of **PCF**:

Exp
$$e$$
 ::= fail fail signal a failure catch $(e_1;e_2)$ catch e_1 ow e_2 catch a failure

The expression fail aborts the current evaluation, and the expression $catch(e_1; e_2)$ catches any failure in e_1 by evaluating e_2 instead. Either e_1 or e_2 may themselves abort, or they may diverge or return a value as usual in **PCF**.

The statics of **FPCF** is given by these rules:

$$\Gamma \vdash \mathsf{fail} : \tau$$
 (29.1a)

268 29.1 Failures

$$\frac{\Gamma \vdash e_1 : \tau \quad \Gamma \vdash e_2 : \tau}{\Gamma \vdash \mathsf{catch}(e_1; e_2) : \tau}$$
 (29.1b)

A failure can have any type, because it never returns. The two expressions in a catch expression must have the same type, because either might determine the value of that expression.

The dynamics of **FPCF** is given using a technique called *stack unwinding*. Evaluation of a catch pushes a frame of the form $\operatorname{catch}(-;e)$ onto the control stack that awaits the arrival of a failure. Evaluation of a fail expression pops frames from the control stack until it reaches a frame of the form $\operatorname{catch}(-;e)$, at which point the frame is removed from the stack and the expression e is evaluated. Failure propagation is expressed by a state of the form $k \triangleleft$, which extends the two forms of state considered in Chapter 28 to express failure propagation.

The FPCF machine extends the PCF machine with the following additional rules:

$$\overline{k \triangleright \text{fail} \longmapsto k \blacktriangleleft}$$
 (29.2a)

$$\frac{1}{k \triangleright \operatorname{catch}(e_1; e_2) \longmapsto k ; \operatorname{catch}(-; e_2) \triangleright e_1}$$
 (29.2b)

$$\frac{1}{k : \operatorname{catch}(-; e_2) \triangleleft v \longmapsto k \triangleleft v} \tag{29.2c}$$

$$\frac{}{k; \operatorname{catch}(-; e_2) \blacktriangleleft \longmapsto k \triangleright e_2} \tag{29.2d}$$

$$\frac{(f \neq \operatorname{catch}(-;e))}{k; f \blacktriangleleft \longmapsto k \blacktriangleleft}$$
 (29.2e)

Evaluating fail propagates a failure up the stack. The act of failing itself, fail, will, of course, give rise to a failure. Evaluating $\operatorname{catch}(e_1;e_2)$ consists of pushing the handler on the control stack and evaluating e_1 . If a value reaches to the handler, the handler is removed and the value is passed to the surrounding frame. If a failure reaches the handler, the stored expression is evaluated with the handler removed from the control stack. Failures propagate through all frames other than the catch frame.

The initial and final states of the FPCF machine are defined by the following rules:

$$\frac{}{\epsilon \triangleright e \text{ initial}} \tag{29.3a}$$

$$\frac{e \text{ val}}{\epsilon \triangleleft e \text{ final}} \tag{29.3b}$$

$$\epsilon \blacktriangleleft \text{ final}$$
 (29.3c)

The definition of stack typing given in Chapter 28 can be extended to account for the new forms of frame so that safety can be proved in the same way as before. The only difference is that the statement of progress must be weakened to take account of failure: a well-typed expression is either a value, or may take a step, or may signal failure.

29.2 Exceptions 269

Theorem 29.1 (Safety for **FPCF**). 1. If s ok and $s \mapsto s'$, then s' ok.

2. If s ok, then either s final or there exists s' such that $s \mapsto s'$.

29.2 Exceptions

The language **XPCF** enriches **FPCF** with *exceptions*, failures to which a value is attached. The syntax of **XPCF** extends that of **PCF** with the following forms of expression:

The argument to raise is evaluated to determine the value passed to the handler. The expression $try(e_1; x.e_2)$ binds a variable x in the handler e_2 . The associated value of the exception is bound to that variable within e_2 , should an exception be raised when e_1 is evaluated.

The statics of exceptions extends the statics of failures to account for the type of the value carried with the exception:

$$\frac{\Gamma \vdash e : \tau_{\mathsf{exn}}}{\Gamma \vdash \mathsf{raise}(e) : \tau} \tag{29.4a}$$

$$\frac{\Gamma \vdash e_1 : \tau \quad \Gamma, x : \tau_{\mathsf{exn}} \vdash e_2 : \tau}{\Gamma \vdash \mathsf{try}(e_1; x.e_2) : \tau}$$
 (29.4b)

The type τ_{exn} is some fixed, but as yet unspecified, type of exception values. (The choice of τ_{exn} is discussed in Section 29.3.)

The dynamics of **XPCF** is similar to that of **FPCF**, except that the failure state $k \triangleleft e$ is replaced by the exception state $k \triangleleft e$ which passes an exception value e to the stack k. There is only one notion of exception, but the associated value can be used to identify the source of the exception. We use a by-value interpretation to avoid the problem of *imprecise exceptions* that arises under a by-name interpretation.

The stack frames of the **PCF** machine are extended to include raise(-) and $try(-;x.e_2)$. These are used in the following rules:

$$\frac{1}{k \triangleright \mathsf{raise}(e) \longmapsto k; \mathsf{raise}(-) \triangleright e} \tag{29.5a}$$

$$\frac{1}{k; \mathtt{raise}(-) \triangleleft e \longmapsto k \blacktriangleleft e} \tag{29.5b}$$

$$\frac{1}{k \triangleright \operatorname{try}(e_1; x.e_2) \longmapsto k; \operatorname{try}(-; x.e_2) \triangleright e_1}$$
 (29.5c)

$$\frac{1}{k; \operatorname{try}(-; x.e_2) \triangleleft e \longmapsto k \triangleleft e} \tag{29.5d}$$

$$\frac{1}{k; \operatorname{try}(-; x.e_2)} \blacktriangleleft e \longmapsto k \triangleright [e/x]e_2$$
 (29.5e)

$$\frac{(f \neq \operatorname{try}(-; x.e_2))}{k; f \blacktriangleleft e \longmapsto k \blacktriangleleft e}$$
 (29.5f)

The main difference compared to Rules (29.2) is that an exception passes a values to the stack, whereas a failure does not.

The initial and final states of the **XPCF** machine are defined by the following rules:

$$\frac{}{\epsilon \triangleright e \text{ initial}} \tag{29.6a}$$

$$\frac{e \text{ val}}{\epsilon \triangleleft e \text{ final}} \tag{29.6b}$$

$$\frac{}{\epsilon \blacktriangleleft e \text{ final}} \tag{29.6c}$$

Theorem 29.2 (Safety for **XPCF**). 1. If s ok and $s \mapsto s'$, then s' ok.

2. If s ok, then either s final or there exists s' such that $s \mapsto s'$.

29.3 Exception Values

The statics of **XPCF** is parameterized by the type τ_{exn} of values associated to exceptions. The choice of τ_{exn} is important because it determines how the source of an exception is identified in a program. If τ_{exn} is the one-element type unit, then exceptions degenerate to failures, which are unable to identify their source. Thus τ_{exn} must have more than one value to be useful.

This fact suggests that τ_{exn} should be a finite sum. The classes of the sum identify the sources of exceptions, and the classified value carries information about the particular instance. For example, τ_{exn} might be a sum type of the form

$$[\mathtt{div} \hookrightarrow \mathtt{unit}, \mathtt{fnf} \hookrightarrow \mathtt{string}, \dots].$$

Here the class div might represent an arithmetic fault, with no associated data, and the class fnf might represent a "file not found" error, with associated data being the name of the file that was not found.

Using a sum means that an exception handler can dispatch on the class of the exception value to identify its source and cause. For example, we might write

```
handle e_1 ow x \hookrightarrow
match x \{
div \langle \rangle \hookrightarrow e_{\text{div}}
| fnf s \hookrightarrow e_{\text{fnf}} \}
```

29.4 Notes 271

to handle the exceptions specified by the above sum type. Because the exception and its associated data are coupled in a sum type, there is no possibility of misinterpreting the data associated to one exception as being that of another.

The disadvantage of choosing a finite sum for τ_{exn} is that it specifies a *closed world* of possible exception sources. All sources must be identified for the entire program, which impedes modular development and evolution. A more modular approach admits an *open world* of exception sources that can be introduced as the program evolves and even as it executes. A generalization of finite sums, called *dynamic classification*, defined in Chapter 33, is required for an open world. (See that Chapter for further discussion.)

When τ_{exn} is a type of classified values, its classes are often called *exceptions*, so that one may speak of "the fnf exception" in the above example. This terminology is harmless, and all but unavoidable, but it invites confusion between two separate ideas:

- 1. Exceptions as a *control mechanism* that allows the course of evaluation to be altered by raising and handling exceptions.
- 2. Exceptions as a *data value* associated with such a deviation of control that allows the source of the deviation to be identified.

As a control mechanism exceptions can be eliminated using explicit *exception passing*. A computation of type τ that may raise an exception is interpreted as an exception-free computation of type $\tau + \tau_{\text{exn}}$.

29.4 Notes

Various forms of exceptions were considered in Lisp (Steele, 1990). The original formulation of ML (Gordon et al., 1979) as a metalanguage for mechanized logic used failures to implement backtracking proof search. Most modern languages now have exceptions, but differ in the forms of data that may be associated with them.

Exercises

- **29.1**. Prove Theorem 29.2. Are any properties of τ_{exp} required for the proof?
- **29.2**. Give an evaluation dynamics for **XPCF** using the following judgment forms:
 - Normal evaluation: $e \downarrow v$, where $e : \tau, v : \tau$, and v val.
 - Exceptional evaluation: $e \uparrow v$, where $e : \tau$, and $v : \tau_{exn}$, and v val.

The first states that e evaluates normally to value v, the second that e raises an exception with value v.

29.3. Give a structural operational dynamics to **XPCF** by inductively defining the following judgment forms:

272 29.4 Notes

- $e \mapsto e'$, stating that expression e transitions to expression e';
- *e* val, stating that expression *e* is a value.

Ensure that $e \Downarrow v$ iff $e \longmapsto^* v$, and $e \Uparrow v$ iff $e \longmapsto^* \mathtt{raise}(v)$, where v val in both cases.



Part XIII Symbolic Data





Part XIV Mutable State



Chapter 34

Modernized Algol

Modernized Algol, or MA, is an imperative, block-structured programming language based on the classic language Algol. MA extends PCF with a new syntactic sort of *commands* that act on assignables by retrieving and altering their contents. Assignables are introduced by declaring them for use within a specified scope; this is the essence of block structure. Commands are combined by sequencing, and are iterated using recursion.

MA maintains a careful separation between *pure* expressions, whose meaning does not depend on any assignables, and *impure* commands, whose meaning is given in terms of assignables. The segregation of pure from impure ensures that the evaluation order for expressions is not constrained by the presence of assignables in the language, so that they can be manipulated just as in **PCF**. Commands, on the other hand, have a constrained execution order, because the execution of one may affect the meaning of another.

A distinctive feature of **MA** is that it adheres to the *stack discipline*, which means that assignables are allocated on entry to the scope of their declaration, and deallocated on exit, using a conventional stack discipline. Stack allocation avoids the need for more complex forms of storage management, at the cost of reducing the expressive power of the language.

34.1 Basic Commands

The syntax of the language **MA** of modernized Algol distinguishes pure *expressions* from impure *commands*. The expressions include those of **PCF** (as described in Chapter 19), augmented with one construct, and the commands are those of a simple imperative programming language based on assignment. The language maintains a sharp distinction between *variables* and *assignables*. Variables are introduced by λ -abstraction and are given meaning by substitution. Assignables are introduced by a declaration and are given meaning by assignment and retrieval of their *contents*, which is, for the time being, restricted to natural numbers. Expressions evaluate to values, and have no effect on assignables. Commands are executed for their effect on assignables, and return a value. Composition of commands not only sequences their execution order, but also passes the value returned by the first to the second before it is executed. The returned value of a command

310 34.1 Basic Commands

is, for the time being, restricted to the natural numbers. (But see Section 34.3 for the general case.) The syntax of **MA** is given by the following grammar, from which we have omitted repetition of the expression syntax of **PCF** for the sake of brevity.

The expression $\operatorname{cmd}(m)$ consists of the unevaluated command m thought of as a value of type cmd . The command $\operatorname{ret}(e)$ returns the value of the expression e without having any effect on the assignables. The command $\operatorname{bnd}(e;x.m)$ evaluates e to an encapsulated command, then this command is executed for its effects on assignables, with its value substituted for x in m. The command $\operatorname{dcl}(e;a.m)$ introduces a new assignable, a, for use within the command m whose initial contents is given by the expression e. The command $\operatorname{get}[a]$ returns the current contents of the assignable e and the command $\operatorname{set}[a](e)$ changes the contents of the assignable e to the value of e, and returns that value.

34.1.1 Statics

The statics of **MA** consists of two forms of judgment:

- 1. Expression typing: $\Gamma \vdash_{\Sigma} e : \tau$.
- 2. Command formation: $\Gamma \vdash_{\Sigma} m$ ok.

The context Γ specifies the types of variables, as usual, and the signature Σ consists of a finite set of assignables. As with other uses of symbols, the signature cannot be interpreted as a form of typing hypothesis (it enjoys no structural properties of entailment), but must be considered as an index of a family of judgments, one for each choice of Σ .

The statics of **MA** is inductively defined by the following rules:

$$\frac{\Gamma \vdash_{\Sigma} m \text{ ok}}{\Gamma \vdash_{\Sigma} \text{cmd}(m) : \text{cmd}}$$
 (34.1a)

$$\frac{\Gamma \vdash_{\Sigma} e : \mathtt{nat}}{\Gamma \vdash_{\Sigma} \mathtt{ret}(e) \mathsf{ok}} \tag{34.1b}$$

$$\frac{\Gamma \vdash_{\Sigma} e : \operatorname{cmd} \quad \Gamma, x : \operatorname{nat} \vdash_{\Sigma} m \text{ ok}}{\Gamma \vdash_{\Sigma} \operatorname{bnd}(e; x.m) \text{ ok}}$$
(34.1c)

$$\frac{\Gamma \vdash_{\Sigma} e : \text{nat} \quad \Gamma \vdash_{\Sigma,a} m \text{ ok}}{\Gamma \vdash_{\Sigma} \text{dcl}(e; a.m) \text{ ok}}$$
(34.1d)

34.1 Basic Commands 311

$$\frac{}{\Gamma \vdash_{\Sigma,a} \operatorname{get}[a] \operatorname{ok}} \tag{34.1e}$$

$$\frac{\Gamma \vdash_{\Sigma,a} e : \mathtt{nat}}{\Gamma \vdash_{\Sigma,a} \mathtt{set}[a](e) \mathsf{ok}} \tag{34.1f}$$

Rule (34.1a) is the introduction rule for the type cmd, and rule (34.1c) is the corresponding elimination form. Rule (34.1d) introduces a new assignable for use within a specified command. The name a of the assignable is bound by the declaration, and so may be renamed to satisfy the implicit constraint that it not already occur in Σ . Rule (34.1e) states that the command to retrieve the contents of an assignable a returns a natural number. Rule (34.1f) states that we may assign a natural number to an assignable.

34.1.2 Dynamics

The dynamics of **MA** is defined in terms of a *memory* μ a finite function assigning a numeral to each of a finite set of assignables.

The dynamics of expressions consists of these two judgment forms:

- 1. $e \text{ val}_{\Sigma}$, stating that e is a value relative to Σ .
- 2. $e \xrightarrow{\Sigma} e'$, stating that the expression e steps to the expression e'.

These judgments are inductively defined by the following rules, together with the rules defining the dynamics of **PCF** (see Chapter 19). It is important, however, that the successor operation be given an *eager*, instead of *lazy*, dynamics so that a closed value of type nat is a numeral (for reasons that will be explained in Section 34.3).

$$\frac{}{\operatorname{cmd}(m)\operatorname{val}_{\Sigma}}\tag{34.2a}$$

Rule (34.2a) states that an encapsulated command is a value.

The dynamics of commands is defined in terms of states $\mu \parallel m$, where μ is a memory mapping assignables to values, and m is a command. There are two judgments governing such states:

- 1. $\mu \parallel m$ final_{Σ}. The state $\mu \parallel m$ is complete.
- 2. $\mu \parallel m \mapsto_{\Sigma} \mu' \parallel m'$. The state $\mu \parallel m$ steps to the state $\mu' \parallel m'$; the set of active assignables is given by the signature Σ .

These judgments are inductively defined by the following rules:

$$\frac{e \operatorname{val}_{\Sigma}}{\mu \parallel \operatorname{ret}(e) \operatorname{final}_{\Sigma}} \tag{34.3a}$$

$$\frac{e \underset{\Sigma}{\longmapsto} e'}{\mu \parallel \operatorname{ret}(e) \underset{\Sigma}{\longmapsto} \mu \parallel \operatorname{ret}(e')} \tag{34.3b}$$

312 34.1 Basic Commands

$$\frac{e \underset{\Sigma}{\longmapsto} e'}{\mu \parallel \operatorname{bnd}(e; x.m) \underset{\Sigma}{\longmapsto} \mu \parallel \operatorname{bnd}(e'; x.m)}$$
(34.3c)

$$\frac{e \operatorname{val}_{\Sigma}}{\mu \parallel \operatorname{bnd}(\operatorname{cmd}(\operatorname{ret}(e)); x.m) \xrightarrow{\Sigma} \mu \parallel [e/x]m}$$
(34.3d)

$$\frac{\mu \parallel m_1 \stackrel{\longleftarrow}{\longmapsto} \mu' \parallel m'_1}{\mu \parallel \operatorname{bnd}(\operatorname{cmd}(m_1); x.m_2) \stackrel{\longleftarrow}{\longmapsto} \mu' \parallel \operatorname{bnd}(\operatorname{cmd}(m'_1); x.m_2)}$$
(34.3e)

$$\mu \otimes a \hookrightarrow e \parallel \operatorname{get}[a] \xrightarrow{\sum_{a} \mu \otimes a} \mu \otimes a \hookrightarrow e \parallel \operatorname{ret}(e)$$
 (34.3f)

$$\frac{e \underset{\Sigma,a}{\longmapsto} e'}{\mu \parallel \operatorname{set}[a](e) \underset{\Sigma,a}{\longmapsto} \mu \parallel \operatorname{set}[a](e')}$$
(34.3g)

$$\frac{e \operatorname{\mathsf{val}}_{\Sigma,a}}{\mu \otimes a \hookrightarrow_{-} \| \operatorname{\mathsf{set}}[a](e) \underset{\Sigma,a}{\longmapsto} \mu \otimes a \hookrightarrow e \| \operatorname{\mathsf{ret}}(e)} \tag{34.3h}$$

$$\frac{e \longmapsto_{\Sigma} e'}{\mu \parallel \operatorname{dcl}(e; a.m) \longmapsto_{\Sigma} \mu \parallel \operatorname{dcl}(e'; a.m)}$$
(34.3i)

$$\frac{e \operatorname{val}_{\Sigma} \quad \mu \otimes a \hookrightarrow e \parallel m \underset{\Sigma, a}{\longmapsto} \mu' \otimes a \hookrightarrow e' \parallel m'}{\mu \parallel \operatorname{dcl}(e; a.m) \underset{\Sigma}{\longmapsto} \mu' \parallel \operatorname{dcl}(e'; a.m')}$$
(34.3j)

$$\frac{e \operatorname{val}_{\Sigma} \quad e' \operatorname{val}_{\Sigma,a}}{\mu \parallel \operatorname{dcl}(e; a.\operatorname{ret}(e')) \xrightarrow{\Sigma} \mu \parallel \operatorname{ret}(e')}$$
(34.3k)

Rule (34.3a) specifies that a ret command is final if its argument is a value. Rules (34.3c) to (34.3e) specify the dynamics of sequential composition. The expression e must, by virtue of the type system, evaluate to an encapsulated command, which is executed to find its return value, which is then substituted into the command m before executing it.

Rules (34.3i) to (34.3k) define the concept of *block structure* in a programming language. Declarations adhere to the *stack discipline* in that an assignable is allocated during evaluation of the body of the declaration, and deallocated after evaluation of the body is complete. Therefore the lifetime of an assignable can be identified with its scope, and hence we may visualize the dynamic lifetimes of assignables as being nested inside one another, in the same way as their static scopes are nested inside one another. The stack-like behavior of assignables is a characteristic feature of what are known as *Algol-like languages*.

34.1 Basic Commands 313

34.1.3 **Safety**

The judgment $\mu \parallel m$ ok Σ is defined by the rule

$$\frac{\vdash_{\Sigma} m \text{ ok} \quad \mu : \Sigma}{\mu \parallel m \text{ ok}_{\Sigma}}$$
 (34.4)

where the auxiliary judgment μ : Σ is defined by the rule

$$\frac{\forall a \in \Sigma \quad \exists e \quad \mu(a) = e \text{ and } e \text{ val}_{\emptyset} \text{ and } \vdash_{\emptyset} e : \text{nat}}{\mu : \Sigma}$$
 (34.5)

That is, the memory must bind a number to each assignable in Σ .

Theorem 34.1 (Preservation).

- 1. If $e \mapsto_{\Sigma} e'$ and $\vdash_{\Sigma} e : \tau$, then $\vdash_{\Sigma} e' : \tau$.
- 2. If $\mu \parallel m \xrightarrow{\Sigma} \mu' \parallel m'$, with $\vdash_{\Sigma} m$ ok and $\mu : \Sigma$, then $\vdash_{\Sigma} m'$ ok and $\mu' : \Sigma$.

Proof. Simultaneously, by induction on rules (34.2) and (34.3).

Consider rule (34.3j). Assume that $\vdash_{\Sigma} \mathtt{dcl}(e; a.m)$ ok and $\mu : \Sigma$. By inversion of typing we have $\vdash_{\Sigma} e :$ nat and $\vdash_{\Sigma,a} m$ ok. Because e val $_{\Sigma}$ and $\mu : \Sigma$, we have $\mu \otimes a \hookrightarrow e : \Sigma$, a. By induction we have $\vdash_{\Sigma,a} m'$ ok and $\mu' \otimes a \hookrightarrow e' : \Sigma$, a, from which the result follows immediately.

Consider rule (34.3k). Assume that $\vdash_{\Sigma} \mathtt{dcl}(e; a.\mathtt{ret}(e'))$ ok and $\mu : \Sigma$. By inversion we have $\vdash_{\Sigma} e : \mathtt{nat}$, and $\vdash_{\Sigma,a} \mathtt{ret}(e')$ ok, and so $\vdash_{\Sigma,a} e' : \mathtt{nat}$. But because $e' \mathtt{val}_{\Sigma,a}$, and e' is a numeral, and we also have $\vdash_{\Sigma} e' : \mathtt{nat}$, as required.

Theorem 34.2 (Progress).

- 1. If $\vdash_{\Sigma} e : \tau$, then either $e \ val_{\Sigma}$, or there exists e' such that $e \xrightarrow[\Sigma]{} e'$.
- 2. If $\vdash_{\Sigma} m$ ok and $\mu : \Sigma$, then either $\mu \parallel m$ final Σ or $\mu \parallel m \mapsto_{\Sigma} \mu' \parallel m'$ for some μ' and m'.

Proof. Simultaneously, by induction on rules (34.1). Consider rule (34.1d). By the first inductive hypothesis we have either $e \mapsto e'$ or e val $_{\Sigma}$. In the former case rule (34.3i) applies. In the latter, we have by the second inductive hypothesis,

$$\mu \otimes a \hookrightarrow e \parallel m \text{ final}_{\Sigma,a} \quad \text{or} \quad \mu \otimes a \hookrightarrow e \parallel m \underset{\Sigma,a}{\longmapsto} \mu' \otimes a \hookrightarrow e' \parallel m'.$$

In the former case we apply rule (34.3k), and in the latter, rule (34.3j).

34.2 Some Programming Idioms

The language **MA** is designed to expose the elegant interplay between the execution of an expression for its value and the execution of a command for its effect on assignables. In this section we show how to derive several standard idioms of imperative programming in **MA**.

We define the *sequential composition* of commands, written $\{x \leftarrow m_1 ; m_2\}$, to stand for the command bnd $x \leftarrow \text{cmd}(m_1) ; m_2$. Binary composition readily generalizes to an n-ary form by defining

$$\{x_1 \leftarrow m_1; \ldots x_{n-1} \leftarrow m_{n-1}; m_n\}$$

to stand for the iterated composition

$$\{x_1 \leftarrow m_1; \dots \{x_{n-1} \leftarrow m_{n-1}; m_n\}\}$$

We sometimes write just $\{m_1; m_2\}$ for the composition $\{-\leftarrow m_1; m_2\}$ where the returned value from m_1 is ignored; this generalizes in the obvious way to an n-ary form.

A related idiom, the command do e, executes an encapsulated command and returns its result. By definition do e stands for the command bnd $x \leftarrow e$; ret x.

The *conditional* command if (m) m_1 else m_2 executes either m_1 or m_2 according to whether the result of executing m is zero or not:

$$\{x \leftarrow m : do(ifz x \{z \hookrightarrow cmd m_1 \mid s(_) \hookrightarrow cmd m_2\})\}.$$

The returned value of the conditional is the value returned by the selected command.

The *while loop* command while $(m_1)m_2$ repeatedly executes the command m_2 while the command m_1 yields a non-zero number. It is defined as follows:

$$do(fix loop: cmd is cmd(if(m_1) \{ret z\} else\{m_2; do loop\})).$$

This commands runs the self-referential encapsulated command that, when executed, first executes m_1 , branching on the result. If the result is zero, the loop returns zero (arbitrarily). If the result is non-zero, the command m_2 is executed and the loop is repeated.

A *procedure* is a function of type $\tau \to \text{cmd}$ that takes an argument of some type τ and yields an unexecuted command as result. Many procedures have the form λ ($x:\tau$) cmd m, which we abbreviate to proc ($x:\tau$) m. A *procedure call* is the composition of a function application with the activation of the resulting command. If e_1 is a procedure and e_2 is its argument, then the procedure call call $e_1(e_2)$ is defined to be the command do ($e_1(e_2)$), which immediately runs the result of applying e_1 to e_2 .

As an example, here is a procedure of type $\mathtt{nat} \rightharpoonup \mathtt{cmd}$ that returns the factorial of its argument:

```
proc (x:nat) {
  dcl r := 1 in
  dcl a := x in
  { while ( @ a ) {
     y ← @ r
    ; z ← @ a
    ; r := (x-z+1) × y
    ; a := z-1
  }
  ; x ← @ r
  ; ret x
}
```

The loop maintains the invariant that the contents of r is the factorial of x minus the contents of a. Initialization makes this invariant true, and it is preserved by each iteration of the loop, so that upon completion of the loop the assignable a contains a and a contains the factorial of a, as required.

34.3 Typed Commands and Typed Assignables

So far we have restricted the type of the returned value of a command, and the contents of an assignable, to be nat. Can this restriction be relaxed, while adhering to the stack discipline?

The key to admitting returned and assignable values of other types may be uncovered by a close examination of the proof of Theorem 34.1. For the proof to go through it is crucial that values of type nat, the type of assignables and return values, cannot contain an assignable, for otherwise the embedded assignable would escape the scope of its declaration. This property is self-evidently true for eagerly evaluated natural numbers, but fails when they are evaluated lazily. Thus the safety of **MA** hinges on the evaluation order for the successor operation, in contrast to most other situations where either interpretation is also safe.

When extending **MA** to admit assignables and returned values of other types, it is necessary to pay close attention to whether assignables can be embedded in a value of a candidate type. For example, if return values of procedure type are allowed, then the following command violates safety:

```
dcl a := z in \{ret(proc(x:nat) \{a := x\})\}.
```

This command, when executed, allocates a new assignable a and returns a procedure that, when called, assigns its argument to a. But this makes no sense, because the assignable a is deallocated when the body of the declaration returns, but the returned value still refers to it. If the returned procedure is called, execution will get stuck in the attempt to assign to a.

A similar example shows that admitting assignables of procedure type is also unsound. For example, suppose that b is an assignable whose contents are of type $\mathtt{nat} \rightharpoonup \mathtt{cmd}$, and consider the command

```
dcl a := z in \{b := proc(x : nat) cmd(a := x); ret z\}.
```

We assign to b a procedure that uses a locally declared assignable a and then leaves the scope of the declaration. If we then call the procedure stored in b, execution will get stuck attempting to assign to the non-existent assignable a.

To admit declarations that return values other than nat and to admit assignables with contents of types other than nat, we must rework the statics of **MA** to record the returned type of a command and to record the type of the contents of each assignable. First, we generalize the finite set Σ of active assignables to assign a mobile type to each active assignable so that Σ has the form of a finite set of assumptions of the form $a \sim \tau$, where a is an assignable. Second, we replace the judgment $\Gamma \vdash_{\Sigma} m$ ok by the more general form $\Gamma \vdash_{\Sigma} m \stackrel{.}{\sim} \tau$, stating that m is a well-formed command returning a value of type τ . Third, the type cmd is generalized to cmd(τ), which is written in examples as τ cmd, to specify the return type of the encapsulated command.

The statics given in Section 34.1.1 is generalized to admit typed commands and typed assignables as follows:

$$\frac{\Gamma \vdash_{\Sigma} m \stackrel{.}{\sim} \tau}{\Gamma \vdash_{\Sigma} \operatorname{cmd}(m) : \operatorname{cmd}(\tau)}$$
 (34.6a)

$$\frac{\Gamma \vdash_{\Sigma} e : \tau}{\Gamma \vdash_{\Sigma} \operatorname{ret}(e) \stackrel{>}{\sim} \tau}$$
 (34.6b)

$$\frac{\Gamma \vdash_{\Sigma} e : \operatorname{cmd}(\tau) \quad \Gamma, x : \tau \vdash_{\Sigma} m \stackrel{.}{\sim} \tau'}{\Gamma \vdash_{\Sigma} \operatorname{bnd}(e; x.m) \stackrel{.}{\sim} \tau'}$$
(34.6c)

$$\frac{\Gamma \vdash_{\Sigma} e : \tau \quad \tau \text{ mobile} \quad \Gamma \vdash_{\Sigma, a \sim \tau} m \stackrel{.}{\sim} \tau' \quad \tau' \text{ mobile}}{\Gamma \vdash_{\Sigma} \operatorname{dcl}(e; a.m) \stackrel{.}{\sim} \tau'}$$
(34.6d)

$$\frac{}{\Gamma \vdash_{\Sigma, a \sim \tau} \operatorname{get}[a] \div \tau}$$
 (34.6e)

$$\frac{\Gamma \vdash_{\Sigma, a \sim \tau} e : \tau}{\Gamma \vdash_{\Sigma, a \sim \tau} \operatorname{set}[a](e) \stackrel{.}{\sim} \tau}$$
(34.6f)

Apart from the generalization to track returned types and content types, the most important change is that in Rule (34.6d) both the type of a declared assignable and the return type of the declaration is required to be *mobile*. The definition of the judgment τ mobile is guided by the following *mobility condition*:

if
$$\tau$$
 mobile, $\vdash_{\Sigma} e : \tau$ and $e \text{ val}_{\Sigma}$, then $\vdash_{\emptyset} e : \tau$ and $e \text{ val}_{\emptyset}$. (34.7)

That is, a value of mobile type may not depend on any active assignables.

As long as the successor operation is evaluated eagerly, the type nat is mobile:

Similarly, a product of mobile types may safely be deemed mobile, if pairs are evaluated eagerly:

$$\frac{\tau_1 \text{ mobile} \quad \tau_2 \text{ mobile}}{\tau_1 \times \tau_2 \text{ mobile}}$$
 (34.9)

34.4 Notes 317

And the same goes for sums, if the injections are evaluated eagerly:

$$\frac{\tau_1 \text{ mobile} \quad \tau_2 \text{ mobile}}{\tau_1 + \tau_2 \text{ mobile}}$$
 (34.10)

In each of these cases laziness defeats mobility, because values may contain suspended computations that depend on an assignable. For example, if the successor operation for the natural numbers were evaluated lazily, then s(e) would be a value for any expression e including one that refers to an assignable a.

Because the body of a procedure may involve an assignable, no procedure type is mobile, nor is any command type. What about function types other than procedure types? We may think they are mobile, because a pure expression cannot depend on an assignable. Although this is the case, the mobility condition need not hold. For example, consider the following value of type $\mathtt{nat} \rightharpoonup \mathtt{nat}$:

$$\lambda (x: nat) (\lambda (\underline{\cdot}: \tau cmd) z) (cmd {@a}).$$

Although the assignable *a* is not actually needed to compute the result, it nevertheless occurs in the value, violating the mobility condition.

The mobility restriction on the statics of declarations ensures that the type associated to an assignable is always mobile. We may therefore assume, without loss of generality, that the types associated to the assignables in the signature Σ are mobile.

Theorem 34.3 (Preservation for Typed Commands).

- 1. If $e \xrightarrow{\Sigma} e'$ and $\vdash_{\Sigma} e : \tau$, then $\vdash_{\Sigma} e' : \tau$.
- 2. If $\mu \parallel m \xrightarrow{\Sigma} \mu' \parallel m'$, with $\vdash_{\Sigma} m \stackrel{.}{\sim} \tau$ and $\mu : \Sigma$, then $\vdash_{\Sigma} m' \stackrel{.}{\sim} \tau$ and $\mu' : \Sigma$.

Theorem 34.4 (Progress for Typed Commands).

- 1. If $\vdash_{\Sigma} e : \tau$, then either $e \ val_{\Sigma}$, or there exists e' such that $e \xrightarrow[\Sigma]{} e'$.
- 2. If $\vdash_{\Sigma} m \stackrel{.}{\sim} \tau$ and $\mu : \Sigma$, then either $\mu \parallel m$ final Σ or $\mu \parallel m \stackrel{.}{\longmapsto} \mu' \parallel m'$ for some μ' and m'.

The proofs of Theorems 34.3 and 34.4 follows very closely the proof of Theorems 34.1 and 34.2. The main difference is that we appeal to the mobility condition to ensure that returned values and stored values are independent of the active assignables.

34.4 Notes

Modernized Algol is a derivative of Reynolds's Idealized Algol (Reynolds, 1981). In contrast to Reynolds's formulation, Modernized Algol maintains a separation between computations that depend on the memory and those that do not, and does not rely on call-by-name for function application, but rather has a type of encapsulated commands that can be used where call-by-name

318 34.4 Notes

would otherwise be required. The modal distinction between expressions and commands was present in the original formulation of Algol 60, but is developed here in light of the concept of monadic effects introduced by Moggi (1989). Its role in functional programming was emphasized by Wadler (1992). The modal separation in MA is adapted directly from Pfenning and Davies (2001), which stresses the connection to lax modal logic.

What are called *assignables* here are invariably called *variables* elsewhere. The distinction between variables and assignables is blurred in languages that allow assignables as forms of expression. (Indeed, Reynolds himself (personal communication, 2012) regards this as a defining feature of Algol, in opposition to the formulation given here.) In **MA** we choose to make the distinction between variables, which are given meaning by substitution, and assignables, which are given meaning by mutation. Drawing this distinction requires new terminology; the term *assignable* seems apt for the imperative programming concept.

The concept of mobility of a type was introduced in the ML5 language for distributed computing (Murphy et al., 2004), with the similar meaning that a value of a mobile type cannot depend on local resources. Here the mobility restriction is used to ensure that the language adheres to the stack discipline.

Exercises

- **34.1.** Originally Algol had both *scalar* assignables, whose contents are atomic values, and *array* assignables, which is a finite sequence of scalar assignables. Like scalar assignables, array assignables are stack-allocated. Extend **MA** with array assignables, ensuring that the language remains type safe, but allowing that computation may abort if a non-existent array element is accessed.
- **34.2**. Consider carefully the behavior of assignable declarations within recursive procedures, as in the following expression

$$fix p is \lambda (p:\tau) dcl a := e in cmd(m)$$

of type $\tau \rightharpoonup \rho$ cmd for some ρ . Because p is recursive, the body m of the procedure may call itself during its execution, causing the *same* declaration to be executed more than once. Explain the dynamics of getting and setting a in such a situation.

34.3. Originally Algol considered assignables as expressions that stand for their contents in memory. Thus, if a is an assignable containing a number, one could write expressions such as a + a that would evaluate to twice the contents of a. Moreover, one could write commands such as a := a + a to double the contents of a. These conventions encouraged programmers to think of assignables as variables, quite the opposite of their separation in **MA**. This convention, combined with an over-emphasis on concrete syntax, led to a conundrum about the different roles of a in the above assignment command: its meaning on the left of the assignment is different from its meaning on the right. These came to be called the *left-*, or *l-value*, and the *right-*, or *r-value* of the assignable a, corresponding to its position in the assignment statement. When viewed as abstract syntax, though, there is no ambiguity to be explained:

34.4 Notes 319

the assignment operator is indexed by its target assignable, instead of taking as argument an expression that happens to be an assignable, so that the command is set[a](a+a), not set(a; a+a).

This still leaves the puzzle of how to regard assignables as forms of expression. As a first cut, reformulate the dynamics of **MA** to account for this. Reformulate the dynamics of expressions in terms of the judgments $\mu \parallel e \underset{\Sigma}{\longmapsto} \mu' \parallel e'$ and $\mu \parallel e$ final that allow evaluation of e to depend on the contents of the memory. Each use of an assignable as an expression should require one access to the memory. Then prove *memory invariance*:: if $\mu \parallel e \underset{\Sigma}{\longmapsto} \mu' \parallel e'$, then $\mu' = \mu$.

A natural generalization is to allow any sequence of commands to be considered as an expression, if they are all *passive* in the sense that no assignments are allowed. Write do $\{m\}$, where m is a passive command, for a *passive block* whose evaluation consists of executing the command m on the current memory, using its return value as the value of the expression. Observe that memory invariance holds for passive blocks.

The use of an assignable a as an expression may now be rendered as the passive block do $\{@a\}$. More complex uses of assignables as expressions admit several different interpretations using passive blocks. For example, an expression such as a + a might be rendered in one of two ways:

- (a) do $\{@a\}$ + do $\{@a\}$, or
- (b) let x be do $\{@a\}$ in x + x.

The latter formulation accesses a only once, but uses its value twice. Comment on there being two different interpretations of a + a.

34.4. Recursive procedures in Algol are *declared* using a command of the form $proc p(x : \tau) : \rho$ is m in m', which is governed by the typing rule

$$\frac{\Gamma, p : \tau \rightharpoonup \rho \operatorname{cmd}, x : \tau \vdash_{\Sigma} m \stackrel{.}{\sim} \rho \quad \Gamma, p : \tau \rightharpoonup \rho \operatorname{cmd} \vdash_{\Sigma} m' \stackrel{.}{\sim} \tau'}{\Gamma \vdash_{\Sigma} \operatorname{proc} p(x : \tau) : \rho \operatorname{is} m \operatorname{in} m' \stackrel{.}{\sim} \tau'} . \tag{34.11}$$

From the present viewpoint it is peculiar to insist on declaring procedures at all, because they are simply values of procedure type, and even more peculiar to insist that they be confined for use within a command. One justification for this limitation, though, is that Algol included a peculiar feature, called an *own variable*¹ that was declared for use within the procedure, but whose state persisted across calls to the procedure. One application would be to a procedure that generated pseudo-random numbers based on a stored seed that influenced the behavior of successive calls to it. Give a formulation in **MA** of the extended declaration

$$\operatorname{proc} p(x : \tau) : \rho \operatorname{is} \{\operatorname{own} a := e \operatorname{in} m\} \operatorname{in} m'$$

¹That is to say, an *own assignable*.

320 34.4 Notes

where a is declared as an "own" of the procedure p. Contrast the meaning of the foregoing declaration with the following one:

$$\operatorname{proc} p(x:\tau) : \rho \text{ is } \{\operatorname{dcl} a := e \text{ in } m\} \text{ in } m'.$$

- **34.5**. A natural generalization of own assignables is to allow the creation of many such scenarios for a single procedure (or mutually recursive collection of procedures), with each instance creating its own persistent state. This ability motivated the concept of a *class* in Simula-67 as a collection of procedures, possibly mutually recursive, that shared common persistent state. Each instance of a class is called an *object* of that class; calls to its constituent procedures mutate the private persistent state. Formulate this 1967 precursor of imperative object-oriented programming in the context of **MA**.
- **34.6**. There are several ways to formulate an abstract machine for **MA** that accounts for both the *control stack*, which sequences execution (as described in Chapter 28 for **PCF**), and the *data stack*, which records the contents of the assignables. A *consolidated stack* combines these two separate concepts into one, whereas *separated stacks* keeps the memory separate from the control stack, much as we have done in the structural dynamics given by Rules (34.3). In either case the storage required for an assignable is deallocated when exiting the scope of that assignable, a key benefit of the stack discipline for assignables in **MA**.

With a modal separation between expressions and commands it is natural to use a structural dynamics for expressions (given by the transition and value judgments, $e \mapsto e'$ and e val), and a stack machine dynamics for commands.

- (a) Formulate a consolidated stack machine where both assignables and stack frames are recorded on the same stack. Consider states $k \triangleright_{\Sigma} m$, where $\vdash_{\Sigma} k \div \tau$ and $\vdash_{\Sigma} m \stackrel{.}{\sim} \tau$, and $k \triangleleft_{\Sigma} e$, where $\vdash_{\Sigma} k \div \tau$ and $\vdash_{\Sigma} e : \tau$. Comment on the implementation methods required for a consolidated stack.
- (b) Formulate a separated stack machine where the memory is maintained separately from the control stack. Consider states of the form $k \parallel \mu \triangleright_{\Sigma} m$, where $\mu : \Sigma, \vdash_{\Sigma} k \div \tau$, and $\vdash_{\Sigma} m \div \tau$, and of the form $k \parallel \mu \triangleleft_{\Sigma} e$, where $\vdash_{\Sigma} k \div \tau, \vdash_{\Sigma} e : \tau$, and e val.

Chapter 35

Assignable References

A reference to an assignable *a* is a value, written & *a*, of reference type that determines the assignable *a*. A reference to an assignable provides the *capability* to get or set the contents of that assignable, even if the assignable itself is not in scope when it is used. Two references can be compared for equality to test whether they govern the same underlying assignable. If two references are equal, then setting one will affect the result of getting the other; if they are not equal, then setting one cannot influence the result of getting from the other. Two references that govern the same underlying assignable are *aliases*. Aliasing complicates reasoning about programs that use references, because any two references may refer to the assignable.

Reference types are compatible with both a scoped and a scope-free allocation of assignables. When assignables are scoped, the range of significance of a reference type is limited to the scope of the assignable to which it refers. Reference types are therefore immobile, so that they cannot be returned from the body of a declaration, nor stored in an assignable. Although ensuring adherence to the stack discipline, this restriction precludes using references to create mutable data structures, those whose structure can be altered during execution. Mutable data structures have a number of applications in programming, including improving efficiency (often at the expense of expressiveness) and allowing cyclic (self-referential) structures to be created. Supporting mutability requires that assignables be given a scope-free dynamics, so that their lifetime persists beyond the scope of their declaration. Consequently, all types are mobile, so that a value of any type may be stored in an assignable or returned from a command.

35.1 Capabilities

The commands get[a] and set[a](e) in **MA** operate on statically specified assignable a. Even to write these commands requires that the assignable a be in scope where the command occurs. But suppose that we wish to define a procedure that, say, updates an assignable to double its previous value, and returns the previous value. We can write such a procedure for any given assignable, a, but what if we wish to write a generic procedure that works for any given assignable?

One way to do this is give the procedure the *capability* to get and set the contents of some caller-specified assignable. Such a capability is a pair consisting of a *getter* and a *setter* for that assignable. The getter for an assignable a is a command that, when executed, returns the contents of a. The setter for an assignable a is a procedure that, when applied to a value of suitable type, assigns that value to a. Thus, a capability for an assignable a containing a value of type τ is a value of type

$$\tau \operatorname{cap} \triangleq \tau \operatorname{cmd} \times (\tau \rightharpoonup \tau \operatorname{cmd}).$$

A capability for getting and setting an assignable a containing a value of type τ is given by the pair

$$\langle \operatorname{cmd}(@a), \operatorname{proc}(x:\tau) a := x \rangle$$

of type τ cap. Because a capability type is a product of a command type and a procedure type, no capability type is mobile. Thus, a capability cannot be returned from a command, nor stored into an assignable. This is as it should be, for otherwise we would violate the stack discipline for allocating assignables.

The proposed generic doubling procedure is programmed using capabilities as follows:

```
\texttt{proc}\left(\left\langle get, set\right\rangle : \texttt{nat} \ \texttt{cmd} \ \times \left(\ \texttt{nat} \ \rightharpoonup \ \texttt{nat} \ \texttt{cmd}\ \right)\ \right) \ \{x \leftarrow \texttt{do} \ get \ ; y \leftarrow \texttt{do} \ \left(set(x+x)\right)\ ) \ ; \texttt{ret} \ x\}.
```

The procedure is called with the capability to access an assignable a. When executed, it invokes the getter to obtain the contents of a, and then invokes the setter to assign to a, returning the previous value. Observe that the assignable a need not be accessible by this procedure; the capability given by the caller comprises the commands required to get and set a.

35.2 Scoped Assignables

A weakness of using a capability to give indirect access to an assignable is that there is no guarantee that a given getter/setter pair are in fact the capability for a particular assignable. For example, we might pair the getter for a with the setter for b, leading to unexpected behavior. There is nothing in the type system that prevents creating such mismatched pairs.

To avoid this we introduce the concept of a *reference* to an assignable. A reference is a value from which we may obtain the capability to get and set a particular assignable. Moreover, two references can be tested for equality to see whether they act on the same assignable. The *reference type* $ref(\tau)$ has as values references to assignables of type τ . The introduction and elimination forms for this type are given by the following syntax chart:

¹The getter and setter do not suffice to define equality, because not all types admit a test for equality. When they do, and when there are at least two distinct values of their type, we can determine whether they are aliases by assigning to one and checking whether the contents of the other is changed.

The statics of reference types is defined by the following rules:

$$\frac{}{\Gamma \vdash_{\Sigma, a \sim \tau} \operatorname{ref}[a] : \operatorname{ref}(\tau)} \tag{35.1a}$$

$$\frac{\Gamma \vdash_{\Sigma} e : \text{ref}(\tau)}{\Gamma \vdash_{\Sigma} \text{getref}(e) \stackrel{.}{\sim} \tau}$$
 (35.1b)

$$\frac{\Gamma \vdash_{\Sigma} e_1 : \text{ref}(\tau) \quad \Gamma \vdash_{\Sigma} e_2 : \tau}{\Gamma \vdash_{\Sigma} \text{setref}(e_1; e_2) \stackrel{.}{\sim} \tau}$$
(35.1c)

Rule (35.1a) specifies that the name of any active assignable is an expression of type $ref(\tau)$. The dynamics of reference types defers to the corresponding operations on assignables, and does not alter the underlying dynamics of assignables:

$$\frac{1}{\operatorname{ref}[a]\operatorname{val}_{\Sigma,a\sim\tau}} \tag{35.2a}$$

$$\frac{e \underset{\Sigma}{\longmapsto} e'}{\mu \parallel \operatorname{getref}(e) \underset{\Sigma}{\longmapsto} \mu \parallel \operatorname{getref}(e')}$$
(35.2b)

$$\overline{\mu \parallel \operatorname{getref}(\operatorname{ref}[a])} \xrightarrow{\Sigma, a \sim \tau} \mu \parallel \operatorname{get}[a]$$
(35.2c)

$$\frac{e_1 \underset{\Sigma}{\longmapsto} e_1'}{\mu \parallel \mathsf{setref}(e_1; e_2) \underset{\Sigma}{\longmapsto} \mu \parallel \mathsf{setref}(e_1'; e_2)} \tag{35.2d}$$

$$\frac{1}{\mu \parallel \operatorname{setref}(\operatorname{ref}[a]; e) \xrightarrow{\Sigma, a \sim \tau} \mu \parallel \operatorname{set}[a](e)}$$
(35.2e)

A reference to an assignable is a value. The getref and setref operations on references defer to the corresponding operations on assignables once the referent has been resolved.

Because references give rise to capabilities, the reference type is immobile. As a result references cannot be stored in assignables or returned from commands. The immobility of references ensures safety, as can be seen by extending the safety proof given in Chapter 34.

As an example of using references, the generic doubling procedure discussed in the preceding section is programmed using references as follows:

$$\texttt{proc}(r:\texttt{natref})\{x \leftarrow *r; r*=x+x; \texttt{ret}x\}.$$

Because the argument is a reference, rather than a capability, there is no possibility that the getter and setter refer to different assignables.

The ability to pass references to procedures comes at a price, because any two references might refer to the same assignable (if they have the same type). Consider a procedure that, when given two references x and y, adds twice the contents of y to the contents of x. One way to write this code creates no complications:

$$\lambda$$
 (x :natref) λ (y :natref) cmd { $x' \leftarrow *x; y' \leftarrow *y; x*=x'+y'+y'$ }.

Even if *x* and *y* refer to the same assignable, the effect will be to set the contents of the assignable referenced by *x* to the sum of its original contents and twice the contents of the assignable referenced by *y*.

But now consider the following seemingly equivalent implementation of this procedure:

$$\lambda$$
 (x:natref) λ (y:natref) cmd { $x+=y$; $x+=y$ },

where x + = y is the command

$$\{x' \leftarrow *x; y' \leftarrow *y; x *= x' + y'\}$$

that adds the contents of y to the contents of x. The second implementation works right, as long as x and y do not refer to the same assignable. If they do refere to a common assignable a, with contents n, the result is that a is to set $4 \times n$, instead of the intended $3 \times n$. The second get of y is affected by the first set of x.

In this case it is clear how to avoid the problem: use the first implementation, rather than the second. But the difficulty is not in fixing the problem once it has been discovered, but in noticing the problem in the first place. Wherever references (or capabilities) are used, the problems of interference lurk. Avoiding them requires very careful consideration of all possible aliasing relationships among all of the references in play. The problem is that the number of possible aliasing relationships among n references grows combinatorially in n.

35.3 Free Assignables

Although it is interesting to note that references and capabilities are compatible with the stack discipline, for references to be useful requires that this restriction be relaxed. With immobile references it is impossible to build data structures containing references, or to return references from procedures. To allow this we must arrange that the lifetime of an assignable extend beyond its scope. In other words we must give up stack allocation for heap allocation. Assignables that persist beyond their scope of declaration are called *scope-free*, or just *free*, assignables. When all assignables are free, every type is mobile and so any value, including a reference, may be used in a data structure.

Supporting free assignables amounts to changing the dynamics so that allocation of assignables persists across transitions. We use transition judgments of the form

$$\nu \Sigma \{ \mu \parallel m \} \longmapsto \nu \Sigma' \{ \mu' \parallel m' \}.$$

Execution of a command may allocate new assignables, may alter the contents of existing assignables, and may give rise to a new command to be executed at the next step. The rules defining the dynamics of free assignables are as follows:

$$\frac{e \operatorname{val}_{\Sigma}}{\nu \Sigma \{ \mu \parallel \operatorname{ret}(e) \} \text{ final}}$$
 (35.3a)

$$\frac{e \longmapsto_{\Sigma} e'}{\nu \Sigma \{ \mu \parallel \operatorname{ret}(e) \} \longmapsto_{\nu} \Sigma \{ \mu \parallel \operatorname{ret}(e') \}}$$
(35.3b)

$$\frac{e \longmapsto_{\Sigma} e'}{\nu \Sigma \{ \mu \parallel \operatorname{bnd}(e; x.m) \} \longmapsto_{} \nu \Sigma \{ \mu \parallel \operatorname{bnd}(e'; x.m) \}}$$
(35.3c)

$$\frac{e \operatorname{val}_{\Sigma}}{\nu \Sigma \{\mu \parallel \operatorname{bnd}(\operatorname{cmd}(\operatorname{ret}(e)); x.m)\} \longmapsto \nu \Sigma \{\mu \parallel [e/x]m\}}$$
(35.3d)

$$\frac{\nu \Sigma \{\mu \parallel m_1\} \longmapsto \nu \Sigma' \{\mu' \parallel m_1'\}}{\nu \Sigma \{\mu \parallel \operatorname{bnd}(\operatorname{cmd}(m_1); x.m_2)\} \longmapsto \nu \Sigma' \{\mu' \parallel \operatorname{bnd}(\operatorname{cmd}(m_1'); x.m_2)\}}$$
(35.3e)

$$\frac{1}{\nu \Sigma, a \sim \tau \{\mu \otimes a \hookrightarrow e \parallel get[a]\} \longmapsto \nu \Sigma, a \sim \tau \{\mu \otimes a \hookrightarrow e \parallel ret(e)\}}$$
(35.3f)

$$\frac{e \underset{\Sigma}{\longmapsto} e'}{\nu \Sigma \{\mu \parallel \operatorname{set}[a](e)\} \longmapsto \nu \Sigma \{\mu \parallel \operatorname{set}[a](e')\}}$$
(35.3g)

$$\frac{e \operatorname{val}_{\Sigma, a \sim \tau}}{\nu \, \Sigma, a \sim \tau \, \{ \, \mu \otimes a \hookrightarrow_{-} \, \| \, \operatorname{set}[\, a \,](\, e \,) \, \} \longmapsto \nu \, \Sigma, a \sim \tau \, \{ \, \mu \otimes a \hookrightarrow e \, \| \, \operatorname{ret}(\, e \,) \, \}}$$
(35.3h)

$$\frac{e \longmapsto_{\Sigma} e'}{\nu \sum \{ \mu \parallel \operatorname{dcl}(e; a.m) \} \longmapsto_{\nu} \nu \sum \{ \mu \parallel \operatorname{dcl}(e'; a.m) \}}$$
(35.3i)

$$\frac{e \operatorname{val}_{\Sigma}}{\nu \Sigma \{\mu \parallel \operatorname{dcl}(e; a.m)\} \longmapsto \nu \Sigma, a \sim \tau \{\mu \otimes a \hookrightarrow e \parallel m\}} \tag{35.3j}$$

The language **RMA** extends **MA** with references to free assignables. Its dynamics is similar to that of references to scoped assignables given earlier.

$$\frac{e \longmapsto_{\Sigma} e'}{\nu \sum \{ \mu \parallel \text{getref}(e) \} \longmapsto_{\nu} \nu \sum \{ \mu \parallel \text{getref}(e') \}}$$
(35.4a)

$$\frac{}{\nu \Sigma \{\mu \parallel \operatorname{getref}(\operatorname{ref}[a])\} \longmapsto \nu \Sigma \{\mu \parallel \operatorname{get}[a]\}}$$
(35.4b)

$$\frac{e_1 \underset{\Sigma}{\longmapsto} e'_1}{\nu \sum \{ \mu \parallel \operatorname{setref}(e_1; e_2) \} \longmapsto \nu \sum \{ \mu \parallel \operatorname{setref}(e'_1; e_2) \}}$$
(35.4c)

326 35.4 Safety

$$\frac{}{\nu \Sigma \left\{ \mu \parallel \operatorname{setref}(\operatorname{ref}[a]; e_2) \right\} \longmapsto \nu \Sigma \left\{ \mu \parallel \operatorname{set}[a](e_2) \right\}}$$
 (35.4d)

The expressions cannot alter or extend the memory, only commands may do so.

As an example of using **RMA**, consider the command newref $[\tau](e)$ defined by

$$dcl a := e \operatorname{inret}(\&a). \tag{35.5}$$

This command allocates a fresh assignable, and returns a reference to it. Its static and dynamics are derived from the foregoing rules as follows:

$$\frac{\Gamma \vdash_{\Sigma} e : \tau}{\Gamma \vdash_{\Sigma} \text{newref}[\tau](e) \stackrel{\sim}{\sim} \text{ref}(\tau)}$$
(35.6)

$$\frac{e \longmapsto_{\Sigma} e'}{\nu \Sigma \{ \mu \parallel \text{newref}[\tau](e) \} \longmapsto_{\nu} \Sigma \{ \mu \parallel \text{newref}[\tau](e') \}}$$
(35.7a)

$$\frac{e \operatorname{val}_{\Sigma}}{\nu \Sigma \{\mu \parallel \operatorname{newref}[\tau](e)\} \longmapsto \nu \Sigma, a \sim \tau \{\mu \otimes a \hookrightarrow e \parallel \operatorname{ret}(\operatorname{ref}[a])\}}$$
(35.7b)

Oftentimes the command $newref[\tau](e)$ is taken as primitive, and the declaration command is omitted. In that case all assignables are accessed by reference, and no direct access to assignables is provided.

35.4 Safety

Although the proof of safety for references to scoped assignables presents few difficulties, the safety for free assignables is tricky. The main difficulty is to account for cyclic dependencies within data structures. The contents of one assignable may contain a reference to itself, or a reference to another assignable that contains a reference to it, and so forth. For example, consider the following procedure e of type nat \rightarrow nat cmd:

$$\texttt{proc}(x:\texttt{nat})\{\texttt{if}(x)\texttt{ret}(1)\texttt{else}\{f\leftarrow @a;y\leftarrow f(x-1);\texttt{ret}(x*y)\}\}.$$

Let μ be a memory of the form $\mu' \otimes a \hookrightarrow e$ in which the contents of a contains, via the body of the procedure, a reference to a itself. Indeed, if the procedure e is called with a non-zero argument, it will "call itself" by indirect reference through a.

Cyclic dependencies complicate the definition of the judgment μ : Σ . It is defined by the following rule:

$$\frac{\vdash_{\Sigma} m \stackrel{.}{\sim} \tau \vdash_{\Sigma} \mu : \Sigma}{\nu \Sigma \{ \mu \parallel m \} \text{ ok}}$$
 (35.8)

The first premise of the rule states that the command m is well-formed relative to Σ . The second premise states that the memory μ conforms to Σ , relative to all of Σ so that cyclic dependencies are permitted. The judgment $\vdash_{\Sigma'} \mu : \Sigma$ is defined as follows:

$$\frac{\forall a \sim \tau \in \Sigma \quad \exists e \quad \mu(a) = e \text{ and } \vdash_{\Sigma'} e : \tau}{\vdash_{\Sigma'} \mu : \Sigma}$$
(35.9)

35.4 Safety

327

Theorem 35.1 (Preservation).

- 1. If $\vdash_{\Sigma} e : \tau$ and $e \mapsto_{\Sigma} e'$, then $\vdash_{\Sigma} e' : \tau$.
- 2. If $\nu \Sigma \{ \mu \parallel m \}$ ok and $\nu \Sigma \{ \mu \parallel m \} \longmapsto \nu \Sigma' \{ \mu' \parallel m' \}$, then $\nu \Sigma' \{ \mu' \parallel m' \}$ ok.

Proof. Simultaneously, by induction on transition. We prove the following stronger form of the second statement:

If $\nu \Sigma \{ \mu \parallel m \} \longmapsto \nu \Sigma' \{ \mu' \parallel m' \}$, where $\vdash_{\Sigma} m \stackrel{.}{\sim} \tau$, $\vdash_{\Sigma} \mu : \Sigma$, then Σ' extends Σ , and $\vdash_{\Sigma'} m' \stackrel{.}{\sim} \tau$, and $\vdash_{\Sigma'} \mu' : \Sigma'$.

Consider the transition

$$\nu \Sigma \{ \mu \parallel dcl(e; a.m) \} \longmapsto \nu \Sigma, a \sim \rho \{ \mu \otimes a \hookrightarrow e \parallel m \}$$

where $e \text{ val}_{\Sigma}$. By assumption and inversion of rule (34.6d) we have $\vdash_{\Sigma} e : \rho, \vdash_{\Sigma, a \sim \rho} m \stackrel{.}{\sim} \tau$, and $\vdash_{\Sigma} \mu : \Sigma$. But because extension of Σ with a fresh assignable does not affect typing, we also have $\vdash_{\Sigma, a \sim \rho} \mu : \Sigma$ and $\vdash_{\Sigma, a \sim \rho} e : \rho$, from which it follows by rule (35.9) that $\vdash_{\Sigma, a \sim \rho} \mu \otimes a \hookrightarrow e : \Sigma, a \sim \rho$. The other cases follow a similar pattern, and are left as an exercise for the reader.

Theorem 35.2 (Progress).

- 1. If $\vdash_{\Sigma} e : \tau$, then either e val $_{\Sigma}$ or there exists e' such that $e \xrightarrow[\Sigma]{} e'$.
- 2. If $\nu \Sigma \{ \mu \parallel m \}$ ok then either $\nu \Sigma \{ \mu \parallel m \}$ final or $\nu \Sigma \{ \mu \parallel m \} \mapsto \nu \Sigma' \{ \mu' \parallel m' \}$ for some Σ' , μ' , and m'.

Proof. Simultaneously, by induction on typing. For the second statement we prove the stronger form

If $\vdash_{\Sigma} m \stackrel{.}{\sim} \tau$ and $\vdash_{\Sigma} \mu : \Sigma$, then either $\nu \Sigma \{ \mu \parallel m \}$ final, or $\nu \Sigma \{ \mu \parallel m \} \longmapsto \nu \Sigma' \{ \mu' \parallel m' \}$ for some Σ' , μ' , and m'.

Consider the typing rule

$$\frac{\Gamma \vdash_{\Sigma} e : \rho \quad \Gamma \vdash_{\Sigma, a \sim \rho} m \stackrel{.}{\sim} \tau}{\Gamma \vdash_{\Sigma} \operatorname{dcl}(e; a.m) \stackrel{.}{\sim} \tau}$$

We have by the first inductive hypothesis that either e val $_{\Sigma}$ or $e \mapsto_{\Sigma} e'$ for some e'. In the latter case we have by rule (35.3i)

$$\nu \Sigma \{ \mu \parallel dcl(e; a.m) \} \longmapsto \nu \Sigma \{ \mu \parallel dcl(e'; a.m) \}.$$

In the former case we have by rule (35.3j) that

$$\nu \Sigma \{ \mu \parallel dcl(e; a.m) \} \longmapsto \nu \Sigma, a \sim \rho \{ \mu \otimes a \hookrightarrow e \parallel m \}.$$

328 35.5 Benign Effects

Now consider the typing rule

$$\overline{\Gamma \vdash_{\Sigma,a \sim \tau} \operatorname{get}[a] \stackrel{.}{\sim} \tau}$$

By assumption $\vdash_{\Sigma, a \sim \tau} \mu : \Sigma, a \sim \tau$, and hence there exists e val $_{\Sigma, a \sim \tau}$ such that $\mu = \mu' \otimes a \hookrightarrow e$ and $\vdash_{\Sigma, a \sim \tau} e : \tau$. By rule (35.3f)

$$\nu \Sigma, a \sim \tau \{ \mu' \otimes a \hookrightarrow e \mid | get[a] \} \longmapsto \nu \Sigma, a \sim \tau \{ \mu' \otimes a \hookrightarrow e \mid | ret(e) \},$$

as required. The other cases are handled similarly.

35.5 Benign Effects

The modal separation between commands and expressions ensures that the meaning of an expression does not depend on the (ever-changing) contents of assignables. Although this is helpful in many, perhaps most, situations, it also precludes programming techniques that use storage effects to implement purely functional behavior. A prime example is memoization. Externally, a suspended computation behaves exactly like the underlying computation; internally, an assignable is associated with the computation that stores the result of any evaluation of the computation for future use. Other examples are self-adjusting data structures, which use state to improve their efficiency without changing their functional behavior. For example, a splay tree is a binary search tree that uses mutation internally to re-balance the tree as elements are inserted, deleted, and retrieved, so that lookup takes time proportional to the logarithm of the number of elements.

These are examples of *benign storage effects*, uses of mutation in a data structure to improve efficiency without disrupting its functional behavior. One class of examples are self-adjusting data structures that reorganize themselves during one use to improve efficiency of later uses. Another class of examples are memoized, or lazy, data structures, which are discussed in Chapter 36. Benign effects such as these are impossible to implement if a strict separation between expressions and commands is maintained. For example, a self-adjusting tree involves mutation, but is a value just like any other, and this cannot be achieved in **MA**. Although several special-case techniques are known, the most general solution is to do away with the modal distinction, coalescing expressions and commands into a single syntactic category. The penalty is that the type system no longer ensures that an expression of type τ denotes a value of that type; it might also have storage effects during its evaluation. The benefit is that one may freely use benign effects, but it is up to the programmer to ensure that they truly are benign.

The language **RPCF** extends **PCF** with references to free assignables. The following rules define the statics of the distinctive features of **RPCF**:

$$\frac{\Gamma \vdash_{\Sigma} e_1 : \tau_1 \quad \Gamma \vdash_{\Sigma, a \sim \tau_1} e_2 : \tau_2}{\Gamma \vdash_{\Sigma} \operatorname{dcl}(e_1; a.e_2) : \tau_2}$$
(35.10a)

$$\frac{}{\Gamma \vdash_{\Sigma, a \sim \tau} \operatorname{get}[a] : \tau}$$
 (35.10b)

35.6 Notes 329

$$\frac{\Gamma \vdash_{\Sigma,a} \sim_{\tau} e : \tau}{\Gamma \vdash_{\Sigma,a} \sim_{\tau} \operatorname{set}[a](e) : \tau}$$
(35.10c)

Correspondingly, the dynamics of **RPCF** is given by transitions of the form

$$\nu \Sigma \{ \mu \parallel e \} \longmapsto \nu \Sigma' \{ \mu' \parallel e' \},$$

where e is an expression, and not a command. The rules defining the dynamics are very similar to those for **RMA**, but with commands and expressions integrated into a single category.

To illustrate the concept of a benign effect, consider the technique of *back-patching* to implement recursion. Here is an implementation of the factorial function that uses an assignable to implement recursive calls:

```
dcl a := \lambda n : \text{nat.0 in}

{ f \leftarrow a := \lambda n : \text{nat.ifz}(n, 1, n' . n \times (@a)(n'))

; \text{ret}(f)

}
```

This declaration returns a function of type $\mathtt{nat} \rightharpoonup \mathtt{nat}$ that is obtained by (a) allocating a free assignable initialized arbitrarily with a function of this type, (b) defining a λ -abstraction in which each "recursive call" consists of retrieving and applying the function stored in that assignable, (c) assigning this function to the assignable, and (d) returning that function. The result is a function on the natural numbers, even though it uses state in its implementation.

Backpatching is not expressible in **RMA**, because it relies on assignment. Let us attempt to recode the previous example in **RMA**:

```
dcl a := proc(n:nat)\{ret 0\} in \{ f \leftarrow a := \dots \}; ret(f) \},
```

where the elided procedure assigned to a is given by

```
proc(n:nat) {if (ret(n)) {ret(1)} else {f\leftarrow@a; x\leftarrowf(n-1); ret(n\timesx)}}.
```

The difficulty is that what we have is a command, not an expression. Moreover, the result of the command is of the procedure type $\mathtt{nat} \rightharpoonup (\mathtt{nat}\,\mathtt{cmd})$, and not of the function type $\mathtt{nat} \rightharpoonup \mathtt{nat}$. Consequently, we cannot use the factorial procedure in an expression, but have to execute it as a command using code such as this:

```
\{ f \leftarrow fact; x \leftarrow f(n); ret(x) \}.
```

35.6 Notes

Reynolds (1981) uses capabilities to provide indirect access to assignables; references are just an abstract form of capability. References are often permitted only for free assignables, but with

330 35.6 Notes

mobility restrictions one may also have references to scoped assignables. The proof of safety of free references outlined here follows those given by Wright and Felleisen (1994) and Harper (1994).

Benign effects are central to the distinction between Haskell, which provides an Algol-like separation between commands and expressions, and ML, which integrates evaluation with execution. The choice between them is classic trade-off, with neither superior to the other in all respects.

Exercises

- **35.1.** Consider scoped array assignables as described in Exercise **34.1**. Extend the treatment of array assignables in Exercise **34.1**, to account for array assignable references.
- **35.2.** References to scope-free assignables are often used to implement recursive data structures such as mutable lists and trees. Examine such data structures in the context of **RMA** enriched with sum, product, and recursive types.

Give six different types that could be considered a type of linked lists, according to the following characteristics:

- (a) A mutable list may only be updated *in toto* by replacing it with another (immutable) list.
- (b) A mutable list can be altered in one of two ways, to make it empty, or to change both its head and tail element simultaneously. The tail element is any other such mutable list, so circularities may arise.
- (c) A mutable list is, permanently, either empty or non-empty. If not, both its head and tail can be modified simultaneously.
- (d) A mutable list is, permanently, either empty or non-empty. If not, its tail, but not its head, can be set to another such list.
- (e) A mutable list is, permanently, either empty or non-empty. If not, either its head or its tail elements can be modified independently.
- (f) A mutable list can be altered to become either empty or non-empty. If it is non-empty, either it head, or its tail, can be modified independently of one another.

Discuss the merits and deficiencies of each representation.





Chapter 37

Nested Parallelism

Parallel computation seeks to reduce the running times of programs by allowing many computations to be carried out simultaneously. For example, if we wish to add two numbers, each given by a complex computation, we may consider evaluating the addends simultaneously, then computing their sum. The ability to exploit parallelism is limited by the dependencies among parts of a program. Obviously, if one computation depends on the result of another, then we have no choice but to execute them sequentially so that we may propagate the result of the first to the second. Consequently, the fewer dependencies among sub-computations, the greater the opportunities for parallelism.

In this chapter we discuss the language **PPCF**, which is the extension of **PCF** with *nested parallelism*. Nested parallelism has a hierarchical structure arising from *forking* two (or more) parallel computations, then *joining* these computations to combine their results before proceeding. Nested parallelism is also known as *fork-join parallelism*. We will consider two forms of dynamics for nested parallelism. The first is a structural dynamics in which a single transition on a compound expression may involve multiple transitions on its constituent expressions. The second is a cost dynamics (introduced in Chapter 7) that focuses attention on the sequential and parallel complexity (also known as the *work* and the *depth*, or *span*) of a parallel program by associating a *series-parallel graph* with each computation.

37.1 Binary Fork-Join

The syntax of **PPCF** extends that of **PCF** with the following construct:

$$\mathsf{Exp} \quad e \quad ::= \quad \mathsf{par}\big(\,e_1;e_2;x_1.x_2.e\,\big) \quad \mathsf{par}\,x_1 = e_1\,\mathsf{and}\,x_2 = e_2\,\mathsf{in}\,e \quad \mathsf{parallel}\,\mathsf{let}$$

The variables x_1 and x_2 are bound only within e, and not within e_1 or e_2 , which ensures that they are not mutually dependent and hence can be evaluated simultaneously. The variable bindings represent a fork of two parallel computations e_1 and e_2 , and the body e represents their join.

The static of **PPCF** enriches that of **PCF** with the following rule for parallel let:

$$\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2 \quad \Gamma, x_1 : \tau_1, x_2 : \tau_2 \vdash e : \tau}{\Gamma \vdash \operatorname{par}(e_1; e_2; x_1. x_2. e) : \tau}$$
(37.1)

The sequential structural dynamics of **PPCF** is defined by a transition judgment of the form $e \mapsto_{\text{seq}} e'$ defined by these rules:

$$\frac{e_1 \underset{\text{seq}}{\longmapsto} e'_1}{\underset{\text{par}(e_1; e_2; x_1. x_2. e)}{\longmapsto} \underset{\text{seq}}{\longmapsto} \text{par}(e'_1; e_2; x_1. x_2. e)}$$
(37.2a)

$$\frac{e_1 \text{ val} \quad e_2 \underset{\text{seq}}{\longmapsto} e_2'}{\operatorname{par}(e_1; e_2; x_1.x_2.e)} \xrightarrow{\text{seq}} \operatorname{par}(e_1; e_2'; x_1.x_2.e)$$
(37.2b)

$$\frac{e_1 \text{ val} \quad e_2 \text{ val}}{\operatorname{par}(e_1; e_2; x_1. x_2. e) \underset{\text{seq}}{\longmapsto} [e_1, e_2/x_1, x_2]e}$$
(37.2c)

The parallel structural dynamics of **PPCF** is given by a transition judgment of the form $e \mapsto_{par} e'$, defined as follows:

$$\frac{e_1 \underset{\mathsf{par}}{\longmapsto} e_1' \quad e_2 \underset{\mathsf{par}}{\longmapsto} e_2'}{\operatorname{par}(e_1; e_2; x_1. x_2. e) \underset{\mathsf{par}}{\longmapsto} \operatorname{par}(e_1'; e_2'; x_1. x_2. e)}$$
(37.3a)

$$\frac{e_1 \underset{\mathsf{par}}{\longmapsto} e_1' \quad e_2 \text{ val}}{\operatorname{par}(e_1; e_2; x_1. x_2. e) \underset{\mathsf{par}}{\longmapsto} \operatorname{par}(e_1'; e_2; x_1. x_2. e)}$$
(37.3b)

$$\frac{e_1 \text{ val} \quad e_2 \underset{\text{par}}{\longmapsto} e_2'}{\underset{\text{par}}{\longmapsto} e_2'}$$

$$\frac{e_1 \text{ val} \quad e_2 \underset{\text{par}}{\longmapsto} e_2'}{\underset{\text{par}}{\longmapsto} par(e_1; e_2'; x_1.x_2.e)}$$
(37.3c)

$$\frac{e_1 \text{ val} \quad e_2 \text{ val}}{\operatorname{par}(e_1; e_2; x_1. x_2. e) \underset{\mathsf{par}}{\longmapsto} [e_1, e_2/x_1, x_2] e}$$
(37.3d)

The parallel dynamics abstracts away from any limitations on processing capacity; such limitations are considered in Section 37.4.

The *implicit parallelism theorem* states that the sequential and the parallel dynamics coincide. Consequently, we need never be concerned with the *meaning* of a parallel program (its meaning is given by the sequential dynamics), but only with its *efficiency*. As a practical matter, this means that a program can be developed on a sequential platform, even if it is meant to run on a parallel platform, because the behavior is not affected by whether we execute it using a sequential or a parallel dynamics. Because the sequential dynamics is deterministic (every expression has at most

one value), the implicit parallelism theorem implies that the parallel dynamics is also deterministic. For this reason the implicit parallelism theorem is also known as the *deterministic parallelism theorem*. This terminology emphasizes the distinction between *deterministic parallelism*, the subject of this chapter, from *non-deterministic concurrency*, the subject of Chapters 39 and 40.

A proof of the implicit parallelism theorem can be given by giving an evaluation dynamics $e \downarrow v$ in the style of Chapter 7, and showing that

$$e \mapsto^* v$$
 iff $e \Downarrow v$ iff $e \mapsto^* v$

(where v is a closed expression such that v val). The most important rule of the evaluation dynamics is for the evaluation of a parallel let:

$$\frac{e_1 \Downarrow v_1 \quad e_2 \Downarrow v_2 \quad [v_1, v_2/x_1, x_2]e \Downarrow v}{\operatorname{par}(e_1; e_2; x_1. x_2.e) \Downarrow v}$$
(37.4)

The other rules are easily derived from the structural dynamics of **PCF** as in Chapter 7.

It is possible to show that the sequential dynamics of **PPCF** agrees with its evaluation dynamics by extending the proof of Theorem 7.2.

Lemma 37.1. For all v val, $e \mapsto_{seq}^* v$ if, and only if, $e \Downarrow v$.

Proof. It suffices to show that if $e \mapsto_{\text{seq}} e'$ and $e' \Downarrow v$, then $e \Downarrow v$, and that if $e_1 \mapsto_{\text{seq}}^* v_1$ and $e_2 \mapsto_{\text{seq}}^* v_2$ and $[v_1, v_2/x_1, x_2]e \mapsto_{\text{seq}}^* v$, then

$$\operatorname{par} x_1 = e_1 \operatorname{and} x_2 = e_2 \operatorname{in} e \xrightarrow[\operatorname{seq}]{}^* v.$$

By a similar argument we may show that the parallel dynamics also agrees with the evaluation dynamics, and hence with the sequential dynamics.

Lemma 37.2. For all v val, $e \mapsto_{par}^* v$ if, and only if, $e \Downarrow v$.

Proof. It suffices to show that if $e \mapsto_{\mathsf{par}} e'$ and $e' \Downarrow v$, then $e \Downarrow v$, and that if $e_1 \mapsto_{\mathsf{par}} v_1$ and $e_2 \mapsto_{\mathsf{par}} v_2$ and $[v_1, v_2/x_1, x_2]e \mapsto_{\mathsf{par}} v$, then

$$\operatorname{par} x_1 = e_1 \operatorname{and} x_2 = e_2 \operatorname{in} e \underset{\operatorname{par}}{\longmapsto}^* v.$$

The proof of the first is by induction on the parallel dynamics. The proof of the second proceeds by simultaneous induction on the derivations of $e_1 \underset{\mathsf{par}}{\longmapsto} v_1$ and $e_2 \underset{\mathsf{par}}{\longmapsto} v_2$. If $e_1 = v_1$ with v_1 val and $e_2 = v_2$ with v_2 val, then the result follows immediately from the third premise. If $e_2 = v_2$ but $e_1 \underset{\mathsf{par}}{\longmapsto} e_1' \underset{\mathsf{par}}{\longmapsto} v_1$, then by induction we have that $\mathsf{par} \ x_1 = e_1' \ \mathsf{and} \ x_2 = v_2 \ \mathsf{in} \ e_1 \underset{\mathsf{par}}{\longmapsto} v_1$, and hence the result follows by an application of rule (37.3b). The symmetric case follows similarly by an application of rule (37.3c), and in the case that both e_1 and e_2 transition, the result follows by induction and rule (37.3a).

Theorem 37.3 (Implicit Parallelism). The sequential and parallel dynamics coincide: for all v val, $e \mapsto_{seq}^* v$ iff $e \mapsto_{par}^* v$.

Proof. By Lemmas 37.1 and 37.2. □

The implicit parallelism theorem states that parallelism does not affect the meaning of a program, only the efficiency of its execution. Correctness is not affected by parallelism, only efficiency.

37.2 Cost Dynamics

In this section we define a *parallel cost dynamics* that assigns a *cost graph* to the evaluation of a **PPCF** expression. Cost graphs are defined by the following grammar:

$$\begin{array}{ccccc} \mathsf{Cost} & c & ::= & \mathbf{0} & \mathsf{zero} \ \mathsf{cost} \\ & & \mathbf{1} & \mathsf{unit} \ \mathsf{cost} \\ & & c_1 \otimes c_2 & \mathsf{parallel} \ \mathsf{combination} \\ & & & c_1 \oplus c_2 & \mathsf{sequential} \ \mathsf{combination} \end{array}$$

A cost graph is a *series-parallel* ordered directed acyclic graph, with a designated *source* node and *sink* node. For **0** the graph consists of one node and no edges, with the source and sink both being the node itself. For **1** the graph consists of two nodes and one edge directed from the source to the sink. For $c_1 \otimes c_2$, if g_1 and g_2 are the graphs of g_1 and g_2 , respectively, then the graph has two extra nodes, a source node with two edges to the source nodes of g_1 and g_2 , and a sink node, with edges from the sink nodes of g_1 and g_2 to it. The children of the source are ordered according to the sequential evaluation order. Finally, for $g_1 \oplus g_2$, where g_1 and g_2 are the graphs of g_1 and g_2 to the source node the source of g_1 , as sink node the sink of g_2 , and an edge from the sink of g_1 to the source of g_2 .

The intuition behind a cost graph is that nodes represent subcomputations of an overall computation, and edges represent *sequentiality constraints* stating that one computation depends on the result of another, and hence cannot be started before the one on which it depends completes. The product of two graphs represents *parallelism opportunities* in which there are no sequentiality constraints between the two computations. The assignment of source and sink nodes reflects the overhead of *forking* two parallel computations and *joining* them after they have both completed. At the structural level, we note that only the root has no ancestors, and only the final node of the cost graph has no descendents. Interior nodes may have one or two descendents, the former representing a sequential dependency, and the latter representing a *fork point*. Such nodes may have one or two ancestors, the former corresponding to a sequential dependency and the latter representing a *join point*.

We associate with each cost graph two numeric measures, the work, wk(c), and the depth, dp(c).

37.2 Cost Dynamics 347

The work is defined by the following equations:

$$wk(c) = \begin{cases} 0 & \text{if } c = \mathbf{0} \\ 1 & \text{if } c = \mathbf{1} \\ wk(c_1) + wk(c_2) & \text{if } c = c_1 \otimes c_2 \\ wk(c_1) + wk(c_2) & \text{if } c = c_1 \oplus c_2 \end{cases}$$
(37.5)

The depth is defined by the following equations:

$$dp(c) = \begin{cases} 0 & \text{if } c = \mathbf{0} \\ 1 & \text{if } c = \mathbf{1} \\ \max(dp(c_1), dp(c_2)) & \text{if } c = c_1 \otimes c_2 \\ dp(c_1) + dp(c_2) & \text{if } c = c_1 \oplus c_2 \end{cases}$$
(37.6)

Informally, the work of a cost graph determines the total number of computation steps represented by the cost graph, and thus corresponds to the sequential complexity of the computation. The depth of the cost graph determines the critical path length, the length of the longest dependency chain within the computation, which imposes a lower bound on the parallel complexity of a computation. The critical path length is a lower bound on the number of steps required to complete the computation.

In Chapter 7 we introduced cost dynamics to assign time complexity to computations. The proof of Theorem 7.7 shows that $e \downarrow ^k v$ iff $e \mapsto ^k v$. That is, the step complexity of an evaluation of e to a value v is just the number of transitions required to derive $e \mapsto^* v$. Here we use cost graphs as the measure of complexity, then relate these cost graphs to the structural dynamics given in Section 37.1.

The judgment $e \Downarrow^c v$, where e is a closed expression, v is a closed value, and c is a cost graph specifies the cost dynamics. By definition let is given by the following rule: $\frac{e_1 \Downarrow^{c_1} v_1 \quad e_2 \Downarrow^{c_2} v_2 \quad [v_1,v_2/x_1,x_2]e \Downarrow^c v}{\text{par}(e_1;e_2;x_1.x_2.e) \Downarrow^{(c_1\otimes c_2)\oplus 1\oplus c} v}$ specifies the cost dynamics. By definition we arrange that $e \downarrow 0$ e when e val. The cost assignment

$$\frac{e_1 \Downarrow^{c_1} v_1 \quad e_2 \Downarrow^{c_2} v_2 \quad [v_1, v_2/x_1, x_2]e \Downarrow^{c} v}{\operatorname{par}(e_1; e_2; x_1, x_2, e) \Downarrow^{(c_1 \otimes c_2) \oplus \mathbf{1} \oplus c} v}$$
(37.7)

The cost assignment specifies that, under ideal conditions, e_1 and e_2 are evaluated in parallel, and that their results are passed to e. The cost of fork and join is implicit in the parallel combination of costs, and assign unit cost to the substitution because we expect it to be implemented by a constant-time mechanism for updating an environment. The cost dynamics of other language constructs is specified in a similar way, using only sequential combination to isolate the source of parallelism to the let construct.

Two simple facts about the cost dynamics are important to keep in mind. First, the cost assignment does not influence the outcome.

Lemma 37.4. $e \Downarrow v \text{ iff } e \Downarrow^c v \text{ for some } c.$

Proof. From right to left, erase the cost assignments to construct an evaluation derivation. From left to right, decorate the evaluation derivations with costs as determined by the rules defining the cost dynamics.

Second, the cost of evaluating an expression is uniquely determined.

Lemma 37.5. If $e \Downarrow^c v$ and $e \Downarrow^{c'} v$, then c is c'.

Proof. By induction on the derivation of $e \downarrow^c v$.

The link between the cost dynamics and the structural dynamics is given by the following theorem, which states that the work cost is the sequential complexity, and the depth cost is the parallel complexity, of the computation.

Theorem 37.6. If $e \Downarrow^c v$, then $e \underset{seq}{\longmapsto^w} v$ and $e \underset{par}{\longmapsto^d} v$, where w = wk(c) and d = dp(c). Conversely, if $e \underset{seq}{\longmapsto^w} v$, then there exists c such that $e \Downarrow^c v$ with wk(c) = w, and if $e \underset{par}{\longmapsto^d} v'$, then there exists c' such that $e \Downarrow^{c'} v'$ with dp(c') = d.

Proof. The first part is proved by induction on the derivation of $e \Downarrow^c v$, the interesting case being rule (37.7). By induction we have $e_1 \mapsto_{\text{seq}}^{w_1} v_1$, $e_2 \mapsto_{\text{seq}}^{w_2} v_2$, and $[v_1, v_2/x_1, x_2]e \mapsto_{\text{seq}}^w v$, where $w_1 = wk(c_1)$, $w_2 = wk(c_2)$, and w = wk(c). By pasting together derivations we get a derivation

$$\begin{array}{c} \mathtt{par}(\,e_1;e_2;x_1.x_2.e\,) & \underset{\mathsf{seq}}{\longmapsto}^{w_1} \,\mathtt{par}(\,v_1;e_2;x_1.x_2.e\,) \\ & \underset{\mathsf{seq}}{\longmapsto}^{w_2} \,\mathtt{par}(\,v_1;v_2;x_1.x_2.e\,) \\ & \underset{\mathsf{seq}}{\longmapsto}^{w} \,[v_1,v_2/x_1,x_2]e \\ & \underset{\mathsf{seq}}{\longmapsto}^{w} \,v. \end{array}$$

Noting that $wk((c_1 \otimes c_2) \oplus \mathbf{1} \oplus c) = w_1 + w_2 + 1 + w$ completes the proof. Similarly, we have by induction that $e_1 \underset{\mathsf{par}}{\longmapsto} d_1 v_1$, $e_2 \underset{\mathsf{par}}{\longmapsto} d_2 v_2$, and $[v_1, v_2/x_1, x_2]e \underset{\mathsf{par}}{\longmapsto} d_1 v$, where $d_1 = dp(c_1)$, $d_2 = dp(c_2)$, and d = dp(c). Assume, without loss of generality, that $d_1 \leq d_2$ (otherwise simply swap the roles of d_1 and d_2 in what follows). We may paste together derivations as follows:

$$\begin{array}{c} \operatorname{par}(e_1;e_2;x_1.x_2.e) & \underset{\operatorname{par}}{\longmapsto}^{d_1} \operatorname{par}(v_1;e_2';x_1.x_2.e) \\ & \underset{\operatorname{par}}{\longmapsto}^{d_2-d_1} \operatorname{par}(v_1;v_2;x_1.x_2.e) \\ & \underset{\operatorname{par}}{\longmapsto} [v_1,v_2/x_1,x_2]e \\ & \underset{\operatorname{par}}{\longmapsto}^{d} v. \end{array}$$

Calculating $dp((c_1 \otimes c_2) \oplus \mathbf{1} \oplus c) = \max(d_1, d_2) + 1 + d$ completes the proof.

Turning to the second part, it suffices to show that if $e \mapsto_{\text{seq}} e'$ with $e' \Downarrow^{c'} v$, then $e \Downarrow^{c} v$ with wk(c) = wk(c') + 1, and if $e \mapsto_{\text{par}} e'$ with $e' \Downarrow^{c'} v$, then $e \Downarrow^{c} v$ with dp(c) = dp(c') + 1.

Suppose that $e = par(e_1; e_2; x_1.x_2.e_0)$ with e_1 val and e_2 val. Then $e \mapsto_{seq} e'$, where e = e' $[e_1, e_2/x_1, x_2]e_0$ and there exists c' such that $e' \Downarrow^{c'} v$. But then $e \Downarrow^c v$, where $c = (\mathbf{0} \otimes \mathbf{0}) \oplus \mathbf{1} \oplus c'$, and a simple calculation shows that wk(c) = wk(c') + 1, as required. Similarly, $e \bowtie_{\mathsf{par}} e'$ for e' as

above, and hence $e \Downarrow^c v$ for some e such that dp(e) = dp(e') + 1, as required. Suppose that $e = par(e_1; e_2; x_1.x_2.e_0)$ and $e \mapsto_{seq} e'$, where $e' = par(e'_1; e_2; x_1.x_2.e_0)$ and $e_1 \longmapsto e_1'$. From the assumption that $e' \Downarrow^{c'} v$, we have by inversion that $e'_1 \Downarrow^{c'_1} v_1$, $e_2 \Downarrow^{c'_2} v_2$, and $[v_1, v_2/x_1, x_2]e_0 \Downarrow^{c'_0} v$, with $c' = (c'_1 \otimes c'_2) \oplus \mathbf{1} \oplus c'_0$. By induction there exists c_1 such that $wk(c_1) = 1 + wk(c'_1)$ and $e_1 \Downarrow^{c_1} v_1$. But then $e \Downarrow^{c} v$, with $c = (c_1 \otimes c'_2) \oplus \mathbf{1} \oplus c'_0$. By a similar argument, suppose that $e = \text{par}(e_1; e_2; x_1.x_2.e_0)$ and $e \mapsto_{\text{par}} e'$, where $e' = \text{par}(e'_1; e'_2; x_1.x_2.e_0)$

and $e_1 \underset{\mathsf{par}}{\longmapsto} e_1'$, $e_2 \underset{\mathsf{par}}{\longmapsto} e_2'$, and $e' \Downarrow^{c'} v$. Then by inversion $e_1' \Downarrow^{c_1'} v_1$, $e_2' \Downarrow^{c_2'} v_2$, $[v_1, v_2/x_1, x_2]e_0 \Downarrow^{c_0} v$. But then $e \Downarrow^c v$, where $c = (c_1 \otimes c_2) \oplus \mathbf{1} \oplus c_0$, $e_1 \Downarrow^{c_1} v_1$ with $dp(c_1) = 1 + dp(c_1')$, $e_2 \Downarrow^{c_2} v_2$ with $dp(c_2) = 1 + dp(c_2')$, and $[v_1, v_2/x_1, x_2]e_0 \Downarrow^{c_0} v$. Calculating, we get

$$dp(c) = \max(dp(c'_1) + 1, dp(c'_2) + 1) + 1 + dp(c_0)$$

$$= \max(dp(c'_1), dp(c'_2)) + 1 + 1 + dp(c_0)$$

$$= dp((c'_1 \otimes c'_2) \oplus \mathbf{1} \oplus c_0) + 1$$

$$= dp(c') + 1,$$

which completes the proof.

Corollary 37.7. If $e \mapsto_{seq}^w v$ and $e \mapsto_{par}^d v'$, then v is v' and $e \downarrow ^c v$ for some c such that wk(c) = w and dp(c) = d.

Multiple Fork-Join 37.3

So far we have confined attention to binary fork/join parallelism induced by the parallel let construct. A generalizaton, called data parallelism, allows the simultaneous creation of any number of tasks that compute on the components of a data structure. The main example is a sequence of values of a specified type. The primitive operations on sequences are a natural source of unbounded parallelism. For example, we may consider a parallel map construct that applies a given function to every element of a sequence simultaneously, forming a sequence of the results.

We will consider here a simple language of sequence operations to illustrate the main ideas.

The expression $seq\{n\}(e_0,\ldots,e_{n-1})$ evaluates to a sequence whose length is n and whose elements are given by the expressions e_0,\ldots,e_{n-1} . The operation len(e) returns the number of elements in the sequence given by e. The operation $sub(e_1;e_2)$ retrieves the element of the sequence given by e_1 at the index given by e_2 . The tabulate operation, $tab(x.e_1;e_2)$, yields the sequence of length given by e_2 whose ith element is given by $[i/x]e_1$. The operation $map(x.e_1;e_2)$ computes the sequence whose ith element is given by $[e/x]e_1$, where e is the ith element of the sequence given by e_2 . The operation $cat(e_1;e_2)$ concatenates two sequences of the same type.

The statics of these operations is given by the following typing rules:

$$\frac{\Gamma \vdash e_0 : \tau \dots \Gamma \vdash e_{n-1} : \tau}{\Gamma \vdash \mathsf{seq}\{n\}(e_0, \dots, e_{n-1}) : \mathsf{seq}(\tau)}$$
(37.8a)

$$\frac{\Gamma \vdash e : \operatorname{seq}(\tau)}{\Gamma \vdash \operatorname{len}(e) : \operatorname{nat}} \tag{37.8b}$$

$$\frac{\Gamma \vdash e_1 : \operatorname{seq}(\tau) \quad \Gamma \vdash e_2 : \operatorname{nat}}{\Gamma \vdash \operatorname{sub}(e_1 : e_2) : \tau}$$
(37.8c)

$$\frac{\Gamma, x : \text{nat} \vdash e_1 : \tau \quad \Gamma \vdash e_2 : \text{nat}}{\Gamma \vdash \text{tab}(x.e_1; e_2) : \text{seq}(\tau)}$$
(37.8d)

$$\frac{\Gamma \vdash e_2 : \operatorname{seq}(\tau) \quad \Gamma, x : \tau \vdash e_1 : \tau'}{\Gamma \vdash \operatorname{map}(x.e_1; e_2) : \operatorname{seq}(\tau')}$$
(37.8e)

$$\frac{\Gamma \vdash e_1 : \operatorname{seq}(\tau) \quad \Gamma \vdash e_2 : \operatorname{seq}(\tau)}{\Gamma \vdash \operatorname{cat}(e_1; e_2) : \operatorname{seq}(\tau)}$$
(37.8f)

The cost dynamics of these constructs is defined by the following rules:

$$\frac{e_0 \Downarrow^{c_0} v_0 \dots e_{n-1} \Downarrow^{c_{n-1}} v_{n-1}}{\operatorname{seq}\{n\}(e_0, \dots, e_{n-1}) \Downarrow^{\bigotimes_{i=0}^{n-1} c_i} \operatorname{seq}\{n\}(v_0, \dots, v_{n-1})}$$
(37.9a)

$$\frac{e \Downarrow^{c} \operatorname{seq}\{n\}(v_{0},\ldots,v_{n-1})}{\operatorname{len}(e) \Downarrow^{c \oplus 1} \operatorname{num}[n]}$$
(37.9b)

$$\frac{e_1 \, \Downarrow^{c_1} \, \operatorname{seq}\{n\}(\, v_0, \dots, v_{n-1} \,) \quad e_2 \, \Downarrow^{c_2} \, \operatorname{num}[\, i \,] \quad (0 \leq i < n)}{\operatorname{sub}(\, e_1; e_2 \,) \, \Downarrow^{c_1 \oplus c_2 \oplus \mathbf{1}} v_i} \tag{37.9c}$$

$$\frac{e_2 \Downarrow^c \text{num}[n] \quad [\text{num}[0]/x] e_1 \Downarrow^{c_0} v_0 \quad \dots \quad [\text{num}[n-1]/x] e_1 \Downarrow^{c_{n-1}} v_{n-1}}{\text{tab}(x.e_1; e_2) \Downarrow^{c \oplus \bigotimes_{i=0}^{n-1} c_i} \text{seq}\{n\}(v_0, \dots, v_{n-1})}$$
(37.9d)

$$\frac{e_2 \Downarrow^c \operatorname{seq}\{n\}(v_0, \dots, v_{n-1})}{[v_0/x]e_1 \Downarrow^{c_0} v'_0 \dots [v_{n-1}/x]e_1 \Downarrow^{c_{n-1}} v'_{n-1}}$$

$$\frac{[v_0/x]e_1 \Downarrow^{c_0} v'_0 \dots [v_{n-1}/x]e_1 \Downarrow^{c_{n-1}} v'_{n-1}}{\operatorname{map}(x.e_1; e_2) \Downarrow^{c \oplus \bigotimes_{i=0}^{n-1} c_i} \operatorname{seq}\{n\}(v'_0, \dots, v'_{n-1})}$$
(37.9e)

$$\frac{e_1 \Downarrow^{c_1} \operatorname{seq}\{m\}(v_0, \dots, v_{m-1}) \quad e_2 \Downarrow^{c_2} \operatorname{seq}\{n\}(v'_0, \dots, v'_{n-1}) \quad p = m+n}{\operatorname{cat}(e_1; e_2) \Downarrow^{c_1 \oplus c_2 \oplus \bigotimes_{i=0}^{m+n} 1} \operatorname{seq}\{p\}(v_0, \dots, v_{m-1}, v'_0, \dots, v'_{n-1})}$$
(37.9f)

The cost dynamics for sequence operations is validated by introducing a sequential and parallel cost dynamics and extending the proof of Theorem 37.6 to cover this extension.

37.4 Bounded Implementations

Theorem 37.6 states that the cost dynamics accurately models the dynamics of the parallel 1et construct, whether executed sequentially or in parallel. The theorem validates the cost dynamics from the point of view of the dynamics of the language, and permits us to draw conclusions about the asymptotic complexity of a parallel program that abstracts away from the limitations imposed by a concrete implementation. Chief among these is the restriction to a fixed number, p > 0, of processors on which to schedule the workload. Besides limiting the available parallelism this also imposes some synchronization overhead that must be taken into account. A *bounded implementation* is one for which we may establish an asymptotic bound on the execution time once these overheads are taken into account.

A bounded implementation must take account of the limitations and capabilities of the hardware on which the program is run. Because we are only interested in asymptotic upper bounds, it is convenient to formulate an abstract machine model, and to show that the primitives of the language can be implemented on this model with guaranteed time (and space) bounds. One example of such a model is the *shared-memory multiprocessor*, or *SMP*, model. The basic assumption of the SMP model is that there are some fixed p>0 processors coordinated by an interconnect that permits constant-time access to any object in memory shared by all p processors. An SMP is assumed to provide a constant-time synchronization primitive with which to control simultaneous access to a memory cell. There are a variety of such primitives, any of which are enough to provide a parallel fetch-and-add instruction that allows each processor to get the current contents of a memory cell and update it by adding a fixed constant in a single atomic operation—the interconnect serializes any simultaneous accesses by more than one processor.

Building a bounded implementation of parallelism involves two major tasks. First, we must show that the primitives of the language can be implemented efficiently on the abstract machine model. Second, we must show how to schedule the workload across the processors to minimize execution time by maximizing parallelism. When working with a low-level machine model such as an SMP, both tasks involve a fair bit of technical detail to show how to use low-level machine instructions, including a synchronization primitive, to implement the language primitives and to schedule the workload. Collecting these together, we may then give an asymptotic bound on the time complexity of the implementation that relates the abstract cost of the computation to cost of implementing the workload on a *p*-way multiprocessor. The prototypical result of this kind is *Brent's Theorem*.

Theorem 37.8. If $e \downarrow^c v$ with wk(c) = w and dp(c) = d, then e can be evaluated on a p-processor SMP in time $O(\max(w/p,d))$.

 $^{^{1}}$ A slightly weaker assumption is that each access may require up to $\lg p$ time to account for the overhead of synchronization, but we shall neglect this refinement in the present, simplified account.

The theorem tells us that we can never execute a program in fewer steps than its depth d and that, at best, we can divide the work up evenly into w/p rounds of execution by the p processors. Note that if p=1 then the theorem establishes an upper bound of O(w) steps, the sequential complexity of the computation. Moreover, if the work is proportional to the depth, then we are unable to exploit parallelism, and the overall time is proportional to the work alone.

Theorem 37.8 motivates consideration of a useful figure of merit, the *parallelizability ratio*, which is the ratio w/d of work to depth. If $w/d \gg p$, then the program is *parallelizable*, because then $w/p \gg d$, and we may therefore reduce running time by using p processors at each step. If the parallelizability ratio is a constant, then d will dominate w/p, and we will have little opportunity to exploit parallelizable solution. The best we can say, on present knowledge, is that there are algorithms for some problems that have a high degree of parallelizability, and there are problems for which no such algorithm is known. It is a difficult problem in complexity theory to analyze which problems are parallelizable, and which are not.

Proving Brent's Theorem for an SMP would take us much too far afield for the present purposes. Instead we shall prove a Brent-type Theorem for an abstract machine, the **P** machine. The machine is unrealistic in that it is defined at a very high level of abstraction. But it is designed to match well the cost dynamics given earlier in this chapter. In particular, there are mechanisms that account for both sequential and parallel dependencies in a computation.

At the highest level, the state of the **P** machine consists of a global task graph whose structure corresponds to a "diagonal cut" through the cost graph of the overall computation. Nodes immediately above the cut are eligible to be executed, higher ancestors having already been completed, and whose immediate descendents are waiting for their ancestors to complete. Further descendents in the full task graph are tasks yet to be created, once the immediate descendents are finished. The **P** machine discards completed tasks, and future tasks beyond the immediate dependents are only created as execution proceeds. Thus it is only those nodes next to the cut line through the cost graph that are represented in the **P** machine state.

The *global state* of the **P** machine is a configuration of the form $v \Sigma \{ \mu \}$, where Σ is degenerated to just a finite set of (pairwise distinct) *task names* and μ is a finite mapping of the task names in Σ to *local states*, representing the state of an individual task. A *local state* is either a closed **PCF** expression, or one of two special *join points* that implement the sequential and parallel dependencies of a task on one or two ancestors, respectively.² Thus, when expanded out, a global state has the form

$$v a_1, \ldots, a_n \{ a_1 \hookrightarrow s_1 \otimes \ldots \otimes a_n \hookrightarrow s_n \},$$

where $n \ge 1$, and each s_i is a local state. The ordering of the tasks in a state, like the order of declarations in the signature, is not significant.

A **P** machine state transition has the form $v \Sigma \{ \mu \} \mapsto v \Sigma' \{ \mu' \}$. There are two forms of such transitions, the *global* and the *local*. A global step selects as many tasks as are available, up to a pre-specified parameter p > 0, which represents the number of processors available at each round. (Such a scheduler is *greedy* in the sense that it never fails to execute an available task, up

²The use of join points for each sequential dependency is profligate, but aligns the machine with the cost dynamics. Realistically, individual tasks manage sequential dependencies without synchronization, by using local control stacks as in Chapter 28.

to the specified limit for each round.) A task is *finished* if it consists of a closed **PCF** value, or is a join point whose dependents are not yet finished; otherwise a task is available, or ready. A ready task is always capable of taking a local step consisting of either a step of PCF, expressed in the setting of the **P** machine, or a *synchronization* step that manages the join-point logic. Because the P machine employs a greedy scheduler, it must complete execution in time proportional to $\max(w/p,d)$ steps by doing up to p steps of work at a time, insofar as it is possible within the limits of the depth of the computation. We thus get a Brent-type Theorem for the abstract machine that illustrates more sophisticated Brent-type Theorems for other models, such as the PRAM, that are used in the analysis of parallel algorithms.

The local transitions of the P machine corresponding to the steps of PCF itself are illustrated by the following example rules for application; the others follow a similar pattern.³

$$\frac{\neg(e_1 \text{ val})}{\nu \, a \, \{ \, a \hookrightarrow e_1(e_2) \, \} \, \underset{\text{loc}}{\longmapsto} \nu \, a \, a_1 \, \{ \, a \hookrightarrow \text{join}[\, a_1 \,](\, x_1.x_1(e_2) \,) \otimes a_1 \hookrightarrow e_1 \, \}}$$

$$\frac{e_1 \, \text{val} \, \neg(e_2 \, \text{val})}{\nu \, a \, \{ \, a \hookrightarrow e_1(e_2) \, \} \, \underset{\text{loc}}{\longmapsto} \nu \, a \, a_2 \, \{ \, a \hookrightarrow \text{join}[\, a_2 \,](\, x_2.e_1(\, x_2) \,) \otimes a_2 \hookrightarrow e_2 \, \}}$$
(37.10a)

$$\frac{e_1 \text{ val } \neg (e_2 \text{ val})}{\nu \, a \, \{ \, a \hookrightarrow e_1(e_2) \, \} \underset{\text{loc}}{\longmapsto} \nu \, a \, a_2 \, \{ \, a \hookrightarrow \text{join}[a_2](x_2.e_1(x_2)) \otimes a_2 \hookrightarrow e_2 \, \}}$$
(37.10b)

$$\frac{e_1 \text{ val}}{v \text{ a } a_1 \left\{ a \hookrightarrow \text{join}[a_1](x_1.x_1(e_2)) \otimes a_1 \hookrightarrow e_1 \right\} \underset{\text{loc}}{\longmapsto} v \text{ a } \left\{ a \hookrightarrow e_1(e_2) \right\}}$$
(37.10c)

$$\frac{e_1 \text{ val}}{v \text{ a } a_1 \left\{ a \hookrightarrow \text{join}[a_1](x_1.x_1(e_2)) \otimes a_1 \hookrightarrow e_1 \right\} \longmapsto_{\text{loc}} v \text{ a } \left\{ a \hookrightarrow e_1(e_2) \right\}}$$

$$\frac{e_1 \text{ val}}{v \text{ a } a_2 \left\{ a_1 \hookrightarrow \text{join}[a_2](x_2.e_1(x_2)) \otimes a_2 \hookrightarrow e_2 \right\} \longmapsto_{\text{loc}} v \text{ a } \left\{ a \hookrightarrow e_1(e_2) \right\}}$$

$$\frac{e_2 \text{ val}}{v \text{ a } \left\{ a \hookrightarrow (\lambda(x : \tau_2) e)(e_2) \right\} \longmapsto_{\text{loc}} v \text{ a } \left\{ a \hookrightarrow [e_2/x]e \right\}}$$
(37.10e)

$$\frac{e_2 \text{ val}}{v \, a \, \{ \, a \hookrightarrow (\lambda \, (x : \tau_2) \, e \,)(e_2) \,\} \underset{\mathsf{loc}}{\longmapsto} v \, a \, \{ \, a \hookrightarrow [e_2/x]e \,\}} \tag{37.10e}$$

Rules (37.10a) and (37.10b) create create tasks for the evaluation of the function and argument of an expression. Rules (37.10c) and (37.10d) propagate the result of evaluation of the function or argument of an application to the appropriate application expression. This rule mediates between the first two rules and Rule (37.10e), which effects a β -reduction in-place.

The local transitions of the P machine corresponding to binary fork and join are as follows:

$$\left\{
\begin{array}{c}
\nu \, a \, \left\{ \, a \hookrightarrow \operatorname{par}(\, e_1; e_2; x_1.x_2.e \,) \,\right\} \\
& \stackrel{\longmapsto}{\operatorname{loc}} \\
\nu \, a_1, a_2, a \, \left\{ \, a_1 \hookrightarrow e_1 \otimes a_2 \hookrightarrow e_2 \otimes a \hookrightarrow \operatorname{join}[\, a_1; a_2 \,](\, x_1; x_2.e \,) \,\right\}
\end{array}
\right\} (37.11a)$$

$$\frac{e_1 \text{ val} \quad e_2 \text{ val}}{\left\{\begin{array}{c}
\nu \, a_1, a_2, a \, \{ \, a_1 \hookrightarrow e_1 \otimes a_2 \hookrightarrow e_2 \otimes a \hookrightarrow \text{join}[\, a_1; a_2 \,](\, x_1; x_2.e \,) \, \} \\
& \qquad \qquad \mapsto \\
\text{loc} \\
\nu \, a \, \{ \, a \hookrightarrow [e_1, e_2/x_1, x_2]e \, \}
\end{array}\right\} \tag{37.11b}$$

³Here and elsewhere typing information is omitted from Σ , because it is not relevant to the dynamics.

35.5 Scheduling

Rule (37.11a) creates two parallel tasks on which the executing task depends. The expression $join[a_1; a_2](x_1; x_2.e)$ is blocked on tasks a_1 and a_2 , so that no local step applies to it. Rule (37.11b) synchronizes a task with the tasks on which it depends once their execution has completed; those tasks are no longer required, and are eliminated from the state.

Each global transition is the simultaneous execution of one step of computation on as many as $p \ge 1$ processors.

$$\nu \Sigma_{1} a_{1} \{ \mu_{1} \otimes a_{1} \hookrightarrow s_{1} \} \xrightarrow{\text{loc}} \nu \Sigma_{1}' a_{1} \{ \mu_{1}' \otimes a_{1} \hookrightarrow s_{1}' \} \\
\dots \\
\nu \Sigma_{n} a_{n} \{ \mu_{n} \otimes a_{n} \hookrightarrow s_{n} \} \xrightarrow{\text{loc}} \nu \Sigma_{n}' a_{n} \{ \mu_{n}' \otimes a_{n} \hookrightarrow s_{n}' \} \\
\frac{\nu \Sigma_{0} \Sigma_{1} a_{1} \dots \Sigma_{n} a_{n} \{ \mu_{0} \otimes \mu_{1} \otimes a_{1} \hookrightarrow s_{1} \otimes \dots \otimes \mu_{n} \otimes a_{n} \hookrightarrow s_{n} \} \xrightarrow{\text{loc}} \\
\nu \Sigma_{0} \Sigma_{1}' a_{1} \dots \Sigma_{n}' a_{n} \{ \mu_{0} \otimes \mu_{1}' \otimes a_{1} \hookrightarrow s_{1}' \otimes \dots \otimes \mu_{n}' \otimes a_{n} \hookrightarrow s_{n}' \} \\
\downarrow \nu \Sigma_{0} \Sigma_{1}' a_{1} \dots \Sigma_{n}' a_{n} \{ \mu_{0} \otimes \mu_{1}' \otimes a_{1} \hookrightarrow s_{1}' \otimes \dots \otimes \mu_{n}' \otimes a_{n} \hookrightarrow s_{n}' \} \\
\end{pmatrix}$$

At each global step some number $1 \le n \le p$ of ready tasks are scheduled for execution, where n is maximal among the number of ready tasks. Because no two distinct tasks may depend on the same task, we may partition the n tasks so that each scheduled task is grouped with the tasks on which it depends as necessary for any local join step. Any local fork step adds two fresh tasks to the state resulting from the global transition; any local join step eliminates two tasks whose execution has completed. A subtle point is that it is implicit in our name binding conventions that the names of any created tasks are *globally unique*, even though they are *locally created*. In implementation terms this requires a synchronization step among the processors to ensure that task names are not accidentally reused among the parallel tasks.

The proof of a Brent-type Theorem for the **P** machine is now obvious. We need only ensure that the parameter n of Rule (37.12) is chosen as large as possible at each step, limited only by the parameter p and the number of ready tasks. A scheduler with this property is greedy; it never allows a processor to go idle if work remains to be done. Consequently, if there are always p available tasks at each global step, then the evaluation will complete in w/p steps, where w is the work complexity of the program. If, at some stage, fewer than p tasks are available, then performance degrades according to the sequential dependencies among the sub-computations. In the limiting case the **P** machine must take at least d steps, where d is the depth of the computation.

37.5 Scheduling

The global transition relation of the **P** machine defined in Section 37.4 affords wide latitude in the choice of tasks that are advanced by taking a local transition. Doing so abstracts from implementation details that are irrelevant to the proof of the Brent-type Theorem given later in that section, the only requirement being that the number of tasks chosen be as large as possible up to the specified bound p, representing the number of available processors. When taking into account factors not considered here, it is necessary to specify the scheduling policy more precisely—for example, different scheduling policies may have asymptotically different space requirements. The overall

37.5 Scheduling 355

idea is to consider scheduling a computation on p processors as a p-way parallel traversal of its cost graph, visiting up to p nodes at a time in an order consistent with the dependency ordering. In this section we will consider one such traversal, p-way parallel depth-first-search, or p-DFS, which specializes to the familiar depth-first traversal in the case that p = 1.

Recall that the depth first-search of a directed graph maintain a stack of unvisited nodes, which is initialized with the start node. At each round, a node is popped from the stack and visited, and then its unvisited children are pushed on the stack (in reverse order in the case of ordered graphs), completing that round. The traversal terminates when the stack is empty. When viewed as a scheduling strategy, visiting a node of a cost graph consists of scheduling the work associated with that node on a processor. The job of such as scheduler is to do the work of the computation in depth-first order, visiting the children of a node from left to right, consistently with the sequential dynamics (which would, in particular, treat a parallel binding as two sequential bindings). Notice that because a cost graph is directed acyclic, there are no "back edges" arising from the traversal, and because it is series-parallel in structure, there are no "cross edges". Thus, all children of a node are unvisited, and no task is considered more than once.

Although evocative, viewing scheduling as graph traversal invites one to imagine that the cost graph is given explicitly as a data structure, which is not at all the case. Instead the graph is created dynamically as the sub-computations are executed. At each round the computation associated with a node may *complete* (when it has achieved its value), *continue* (when more work is yet to be done), or *fork* (when it generates parallel sub-computations with a specified join point). Once a computation has completed and its value has been passed to the associated join point, its node in the cost graph is discarded. Furthermore, the children of a node only come into existence as a result of its execution, according to whether it completes (no children), continues (one child), or forks (two children). Thus one may envision that the cost graph "exists" as a cut through the abstract cost graph representing pending tasks that have not yet been activated by the traversal.

A parallel depth-first search works much the same way, except that as many as p nodes are visited at each round, constrained only by the presence of unvisited (yet-to-be-scheduled) nodes. One might naively think that this simply means popping up to p nodes from the stack on each round, visiting them all simultaneously, and pushing their dependents on the stack in reverse order, just as for conventional depth-first search. But a moment's thought reveals that this is not correct. Because the cost graphs are ordered, the visited nodes form a sequence determined by the left-to-right ordering of the children of a node. If a node completes, it has no children and is removed from its position in the sequence in the next round. If a node continues, it has one child that occupies the same relative position as its parent in the next round. And if a node forks two children, they are inserted into the sequence after the predecessor, and immediately prior to that node, related to each other by the left-to-right ordering of the children. The task associated to the visited node itself becomes the join point of the immediately preceding pair of tasks, with which it will synchronize when they complete. Thus the visited sequence of $k \leq p$ nodes becomes, on the next round, anywhere from 0 (if all nodes completes) to $3 \times k$ nodes (if each node forks). These are placed into consideration, in the specified order, for the next round to ensure that they are processed in depth-first order. Importantly, the data structure maintaining the unvisited nodes of the graph is not a simple pushdown stack, because of the "in-place" replacement of each visited node by zero, one, or two nodes in between its predecessor and successor in the sequential ordering of the visited nodes.

356 37.6 Notes

Consider a variant of the **P** machine in which the order of the tasks is significant. A task is *finished* if it is a value, *blocked* if it is a join, and *ready* otherwise. Local transitions remain the same as in Section 37.4, bearing in mind that the ordering is significant. A global transition, however, consists of making a local transition on each of the first $k \le p$ ready tasks.⁴ After this selection the global state is depicted as follows:

$$\nu \Sigma_0 a_1 \Sigma_1 \dots a_k \Sigma_k \Sigma \{ \mu_0 \otimes a_1 \hookrightarrow e_1 \otimes \mu_1 \otimes \dots a_k \hookrightarrow e_k \otimes \mu \}$$

where each μ_i consists of finished or blocked tasks, and each e_i is ready. A schedule is greedy If k < p only when no task in μ is ready.

After a local transition is made on each of the *k* selected tasks, the resulting global state has the form

$$\nu \Sigma_0 \Sigma_1' a_1 \Sigma_1 \dots \Sigma_k' a_k \Sigma_k \Sigma \{ \mu_0 \otimes \mu_1' \otimes a_1 \hookrightarrow e_1' \otimes \mu_1 \otimes \dots \mu_k' \otimes a_k \hookrightarrow e_k' \otimes \mu \}$$

where each μ'_i represents the newly created task(s) of the local transition on task $a_i \hookrightarrow e_i$, and each e'_i is the expression resulting from the transition on that task. Next, all possible synchronizations are made by replacing each occurrence of an adjacent triple of the form

$$a_{i,1} \hookrightarrow e_1 \otimes a_{i,2} \hookrightarrow e_2 \otimes a_i \hookrightarrow \text{join}[a_{i,1}; a_{i,2}](x_1; x_2.e)$$

(with e_1 and e_2 finished) by the task $a_i \hookrightarrow [e_1, e_2/x_1, x_2]e$. Doing so propagates the values of tasks $a_{i,1}$ and $a_{i,2}$ to the join point, enabling the computation to continue. The two finished tasks are removed from the state, and the join point is no longer blocked.

37.6 Notes

Parallelism is a high-level programming concept that increases efficiency by carrying out multiple computations simultaneously when they are mutually independent. Parallelism does not change the meaning of a program, but only how fast it is executed. The cost dynamics specifies the number of steps required to execute a program sequentially and with maximal parallelism. A bounded implementation provides a bound on the number of steps when the number of processors is limited, limiting the degree of parallelism that can be realized. This formulation of parallelism was introduced by Blelloch (1990). The concept of a cost dynamics and the idea of a bounded implementation studied here are derived from Blelloch and Greiner (1995, 1996).

Exercises

- **37.1.** Consider extending **PPCF** with exceptions, as described in Chapter 29, under the assumption that τ_{exn} has at least two exception values. Give a sequential and a parallel structural dynamics to parallel let in such a way that determinacy continues to hold.
- **37.2**. Give a matching cost dynamics to **PPCF** extended with exceptions (descibed in Exercise **37.1**) by inductively defining the following two judgments:

⁴Thus the local transition given by Rule (37.11b) is never applicable; the dynamics of joins will be described shortly.

37.6 Notes 357

- (a) $e \downarrow^c v$, stating that e evaluates to value v with cost c;
- (b) $e \uparrow^c v$, stating that e raises the value v with cost c.

The analog of Theorem 37.6 remains valid for the dynamics. In particular, if $e \uparrow^c v$, then both $e \mapsto^w_{\text{seq}} \text{raise}(v)$, where w = wk(c), and $e \mapsto^d_{\text{par}} \text{raise}(v)$, where d = dp(c), and conversely.

- **37.3**. Extend the **P** machine to admit exceptions to match your solution to Exercise **37.2**. Argue that the revised machine supports a Brent-type validation of the cost dynamics.
- **37.4.** Another way to express the dynamics of **PPCF** enriched with exceptions is by rewriting $par(e_1; e_2; x_1.x_2.e)$ into another such parallel binding, $par(e_1'; e_2'; x_1'.x_2'.e')$, which implements the correct dynamics to ensure determinacy. *Hint*: Extend **XPCF** with sums (Chapter 11), using them to record the outcome of each parallel sub-computation (e_1') derived from e_1 , and e_2' derived from e_2), then check the outcomes (e_1') derived from e_2 in such a way to ensure determinacy.

358 37.6 Notes



Chapter 38

Futures and Speculations

A *future* is a computation that is performed before it is value is needed. Like a suspension, a future represents a value that is to be determined later. Unlike a suspension, a future is always evaluated, regardless of whether its value is required. In a sequential setting futures are of little interest; a future of type τ is just an expression of type τ . In a parallel setting, however, futures are of interest because they provide a means of initiating a parallel computation whose result is not needed until later, by which time it will have been completed.

The prototypical example of the use of futures is to implementing *pipelining*, a method for overlapping the stages of a multistage computation to the fullest extent possible. Pipelining minimizes the latency caused by one stage waiting for a previous stage to complete by allowing the two stages to execute in parallel until an explicit dependency arises. Ideally, the computation of the result of an earlier stage is finished by the time a later stage needs it. At worst the later stage is delayed until the earlier stage completes, incurring what is known as a *pipeline stall*.

A *speculation* is a delayed computation whose result might be needed for the overall computation to finish. The dynamics for speculations executes suspended computations in parallel with the main thread of computation, without regard to whether the value of the speculation is needed by the main thread. If the value of the speculation is needed, then such a dynamics pays off, but if not, the effort to compute it is wasted.

Futures are *work efficient* in that the overall work done by a computation involving futures is no more than the work done by a sequential execution. Speculations, in contrast, are *work inefficient* in that speculative execution might be in vain—the overall computation may involve more steps than the work needed to compute the result. For this reason speculation is a risky strategy for exploiting parallelism. It can make use available resources, but perhaps only at the expense of doing more work than necessary!

360 38.1 Futures

38.1 Futures

The syntax of futures is given by the following grammar:

The type τ fut is the type of futures of type τ . Futures are introduced by the expression $\mathtt{fut}(e)$, which schedules e for evaluation and returns a reference to it. Futures are eliminated by the expression $\mathtt{fsyn}(e)$, which synchronizes with the future referred to by e, returning its value. Indirect references to future values are represented by $\mathtt{fcell}[a]$, indicating a future value to be stored at a.

38.1.1 Statics

The statics of futures is given by the following rules:

$$\frac{\Gamma \vdash e : \tau}{\Gamma \vdash \text{fut}(e) : \text{fut}(\tau)}$$
(38.1a)

$$\frac{\Gamma \vdash e : \mathsf{fut}(\tau)}{\Gamma \vdash \mathsf{fsyn}(e) : \tau} \tag{38.1b}$$

These rules are unsurprising, because futures add no new capabilities to the language beyond providing an opportunity for parallel evaluation.

38.1.2 Sequential Dynamics

The sequential dynamics of futures is easily defined. Futures are evaluated eagerly; synchronization returns the value of the future.

$$\frac{e \text{ val}}{\text{fut}(e) \text{ val}} \tag{38.2a}$$

$$\frac{e \longmapsto e'}{\operatorname{fut}(e) \longmapsto \operatorname{fut}(e')} \tag{38.2b}$$

$$\frac{e \longmapsto e'}{\mathtt{fsyn}(e) \longmapsto \mathtt{fsyn}(e')} \tag{38.2c}$$

$$\frac{e \, \mathsf{val}}{\mathsf{fsyn}(\,\mathsf{fut}(\,e\,)\,) \longmapsto e} \tag{38.2d}$$

Under a sequential dynamics futures have little purpose: they introduce a pointless level of indirection.

38.2 Speculations 361

38.2 Speculations

The syntax of (non-recursive) speculations is given by the following grammar:¹

The type τ spec is the type of speculations of type τ . The introduction form $\operatorname{spec}(e)$ creates a computation that can be speculatively evaluated, and the elimination form $\operatorname{ssyn}(e)$ synchronizes with a speculation. A reference to the result of a speculative computation stored at a is written $\operatorname{scell}[a]$.

38.2.1 Statics

The statics of speculations is given by the following rules:

$$\frac{\Gamma \vdash e : \tau}{\Gamma \vdash \operatorname{spec}(e) : \operatorname{spec}(\tau)}$$
 (38.3a)

$$\frac{\Gamma \vdash e : \operatorname{spec}(\tau)}{\Gamma \vdash \operatorname{ssyn}(e) : \tau}$$
(38.3b)

Thus, the statics for speculations as given by rules (38.3) is equivalent to the statics for futures given by rules (38.1).

38.2.2 Sequential Dynamics

The definition of the sequential dynamics of speculations is like that of futures, except that speculations are values.

$$\frac{}{\operatorname{spec}(e)\operatorname{val}}\tag{38.4a}$$

$$\frac{e \longmapsto e'}{\operatorname{ssyn}(e) \longmapsto \operatorname{ssyn}(e')} \tag{38.4b}$$

$$\frac{}{\operatorname{ssyn}(\operatorname{spec}(e)) \longmapsto e} \tag{38.4c}$$

Under a sequential dynamics speculations are simply a re-formulation of suspensions.

 $^{^1\}mbox{We}$ confine ourselves to the non-recursive case to ease the comparison with futures.

38.3 Parallel Dynamics

Futures are only interesting insofar as they admit a parallel dynamics that allows the computation of the future to go ahead concurrently with some other computation. In this section we give a parallel dynamics of futures and speculation in which the creation, execution, and synchronization of tasks is made explicit. The parallel dynamics of futures and speculations is *identical*, except for the termination condition. Whereas futures require that all tasks are completed before termination, speculations may be abandoned before they are completed. For the sake of concision we will give the parallel dynamics of futures, remarking only where alterations are made for the parallel dynamics of speculations.

The parallel dynamics of futures relies on a modest extension to the language given in Section 38.1 to introduce *names* for tasks. Let Σ be a finite mapping assigning types to names. As mentioned earlier, the expression fcell[a] is a value referring to the outcome of task a. The statics of this expression is given by the following rule:²

$$\frac{}{\Gamma \vdash_{\Sigma, a \sim \tau} \text{fcell}[a] : \text{fut}(\tau)} \tag{38.5}$$

Rules (38.1) carry over in the obvious way with Σ recording the types of the task names.

States of the parallel dynamics have the form $v \Sigma \{e \mid \mu\}$, where e is the *focus* of evaluation, and μ records the active parallel futures (or speculations). Formally, μ is a finite mapping assigning expressions to the task names declared in Σ . A state is well-formed according to the following rule:

$$\frac{\vdash_{\Sigma} e : \tau \quad (\forall a \in dom(\Sigma)) \vdash_{\Sigma} \mu(a) : \Sigma(a)}{\nu \Sigma \{e \mid |\mu\} \text{ ok}}$$
(38.6)

As discussed in Chapter 35 this rule admits self-referential and mutually referential futures. A more refined condition could as well be given that avoids circularities; we leave this as an exercise for the reader.

The parallel dynamics is divided into two phases, the *local* phase, which defines the basic steps of evaluation of an expression, and the *global* phase, which executes all possible local steps in parallel. The local dynamics of futures is defined by the following rules:³

$$\frac{1}{\text{fcell}[a] \text{val}_{\Sigma,a\sim\tau}} \tag{38.7a}$$

$$\overline{\nu \, \Sigma \, \{ \, \text{fut}(e) \parallel \mu \, \} \,} \underset{\text{loc}}{\longmapsto} \nu \, \Sigma, a \sim \tau \, \{ \, \text{fcell}[a] \parallel \mu \otimes a \hookrightarrow e \, \}$$
 (38.7b)

$$\frac{\nu \Sigma \{e \parallel \mu\} \underset{\mathsf{loc}}{\longmapsto} \nu \Sigma' \{e' \parallel \mu'\}}{\nu \Sigma \{\mathsf{fsyn}(e) \parallel \mu\} \underset{\mathsf{loc}}{\longmapsto} \nu \Sigma' \{\mathsf{fsyn}(e') \parallel \mu'\}}$$
(38.7c)

²A similar rule applies to scell[a] in the case of speculations.

³These rules are augmented by a reformulation of the dynamics of the other constructs of the language phrased in terms of the present notion of state.

$$\frac{e' \operatorname{val}_{\Sigma, a \sim \tau}}{\left\{ \begin{array}{c} \nu \, \Sigma, a \sim \tau \, \{ \, \operatorname{fsyn}(\, \operatorname{fcell}[\, a \,] \,) \parallel \mu \otimes a \hookrightarrow e' \, \} \\ & \underset{\operatorname{loc}}{\longmapsto} \\ \nu \, \Sigma, a \sim \tau \, \{ \, e' \parallel \mu \otimes a \hookrightarrow e' \, \} \end{array} \right\}} \tag{38.7d}$$

Rule (38.7b) activates a future named a executing the expression e and returns a reference to it. Rule (38.7d) synchronizes with a future whose value has been determined. Note that a local transition always has the form

$$\nu \Sigma \{e \parallel \mu\} \underset{\mathsf{loc}}{\longmapsto} \nu \Sigma \Sigma' \{e' \parallel \mu \otimes \mu'\}$$

where Σ' is either empty or declares the type of a single symbol, and μ' is either empty or of the form $a \hookrightarrow e'$ for some expression e'.

A global step of the parallel dynamics consists of at most one local step for the focal expression and one local step for each of up to p futures, where p > 0 is a fixed parameter representing the number of processors.

$$\mu = \mu_{0} \otimes a_{1} \hookrightarrow e_{1} \otimes \ldots \otimes a_{n} \hookrightarrow e_{n}$$

$$\mu'' = \mu_{0} \otimes a_{1} \hookrightarrow e'_{1} \otimes \ldots \otimes a_{n} \hookrightarrow e'_{n}$$

$$\nu \Sigma \{e \parallel \mu\} \xrightarrow{\bowtie} {}^{?} \nu \Sigma \Sigma' \{e' \parallel \mu \otimes \mu' \}$$

$$\frac{(\forall 1 \leq i \leq n \leq p) \quad \nu \Sigma \{e_{i} \parallel \mu\} \xrightarrow{\bowtie} \nu \Sigma \Sigma'_{i} \{e'_{i} \parallel \mu \otimes \mu'_{i} \}}{}$$

$$\frac{\nu \Sigma \{e \parallel \mu\}}{}
\xrightarrow{\text{glo}}$$

$$\nu \Sigma \Sigma' \Sigma'_{1} \ldots \Sigma'_{n} \{e' \parallel \mu'' \otimes \mu' \otimes \mu'_{1} \otimes \ldots \otimes \mu'_{n} \}$$
(38.8a)

Rule (38.8a) allows the focus expression to take either zero or one steps because it might be blocked awaiting the completion of evaluation of a parallel future (or synchronizing with a speculation). The futures allocated by the local steps of execution are consolidated in the result of the global step. We assume without loss of generality that the names of the new futures in each local step are pairwise disjoint so that the combination makes sense. In implementation terms satisfying this disjointness assumption means that the processors must synchronize their access to memory.

The initial state of a computation, for futures or speculations, is defined by the rule

$$\frac{}{\nu \oslash \{e \parallel \oslash\} \text{ initial}} \tag{38.9}$$

For futures a state is final only if the focus and all parallel futures have completed evaluation:

$$\frac{e \operatorname{val}_{\Sigma} \quad \mu \operatorname{val}_{\Sigma}}{\nu \Sigma \left\{ e \parallel \mu \right\} \operatorname{final}} \tag{38.10a}$$

$$\frac{(\forall a \in dom(\Sigma)) \ \mu(a) \ \mathsf{val}_{\Sigma}}{\mu \ \mathsf{val}_{\Sigma}} \tag{38.10b}$$

For speculations a state is final only if the focus is a value, regardless of whether any other speculations have completed:

$$\frac{e \operatorname{val}_{\Sigma}}{\nu \Sigma \left\{e \parallel \mu\right\} \operatorname{final}} \tag{38.11}$$

All futures must terminate to ensure that the work performed in parallel matches that performed sequentially; no future is created whose value is not needed according to the sequential dynamics. In contrast, speculations can be abandoned when their values are not needed.

38.4 Pipelining With Futures

Pipelining is an interesting example of the use of parallel futures. Consider a situation in which a *producer* builds a list whose elements represent units of work, and a *consumer* traverses the work list and acts on each element of that list. The elements of the work list can be thought of as "instructions" to the consumer, which maps a function over that list to carry out those instructions. An obvious sequential implementation first builds the work list, then traverses it to perform the work indicated by the list. This strategy works well provided that the elements of the list can be produced quickly, but if each element needs a lot of computation, it would be preferable to overlap production of the next list element with execution of the previous unit of work, which can be programmed using futures.

Let flist be the recursive type rec t is unit + (nat \times t fut), whose elements are nil, defined to be fold(1 \cdot $\langle \rangle$), and cons(e_1,e_2), defined to be fold(r \cdot $\langle e_1$, fut(e_2) \rangle). The producer is a recursive function that generates a value of type flist:

```
fix produce : (nat \rightarrow nat opt) \rightarrow nat \rightarrow flist is \lambda f. \lambda i. case f(i) {
    null \hookrightarrow nil
    just x \hookrightarrow cons(x, fut (produce f (i+1)))
}
```

On each iteration the producer generates a parallel future to produce the tail. The future continues to execute after the producer returns so that its evaluation overlaps with subsequent computation.

The consumer folds an operation over the work list as follows:

```
fix consume : ((\text{nat} \times \text{nat}) \rightarrow \text{nat} \rightarrow \text{flist} \rightarrow \text{nat} \text{ is } \lambda \text{ g. } \lambda \text{ a. } \lambda \text{ xs.}
case xs {
    nil \hookrightarrow a
    | cons (x, xs) \hookrightarrow consume g (g (x, a)) (fsyn xs)
}
```

38.5 Notes 365

The consumer synchronizes with the tail of the work list just at the point where it makes a recursive call and hence needs the head element of the tail to continue processing. At this point the consumer will block, if necessary, to await computation of the tail before continuing the recursion.

Speculations arise naturally in lazy languages. But although they provide opportunities for parallelism, they are not, in general, *work efficient*: a speculation might be evaluated even though its value is never needed. An alternative is to combine suspensions (see Chapter 36) with futures so that the programmer may specify which suspensions ought to be evaluated in parallel. The notion of a *spark* is designed to achieve this. A spark evaluates a computation in parallel only for its effect on suspensions that are likely to be needed later. Specifically, we may define

$$\operatorname{spark}(e_1; e_2) \triangleq \operatorname{letfut_beforce}(e_1) \operatorname{in} e_2$$
,

where $e_1 : \tau_1$ susp and $e_2 : \tau_2$.⁴ The expression force(e_1) is evaluated in parallel, forcing the evaluation of e_1 , in hopes that it will have completed evaluation before its value is needed by e_2 .

As an example, consider the type strm of streams of numbers defined by the recursive type $\operatorname{rec} t$ is $(\operatorname{unit} + (\operatorname{nat} \times t))$ susp. Elements of this type are suspended computations that, when forced, either signals the end of stream, or produces a number and another such stream. Suppose that s is such a stream, and assume that we know, for reasons of its construction, that it is finite. We wish to compute $\operatorname{map}(f)(s)$ for some function f, and to overlap this computation with the production of the stream elements. We will make use of a function $\operatorname{mapforce}$ that forces successive elements of the input stream, but yields no useful output. The computation

letfut_bemap(force)(
$$s$$
)inmap(f)(s)

forces the elements of the stream in parallel with the computation of map(f)(s), with the intention that all suspensions in s are forced before their values are needed by the main computation.

38.5 Notes

Futures were introduced by Friedman and Wise (1976), and featured in the MultiLisp language (Halstead, 1985) for parallel programming. A similar concept is proposed by Arvind et al. (1986) under the name "I-structures." The formulation given here is derived from Greiner and Blelloch (1999). Sparks were introduced by Trinder et al. (1998).

Exercises

- **38.1**. Use futures to define letfut x be e_1 in e_2 , a parallel let in which e_2 is evaluated in parallel with e_1 up to the point that e_2 needs the value of x.
- **38.2.** Use futures to encode binary nested parallelism by giving a definition of $par(e_1; e_2; x_1.x_2.e)$. *Hint*: Only one future is needed if you are careful.

⁴The expression evaluates e_1 simultaneously with e_2 , up to the point that the value of x is needed. Its definition in terms of futures is the subject of Exercise 38.1.

366 38.5 Notes



Part XVI Concurrency and Distribution



Chapter 40

Concurrent Algol

In this chapter we integrate concurrency into the framework of Modernized Algol described in Chapter 34. The resulting language, called Concurrent Algol, or **CA**, illustrates the integration of the mechanisms of the process calculus described in Chapter 39 into a practical programming language. To avoid distracting complications, we drop assignables from Modernized Algol entirely. (There is no loss of generality, however, because free assignables are definable in Concurrent Algol using processes as cells.)

The process calculus described in Chapter 39 is intended as a self-standing model of concurrent computation. When viewed in the context of a programming language, however, it is possible to streamline the machinery to take full advantage of types that are in any case required for other purposes. In particular the concept of a *channel*, which features prominently in Chapter 39, is identified with the concept of a *dynamic class* as described in Chapter 33. More precisely, we take *broadcast communication* of dynamically classified values as the basic synchronization mechanism of the language. Being dynamically classified, messages consist of a *payload* tagged with a *class*, or *channel*. The type of the channel determines the type of the payload. Importantly, only those processes that have access to the channel may decode the message; all others must treat it as inscrutable data that can be passed around but not examined. In this way we can model not only the mechanisms described in Chapter 39, but also formulate an abstract account of encryption and decryption in a network using the methods described in Chapter 39.

Concurrent Algol features a modal separation between commands and expressions like in Modernized Algol. It is also possible to combine these two levels (so as to allow benign concurrency effects), but we do not develop this approach in detail here.

40.1 Concurrent Algol

The syntax of **CA** is obtained by removing assignables from **MA**, and adding a syntactic level of *processes* to represent the global state of a program:

```
commands
Тур
                 cmd(\tau)
                                     	au cmd
Exp
                 cmd(m)
                                                      command
           ::=
                                     {\tt cmd}\, m
Cmd m
                 rete
                                     rete
                                                      return
                  bnd(e; x.m)
                                                      sequence
                                     bnd x \leftarrow e; m
Proc p
                 stop
                                                      idle
                                                      atomic
                 run(m)
                                     run(m)
                  conc(p_1; p_2)
                                                      concurrent
                                     p_1 \otimes p_2
                 newch\{\tau\}(a.p) \quad \nu \, a \sim \tau.p
                                                      new channel
```

The process run(m) is an atomic process executing the command m. The other forms of process are adapted from Chapter 39. If Σ has the form $a_1 \sim \tau_1, \ldots, a_n \sim \tau_n$, then we sometimes write $\nu \Sigma \{p\}$ for the iterated form $\nu a_1 \sim \tau_1, \ldots, \nu a_n \sim \tau_n, p$.

The statics of **CA** is given by these judgments:

 $\Gamma \vdash_{\Sigma} e : \tau$ expression typing $\Gamma \vdash_{\Sigma} m \stackrel{.}{\sim} \tau$ command typing $\Gamma \vdash_{\Sigma} p$ proc process formation $\Gamma \vdash_{\Sigma} \alpha$ action action formation

The expression and command typing judgments are essentially those of **MA**, augmented with the constructs described below.

Process formation is defined by the following rules:

$$-\frac{}{\vdash_{\Sigma} \mathbf{1} \text{ proc}}$$
 (40.1a)

$$\frac{\vdash_{\Sigma} m \stackrel{.}{\sim} \tau}{\vdash_{\Sigma} \text{run}(m) \text{ proc}} \tag{40.1b}$$

$$\frac{\vdash_{\Sigma} p_1 \text{ proc} \vdash_{\Sigma} p_2 \text{ proc}}{\vdash_{\Sigma} p_1 \otimes p_2 \text{ proc}}$$
(40.1c)

$$\frac{\vdash_{\Sigma,a\sim\tau} p \text{ proc}}{\vdash_{\Sigma} \nu \, a\sim\tau. p \text{ proc}}$$
(40.1d)

Processes are identified up to structural congruence, as described in Chapter 39.

Action formation is defined by the following rules:

$$\frac{}{\vdash_{\Sigma} \varepsilon \text{ action}}$$
 (40.2a)

$$\frac{\vdash_{\Sigma} e : \mathsf{clsfd} \quad e \, \mathsf{val}_{\Sigma}}{\vdash_{\Sigma} e \, ! \, \mathsf{action}} \tag{40.2b}$$

$$\frac{\vdash_{\Sigma} e : \mathsf{clsfd} \quad e \, \mathsf{val}_{\Sigma}}{\vdash_{\Sigma} e ? \, \mathsf{action}} \tag{40.2c}$$

Messages are values of the type clsfd defined in Chapter 33.

The dynamics of **CA** is defined by transitions between processes, which represent the state of the computation. More precisely, the judgment $p \mapsto p'$ states that the process p evolves in one step to the process p' while undertaking action α .

$$\frac{m \stackrel{\alpha}{\Longrightarrow} \nu \Sigma' \{ m' \otimes p \}}{\operatorname{run}(m) \stackrel{\alpha}{\longmapsto} \nu \Sigma' \{ \operatorname{run}(m') \otimes p \}}$$
(40.3a)

$$\frac{e \operatorname{val}_{\Sigma}}{\operatorname{run}(\operatorname{ret} e) \stackrel{e}{\underset{\Sigma}{\vdash}} 1} \tag{40.3b}$$

$$\frac{p_1 \stackrel{\alpha}{\stackrel{\leftarrow}{\Sigma}} p_1'}{p_1 \otimes p_2 \stackrel{\alpha}{\stackrel{\leftarrow}{\Sigma}} p_1' \otimes p_2}$$
(40.3c)

$$\frac{p_1 \stackrel{\alpha}{\rightarrowtail} p'_1}{p_1 \otimes p_2 \stackrel{\varepsilon}{\rightarrowtail} p'_1 \otimes p'_2} \frac{p'_2}{p'_1 \otimes p'_2}$$

$$(40.3d)$$

$$\frac{p \underset{\Sigma, a \sim \tau}{\longleftrightarrow} p' \vdash_{\Sigma} \alpha \text{ action}}{\nu \, a \sim \tau. p \underset{\Sigma}{\longleftrightarrow} \nu \, a \sim \tau. p'}$$

$$(40.3e)$$

Rule (40.3a) states that a step of execution of the atomic process run(m) consists of a step of execution of the command m, which may allocate some set Σ' of symbols or create a concurrent process p. This rule implements scope extrusion for classes (channels) by expanding the scope of the channel declaration to the context in which the command m occurs. Rule (40.3b) states that a completed command evolves to the inert (stopped) process; processes are executed solely for their effect, and not for their value.

Executing a command in **CA** may, in addition to evolving to another command, allocate a new channel or may spawn a new process. More precisely, the judgment¹

$$m \stackrel{\alpha}{\underset{\Sigma}{\Longrightarrow}} \nu \Sigma' \{ m' \otimes p' \}$$

¹The right-hand side of this judgment is a triple consisting of Σ' , m', and p', not a process expression comprising these parts.

states that the command m transitions to the command m' while creating new channels Σ' and new processes p'. The action α specifies the interactions of which m is capable when executed. As a notational convenience we drop mention of the new channels or processes when either are trivial.

The following rules define the execution of the basic forms of command inherited from MA:

$$\frac{e \xrightarrow{\Sigma} e'}{\operatorname{ret} e \stackrel{\varepsilon}{\Rightarrow} \operatorname{ret} e'} \tag{40.4a}$$

$$\frac{m_1 \stackrel{\alpha}{\Longrightarrow} \nu \Sigma' \{ m_1' \otimes p' \}}{\operatorname{bnd} x \leftarrow \operatorname{cmd} m_1 ; m_2 \stackrel{\alpha}{\Longrightarrow} \nu \Sigma' \{ \operatorname{bnd} x \leftarrow \operatorname{cmd} m_1' ; m_2 \otimes p' \}}$$
(40.4b)

$$\frac{e \operatorname{val}_{\Sigma}}{\operatorname{bnd} x \leftarrow \operatorname{cmd} (\operatorname{ret} e); m_2 \stackrel{\varepsilon}{\xrightarrow{\Sigma}} [e/x] m_2}$$
(40.4c)

$$\frac{e_1 \underset{\Sigma}{\longmapsto} e_1'}{\text{bnd } x \leftarrow e_1 ; m_2 \stackrel{\varepsilon}{\Longrightarrow} \text{bnd } x \leftarrow e_1' ; m_2}$$

$$(40.4d)$$

These rules are supplemented by rules governing communication and synchronization among processes in the next two sections.

40.2 Broadcast Communication

In this section we consider a very general form of process synchronization called *broadcast*. Processes emit and accept messages of type clsfd, the type of dynamically classified values considered in Chapter 33. A message consists of a *channel*, which is its class, and a *payload*, which is a value of the type associated with the channel (class). Recipients may pattern match against a message to determine whether it is of a given class, and, if so, recover the associated payload. No process that lacks access to the class of a message may recover the payload of that message. (See Section 33.4.1 for a discussion of how to enforce confidentiality and integrity restrictions using dynamic classification).

The syntax of the commands pertinent to broadcast communication is given by the following grammar:

$$\begin{array}{ccccc} \mathsf{Cmd} & m & \coloneqq & \mathsf{spawn}(e) & \mathsf{spawn}(e) & \mathsf{spawn} \\ & & \mathsf{emit}(e) & \mathsf{emit}(e) & \mathsf{emit} \, \mathsf{message} \\ & & \mathsf{acc} & \mathsf{acc} & \mathsf{accept} \, \mathsf{message} \\ & & & \mathsf{newch}\{\tau\} & \mathsf{newch} & \mathsf{new} \, \mathsf{channel} \end{array}$$

The command spawn(e) spawns a process that executes the encapsulated command given by e. The commands emit(e) and acc emit and accept messages, which are classified values whose

class is the channel on which the message is sent. The command $newch\{\tau\}$ returns a reference to a fresh class carrying values of type τ .

The statics of broadcast communication is given by the following rules:

$$\frac{\Gamma \vdash_{\Sigma} e : \operatorname{cmd}(\operatorname{unit})}{\Gamma \vdash_{\Sigma} \operatorname{spawn}(e) \stackrel{.}{\sim} \operatorname{unit}}$$

$$(40.5a)$$

$$\frac{\Gamma \vdash_{\Sigma} e : \mathtt{clsfd}}{\Gamma \vdash_{\Sigma} \mathtt{emit}(e) \stackrel{.}{\sim} \mathtt{unit}} \tag{40.5b}$$

$$\Gamma \vdash_{\Sigma} \operatorname{acc} \stackrel{.}{\sim} \operatorname{clsfd}$$
 (40.5c)

$$\frac{}{\Gamma \vdash_{\Sigma} \operatorname{newch}\{\tau\} \div \operatorname{cls}(\tau)} \tag{40.5d}$$

Execution of these commands is defined as follows:

$$\overline{\operatorname{spawn}(\operatorname{cmd}(m)) \stackrel{\mathcal{E}}{\Rightarrow} \operatorname{ret} \langle \rangle \otimes \operatorname{run}(m)}$$
 (40.6a)

$$\frac{e \underset{\Sigma}{\longmapsto} e'}{\operatorname{spawn}(e) \overset{\varepsilon}{\underset{\Sigma}{\Longrightarrow}} \operatorname{spawn}(e')} \tag{40.6b}$$

$$\frac{e \operatorname{val}_{\Sigma}}{\operatorname{emit}(e) \xrightarrow{e!} \operatorname{ret} \langle \rangle}$$
(40.6c)

$$\frac{e \underset{\Sigma}{\longleftarrow} e'}{\frac{e}{\Sigma} \text{ emit}(e')} \tag{40.6d}$$

$$\frac{e \operatorname{val}_{\Sigma}}{\operatorname{acc} \xrightarrow{e?}_{\Sigma} \operatorname{ret} e}$$
 (40.6e)

$$\overline{\operatorname{newch}\{\tau\}} \underset{\Sigma}{\overset{\varepsilon}{\Longrightarrow}} \nu \, a \sim \tau \, \{ \, \operatorname{ret} \, (\,\&\, a\,) \, \} \tag{40.6f}$$

Rule (40.6c) specifies that emit(e) has the effect of emitting the message e. Correspondingly, rule (40.6e) specifies that acc may accept (any) message that is being sent.

As usual, the preservation theorem for **CA** ensures that well-typed programs remain well-typed during execution. The proof of preservation requires a lemma about command execution.

Lemma 40.1. If
$$m \stackrel{\alpha}{\Longrightarrow} \nu \Sigma' \{ m' \otimes p' \}$$
, $\vdash_{\Sigma} m \stackrel{.}{\leadsto} \tau$, then $\vdash_{\Sigma} \alpha$ action, $\vdash_{\Sigma\Sigma'} m' \stackrel{.}{\leadsto} \tau$, and $\vdash_{\Sigma\Sigma'} p'$ proc.

Proof. By induction on rules (40.4).

With this in hand the proof of preservation goes along familiar lines.

Theorem 40.2 (Preservation). *If* $\vdash_{\Sigma} p$ *proc and* $p \vdash_{\Sigma} p'$, then $\vdash_{\Sigma} p'$ *proc.*

Proof. By induction on transition, appealing to Lemma 40.1 for the crucial steps.

Typing does not, however, guarantee progress with respect to unlabeled transition, for the simple reason that there may be no other process with which to communicate. By extending progress to labeled transitions we may state that this is the *only* way for process exceution to get stuck. But some care must be taken to account for allocating new channels.

Theorem 40.3 (Progress). If $\vdash_{\Sigma} p$ proc, then either $p \equiv 1$, or $p \equiv \nu \Sigma' \{p'\}$ such that $p' \stackrel{\alpha}{\vdash_{\Sigma \Sigma'}} p''$ for some $\vdash_{\Sigma \Sigma'} p''$ and some $\vdash_{\Sigma \Sigma'} \alpha$ action.

Proof. By induction on rules (40.1) and (40.5).

The progress theorem says that no process can get stuck for any reason other than ithe nability to communicate with another process. For example, a process that receives on a channel for which there is no sender is "stuck", but this does not violate Theorem 40.3.

40.3 Selective Communication

Broadcast communication provides no means of restricting acceptance to messages of a particular class (that is, of messages on a particular channel). Using broadcast communication we may restrict attention to a particular channel a of type τ by running the following command:

```
fix loop : \tau cmd is cmd \{x \leftarrow acc; match x as a \cdot y \hookrightarrow ret y ow \hookrightarrow emit(x); do loop\}
```

This command is always capable of receiving a broadcast message. When one arrives, it is examined to see whether it is classified by *a*. If so, the underlying classified value is returned; otherwise the message is re-broadcast so that another process may consider it. *Polling* consists of repeatedly executing the above command until a message of channel *a* is successfully accepted, if ever it is.

Polling is evidently impractical in most situations. An alternative is to change the language to allow for *selective communication*. Rather than accept any broadcast message, we may confine attention to messages sent only on certain channels. The type $event(\tau)$ of *events* consists of a finite choice of accepts, all of whose payloads are of type τ .

```
:= event(\tau)
                                	au event
                                               events
               rcv[a]
                                               selective read
                                 ? a
                never\{\tau\}
                                never
                                               null
                or(e_1; e_2)
                                e_1 or e_2
                                               choice
                wrap(e_1; x.e_2) e_1 as x in e_2
                                               post-composition
Cmd m ::= sync(e)
                                sync(e)
                                               synchronize
```

Events in **CA** are similar to those of the asynchronous process calculus described in Chapter 39. The chief difference is that post-composition is considered as a general operation on events, instead of one tied to the receive event itself.

The statics of event expressions is given by the following rules:

$$\frac{\Sigma \vdash a \sim \tau}{\Gamma \vdash_{\Sigma} \text{rcv}[a] : \text{event}(\tau)}$$
(40.7a)

$$\frac{\Gamma \vdash_{\Sigma} \text{never}\{\tau\} : \text{event}(\tau)}{\Gamma \vdash_{\Sigma} \text{never}\{\tau\} : \text{event}(\tau)}$$

$$\frac{\Gamma \vdash_{\Sigma} e_1 : \text{event}(\tau) \quad \Gamma \vdash_{\Sigma} e_2 : \text{event}(\tau)}{\Gamma \vdash_{\Sigma} \text{or}(e_1; e_2) : \text{event}(\tau)}$$
(40.7c)

$$\frac{\Gamma \vdash_{\Sigma} e_1 : \text{event}(\tau_1) \quad \Gamma, x : \tau_1 \vdash_{\Sigma} e_2 : \tau_2}{\Gamma \vdash_{\Sigma} \text{wrap}(e_1; x.e_2) : \text{event}(\tau_2)}$$
(40.7d)

The corresponding dynamics is defined by these rules:

$$\frac{\Sigma \vdash a \sim \tau}{\operatorname{rcv}[a] \operatorname{val}_{\Sigma}} \tag{40.8a}$$

$$\frac{}{\mathsf{never}\{\tau\}\,\mathsf{val}_\Sigma} \tag{40.8b}$$

$$\frac{e_1 \operatorname{val}_{\Sigma} \quad e_2 \operatorname{val}_{\Sigma}}{\operatorname{or}(e_1; e_2) \operatorname{val}_{\Sigma}} \tag{40.8c}$$

$$\frac{e_1 \operatorname{val}_{\Sigma}}{\operatorname{wrap}(e_1; x.e_2) \operatorname{val}_{\Sigma}} \tag{40.8d}$$

$$\frac{e_1 \longmapsto_{\Sigma} e'_1}{\operatorname{or}(e_1; e_2) \longmapsto_{\Sigma} \operatorname{or}(e'_1; e_2)}$$

$$(40.8e)$$

$$\frac{e_1 \operatorname{val}_{\Sigma} \quad e_2 \underset{\Sigma}{\longmapsto} e_2'}{\operatorname{or}(e_1; e_2) \underset{\Sigma}{\longmapsto} \operatorname{or}(e_1; e_2')} \tag{40.8f}$$

$$\frac{e_1 \underset{\Sigma}{\longmapsto} e_1'}{\underset{\text{wrap}(e_1; x.e_2)}{\longmapsto} \text{wrap}(e_1'; x.e_2')} \tag{40.8g}$$

Event values are identified up to structural congruence as described in Chapter 39.

The statics of the synchronization command is given by the following rule:

$$\frac{\Gamma \vdash_{\Sigma} e : \text{event}(\tau)}{\Gamma \vdash_{\Sigma} \text{sync}(e) \stackrel{.}{\sim} \tau}$$

$$(40.9a)$$

The type of the event determines the type of value returned by the synchronization command. Execution of a synchronization command depends on the event.

$$\frac{e \xrightarrow{\Sigma} e'}{\operatorname{sync}(e) \stackrel{\varepsilon}{\Rightarrow} \operatorname{sync}(e')}$$

$$(40.10a)$$

$$\frac{e \operatorname{val}_{\Sigma} \vdash_{\Sigma} e : \tau \quad \Sigma \vdash a \sim \tau}{\operatorname{sync}(\operatorname{rcv}[a]) \xrightarrow{\underline{a \cdot e}?} \operatorname{ret}(e)}$$
(40.10b)

$$\frac{\operatorname{sync}(e_1) \stackrel{\alpha}{\Longrightarrow} m_1}{\operatorname{sync}(\operatorname{or}(e_1; e_2)) \stackrel{\alpha}{\Longrightarrow} m_1}$$

$$(40.10c)$$

$$\frac{\operatorname{sync}(e_2) \stackrel{\alpha}{\Longrightarrow} m_2}{\operatorname{sync}(\operatorname{or}(e_1; e_2)) \stackrel{\alpha}{\Longrightarrow} m_2}$$

$$(40.10d)$$

$$\frac{\operatorname{sync}(e_1) \stackrel{\alpha}{\Longrightarrow} m_1}{\operatorname{sync}(\operatorname{wrap}(e_1; x.e_2)) \stackrel{\alpha}{\Longrightarrow} \operatorname{bnd}(\operatorname{cmd}(m_1); x.\operatorname{ret}(e_2))} \tag{40.10e}$$

Rule (40.10b) states that an acceptance on a channel a may synchronize only with messages classified by a. When combined with structural congruence, Rules (40.10c) and (40.10d) state that either event between two choices may engender an action. Rule (40.10e) yields the command that performs the command m_1 resulting from the action α taken by the event e_1 , then returns e_2 with x bound to the return value of m_1 .

Selective communication and dynamic events can be used together to implement a communication protocol in which a channel reference is passed on a channel in order to establish a communication path with the recipient. Let a be a channel carrying values of type $\mathtt{cls}(\tau)$, and let b be a channel carrying values of type τ , so that &b can be passed as a message along channel a. A process that wishes to accept a channel reference on a and then accept on that channel has the form

$$\{x \leftarrow \mathtt{sync}(?a); y \leftarrow \mathtt{sync}(??x); \ldots\}.$$

The event ? a specifies a selective receipt on channel a. Once the value x is accepted, the event ?? x specifies a selective receipt on the channel referenced by x. So, if & b is sent along a, then the event ?? & b evaluates to ? b, which accepts selectively on channel b, even though the receiving process may have no direct access to the channel b itself.

40.4 Free Assignables as Processes

Scope-free assignables are definable in **CA** by associating to each assignable a server process that sets and gets the contents of the assignable. To each assignable a of type τ is associated a server that selectively accepts a message on channel a with one of two forms:

- 1. get \cdot (& b), where b is a channel of type τ . This message requests that the contents of a be sent on channel b.
- 2. set \cdot ($\langle e, \& b \rangle$), where e is a value of type τ , and b is a channel of type τ . This message requests that the contents of a be set to e, and that the new contents be transmitted on channel b.

In other words a is a channel of type τ_{srvr} given by

$$[\,\mathtt{get} \hookrightarrow \tau\,\mathtt{cls},\mathtt{set} \hookrightarrow \tau \times \tau\,\mathtt{cls}\,].$$

The server selectively accepts on channel *a*, then dispatches on the class of the message to satisfy the request.

The server associated with the assignable a of type τ maintains the contents of a using recursion. When called with the current contents of the assignable, the server selectively accepts on channel a, dispatching on the associated request, and calling itself recursively with the (updated, if necessary) contents:

$$\lambda \left(u : \tau_{srvr} \operatorname{cls} \right) \operatorname{fix} srvr : \tau \rightharpoonup \operatorname{void} \operatorname{cmd} \operatorname{is} \lambda \left(x : \tau \right) \operatorname{cmd} \left\{ y \leftarrow \operatorname{sync} \left(?? u \right) ; e_{(40.12)} \right\}. \tag{40.11}$$

The server is a procedure that takes an argument of type τ , the current contents of the assignable, and yields a command that never terminates, because it restarts the server loop after each request. The server selectively accepts a message on channel a, and dispatches on it as follows:

$$\operatorname{case} y \left\{ \operatorname{get} \cdot z \hookrightarrow e_{(40.13)} \mid \operatorname{set} \cdot \langle x', z \rangle \hookrightarrow e_{(40.14)} \right\}. \tag{40.12}$$

A request to get the contents of the assignable *a* is served as follows:

$$\{ -\leftarrow \operatorname{emit}(\operatorname{inref}(z; x)) ; \operatorname{do} \operatorname{srvr}(x) \}$$
 (40.13)

A request to set the contents of the assignable *a* is served as follows:

$$\{_\leftarrow \texttt{emit}(\texttt{inref}(z; x')); \texttt{do} \textit{srvr}(x')\} \tag{40.14}$$

The type τ ref is defined to be τ_{srvr} cls, the type of channels (classes) to servers providing a cell containing a value of type τ . A new free assignable is created by the command ref e_0 , which is defined to be

$$\{x \leftarrow \mathtt{newch}; _ \leftarrow \mathtt{spawn}(e_{(40.11)}(x)(e_0)); \mathtt{ret} x\}.$$
 (40.15)

A channel carrying a value of type τ_{srvr} is allocated to serve as the name of the assignable, and a new server is spawned that accepts requests on that channel, with initial value e_0 of type τ_0 .

The commands $*e_0$ and $e_0 *= e_1$ send a message to the server to get and set the contents of an assignable. The code for $*e_0$ is as follows:

$$\{x \leftarrow \text{newch}; _\leftarrow \text{emit}(\text{inref}(e_0; \text{get} \cdot x)); \text{sync}(??(x))\}$$
 (40.16)

A channel is allocated for the return value, the server is contacted with a get message specifying this channel, and the result of receiving on this channel is returned. Similarly, the code for $e_0 *= e_1$ is as follows:

$$\{x \leftarrow \text{newch}; _\leftarrow \text{emit}(\text{inref}(e_0; \text{set} \cdot \langle e_1, x \rangle)); \text{sync}(??(x))\}$$
 (40.17)

394 40.5 Notes

40.5 Notes

Concurrent Algol is a synthesis of process calculus and Modernized Algol; is essentially an "Algollike" formulation of Concurrent ML (Reppy, 1999). The design is influenced by Parallel Algol (Brookes, 2002). Much work on concurrent interaction takes communication channels as a basic concept, but see Linda (Gelernter, 1985) for an account similar to the one suggested here.

Exercises

- **40.1**. In Section 40.2 channels are allocated using the command newch, which returns a channel reference. Alternatively one may extend **CA** with a means of declaring channels just as assignables are declared in **MA**. Formulate the syntax, statics, and dynamics of such a construct, and derive newch using this extension.
- **40.2**. Extend selective communication (Section 40.3) to account for channel references, which give rise to a new form of event. Give the syntax, statics, and dynamics of this extension.
- 40.3. Adapt the implementation of an RS latch given in Exercise 39.3 to CA.

40.5 Notes



Part XVII Modularity



Chapter 42

Modularity and Linking

Modularity is the most important technique for controlling the complexity of programs. Programs are decomposed into separate *components* with precisely specified, and tightly controlled, interactions. The pathways for interaction among components determine dependencies that constrain the process by which the components are integrated, or *linked*, to form a complete system. Different systems may use the same components, and a single system may use multiple instances of a single component. Sharing of components amortizes the cost of their development across systems, and helps limit errors by limiting coding effort.

Modularity is not limited to programming languages. In mathematics the proof of a theorem is decomposed into a collection of definitions and lemmas. References among the lemmas determine a dependency relation that constrains their integration to form a complete proof of the main theorem. Of course, one person's theorem is another person's lemma; there is no intrinsic limit on the depth and complexity of the hierarchies of results in mathematics. Mathematical structures are themselves composed of separable parts, for example, a ring comprises a group and a monoid structure on the same underlying set.

Modularity arises from the structural properties of the hypothetical and general judgments. Dependencies among components are expressed by free variables whose typing assumptions state the presumed properties of the component. Linking amounts to substitution to discharge the hypothesis.

42.1 Simple Units and Linking

Decomposing a program into units amounts to exploiting the transitivity of the hypothetical judgment (see Chapter 3). The decomposition may be described as an interaction between two parties, the *client* and the *implementor*, mediated by an agreed-upon contract, an *interface*. The client *assumes* that the implementor upholds the contract, and the implementor *guarantees* that the contract will be upheld. The assumption made by the client amounts to a declaration of its dependence on the implementor discharged by *linking* the two parties according to their agreed-upon contract.

The interface that mediates the interaction between a client and an implementor is a *type*. Linking is the implementation of the composite structural rules of substitution and transitivity:

$$\frac{\Gamma \vdash e_{impl} : \tau_{intf} \quad \Gamma, x : \tau_{intf} \vdash e_{client} : \tau_{client}}{\Gamma \vdash [e_{impl} / x] e_{client} : \tau_{client}}$$
(42.1)

The type τ_{intf} is the interface type. It defines the operations provided by the implementor e_{impl} and relied upon by the client e_{client} . The free variable x expresses the dependency of e_{client} on e_{impl} . That is, the client accesses the implementation by using the variable x.

The interface type τ_{intf} is the contract between the client and the implementor. It determines the properties of the implementation on which the client may depend and, at the same time, determines the obligations that the implementor must fulfill. The simplest form of interface type is a finite product type of the form $\langle f_1 \hookrightarrow \tau_1, \ldots, f_n \hookrightarrow \tau_n \rangle$, specifying a component with components f_i of type τ_i . Such a type is an *application program interface*, or *API*, because it determines the operations that the client (application) may expect from the implementor. A more advanced form of interface is one that defines an abstract type of the form $\exists (t.\langle f_1 \hookrightarrow \tau_1, \ldots, f_n \hookrightarrow \tau_n \rangle)$, which defines an abstract type t representing the internal state of an "abstract machine" whose "instruction set" consists of the operations f_1, \ldots, f_n whose types may involve t. Being abstract, the type t is not revealed to the client, but is known only to the implementor.¹

Conceptually, linking is just substitution, but practically this can be implemented in many ways. One method is *separate compilation*. The expressions e_{client} and e_{impl} , the *source modules*, are translated (compiled) into another, lower-level, language, resulting in *object modules*. Linking consists of performing the required substitution at the level of the object language in such a way that the result corresponds to translating $[e_{impl}/x]e_{client}$. Another method, *separate checking*, shifts the requirement for translation to the linker. The client and implementor units are checked for type correctness with respect to the interface, but are not translated into lower-level form. Linking then consists of translating the composite program as a whole, often resulting in a more efficient outcome than would be possible when compiling separately.

The foregoing are all forms of *static linking* because the program is composed before it is executed. Another method, *dynamic linking*, defers program composition until run-time, so that a component is loaded only if it is actually required during execution. This might seem to involve executing programs with free variables, but it does not. Each client implemented by a *stub* that forwards accesses to a stored implementation (typically, in an ambient file system). The difficulty with dynamic linking is that it refers to components by name (say, a path in a file system), and the binding of that name may change at any time, wreaking havoc on program behavior.

42.2 Initialization and Effects

Linking resolves the dependencies among the components of a program by substitution. This view is valid so long as the components are given by pure expressions, those that evaluate to a value without inducing any effects. For in such cases there is no problem with the replication, or

¹See Chapters 17 and 48 for a discussion of type abstraction.

complete omission, of a component arising from repeated, or absent, uses of a variable representing it. But what if the expression defining the implementation of a component has an effect when evaluated? At a minimum replication of the component implies replication of its effects. Worse, effects introduce *implicit dependencies* among components that are not apparent from their types. For example, if each of two components mutates a shared assignable, the order in which they are linked with a client program affects the behavior of the whole.

This may raise doubts about the treatment of linking as substitution, but on closer inspection it becomes clear that implicit dependencies are naturally expressed by the modal distinction between expressions and commands introduced in Chapter 34. Specifically, a component that may have an effect when executed does not have type τ_{intf} of implementations of the interface type, but rather the type τ_{intf} cmd of encapsulated commands that, when executed, have effects and yield implementations. Being encapsulated, a value of this type is itself free of effects, but it may have effects when evaluated.

The distinction between the types τ_{intf} and τ_{intf} cmd is mediated by the sequencing command introduced in Chapter 34. For the sake of generality, let us assume that the client is itself an encapsulated command of type τ_{client} cmd, so that it may itself have effects when executed, and may serve as a component of a yet larger system. Assuming that the client refers to the encapsulated implementation by the variable x, the command

bnd
$$x \leftarrow x$$
; do e_{client}

first determines the implementation of the interface by running the encapsulated command x then running the client code with the result bound to x. The implicit dependencies of the client on the implementor are made explicit by the sequencing command, which ensures that the implementor's effects occur prior to those of the client, precisely because the client depends on the implementor for its execution.

More generally, to manage such interactions in a large program it is common to isolate an *initialization procedure* whose role is to stage the effects engendered by the various components according to some policy or convention. Rather than attempt to survey all possible policies, let us just note that the upshot of such conventions is that the initialization procedure is a command of the form

$$\{x_1 \leftarrow x_1; \dots x_n \leftarrow x_n; m_{main}\},\$$

where $x_1, ..., x_n$ represent the components of the system and m_{main} is the main (startup) routine. After linking the initialization procedure has the form

$$\{x_1 \leftarrow e_1; \dots x_n \leftarrow e_n; m_{main}\},\$$

where e_1, \ldots, e_n are the encapsulated implementations of the linked components. When the initialization procedure is executed, it results in the substitution

$$[v_1,\ldots,v_n/x_1,\ldots,x_n]m_{main},$$

where the expressions v_1, \ldots, v_n represent the values resulting from executing e_1, \ldots, e_n , respectively, and the implicit effects have occurred in the order specified by the initializer.

408 42.3 Notes

42.3 Notes

The relationship between the structural properties of entailment and the practical problem of separate development was implicit in much early work on programming languages, but became explicit once the correspondence between propositions and types was developed. There are many indications of this correspondence in sources such as *Proofs and Types* (Girard, 1989) and *Intuitionistic Type Theory* (Martin-Löf, 1984), but it was first made explicit by Cardelli (1997).

Chapter 44

Type Abstractions and Type Classes

An interface is a contract that specifies the rights of a client and the responsibilities of an implementor. Being a specification of behavior, an interface is a type. In principle any type may serve as an interface, but in practice it is usual to structure code into *modules* consisting of separable and reusable components. An interface specifies the behavior of a module expected by a client and imposed on the implementor. It is the fulcrum balancing the tension between separation and integration. As a rule, a module ought to have a well-defined behavior that can be understood separately, but it is equally important that it be easy to combine modules to form an integrated whole.

A fundamental question is, what is the type of a module? That is, what form should an interface take? One long-standing idea is that an interface is a labeled tuple of functions and procedures with specified types. The types of the fields of the tuple are often called *function headers*, because they summarize the call and return types of each function. Using interfaces of this form is called *procedural abstraction*, because it limits the dependencies between modules to a specified set of procedures. We may think of the fields of the tuple as being the instruction set of a virtual machine. The client makes use of these instructions in its code, and the implementor agrees to provide their implementations.

The problem with procedural abstraction is that it does not provide as much insulation as one might like. For example, a module that implements a dictionary must expose in the types of its operations the exact representation of the tree as, say, a recursive type (or, in more rudimentary languages, a pointer to a structure that itself may contain such pointers). Yet the client ought not depend on this representation: the purpose of abstraction is to get rid of pointers. The solution, as discussed in Chapter 17, is to extend the abstract machine metaphor to allow the internal state of the machine to be hidden from the client. In the case of a dictionary the representation of the dictionary as a binary search tree is hidden by existential quantification. This concept is called *type abstraction*, because the type of the underlying data (state of the abstract machine) is hidden.

Type abstraction is a powerful method for limiting the dependencies among the modules that constitute a program. It is very useful in many circumstances, but is not universally applicable. It is often useful to expose, rather than obscure, type information across a module boundary. A typical example is the implementation of a dictionary, which is a mapping from keys to values. To

use, say, a binary search tree to implement a dictionary, we require that the key type admit a total ordering with which keys can be compared. The dictionary abstraction does not depend on the exact type of the keys, but only requires that the key type be constrained to provide a comparison operation. A *type class* is a specification of such a requirement. The class of comparable types, for example, specifies a type t together with an operation leq of type $(t \times t) \to bool$ with which to compare them. Superficially, such a specification looks like a type abstraction, because it specifies a type and one or more operations on it, but with the important difference that the type t is not hidden from the client. For if it were, the client would only be able to compare keys using leq, but would have no means of obtaining keys to compare. A type class, in contrast to a type abstraction, is not intended to be an exhaustive specification of the operations on a type, but as a constraint on its behavior expressed by demanding that certain operations, such as comparison, be available, without limiting the other operations that might be defined on it.

Type abstractions and type classes are the extremal cases of a general concept of module type that we shall discuss in detail in this chapter. The crucial idea is the *controlled revelation* of type information across module boundaries. Type abstractions are opaque; type classes are transparent. These are both instances of *translucency*, which arises from combining existential types (Chapter 17), subtyping (Chapter 24), and singleton kinds and subkinding (Chapter 43). Unlike in Chapter 17, however, we will distinguish the types of modules, which are called *signatures*, from the types of ordinary values. The distinction is not essential, but it will be helpful to keep the two concepts separate at the outset, deferring discussion of how to ease the segregation once the basic concepts are in place.

44.1 Type Abstraction

Type abstraction is captured by a form of existential type quantification similar to that described in Chapter 17. For example, a dictionary with keys of type τ_{key} and values of type τ_{val} implements the signature σ_{dict} defined by $[t: T; \tau_{\text{dict}}]$, where τ_{dict} is the labeled tuple type

$$\langle \texttt{emp} \hookrightarrow t \,, \texttt{ins} \hookrightarrow \tau_{\mathsf{key}} \times \tau_{\mathsf{val}} \times t \to t \,, \texttt{fnd} \hookrightarrow \tau_{\mathsf{key}} \times t \to \tau_{\mathsf{val}} \, \texttt{opt} \rangle.$$

The type variable t occurring in $\tau_{\rm dict}$ and bound by $\sigma_{\rm dict}$ is the abstract type of dictionaries on which are defined three operations emp, ins, and fnd with the specified types. The type $\tau_{\rm val}$ is immaterial to the discussion, because the dictionary operations impose no restrictions on the values that are associated to keys. However, it is important that the type $\tau_{\rm key}$ be some fixed type, such as str, equipped with a suite of operations such as comparison. Observe that the signature $\sigma_{\rm dict}$ merely specifies that a dictionary is a value of some type that admits the operations emp, ins, and fnd with the types given by $\tau_{\rm dict}$.

An implementation of the signature σ_{dict} is a *structure* M_{dict} of the form $[\![\rho_{\text{dict}}]\!]$, where ρ_{dict} is some concrete representation of dictionaries, and e_{dict} is a labeled tuple of type $[\![\rho_{\text{dict}}]\!]$ of the general form

$$\langle \texttt{emp} \hookrightarrow \ldots, \texttt{ins} \hookrightarrow \ldots, \texttt{fnd} \hookrightarrow \ldots \rangle$$
.

The elided parts implement the dictionary operations in terms of the chosen representation type ρ_{dict} making use of the comparison operation that we assume is available of values of type τ_{key} . For

example, the type ρ_{dict} might be a recursive type defining a balanced binary search tree, such as a red-black tree. The dictionary operations work on the underlying representation of the dictionary as such a tree, just as would a package of existential type (see Chapter 17). The supposition about τ_{key} is temporary, and is lifted in Section 44.2.

To ensure that the representation of the dictionary is hidden from a client, the structure $M_{\rm dict}$ is *sealed* with the signature $\sigma_{\rm dict}$ to obtain the module

$M_{\rm dict}$ 1 $\sigma_{\rm dict}$.

The effect of sealing is to ensure that the *only* information about M_{dict} that propagates to the client is given by σ_{dict} . In particular, because σ_{dict} only specifies that the type t have kind T, no information about the choice of t as ρ_{dict} in M_{dict} is made available to the client.

A module is a *two-phase* object consisting of a *static part* and a *dynamic part*. The static part is a constructor of a specified kind; the dynamic part is a value of a specified type. There are two elimination forms that extract the static and dynamic parts of a module. These are, respectively, a form of constructor and a form of expression. More precisely, the constructor $M \cdot s$ stands for the static part of M, and the expression $M \cdot d$ stands for its dynamic part. According to the inversion principle, if a module M has introduction form, then $M \cdot s$ should be equivalent to the static part of M. So, for example, $M_{\text{dict}} \cdot s$ should be equivalent to ρ_{dict} .

But consider the static part of a sealed module, which has the form $(M_{\text{dict}} \mid \sigma_{\text{dict}}) \cdot s$. Because sealing hides the representation of an abstract type, this constructor should not be equivalent to ρ_{dict} . If M'_{dict} is another implementation of σ_{dict} , should $(M_{\text{dict}} \mid \sigma_{\text{dict}}) \cdot s$ be equivalent to $(M'_{\text{dict}} \mid \sigma_{\text{dict}}) \cdot s$? To ensure reflexivity of type equivalence this equation should hold when M and M' are equivalent modules. But this violates representation independence for abstract types by making equivalence of abstract types sensitive to their implementation.

It would seem, then, that there is a contradiction between two very fundamental concepts, type equivalence and representation independence. The way out of this conundrum is to *disallow* reference to the static part of a sealed module: the type expression $M \mid \sigma \cdot s$ is deemed ill-formed. More generally, the formation of $M \cdot s$ is disallowed unless M is a *module value*, whose static part is always manifest. An explicit structure is a module value, as is any module variable (provided that module variables are bound by-value).

One effect of this restriction is that sealed modules must be bound to a variable before they are used. Because module variables are bound by-value, doing so has the effect of imposing abstraction at the binding site. In fact, we may think of sealing as a kind of computational effect that "occurs" at the binding site, much as the bind operation in Algol discussed in Chapter 34 engenders the effects induced by an encapsulated command. As a consequence two bindings of the same sealed module result in two abstract types. The type system willfully ignores the identity of the two occurrences of the same module in order to ensure that their representations can be changed independently of one another without disrupting the behavior of any client code (because the client cannot rely on their identity, it must regard them as different).

422 44.2 Type Classes

44.2 Type Classes

Type abstraction is an essential tool for limiting dependencies among modules in a program. The signature of a type abstraction determines all that is known about a module by a client; no other uses of the values of an abstract type are permissible. A complementary tool is to use a signature to partially specify the capabilities of a module. Such a signature is a *type class*, or a *view*; an *instance* of the type class is an implementation of it. Because the signature of a type class only constrains the minimum capabilities of an unknown module, there must be some other means of working with values of that type. The way to achieve this is to expose, rather than hide, the identity of the static part of a module. In this sense type classes are the "opposite" of type abstractions, but we shall see below that there is a smooth progression between them, mediated by a subsignature judgment.

Let us consider the implementation of dictionaries as a client of the implementation of its keys. To implement a dictionary using a binary search tree the only requirement is that keys come equipped with a total ordering given by a comparison operation. This requirement can be expressed by a signature σ_{ord} given by

$$\llbracket t :: T; \langle \mathtt{leq} \hookrightarrow (t \times t) \rightarrow \mathtt{bool} \rangle \rrbracket.$$

Because a given type can be ordered in many ways, it is essential that the ordering be packaged with the type to determine a type of keys.

The implementation of dictionaries as binary search trees takes the form

$$X: \sigma_{\mathsf{ord}} \vdash M^{X}_{\mathsf{bstdict}}: \sigma^{X}_{\mathsf{dict}}.$$

Here σ_{dict}^X is the signature $[\![t::\mathtt{T}\,;\tau_{\mathrm{dict}}^X]\!]$, whose body, τ_{dict}^X , is the tuple type

$$\langle \texttt{emp} \hookrightarrow t \,, \texttt{ins} \hookrightarrow X \cdot \texttt{s} \times \tau_{\mathsf{val}} \times t \to t \,, \texttt{fnd} \hookrightarrow X \cdot \texttt{s} \times t \to \tau_{\mathsf{val}} \, \texttt{opt} \rangle,$$

and M_{bstdict}^X is a structure (not given explicitly here) that implements the dictionary operations using binary search trees. Within M_{bstdict}^X , the static and dynamic parts of the module X are accessed by writing $X \cdot s$ and $X \cdot d$, respectively. In particular, the comparison operation on keys is accessed by the projection $X \cdot d \cdot leq$.

The declared signature of the module variable X expresses a constraint on the capabilities of a key type by specifying an upper bound on its signature in the subsignature ordering. So any module bound to X must provide a type of keys and a comparison operation on that type, but nothing else is assumed of it. Because this is all we know about the unknown module X the dictionary implementation is constrained to rely only on these specified capabilities, and no others. When linking with a module defining X, the implementation need not be sealed with this signature, but must instead have a signature that is no larger than it in the subsignature relation. Indeed, the signature σ_{ord} is useless for sealing, as is easily seen by example. Suppose that M_{natord} : σ_{ord} is an instance of the class of ordered types under the usual ordering. If we seal M_{natord} with σ_{ord} by writing

$$M_{\rm natord} \mid \sigma_{\rm ord}$$
,

¹Here and elsewhere in this chapter and the next, the superscript X serves as a reminder that the module variable X may occur free in the annotated module or signature.

44.2 Type Classes 423

the resulting module is *useless*, because we would then have no way to create values of the key type.

We see, then, that a type class is a description (or view) of a pre-existing type, and is not a means of introducing a new type. Rather than obscure the identity of the static part of $M_{\rm natord}$, we wish to propagate its identity as nat while specifying a comparison with which to order them. Type identity propagation is achieved using singleton kinds (as described in Chapter 43). Specifically, the most precise, or *principal*, signature of a structure is the one that exposes its static part using a singleton kind. In the case of the module $M_{\rm natord}$, the principal signature is the signature $\sigma_{\rm natord}$ given by

$$\llbracket t :: S(\mathtt{nat}); \mathtt{leq} \hookrightarrow (t \times t) \rightarrow \mathtt{bool} \rrbracket$$

which, by the rules of equivalence (defined formally in Section 44.3), is equivalent to the signature

$$\llbracket_{\text{-}} :: S(\,\mathtt{nat}\,) \, ; \mathtt{leq} \, \hookrightarrow \, (\,\mathtt{nat} \, \times \, \mathtt{nat}\,) \, \rightarrow \, \mathtt{bool} \rrbracket \, .$$

The derivation of such an equivalence is called *equivalence propagation*, because it propagates the identity of the type *t* into its scope.

The dictionary implementation M_{bstdict}^X expects a module X with signature σ_{ord} , but the module M_{natord} provides the signature σ_{natord} . Applying the rules of subkinding given in Chapter 43, together with the covariance principle for signatures, we obtain the subsignature relationship

$$\sigma_{\mathsf{natord}} <: \sigma_{\mathsf{ord}}.$$

By the subsumption principle, a module of signature σ_{natord} may be provided when a module of signature σ_{ord} is required. Therefore M_{natord} may be linked to X in M_{bstdict}^X .

Combining subtyping with sealing provides a smooth gradation between type classes and type abstractions. The principal signature for M_{bstdict}^X is the signature ρ_{dict}^X given by

$$\left[\!\!\left[t :: \mathtt{S}(\,\tau^{X}_{\mathtt{bst}}\,) \: ; \langle \mathtt{emp} \,{\hookrightarrow}\, t \: , \mathtt{ins} \,{\hookrightarrow}\, X \cdot \mathtt{s} \times \tau_{\mathtt{val}} \times t \to t \: , \mathtt{fnd} \,{\hookrightarrow}\, X \cdot \mathtt{s} \times t \to \tau_{\mathtt{val}} \, \mathtt{opt} \rangle \right]\!\!\right] ,$$

where τ_{bst}^X is the type of binary search trees with keys given by the module X of signature σ_{ord} . This signature is a subsignature of σ_{dict}^X given earlier, so that the sealed module

$$M_{
m bstdict}^{X} \mid \sigma_{
m dict}^{X}$$

is well-formed, and has type $\sigma_{
m dict}^X$, which hides the representation type of the dictionary abstraction

After linking X to M_{natord} , the signature of the dictionary is specialized by propagating the identity of the static part of M_{natord} using the subsignature judgment. As remarked earlier, the dictionary implementation satisfies the typing

$$X : \sigma_{\mathsf{ord}} \vdash M_{\mathsf{bstdict}}^X : \sigma_{\mathsf{dict}}^X$$

But because $\sigma_{\text{natord}} <: \sigma_{\text{ord}}$, we have, by contravariance, that

$$X : \sigma_{\mathsf{natord}} \vdash M^{X}_{\mathsf{bstdict}} : \sigma^{X}_{\mathsf{dict}}.$$

is also a valid typing judgment. If $X:\sigma_{\mathsf{natord}}$, then $X\cdot \mathsf{s}$ is equivalent to \mathtt{nat} , because it has kind $\mathtt{S}(\mathtt{nat})$, so that the typing

$$X : \sigma_{\mathsf{natord}} \vdash M^{X}_{\mathsf{bstdict}} : \sigma_{\mathsf{natdict}}$$

is also valid. The closed signature $\sigma_{natdict}$ is given explicitly by

$$[\![t :: \mathtt{T} \, ; \langle \mathtt{emp} \, \hookrightarrow t \, , \mathtt{ins} \, \hookrightarrow \mathtt{nat} \, \times \, \tau_{\mathsf{val}} \, \times \, t \, \to \, t \, , \mathtt{fnd} \, \hookrightarrow \mathtt{nat} \, \times \, t \, \to \, \tau_{\mathsf{val}} \, \mathtt{opt} \rangle]\!] \, .$$

The representation of dictionaries is hidden, but the representation of keys as natural numbers is not. The dependency on X has been eliminated by replacing all occurrences of $X \cdot s$ within σ^X_{dict} by the type nat. Having derived this typing we may link X with M_{natord} as described in Chapter 42 to obtain a composite module, M_{natdict} , of signature σ_{natdict} , in which keys are natural numbers ordered as specified by M_{natord} .

It is convenient to exploit subtyping for labeled tuple types to avoid creating an *ad hoc* module specifying the standard ordering on the natural numbers. Instead we can extract the required module directly from the implementation of the abstract type of numbers using subsumption. As an illustration, let X_{nat} be a module variable of signature σ_{nat} , which has the form

$$\llbracket t :: T; \langle {\tt zero} \hookrightarrow t \,, {\tt succ} \hookrightarrow t \to t \,, {\tt leq} \hookrightarrow (\, t \times t \,) \to {\tt bool} \,, \dots \rangle
bracket$$

The fields of the tuple provide all and only the operations that are available on the abstract type of natural numbers. Among them is the comparison operation 1eq, which is required by the dictionary module. Applying the subtyping rules for labeled tuples given in Chapter 24, together with the covariance of signatures, we obtain the subsignature relationship

$$\sigma_{\sf nat} <: \sigma_{\sf ord}$$

so that by subsumption the variable X_{nat} may be linked to the variable X postulated by the dictionary implementation. Subtyping takes care of extracting the required leq field from the abstract type of natural numbers, demonstrating that the natural numbers are an instance of the class of ordered types. Of course, this approach only works if we wish to order the natural numbers in the natural way provided by the abstract type. If, instead, we wish to use another ordering, then we must construct instances of σ_{ord} "by hand" to define the appropriate ordering.

44.3 A Module Language

The module language **Mod** formalizes theideas outlined in the preceding section. The syntax is divided into five levels: expressions classified by types, constructors classified by kinds, and modules classified by signatures. The expression and type level consists of various language mechanisms described earlier in this book, including at least product, sum, and partial function types. The constructor and kind level is as described in Chapters 18 and 43, with singleton and dependent

kinds. The following grammar summarizes the syntax of modules.

```
Sig
                    sig{\kappa}(t.\tau)
                                                 \llbracket t :: \kappa ; \tau \rrbracket
                                                                                  signature
\mathsf{Mod}\ M ::=
                                                                                  variable
                                                                                  structure
                     str(c;e)
                                                 [c;e]
                     seal\{\sigma\}(M)
                                                 M1\sigma
                                                                                  seal
                     let{\sigma}(M_1; X.M_2)
                                                (let X be M_1 in M_2 ): \sigma
                                                                                  definition
                     s(M)
                                                 M \cdot s
                                                                                  static part
Con
                                                 M \cdot d
Exp
                   d(M)
                                                                                  dynamic part
```

The statics of **Mod** consists of the following forms of judgment:

$\Gamma \vdash \sigma sig$	well-formed signature
$\Gamma \vdash \sigma_1 \equiv \sigma_2$	equivalent signatures
$\Gamma \vdash \sigma_1 <: \sigma_2$	subsignature
$\Gamma \vdash M : \sigma$	well-formed module
$\Gamma \vdash M$ val	module value
$\Gamma \vdash e \text{ val}$	expression value

Rather than segregate hypotheses into zones, we instead admit the following three forms of hypothesis groups:

$X:\sigma$, X val	module value variable
$u :: \kappa$	constructor variable
$x:\tau,x$ val	expression value variable

It is important that module and expression variables are always regarded as values to ensure that type abstraction is properly enforced. Correspondingly, each module and expression variable appears in Γ paired with the hypothesis that it is a value. As a notational convenience we will not explicitly state the value hypotheses associated with module and expression variables, under the convention that all such variables implicitly come paired with such an assumption.

The following rules define the formation, equivalence, and subsignature judgments.

$$\frac{\Gamma \vdash \kappa \text{ kind } \Gamma, u :: \kappa \vdash \tau \text{ type}}{\Gamma \vdash \llbracket u :: \kappa ; \tau \rrbracket \text{ sig}}$$

$$(44.1a)$$

$$\frac{\Gamma \vdash \kappa_1 \equiv \kappa_2 \quad \Gamma, u :: \kappa_1 \vdash \tau_1 \equiv \tau_2}{\Gamma \vdash [u :: \kappa_1 ; \tau_1]] \equiv [u :: \kappa_2 ; \tau_2]}$$

$$(44.1b)$$

$$\frac{\Gamma \vdash \kappa_1 < :: \kappa_2 \quad \Gamma, u :: \kappa_1 \vdash \tau_1 <: \tau_2}{\Gamma \vdash \llbracket u :: \kappa_1 ; \tau_1 \rrbracket <: \llbracket u :: \kappa_2 ; \tau_2 \rrbracket}$$

$$(44.1c)$$

Most importantly, signatures are covariant in both the kind and type positions: subkinding and subtyping are preserved by the formation of a signature. It follows from rule (44.1b) that

$$\llbracket u :: S(c); \tau \rrbracket \equiv \llbracket _ :: S(c); [c/u]\tau \rrbracket$$

and, further, it follows from rule (44.1c) that

$$\llbracket _ :: \mathtt{S}(\, c \,) \, ; [c/u]\tau \rrbracket <: \llbracket _ :: \mathtt{T} \, ; [c/u]\tau \rrbracket$$

and so

$$[u :: S(c); \tau] <: [- :: T; [c/u]\tau].$$

It is also the case that

$$[\![u :: S(c); \tau]\!] <: [\![u :: T; \tau]\!].$$

But the two supersignatures of $[u :: S(c); \tau]$ are *incomparable* with respect to the subsignature judgment.

The statics of expressions of **Mod** is given by the following rules:

$$\overline{\Gamma, X : \sigma \vdash X : \sigma} \tag{44.2a}$$

$$\frac{\Gamma \vdash c :: \kappa \quad \Gamma \vdash e : [c/u]\tau}{\Gamma \vdash \llbracket c \, ; e \rrbracket : \llbracket u :: \kappa \, ; \tau \rrbracket}$$
(44.2b)

$$\frac{\Gamma \vdash \sigma \operatorname{sig} \quad \Gamma \vdash M : \sigma}{\Gamma \vdash M \uparrow \sigma : \sigma}$$
(44.2c)

$$\frac{\Gamma \vdash \sigma \operatorname{sig} \quad \Gamma \vdash M_1 : \sigma_1 \quad \Gamma, X : \sigma_1 \vdash M_2 : \sigma}{\Gamma \vdash (\operatorname{let} X \operatorname{be} M_1 \operatorname{in} M_2) : \sigma : \sigma}$$

$$(44.2d)$$

$$\frac{\Gamma \vdash M : \sigma \quad \Gamma \vdash \sigma <: \sigma'}{\Gamma \vdash M : \sigma'} \tag{44.2e}$$

In rule (44.2b) it is always possible to choose κ to be the most specific kind of c in the subkind ordering, which uniquely determines c up to constructor equivalence. For such a choice, the signature $\llbracket u :: \kappa; \tau \rrbracket$ is equivalent to $\llbracket _{-} :: \kappa; \lceil c/u \rceil \tau \rrbracket$, which propagates the identity of the static part of the module expression into the type of its dynamic part. Rule (44.2c) is used together with the subsumption (rule (44.2e)) to ensure that M has the specified signature.

The need for a signature annotation on a module definition is a manifestation of the *avoidance problem*. Rule (44.2d) would be perfectly sensible were the signature σ omitted from the syntax of the definition. However, omitting this information greatly complicates type checking. If σ were omitted from the syntax of the definition, the type checker would be required to find a signature σ for the body of the definition that *avoids* the module variable X. Inductively, we may suppose that we have found a signature σ_1 for the module M_1 , and a signature σ_2 for the module M_2 , under the assumption that X has signature σ_1 . To find a signature for an unadorned definition, we must find a supersignature σ of σ_2 that avoids σ_2 . To ensure that all possible choices of σ_2 are accounted for, we seek the least (most precise) such signature with respect to the subsignature relation; this is called the *principal signature* of a module. The problem is that there may not be a least supersignature of a given signature that avoids a specified variable. (Consider the example above of a signature with two incomparable supersignatures. The example can be chosen so that the supersignatures avoid a variable σ_2 that occurs in the subsignature.) Consequently, modules do not have principal signatures, a significant complication for type checking. To avoid this problem, we insist that the

avoiding supersignature σ be given by the programmer so that the type checker is not required to find one.

Modules give rise to a new form of constructor expression, $M \cdot s$, and a new form of value expression, $M \cdot d$. These operations, respectively, extract the static and dynamic parts of the module M. Their formation rules are as follows:

$$\frac{\Gamma \vdash M \text{ val} \quad \Gamma \vdash M : \llbracket u :: \kappa; \tau \rrbracket}{\Gamma \vdash M \cdot \mathbf{s} :: \kappa}$$
 (44.3a)

$$\frac{\Gamma \vdash M : \llbracket_{-} :: \kappa; \tau\rrbracket}{\Gamma \vdash M \cdot d : \tau}$$
(44.3b)

Rule (44.3a) requires that the module expression *M* be a value according to the following rules:

$$\frac{}{\Gamma, X : \sigma, X \text{ val} \vdash X \text{ val}}$$
 (44.4a)

$$\frac{\Gamma \vdash e \text{ val}}{\Gamma \vdash \llbracket c ; e \rrbracket \text{ val}}$$
 (44.4b)

(It is not strictly necessary to insist that the dynamic part of a structure be a value for the structure to itself be a value.)

Rule (44.3a) specifies that only structure values have well-defined static parts, and hence precludes reference to the static part of a sealed structure, which is not a value. This property ensures representation independence for abstract types, as discussed in Section 44.1. For if $M \cdot s$ were admissible when M is a sealed module, it would be a type whose identity depends on the underlying implementation, in violation of the abstraction principle. Module variables are, on the other hand, values, so that if $X : [t :: T; \tau]$ is a module variable, then $X \cdot s$ is a well-formed type. What this means in practice is that sealed modules must be bound to variables before they can be used. It is for this reason that we include definitions among module expressions.

Rule (44.3b) requires that the signature of the module, M, be non-dependent, so that the result type, τ , does not depend on the static part of the module. This independence may not always be the case. For example, if M is a sealed module, say $N \upharpoonright \llbracket t :: T;t \rrbracket$ for some module N, then projection $M \cdot d$ is ill-formed. For if it were well-formed, its type would be $M \cdot s$, which would violate representation independence for abstract types. But if M is a module value, then it is always possible to derive a non-dependent signature for it, provided that we include the following rule of *self-recognition*:

$$\frac{\Gamma \vdash M : \llbracket u :: \kappa; \tau \rrbracket \quad \Gamma \vdash M \text{ val}}{\Gamma \vdash M : \llbracket u :: S(M \cdot s :: \kappa); \tau \rrbracket}$$
(44.5)

This rule propagates the identity of the static part of a module value into its signature. The dependency of the type of the dynamic part on the static part is then eliminable by sharing propagation.

The following rule of constructor equivalence states that a type projection from a module value is eliminable:

$$\frac{\Gamma \vdash \llbracket c ; e \rrbracket : \llbracket t :: \kappa ; \tau \rrbracket \quad \Gamma \vdash \llbracket c ; e \rrbracket \text{ val}}{\Gamma \vdash \llbracket c ; e \rrbracket \cdot \mathbf{s} \equiv c :: \kappa}$$

$$(44.6)$$

The requirement that the expression e be a value, which is implicit in the second premise of the rule, is not strictly necessary, but does no harm. A consequence is that apparent dependencies of closed constructors (or kinds) on modules may always be eliminated. In particular the identity of the constructor [c;e] · s is independent of e, as would be expected if representation independence is to be assured.

The dynamics of modules is given as follows:

$$\frac{e \longmapsto e'}{\llbracket c ; e \rrbracket \longmapsto \llbracket c ; e' \rrbracket} \tag{44.7a}$$

$$\frac{e \text{ val}}{\llbracket c ; e \rrbracket \cdot d \longmapsto e} \tag{44.7b}$$

There is no need to evaluate constructors at run-time, because the dynamics of expressions does not depend on their types. It is not difficult to prove type safety for this dynamics relative to the foregoing statics.

44.4 First- and Second-Class

It is common to draw a distinction between *first-class* and *second-class* modules based on whether signatures are types, and hence whether modules are just a form of expression like any other. When modules are first-class their values can depend on the state of the world at run-time. When modules are second-class signatures are a separate form of classifier from types, and module expressions may not be used in the same way as ordinary expressions. For example, it may not be possible to compute a module based on the phase of the moon.

Superficially, it seems as though first-class modules are uniformly superior to second-class modules, because you can do more with them. But on closer examination we see that the "less is more" principle applies here as well, much as in the distinction between dynamic and static languages discussed in Chapters 22 and 23. In particular if modules are first-class, then one must adopt a "pessimistic" attitude towards expressions that compute them, precisely because they represent fully general, even state-dependent, computations. One consequence is that it is difficult, or even impossible, to track the identity of the static part of a module during type checking. A general module expression need not have a well-defined static component, precluding its use in type expressions. Second-class modules, on the other hand, can be permissive with the use of the static components of modules in types, precisely because the range of possible computations is reduced. In this respect second-class modules are more powerful than first-class, despite initial impressions. More importantly, a second-class module system can always be enriched to allow first-class modules, without requiring that they be first-class. Thus we have the best of both worlds: the flexibility of first-class modules and the precision of second-class modules. In short you pay for only what you use: if you use first-class capabilities, you should expect to pay a cost, but if you do not, you should not be taxed on the unrealized gain.

First-class modules are added to **Mod** in the following way. First, enrich the type system with existential types, as described in Chapter 17, so that "first-class modules" are just packages of existential type. A second-class module M of signature $[t :: \kappa; \tau]$ is made first-class by forming

44.5 Notes 429

the package pack $M \cdot s$ with $M \cdot d$ as $\exists (t.\tau)$ of type $\exists t :: \kappa.\tau$ consisting of the static and dynamic parts of M. Second, to allow packages to act like modules, we introduce the module expression open e that opens the contents of a package as a module:

$$\frac{\Gamma \vdash e : \exists t :: \kappa.\tau}{\Gamma \vdash \text{open } e : \llbracket t :: \kappa ; \tau \rrbracket}$$
(44.8)

Because the package e is an arbitrary expression of existential type, the module expression open e may not be regarded as a value, and hence does not have a well-defined static part. Instead we must generally bind it to a variable before it is used, mimicking the composite behavior of the existential elimination form given in Chapter 17.

44.5 Notes

The use of dependent types to express modularity was first proposed by MacQueen (1986). Later studies extended this proposal to model the *phase distinction* between compile- and run-time (Harper et al., 1990), and to account for type abstraction as well as type classes (Harper and Lillibridge, 1994; Leroy, 1994). The avoidance problem was first isolated by Castagna and Pierce (1994) and by Harper and Lillibridge (1994). It has come to play a central role in subsequent work on modules, such as Lillibridge (1997) and Dreyer (2005). The self-recognition rule was introduced by Harper and Lillibridge (1994) and by Leroy (1994). That rule was later identified as a manifestation of higher-order singletons (Stone and Harper, 2006). A consolidation of these ideas is used as the foundation for a mechanization of the meta-theory of modules (Lee et al., 2007). A thorough summary of the main issues in module system design is given in Dreyer (2005).

The presentation given here focuses attention on the type structure required to support modularity. An alternative formulation uses *elaboration*, a translation of modularity constructs into more primitive notions such as polymorphism and higher-order functions. *The Definition of Standard ML* (Milner et al., 1997) pioneered the elaboration approach. Building on earlier work of Russo, a more rigorous type-theoretic formulation was given by Rossberg et al. (2010). The advantage of the elaboration-based approach is that it can make do with a simpler type theory as the target language, but at the expense of making the explanation of modularity more complex.

Exercises

44.1. Consider the type abstraction σ_{set} of *finite sets* of elements of type τ_{elt} given by the following equations:

$$\begin{split} \sigma_{\mathsf{set}} &\triangleq \llbracket t :: \mathtt{T} \, ; \tau_{\mathsf{set}} \rrbracket \\ \tau_{\mathsf{set}} &\triangleq \langle \mathsf{emp} \hookrightarrow t \, , \mathsf{ins} \hookrightarrow \tau_{\mathsf{elt}} \times t \to t \, , \mathsf{mem} \hookrightarrow \tau_{\mathsf{elt}} \times t \to \mathsf{bool} \rangle. \end{split}$$

Define an implementation

$$\Gamma$$
, $D: \sigma_{\mathsf{dict}} \vdash M_{\mathsf{set}}: \sigma_{\mathsf{set}}$

430 44.5 Notes

of finite sets of elements in terms of a dictionary whose key and value types are chosen appropriately.

44.2. Fix an ordered type τ_{nod} of *nodes*, and consider the type abstraction σ_{grph} of *finite graphs* given by the following equations:

$$\begin{split} &\sigma_{\text{grph}} \triangleq \ \big[\!\!\big[t_{\text{grph}} :: T \, ; \, \big[\!\!\big[t_{\text{edg}} :: S(\tau_{\text{edg}}) \, ; \tau_{\text{grph}} \big]\!\!\big] \big] \\ &\tau_{\text{edg}} \triangleq \tau_{\text{nod}} \times \tau_{\text{nod}} \\ &\tau_{\text{grph}} \triangleq \big\langle \text{emp} \hookrightarrow t_{\text{grph}} \, , \text{ins} \hookrightarrow \tau_{\text{edg}} \times t_{\text{grph}} \to t_{\text{grph}} \, , \text{mem} \hookrightarrow \tau_{\text{edg}} \times t_{\text{grph}} \to \text{bool} \big\rangle. \end{split}$$

The signature σ_{grph} is translucent, with both opaque and transparent type components: graphs themselves are abstract, but edges are pairs of nodes.

Define an implementation

$$N: \sigma_{\mathsf{ord}}, S: \sigma_{\mathsf{nodset}}, D: \sigma_{\mathsf{nodsetdict}} \vdash M_{\mathsf{grph}}: \sigma_{\mathsf{grph}}$$

in terms of an implementation of nodes, sets of nodes, and a dictionary mapping nodes to sets of nodes. Represent the graph by a dictionary assigning to each node the set of nodes incident upon it. Define the node type τ_{nod} to be the type $N \cdot s$, and choose the signatures of the set and dictionary abstractions appropriately in terms of this choice of node type.

- **44.3.** Define *signature modification*, a variant of kind modification defined in Exercise **43.3**, in which a definition of a constructor component can be imposed on a signature. Let P stand for a composition of static and dynamic projections of the form $\cdot d \dots \cdot d \cdot s$, so that $X \cdot P$ stands for $X \cdot d \dots \cdot d \cdot s$. Assume that $\Gamma \vdash \sigma \operatorname{sig}$, Γ , $X : \sigma \vdash X \cdot P :: \kappa$, and $\Gamma \vdash c :: \kappa$. Define signature $\sigma\{P := c\}$ such that $\Gamma \vdash \sigma\{P := c\} <: \sigma \text{ and } \Gamma$, $X : \sigma\{P := c\} \vdash X \cdot P \equiv c :: \kappa$.
- **44.4**. The signature σ_{grph} is a subsignature (instance) of the type class

$$\sigma_{\mathsf{grphcls}} \triangleq \llbracket t_{\mathsf{grph}} :: \mathtt{T} \, ; \, \llbracket t_{\mathsf{edg}} :: \mathtt{T} \, ; \tau_{\mathsf{grph}} \rrbracket \rrbracket$$

in which the definition of t_{edg} has been made explicit as the product of two nodes.

Check that $\Gamma \vdash \sigma_{grph} \equiv \sigma_{grphcls} \{ \cdot d \cdot s := \tau_{nod} \times \tau_{nod} \}$, so that the former can be defined as the latter.

Part XVIII Equational Reasoning



Chapter 46

Equality for System T

The beauty of functional programming is that equality of expressions in a functional language follows familiar patterns of mathematical reasoning. For example, in the language **T** of Chapter 9 in which we can express addition as the function plus, the expressions

$$\lambda(x: \mathtt{nat})\lambda(y: \mathtt{nat})$$
 plus $(x)(y)$

and

$$\lambda(x:nat)\lambda(y:nat)plus(y)(x)$$

are equal. In other words, the addition function as programmed in **T** is commutative.

Commutativity of addition may seem self-evident, but *why* is it true? What does it mean for two expressions to be equal? These two expressions are not *definitionally* equal; their equality requires proof, and is not merely a matter of calculation. Yet the two expressions are interchangeable because they give the same result when applied to the same number. In general two functions are *logically equivalent* if they give equal results for equal arguments. Because this is all that matters about a function, we may expect that logically equivalent functions are interchangeable in any program. Thinking of the programs in which these functions occur as *observations* of their behavior, these functions are said to be *observationally equivalent*. The main result of this chapter is that observational and logical equivalence coincide for a variant of **T** in which the successor is evaluated eagerly, so that a value of type nat is a numeral.

46.1 Observational Equivalence

When are two expressions equal? Whenever we cannot tell them apart! It may seem tautological to say so, but it is not, because it all depends on what we consider to be a means of telling expressions apart. What "experiment" are we permitted to perform on expressions in order to distinguish them? What counts as an observation that, if different for two expressions, is a sure sign that they are different?

If we permit ourselves to consider the syntactic details of the expressions, then very few expressions could be considered equal. For example, if it is significant that an expression contains, say, more than one function application, or that it has an occurrence of λ -abstraction, then very few expressions would come out as equivalent. But such considerations seem silly, because they conflict with the intuition that the significance of an expression lies in its contribution to the *outcome* of a computation, and not to the process of obtaining that outcome. In short, if two expressions make the same contribution to the outcome of a complete program, then they ought to be regarded as equal.

We must fix what we mean by a complete program. Two considerations inform the definition. First, the dynamics of **T** is defined only for expressions without free variables, so a complete program should clearly be a *closed* expression. Second, the outcome of a computation should be *observable*, so that it is evident whether the outcome of two computations differs or not. We define a *complete program* to be a closed expression of type nat, and define the *observable behavior* of the program to be the numeral to which it evaluates.

An *experiment* on, or *observation* about, an expression is any means of using that expression within a complete program. We define an *expression context* to be an expression with a "hole" in it serving as a place-holder for another expression. The hole is permitted to occur anywhere, including within the scope of a binder. The bound variables within whose scope the hole lies are *exposed to capture* by the expression context. A *program context* is a closed expression context of type nat—that is, it is a complete program with a hole in it. The meta-variable $\mathcal C$ stands for any expression context.

Replacement is the process of filling a hole in an expression context \mathcal{C} with an expression e which is written $\mathcal{C}\{e\}$. Importantly, the free variables of e that are exposed by \mathcal{C} are *captured* by replacement (which is why replacement is not a form of substitution, which is defined so as to avoid capture). If \mathcal{C} is a program context, then $\mathcal{C}\{e\}$ is a complete program iff all free variables of e are captured by the replacement. For example, if $\mathcal{C} = \lambda$ (x: nat) \circ , and e = x + x, then

$$C\{e\} = \lambda (x : \mathtt{nat}) x + x.$$

The free occurrences of x in e are captured by the λ -abstraction as a result of the replacement of the hole in \mathcal{C} by e.

We sometimes write $C\{\circ\}$ to emphasize the occurrence of the hole in C. Expression contexts are closed under *composition* in that if C_1 and C_2 are expression contexts, then so is

$$\mathcal{C}\{\circ\} \triangleq \mathcal{C}_1\{\mathcal{C}_2\{\circ\}\},\,$$

and we have $C\{e\} = C_1\{C_2\{e\}\}\$. The *trivial*, or *identity*, expression context is the "bare hole", written \circ , for which $\circ\{e\} = e$.

The statics of expressions of ${\bf T}$ is extended to expression contexts by defining the typing judgment

$$\mathcal{C}: (\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \tau')$$

so that if $\Gamma \vdash e : \tau$, then $\Gamma' \vdash \mathcal{C}\{e\} : \tau'$. This judgment is inductively defined by a collection of rules derived from the statics of **T** (see rules (9.1)). Some representative rules are as follows:

$$\overline{\circ : (\Gamma \triangleright \tau) \leadsto (\Gamma \triangleright \tau)} \tag{46.1a}$$

$$\frac{\mathcal{C}: (\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \mathtt{nat})}{\mathtt{s}(\mathcal{C}): (\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \mathtt{nat})} \tag{46.1b}$$

$$\frac{\mathcal{C}: (\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \mathtt{nat}) \quad \Gamma' \vdash e_0 : \tau' \quad \Gamma', x : \mathtt{nat}, y : \tau' \vdash e_1 : \tau'}{\mathtt{rec} \, \mathcal{C} \, \{ \mathtt{z} \hookrightarrow e_0 \, | \, \mathtt{s}(\, x \,) \, \mathtt{with} \, y \hookrightarrow e_1 \} : (\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \tau')} \tag{46.1c}$$

$$\frac{\Gamma' \vdash e : \mathtt{nat} \quad \mathcal{C}_0 : (\Gamma \rhd \tau) \leadsto (\Gamma' \rhd \tau') \quad \Gamma', x : \mathtt{nat}, y : \tau' \vdash e_1 : \tau'}{\mathtt{rec} \, e \, \{ z \hookrightarrow \mathcal{C}_0 \, | \, \mathsf{s}(x) \, \mathtt{with} \, y \hookrightarrow e_1 \} : (\Gamma \rhd \tau) \leadsto (\Gamma' \rhd \tau')} \tag{46.1d}$$

$$\frac{\Gamma' \vdash e : \mathtt{nat} \quad \Gamma' \vdash e_0 : \tau' \quad \mathcal{C}_1 : (\Gamma \rhd \tau) \leadsto (\Gamma', x : \mathtt{nat}, y : \tau' \rhd \tau')}{\mathtt{rec} \, e \, \{ \mathbf{z} \hookrightarrow e_0 \, | \, \mathbf{s}(x) \, \mathtt{with} \, y \hookrightarrow \mathcal{C}_1 \} : (\Gamma \rhd \tau) \leadsto (\Gamma' \rhd \tau')} \tag{46.1e}$$

$$\frac{C_2: (\Gamma \triangleright \tau) \leadsto (\Gamma', x: \tau_1 \triangleright \tau_2)}{\lambda (x: \tau_1) C_2: (\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \tau_1 \to \tau_2)}$$

$$(46.1f)$$

$$\frac{\mathcal{C}_1: (\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \tau_2 \to \tau') \quad \Gamma' \vdash e_2 : \tau_2}{\mathcal{C}_1(e_2): (\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \tau')}$$
(46.1g)

$$\frac{\Gamma' \vdash e_1 : \tau_2 \to \tau' \quad \mathcal{C}_2 : (\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \tau_2)}{e_1(\mathcal{C}_2) : (\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \tau')}$$

$$(46.1h)$$

Lemma 46.1. *If* $C : (\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \tau')$, then $\Gamma' \subseteq \Gamma$, and if $\Gamma \vdash e : \tau$, then $\Gamma' \vdash C\{e\} : \tau'$.

Contexts are closed under composition, with the trivial context acting as an identity for it.

Lemma 46.2. *If*
$$C : (\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \tau')$$
, and $C' : (\Gamma' \triangleright \tau') \leadsto (\Gamma'' \triangleright \tau'')$, then $C'\{C\{\circ\}\} : (\Gamma \triangleright \tau) \leadsto (\Gamma'' \triangleright \tau'')$. **Lemma 46.3.** *If* $C : (\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \tau')$ and $x \notin dom(\Gamma)$, then $C : (\Gamma, x : \tau'' \triangleright \tau) \leadsto (\Gamma', x : \tau'' \triangleright \tau')$. *Proof.* By induction on rules (46.1).

A *complete program* is a closed expression of type nat.

Definition 46.4. Two complete programs, e and e', are Kleene equal, written $e \simeq e'$, iff there exists $n \ge 0$ such that $e \longmapsto^* \overline{n}$ and $e' \longmapsto^* \overline{n}$.

Kleene equality is obviously symmetric; its transitivity follows from determinacy of evaluation. Closure under converse evaluation follows similarly. It is immediate from the definition that $\bar{0} \not\simeq \bar{1}$.

Definition 46.5. Suppose that $\Gamma \vdash e : \tau$ and $\Gamma \vdash e' : \tau$ are two expressions of the same type. Two such expressions are observationally equivalent, written $\Gamma \vdash e \cong e' : \tau$, iff $C\{e\} \simeq C\{e'\}$ for every program context $C : (\Gamma \triangleright \tau) \leadsto (\emptyset \triangleright \mathtt{nat})$.

In other words, for all possible experiments, the outcome of an experiment on e is the same as the outcome on e', which is an equivalence relation. For the sake of brevity, we often write $e \cong_{\tau} e'$ for $\emptyset \vdash e \cong e' : \tau$.

A family of relations $\Gamma \vdash e_1 \mathcal{E} e_2 : \tau$ is a *congruence* iff it is preserved by all contexts. That is,

if
$$\Gamma \vdash e \ \mathcal{E} \ e' : \tau$$
, then $\Gamma' \vdash \mathcal{C}\{e\} \ \mathcal{E} \ \mathcal{C}\{e'\} : \tau'$

for every expression context $\mathcal{C}:(\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \tau')$. Such a family of relations is *consistent* iff $\emptyset \vdash e \ \mathcal{E} \ e'$: nat implies $e \simeq e'$.

Theorem 46.6. Observational equivalence is the coarsest consistent congruence on expressions.

Proof. Consistency follows from the definition by noting that the trivial context is a program context. Observational equivalence is clearly an equivalence relation. To show that it is a congruence, we need only observe that type-correct composition of a program context with an arbitrary expression context is again a program context. Finally, it is the coarsest such equivalence relation, for if $\Gamma \vdash e \ \mathcal{E} \ e' : \tau$ for some consistent congruence \mathcal{E} , and if $\mathcal{C} : (\Gamma \triangleright \tau) \leadsto (\emptyset \triangleright \text{nat})$, then by congruence $\emptyset \vdash \mathcal{C}\{e\} \ \mathcal{E} \ \mathcal{C}\{e'\}$: nat, and hence by consistency $\mathcal{C}\{e\} \simeq \mathcal{C}\{e'\}$.

A closing substitution γ for the typing context $\Gamma = x_1 : \tau_1, \ldots, x_n : \tau_n$ is a finite function assigning closed expressions $e_1 : \tau_1, \ldots, e_n : \tau_n$ to x_1, \ldots, x_n , respectively. We write $\hat{\gamma}(e)$ for the substitution $[e_1, \ldots, e_n/x_1, \ldots, x_n]e$, and write $\gamma : \Gamma$ to mean that if $x : \tau$ occurs in Γ , then there exists a closed expression e such that $\gamma(x) = e$ and $e : \tau$. We write $\gamma \cong_{\Gamma} \gamma'$, where $\gamma : \Gamma$ and $\gamma' : \Gamma$, to express that $\gamma(x) \cong_{\Gamma(x)} \gamma'(x)$ for each x declared in Γ .

Lemma 46.7. If $\Gamma \vdash e \cong e' : \tau$ and $\gamma : \Gamma$, then $\hat{\gamma}(e) \cong_{\tau} \hat{\gamma}(e')$. Moreover, if $\gamma \cong_{\Gamma} \gamma'$, then $\hat{\gamma}(e) \cong_{\tau} \widehat{\gamma'}(e)$ and $\hat{\gamma}(e') \cong_{\tau} \widehat{\gamma'}(e')$.

Proof. Let $\mathcal{C}: (\emptyset \triangleright \tau) \leadsto (\emptyset \triangleright \text{nat})$ be a program context; we are to show that $\mathcal{C}\{\hat{\gamma}(e)\} \simeq \mathcal{C}\{\hat{\gamma}(e')\}$. Because \mathcal{C} has no free variables, this is equivalent to showing that $\hat{\gamma}(\mathcal{C}\{e\}) \simeq \hat{\gamma}(\mathcal{C}\{e'\})$. Let \mathcal{D} be the context

$$\lambda(x_1:\tau_1)\ldots\lambda(x_n:\tau_n)\mathcal{C}\{\circ\}(e_1)\ldots(e_n),$$

where $\Gamma = x_1 : \tau_1, \ldots, x_n : \tau_n$ and $\gamma(x_1) = e_1, \ldots, \gamma(x_n) = e_n$. By Lemma 46.3 we have $\mathcal{C} : (\Gamma \triangleright \tau) \leadsto (\Gamma \triangleright \text{nat})$, from which it follows that $\mathcal{D} : (\Gamma \triangleright \tau) \leadsto (\emptyset \triangleright \text{nat})$. Because $\Gamma \vdash e \cong e' : \tau$, we have $\mathcal{D}\{e\} \simeq \mathcal{D}\{e'\}$. But by construction $\mathcal{D}\{e\} \simeq \hat{\gamma}(\mathcal{C}\{e\})$, and $\mathcal{D}\{e'\} \simeq \hat{\gamma}(\mathcal{C}\{e'\})$, so $\hat{\gamma}(\mathcal{C}\{e\}) \simeq \hat{\gamma}(\mathcal{C}\{e'\})$. Because \mathcal{C} is arbitrary, it follows that $\hat{\gamma}(e) \cong_{\tau} \hat{\gamma}(e')$.

Defining \mathcal{D}' like \mathcal{D} , but based on γ' , rather than γ , we may also show that $\mathcal{D}'\{e\} \simeq \mathcal{D}'\{e'\}$, and hence $\widehat{\gamma}'(e) \cong_{\tau} \widehat{\gamma}'(e')$. Now if $\gamma \cong_{\Gamma} \gamma'$, then by congruence we have $\mathcal{D}\{e\} \cong_{\mathtt{nat}} \mathcal{D}'\{e\}$, and $\mathcal{D}\{e'\} \cong_{\mathtt{nat}} \mathcal{D}'\{e'\}$. It follows that $\mathcal{D}\{e\} \cong_{\mathtt{nat}} \mathcal{D}'\{e'\}$, and so, by consistency of observational equivalence, we have $\mathcal{D}\{e\} \simeq \mathcal{D}'\{e'\}$, which is to say that $\widehat{\gamma}(e) \cong_{\tau} \widehat{\gamma}'(e')$.

Theorem 46.6 licenses the principle of *proof by coinduction*: to show that $\Gamma \vdash e \cong e' : \tau$, it is enough to exhibit a consistent congruence \mathcal{E} such that $\Gamma \vdash e \mathcal{E} e' : \tau$. It can be difficult to construct such a relation. In the next section we will provide a general method for doing so that exploits types.

46.2 Logical Equivalence

The key to simplifying reasoning about observational equivalence is to exploit types. Informally, we may classify the uses of expressions of a type into two broad categories, the *passive* and the *active* uses. The passive uses are those that manipulate expressions without inspecting them. For example, we may pass an expression of type τ to a function that simply returns it. The active uses are those that operate on the expression itself; these are the elimination forms associated with the type of that expression. For the purposes of distinguishing two expressions, it is only

the active uses that matter; the passive uses manipulate expressions at arm's length, affording no opportunities to distinguish one from another.

Logical equivalence is therefore defined as follows.

Definition 46.8. Logical equivalence is a family of relations $e \sim_{\tau} e'$ between closed expressions of type τ . It is defined by induction on τ as follows:

$$e \sim_{\mathtt{nat}} e'$$
 iff $e \simeq e'$
$$e \sim_{\tau_1 \to \tau_2} e'$$
 iff if $e_1 \sim_{\tau_1} e'_1$, then $e(e_1) \sim_{\tau_2} e'(e'_1)$

The definition of logical equivalence at type nat licenses the following principle of *proof by* nat-*induction*. To show that $\mathcal{E}(e,e')$ whenever $e \sim_{\mathtt{nat}} e'$, it is enough to show that

- 1. $\mathcal{E}(\overline{0},\overline{0})$, and
- 2. if $\mathcal{E}(\overline{n}, \overline{n})$, then $\mathcal{E}(\overline{n+1}, \overline{n+1})$.

This assertion is justified by mathematical induction on $n \ge 0$, where $e \mapsto^* \overline{n}$ and $e' \mapsto^* \overline{n}$ by the definition of Kleene equivalence.

Lemma 46.9. Logical equivalence is symmetric and transitive: if $e \sim_{\tau} e'$, then $e' \sim_{\tau} e$, and if $e \sim_{\tau} e'$ and $e' \sim_{\tau} e''$, then $e \sim_{\tau} e''$.

Proof. Simultaneously, by induction on the structure of τ . If $\tau=$ nat, the result is immediate. If $\tau=\tau_1\to\tau_2$, then we may assume that logical equivalence is symmetric and transitive at types τ_1 and τ_2 . For symmetry, assume that $e\sim_{\tau}e'$; we wish to show $e'\sim_{\tau}e$. Assume that $e'_1\sim_{\tau_1}e_1$; it suffices to show that $e'(e'_1)\sim_{\tau_2}e(e_1)$. By induction we have that $e_1\sim_{\tau_1}e'_1$. Therefore by assumption $e(e_1)\sim_{\tau_2}e'(e'_1)$, and hence by induction $e'(e'_1)\sim_{\tau_2}e(e_1)$. For transitivity, assume that $e\sim_{\tau}e'$ and $e'\sim_{\tau}e''$; we are to show $e\sim_{\tau}e''$. Suppose that $e_1\sim_{\tau_1}e''_1$; it is enough to show that $e(e_1)\sim_{\tau}e''(e''_1)$. By symmetry and transitivity we have $e_1\sim_{\tau_1}e_1$, so by assumption $e(e_1)\sim_{\tau_2}e''(e'_1)$. We also have by assumption $e'(e_1)\sim_{\tau_2}e''(e''_1)$. By transitivity we have $e'(e_1)\sim_{\tau_2}e''(e''_1)$, which suffices for the result.

Logical equivalence is extended to open terms by substitution of related closed terms to obtain related results. If γ and γ' are two substitutions for Γ , we define $\gamma \sim_{\Gamma} \gamma'$ to hold iff $\gamma(x) \sim_{\Gamma(x)} \gamma'(x)$ for every variable, x, such that $\Gamma \vdash x : \tau$. Open logical equivalence, written $\Gamma \vdash e \sim e' : \tau$, is defined to mean that $\hat{\gamma}(e) \sim_{\tau} \hat{\gamma'}(e')$ whenever $\gamma \sim_{\Gamma} \gamma'$.

Lemma 46.10. Open logical equivalence is symmetric and transitive.

Proof. Follows from Lemma 46.9 and the definition of open logical equivalence.

At this point we are "two thirds of the way" to justifying the use of the name "open logical equivalence." The remaining third, reflexivity, is established in the next section.

46.3 Logical and Observational Equivalence Coincide

In this section we prove the coincidence of observational and logical equivalence.

Lemma 46.11 (Converse Evaluation). *Suppose that* $e \sim_{\tau} e'$. *If* $d \mapsto e$, then $d \sim_{\tau} e'$, and if $d' \mapsto e'$, then $e \sim_{\tau} d'$.

Proof. By induction on the structure of τ . If $\tau=$ nat, then the result follows from the closure of Kleene equivalence under converse evaluation. If $\tau=\tau_1\to\tau_2$, then suppose that $e\sim_{\tau}e'$, and $d\longmapsto e$. To show that $d\sim_{\tau}e'$, we assume $e_1\sim_{\tau_1}e'_1$ and show $d(e_1)\sim_{\tau_2}e'(e'_1)$. It follows from the assumption that $e(e_1)\sim_{\tau_2}e'(e'_1)$. Noting that $d(e_1)\longmapsto e(e_1)$, the result follows by induction.

Lemma 46.12 (Consistency). *If* $e \sim_{nat} e'$, then $e \simeq e'$.

Proof. Immediate, from Definition 46.8.

Theorem 46.13 (Reflexivity). *If* $\Gamma \vdash e : \tau$, then $\Gamma \vdash e \sim e : \tau$.

Proof. We are to show that if $\Gamma \vdash e : \tau$ and $\gamma \sim_{\Gamma} \gamma'$, then $\hat{\gamma}(e) \sim_{\tau} \hat{\gamma'}(e)$. The proof proceeds by induction on typing derivations; we consider two representative cases.

Consider the case of rule (8.4a), in which $\tau = \tau_1 \to \tau_2$ and $e = \lambda (x : \tau_1) e_2$. We are to show that

$$\lambda(x:\tau_1) \hat{\gamma}(e_2) \sim_{\tau_1 \to \tau_2} \lambda(x:\tau_1) \hat{\gamma}'(e_2).$$

Assume that $e_1 \sim_{\tau_1} e_1'$; by Lemma 46.11, it is enough to show that $[e_1/x]\hat{\gamma}(e_2) \sim_{\tau_2} [e_1'/x]\hat{\gamma}'(e_2)$. Let $\gamma_2 = \gamma \otimes x \hookrightarrow e_1$ and $\gamma_2' = \gamma' \otimes x \hookrightarrow e_1'$, and observe that $\gamma_2 \sim_{\Gamma,x:\tau_1} \gamma_2'$. Therefore, by induction we have $\hat{\gamma}_2(e_2) \sim_{\tau_2} \hat{\gamma}_2'(e_2)$, from which the result follows easily.

Now consider the case of rule (9.1d), for which we are to show that

$$\operatorname{\mathtt{rec}}\{\hat{\gamma}(e_0); x.y.\hat{\gamma}(e_1)\}(\,\hat{\gamma}(e)\,) \sim_{\tau} \operatorname{\mathtt{rec}}\{\hat{\gamma'}(e_0); x.y.\hat{\gamma'}(e_1)\}(\,\hat{\gamma'}(e)\,).$$

By the induction hypothesis applied to the first premise of rule (9.1d), we have

$$\hat{\gamma}(e) \sim_{\mathtt{nat}} \widehat{\gamma'}(e).$$

We proceed by nat-induction. It suffices to show that

$$rec\{\hat{\gamma}(e_0); x.y.\hat{\gamma}(e_1)\}(z) \sim_{\tau} rec\{\hat{\gamma'}(e_0); x.y.\hat{\gamma'}(e_1)\}(z), \tag{46.2}$$

and that

$$\operatorname{rec}\{\hat{\gamma}(e_0); x.y. \hat{\gamma}(e_1)\}(\operatorname{s}(\overline{n})) \sim_{\tau} \operatorname{rec}\{\hat{\gamma'}(e_0); x.y. \hat{\gamma'}(e_1)\}(\operatorname{s}(\overline{n})), \tag{46.3}$$

assuming

$$\operatorname{rec}\{\hat{\gamma}(e_0); x.y.\hat{\gamma}(e_1)\}(\overline{n}) \sim_{\tau} \operatorname{rec}\{\hat{\gamma'}(e_0); x.y.\hat{\gamma'}(e_1)\}(\overline{n}). \tag{46.4}$$

To show (46.2), by Lemma 46.11 it is enough to show that $\hat{\gamma}(e_0) \sim_{\tau} \hat{\gamma'}(e_0)$. This condition is assured by the outer inductive hypothesis applied to the second premise of rule (9.1d).

To show (46.3), define

$$\delta = \gamma \otimes x \hookrightarrow \overline{n} \otimes y \hookrightarrow \operatorname{rec}\{\hat{\gamma}(e_0); x.y.\hat{\gamma}(e_1)\}(\overline{n})$$

and

$$\delta' = \gamma' \otimes x \hookrightarrow \overline{n} \otimes y \hookrightarrow \operatorname{rec}\{\widehat{\gamma'}(e_0); x.y.\widehat{\gamma'}(e_1)\}(\overline{n}).$$

By (46.4) we have $\delta \sim_{\Gamma,x:\mathtt{nat},y:\tau} \delta'$. Consequently, by the outer inductive hypothesis applied to the third premise of rule (9.1d), and Lemma 46.11, the required follows.

Corollary 46.14 (Equivalence). Open logical equivalence is an equivalence relation.

Corollary 46.15 (Termination). *If* e : nat, then $e \mapsto^* e'$ for some e' val.

Lemma 46.16 (Congruence). *If* $C_0 : (\Gamma \triangleright \tau) \rightsquigarrow (\Gamma_0 \triangleright \tau_0)$, and $\Gamma \vdash e \sim e' : \tau$, then $\Gamma_0 \vdash C_0\{e\} \sim C_0\{e'\} : \tau_0$.

Proof. By induction on the derivation of the typing of \mathcal{C}_0 . We consider a representative case in which $\mathcal{C}_0 = \lambda \ (x : \tau_1) \ \mathcal{C}_2$ so that $\mathcal{C}_0 : (\Gamma \triangleright \tau) \leadsto (\Gamma_0 \triangleright \tau_1 \to \tau_2)$ and $\mathcal{C}_2 : (\Gamma \triangleright \tau) \leadsto (\Gamma_0, x : \tau_1 \triangleright \tau_2)$. Assuming $\Gamma \vdash e \sim e' : \tau$, we are to show that

$$\Gamma_0 \vdash \mathcal{C}_0\{e\} \sim \mathcal{C}_0\{e'\} : \tau_1 \to \tau_2$$
,

which is to say

$$\Gamma_0 \vdash \lambda (x : \tau_1) C_2 \{e\} \sim \lambda (x : \tau_1) C_2 \{e'\} : \tau_1 \rightarrow \tau_2.$$

We know, by induction, that

$$\Gamma_0, x : \tau_1 \vdash C_2\{e\} \sim C_2\{e'\} : \tau_2.$$

Suppose that $\gamma_0 \sim_{\Gamma_0} \gamma_0'$, and that $e_1 \sim_{\tau_1} e_1'$. Let $\gamma_1 = \gamma_0 \otimes x \hookrightarrow e_1$, $\gamma_1' = \gamma_0' \otimes x \hookrightarrow e_1'$, and observe that $\gamma_1 \sim_{\Gamma_0, x: \tau_1} \gamma_1'$. By Definition 46.8 and Lemma 46.11 it is enough to show that

$$\hat{\gamma}_1(\mathcal{C}_2\{e\}) \sim_{\tau_2} \hat{\gamma}'_1(\mathcal{C}_2\{e'\}),$$

which follows from the inductive hypothesis.

Theorem 46.17. *If* $\Gamma \vdash e \sim e' : \tau$, then $\Gamma \vdash e \cong e' : \tau$.

Proof. By Lemmas 46.12 and 46.16, and Theorem 46.6.

Corollary 46.18. *If* e : nat, then $e \cong_{\text{nat}} \overline{n}$, for some $n \geq 0$.

Proof. By Theorem 46.13 we have $e \sim_{\mathtt{nat}} e$. Hence for some $n \geq 0$, we have $e \sim_{\mathtt{nat}} \overline{n}$, and so by Theorem 46.17, $e \cong_{\mathtt{nat}} \overline{n}$.

Lemma 46.19. For closed expressions $e : \tau$ and $e' : \tau$, if $e \cong_{\tau} e'$, then $e \sim_{\tau} e'$.

Proof. We proceed by induction on the structure of τ . If $\tau = \text{nat}$, consider the empty context to obtain $e \simeq e'$, and hence $e \sim_{\text{nat}} e'$. If $\tau = \tau_1 \to \tau_2$, then we are to show that whenever $e_1 \sim_{\tau_1} e'_1$, we have $e(e_1) \sim_{\tau_2} e'(e'_1)$. By Theorem 46.17 we have $e_1 \cong_{\tau_1} e'_1$, and hence by congruence of observational equivalence it follows that $e(e_1) \cong_{\tau_2} e'(e'_1)$, from which the result follows by induction

Theorem 46.20. *If* $\Gamma \vdash e \cong e' : \tau$, then $\Gamma \vdash e \sim e' : \tau$.

Proof. Assume that $\Gamma \vdash e \cong e' : \tau$, and that $\gamma \sim_{\Gamma} \gamma'$. By Theorem 46.17 we have $\gamma \cong_{\Gamma} \gamma'$, so by Lemma 46.7 $\hat{\gamma}(e) \cong_{\tau} \hat{\gamma'}(e')$. Therefore, by Lemma 46.19, $\hat{\gamma}(e) \sim_{\tau} \hat{\gamma}(e')$.

Corollary 46.21. $\Gamma \vdash e \cong e' : \tau \text{ iff } \Gamma \vdash e \sim e' : \tau.$

Definitional equality is sufficient for observational equivalence:

Theorem 46.22. *If* $\Gamma \vdash e \equiv e' : \tau$, then $\Gamma \vdash e \sim e' : \tau$, and hence $\Gamma \vdash e \cong e' : \tau$.

Proof. By an argument similar to that used in the proof of Theorem 46.13 and Lemma 46.16, then appealing to Theorem 46.17.

Corollary 46.23. *If* $e \equiv e'$: nat, then there exists $n \geq 0$ such that $e \mapsto^* \overline{n}$ and $e' \mapsto^* \overline{n}$.

Proof. By Theorem 46.22 we have $e \sim_{\mathtt{nat}} e'$ and hence $e \simeq e'$.

46.4 Some Laws of Equality

In this section we summarize some useful principles of observational equivalence for \mathbf{T} . For the most part these are laws of logical equivalence, and then transferred to observational equivalence by appeal to Corollary 46.21. The laws are presented as inference rules with the meaning that if all of the premises are true judgments about observational equivalence, then so are the conclusions. In other words each rule is admissible as a principle of observational equivalence.

46.4.1 General Laws

Logical equivalence is indeed an equivalence relation: it is reflexive, symmetric, and transitive.

$$\overline{\Gamma \vdash e \cong e : \tau} \tag{46.5a}$$

$$\frac{\Gamma \vdash e' \cong e : \tau}{\Gamma \vdash e \cong e' : \tau} \tag{46.5b}$$

$$\frac{\Gamma \vdash e \cong e' : \tau \quad \Gamma \vdash e' \cong e'' : \tau}{\Gamma \vdash e \cong e'' : \tau}$$
(46.5c)

Reflexivity is an instance of a more general principle, that all definitional equalities are observational equivalences.

$$\frac{\Gamma \vdash e \equiv e' : \tau}{\Gamma \vdash e \cong e' : \tau} \tag{46.6a}$$

Observational equivalence is a congruence: we may replace equals by equals anywhere in an expression.

$$\frac{\Gamma \vdash e \cong e' : \tau \quad \mathcal{C} : (\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \tau')}{\Gamma' \vdash \mathcal{C}\{e\} \cong \mathcal{C}\{e'\} : \tau'}$$
(46.7a)

Equivalence is stable under substitution for free variables, and substituting equivalent expressions in an expression gives equivalent results.

$$\frac{\Gamma \vdash e : \tau \quad \Gamma, x : \tau \vdash e_2 \cong e'_2 : \tau'}{\Gamma \vdash [e/x]e_2 \cong [e/x]e'_2 : \tau'}$$
(46.8a)

$$\frac{\Gamma \vdash e_1 \cong e'_1 : \tau \quad \Gamma, x : \tau \vdash e_2 \cong e'_2 : \tau'}{\Gamma \vdash [e_1/x]e_2 \cong [e'_1/x]e'_2 : \tau'}$$

$$(46.8b)$$

46.4.2 Equality Laws

Two functions are equal if they are equal on all arguments.

$$\frac{\Gamma, x : \tau_1 \vdash e(x) \cong e'(x) : \tau_2}{\Gamma \vdash e \cong e' : \tau_1 \to \tau_2}$$

$$(46.9)$$

Consequently, every expression of function type is equal to a λ -abstraction:

$$\Gamma \vdash e \cong \lambda (x : \tau_1) e(x) : \tau_1 \to \tau_2$$
(46.10)

46.4.3 Induction Law

An equation involving a free variable x of type nat can be proved by induction on x.

$$\frac{\Gamma \vdash [\overline{n}/x]e \cong [\overline{n}/x]e' : \tau \text{ (for every } n \in \mathbb{N})}{\Gamma, x : \text{nat} \vdash e \cong e' : \tau}$$
(46.11a)

To apply the induction rule, we proceed by mathematical induction on $n \in \mathbb{N}$, which reduces to showing:

- 1. $\Gamma \vdash [\mathbf{z}/x]e \cong [\mathbf{z}/x]e' : \tau$, and
- 2. $\Gamma \vdash [s(\overline{n})/x]e \cong [s(\overline{n})/x]e' : \tau$, if $\Gamma \vdash [\overline{n}/x]e \cong [\overline{n}/x]e' : \tau$.

452 46.5 Notes

46.5 Notes

The method of *logical relations* interprets types as relations (here, equivalence relations) by associating with each type constructor a relational action that transforms the relation interpreting its arguments to the relation interpreting the constructed type. Logical relations (Statman, 1985) are a fundamental tool in proof theory and provide the foundation for the semantics of the NuPRL type theory (Constable, 1986; Allen, 1987; Harper, 1992). The use of logical relations to characterize observational equivalence is an adaptation of the NuPRL semantics to the simpler setting of Gödel's System T.



Chapter 47

Equality for System PCF

In this Chapter we develop the theory of observational equivalence for **PCF**, with an eager interpretation of the type of natural numbers. The development proceeds along lines similar to those in Chapter 46, but is complicated by the presence of general recursion. The proof depends on the concept of an *admissible relation*, one that admits the principle of *proof by fixed point induction*.

47.1 Observational Equivalence

The definition of observational equivalence, along with the auxiliary notion of Kleene equivalence, are defined similarly to Chapter 46, but modified to account for the possibility of non-termination.

The collection of well-formed **PCF** contexts is inductively defined in a manner directly analogous to that in Chapter 46. Specifically, we define the judgment $\mathcal{C}: (\Gamma \triangleright \tau) \leadsto (\Gamma' \triangleright \tau')$ by rules similar to rules (46.1), modified for **PCF**. (We leave the precise definition as an exercise for the reader.) When Γ and Γ' are empty, we write just $\mathcal{C}: \tau \leadsto \tau'$.

A complete program is a closed expression of type nat.

Definition 47.1. We say that two complete programs, e and e', are Kleene equal, written $e \simeq e'$, iff for every $n \geq 0$, $e \longmapsto^* \overline{n}$ iff $e' \longmapsto^* \overline{n}$.

Kleene equality is clearly an equivalence relation and is closed under converse evaluation. Moreover, $\bar{0} \not\simeq \bar{1}$ and, if e and e' are both divergent, then $e \simeq e'$.

Observational equivalence is defined just as it is in Chapter 46.

Definition 47.2. *We say that* $\Gamma \vdash e : \tau$ *and* $\Gamma \vdash e' : \tau$ *are* observationally, *or* contextually, equivalent *iff for every program context* $C : (\Gamma \triangleright \tau) \leadsto (\emptyset \triangleright \text{nat}), C\{e\} \simeq C\{e'\}.$

Theorem 47.3. *Observational equivalence is the coarsest consistent congruence.*

Proof. See the proof of Theorem 46.6.

Lemma 47.4 (Substitution and Functionality). *If* $\Gamma \vdash e \cong e' : \tau$ *and* $\gamma : \Gamma$, *then* $\hat{\gamma}(e) \cong_{\tau} \hat{\gamma}(e')$. *Moreover, if* $\gamma \cong_{\Gamma} \gamma'$, *then* $\hat{\gamma}(e) \cong_{\tau} \hat{\gamma}'(e)$ *and* $\hat{\gamma}(e') \cong_{\tau} \hat{\gamma}'(e')$.

Proof. See Lemma 46.7.

47.2 Logical Equivalence

Definition 47.5. Logical equivalence, $e \sim_{\tau} e'$, between closed expressions of type τ is defined by induction on τ as follows:

$$e \sim_{\mathtt{nat}} e'$$
 iff $e \simeq e'$
$$e \sim_{\tau_1 \rightharpoonup \tau_2} e'$$
 iff $e_1 \sim_{\tau_1} e'_1$ implies $e(e_1) \sim_{\tau_2} e'(e'_1)$

Formally, logical equivalence is defined as in Chapter 46, except that the definition of Kleene equivalence is altered to account for non-termination. Logical equivalence is extended to open terms by substitution. Specifically, we define $\Gamma \vdash e \sim e'$: τ to mean that $\hat{\gamma}(e) \sim_{\tau} \widehat{\gamma'}(e')$ whenever $\gamma \sim_{\Gamma} \gamma'$.

By the same argument as given in the proof of Lemma 46.9 logical equivalence is symmetric and transitive, as is its open extension.

Lemma 47.6 (Strictness). *If* $e : \tau$ *and* $e' : \tau$ *are both divergent, then* $e \sim_{\tau} e'$.

Proof. By induction on the structure of τ . If $\tau = \text{nat}$, then the result follows immediately from the definition of Kleene equivalence. If $\tau = \tau_1 \rightharpoonup \tau_2$, then $e(e_1)$ and $e'(e'_1)$ diverge, so by induction $e(e_1) \sim_{\tau_2} e'(e'_1)$, as required.

Lemma 47.7 (Converse Evaluation). *Suppose that* $e \sim_{\tau} e'$. *If* $d \mapsto e$, then $d \sim_{\tau} e'$, and if $d' \mapsto e'$, then $e \sim_{\tau} d'$.

47.3 Logical and Observational Equivalence Coincide

The proof of coincidence of logical and observational equivalence relies on the concept of *bounded* recursion, which we define by induction on $m \ge 0$ as follows:

$$\begin{array}{l}
0 \\
\text{fix } x : \tau \text{ is } e \triangleq \text{ fix } x : \tau \text{ is } x \\
m+1 \\
\text{fix } x : \tau \text{ is } e \triangleq \begin{bmatrix} m \\ \text{fix } x : \tau \text{ is } e/x \end{bmatrix} e
\end{array}$$

When m=0, bounded recursion is defined to be a divergent expression of type τ . When m>0, bounded recursion is defined by unrolling the recursion m times by iterated substitution. Intuitively, the bounded recursive expression $\text{fix}^m x : \tau \text{ is } e$ is as good as $\text{fix} x : \tau \text{ is } e$ for up to m unrollings, after which it is divergent.

It is easy to check that the following rule is derivable for each $m \ge 0$:

$$\frac{\Gamma, x : \tau \vdash e : \tau}{\Gamma \vdash \text{fix}^m \{\tau\}(x.e) : \tau}$$
 (47.1a)

The proof is by induction on $m \ge 0$, and amounts to an iteration of the substitution lemma for the statics of **PCF**.

The key property of bounded recursion is the principle of fixed point induction, which permits reasoning about a recursive computation by induction on the number of unrollings required to reach a value. The proof relies on *compactness*, which will be stated and proved in Section 47.4 below.

Theorem 47.8 (Fixed Point Induction). *Suppose that* $x : \tau \vdash e : \tau$. *If*

$$(\forall m \geq 0) \operatorname{fix}^m x : \tau \operatorname{is} e \sim_{\tau} \operatorname{fix}^m x : \tau \operatorname{is} e',$$

then $fix x : \tau is e \sim_{\tau} fix x : \tau is e'$.

Proof. Define an *applicative context* \mathcal{A} to be either a hole, \circ , or an application of the form $\mathcal{A}(e)$, where \mathcal{A} is an applicative context. The typing judgment for applicative contexts, $\mathcal{A}: \tau_0 \leadsto \tau$, is a special case of the general typing judgment for contexts. Define logical equivalence of applicative contexts, $\mathcal{A} \sim \mathcal{A}': \tau_0 \leadsto \tau$, by induction on the structure of \mathcal{A} as follows:

1.
$$\circ \sim \circ : \tau_0 \leadsto \tau_0$$
;

2. if
$$A \sim A' : \tau_0 \leadsto \tau_2 \rightharpoonup \tau$$
 and $e_2 \sim_{\tau_2} e_2'$, then $A(e_2) \sim A'(e_2') : \tau_0 \leadsto \tau$.

We prove by induction on the structure of τ , if $A \sim A' : \tau_0 \leadsto \tau$ and

for every
$$m \ge 0$$
, $\mathcal{A}\{\text{fix } x : \tau_0 \text{ is } e\} \sim_{\tau} \mathcal{A}'\{\text{fix } x : \tau_0 \text{ is } e'\}$, (47.2)

then

$$\mathcal{A}\{\text{fix}\,x:\tau_0\,\text{is}\,e\}\sim_\tau\mathcal{A}'\{\text{fix}\,x:\tau_0\,\text{is}\,e'\}. \tag{47.3}$$

Choosing A = A' = 0 with $\tau_0 = \tau$ completes the proof.

If $\tau = \text{nat}$, then assume that $A \sim A^{\tau} : \tau_0 \rightsquigarrow \text{nat}$ and (47.2). By Definition 47.5, we are to show

$$\mathcal{A}\{\text{fix}\,x:\tau_0\,\text{is}\,e\}\simeq\mathcal{A}'\{\text{fix}\,x:\tau_0\,\text{is}\,e'\}.$$

By Corollary 47.17 there exists $m \ge 0$ such that

$$\mathcal{A}\{\mathtt{fix}\,x\!:\! au_0\,\mathtt{is}\,e\}\simeq\mathcal{A}\{\mathtt{fix}^m\,x\!:\! au_0\,\mathtt{is}\,e\}.$$

By (47.2) we have

$$\mathcal{A}\{ egin{aligned} egin{aligned} & \mathcal{A}\{ egin{aligned} & \mathbf{fix} \, x \colon au_0 \, f is \, e' \}. \end{aligned}$$

By Corollary 47.17

$$\mathcal{A}'\{ extstyle{fix}\ x: au_0\ extstyle{is}\ e'\}\simeq \mathcal{A}'\{ extstyle{fix}\ x: au_0\ extstyle{is}\ e'\}.$$

The result follows by transitivity of Kleene equivalence.

If $\tau = \tau_1 \rightharpoonup \tau_2$, then by Definition 47.5, it is enough to show

$$\mathcal{A}\{\operatorname{fix} x : \tau_0 \operatorname{is} e\}(e_1) \sim_{\tau_2} \mathcal{A}'\{\operatorname{fix} x : \tau_0 \operatorname{is} e'\}(e'_1)$$

whenever $e_1 \sim_{\tau_1} e_1'$. Let $\mathcal{A}_2 = \mathcal{A}(e_1)$ and $\mathcal{A}_2' = \mathcal{A}'(e_1')$. It follows from (47.2) that for every $m \geq 0$

$$\mathcal{A}_2\{\operatorname{fix} x: \tau_0 \operatorname{is} e\} \sim_{\tau_2} \mathcal{A}'_2\{\operatorname{fix} x: \tau_0 \operatorname{is} e'\}.$$

Noting that $A_2 \sim A_2' : \tau_0 \leadsto \tau_2$, we have by induction

$$\mathcal{A}_2\{\operatorname{fix} x: \tau_0 \operatorname{is} e\} \sim_{\tau_2} \mathcal{A}'_2\{\operatorname{fix} x: \tau_0 \operatorname{is} e'\},$$

as required.

Lemma 47.9 (Reflexivity). *If* $\Gamma \vdash e : \tau$, then $\Gamma \vdash e \sim e : \tau$.

Proof. The proof proceeds along the same lines as the proof of Theorem 46.13. The main difference is the treatment of general recursion, which is proved by fixed point induction. Consider rule (19.1g). Assuming $\gamma \sim_{\Gamma} \gamma'$, we are to show that

$$\operatorname{fix} x : \tau \operatorname{is} \widehat{\gamma}(e) \sim_{\tau} \operatorname{fix} x : \tau \operatorname{is} \widehat{\gamma'}(e).$$

By Theorem 47.8 it is enough to show that, for every $m \ge 0$,

$$\operatorname{fix}^{m} x : \tau \operatorname{is} \widehat{\gamma}(e) \sim_{\tau} \operatorname{fix}^{m} x : \tau \operatorname{is} \widehat{\gamma'}(e).$$

We proceed by an inner induction on m. When m=0 the result is immediate, because both sides of the desired equivalence diverge. Assuming the result for m, and applying Lemma 47.7, it is enough to show that $\hat{\gamma}(e_1) \sim_{\tau} \hat{\gamma'}(e'_1)$, where

$$e_1 = [fix x : \tau is \hat{\gamma}(e)/x] \hat{\gamma}(e), \text{ and}$$
 (47.4)

$$e_1' = [\widehat{\operatorname{fix}} x : \tau \operatorname{is} \widehat{\gamma'}(e)/x] \widehat{\gamma'}(e). \tag{47.5}$$

But this follows directly from the inner and outer inductive hypotheses. For by the outer inductive hypothesis, if

then

$$[\overset{m}{\text{fix}}\,x\,:\,\tau\,\text{is}\,\hat{\gamma}(e)/x]\hat{\gamma}(e)\sim_{\tau}[\overset{m}{\text{fix}}\,x\,:\,\tau\,\text{is}\,\hat{\gamma'}(e)/x]\hat{\gamma'}(e).$$

But the hypothesis holds by the inner inductive hypothesis, from which the result follows.

To handle the conditional if $z \in \{z \hookrightarrow e_0 \mid s(x) \hookrightarrow e_1\}$, we proceed by cases on whether e diverges, in which case the conditional is divergent and therefore self-related by Lemma 47.6, or e converges, in which case we can proceed by an inner mathematical induction on its value, appealing to the inductive hypotheses governing the branches of the conditional to complete the argument.

Symmetry and transitivity of eager logical equivalence are easily established by induction on types, noting that Kleene equivalence is symmetric and transitive. Eager logical equivalence is therefore an equivalence relation.

47.4 Compactness 457

Lemma 47.10 (Congruence). *If* $C_0 : (\Gamma \triangleright \tau) \leadsto (\Gamma_0 \triangleright \tau_0)$, and $\Gamma \vdash e \sim e' : \tau$, then $\Gamma_0 \vdash C_0\{e\} \sim C_0\{e'\} : \tau_0$.

Proof. By induction on the derivation of the typing of C_0 , following along similar lines to the proof of Lemma 47.9.

Logical equivalence is consistent, by definition. Consequently, it is contained in observational equivalence.

Theorem 47.11. *If* $\Gamma \vdash e \sim e' : \tau$, then $\Gamma \vdash e \cong e' : \tau$.

Proof. By consistency and congruence of logical equivalence.

Lemma 47.12. *If* $e \cong_{\tau} e'$, then $e \sim_{\tau} e'$.

Proof. By induction on the structure of τ . If $\tau=$ nat, then the result is immediate, because the empty expression context is a program context. If $\tau=\tau_1 \rightharpoonup \tau_2$, then suppose that $e_1 \sim_{\tau_1} e_1'$. We are to show that $e(e_1) \sim_{\tau_2} e'(e_1')$. By Theorem 47.11 $e_1 \cong_{\tau_1} e_1'$, and hence by Lemma 47.4 $e(e_1) \cong_{\tau_2} e'(e_1')$, from which the result follows by induction.

Theorem 47.13. *If* $\Gamma \vdash e \cong e' : \tau$, then $\Gamma \vdash e \sim e' : \tau$.

Proof. Assume that $\Gamma \vdash e \cong e' : \tau$. Suppose that $\gamma \sim_{\Gamma} \gamma'$. By Theorem 47.11 we have $\gamma \cong_{\Gamma} \gamma'$, and so by Lemma 47.4 we have

$$\hat{\gamma}(e) \cong_{\tau} \hat{\gamma'}(e').$$

Therefore by Lemma 47.12 we have

$$\hat{\gamma}(e) \sim_{\tau} \hat{\gamma'}(e')$$

Corollary 47.14. $\Gamma \vdash e \cong e' : \tau \text{ iff } \Gamma \vdash e \sim e' : \tau.$

47.4 Compactness

The principle of fixed point induction is derived from a critical property of **PCF**, called *compactness*. This property states that only finitely many unwindings of a fixed point expression are needed in a complete evaluation of a program. Although intuitively obvious (one cannot complete infinitely many recursive calls in a finite computation), it is rather tricky to state and prove rigorously.

The proof of compactness (Theorem 47.16) makes use of the stack machine for **PCF** defined in Chapter 28, augmented with the following transitions for bounded recursive expressions:

$$k \triangleright \text{fix}^0 x : \tau \text{ is } e \longmapsto k \triangleright \text{fix}^0 x : \tau \text{ is } e$$
 (47.6a)

$$\overline{k \triangleright \operatorname{fix}^{m+1} x : \tau \operatorname{is} e \longmapsto k \triangleright [\operatorname{fix}^m x : \tau \operatorname{is} e/x]e}$$
(47.6b)

It is not difficult to extend the proof of Corollary 28.4 to account for bounded recursion.

458 47.4 Compactness

To get a feel for what is involved in the compactness proof, consider first the factorial function *f* in **PCF**:

$$\mathtt{fix}\,f:\mathtt{nat} \rightharpoonup \mathtt{nat}\,\mathtt{is}\,\lambda\,(\,x:\mathtt{nat}\,)\,\mathtt{ifz}\,x\,\{\mathtt{z} \hookrightarrow \mathtt{s}(\,\mathtt{z}\,)\,|\,\mathtt{s}(\,x'\,) \hookrightarrow x*f(\,x'\,)\}.$$

Obviously evaluation of $f(\overline{n})$ requires n recursive calls to the function itself. That is, for a given input n we may place a *bound* m on the recursion that is sufficient to ensure termination of the computation. This property can be expressed formally using the m-bounded form of general recursion,

$$\inf^m f : \mathtt{nat} \rightharpoonup \mathtt{nat} \text{ is } \lambda \ (x : \mathtt{nat}) \text{ ifz } x \ \{ \mathtt{z} \hookrightarrow \mathtt{s}(\mathtt{z}) \ | \ \mathtt{s}(x') \hookrightarrow x * f(x') \}.$$

Call this expression $f^{(m)}$. It follows from the definition of f that if $f(\overline{n}) \mapsto^* \overline{p}$, then $f^{(m)}(\overline{n}) \mapsto^* \overline{p}$ for some m > 0 (in fact, m = n suffices).

When considering expressions of higher type, we cannot expect to get the *same* result from the bounded recursion as from the unbounded. For example, consider the addition function a of type $\tau = \mathtt{nat} \rightharpoonup (\mathtt{nat} \rightharpoonup \mathtt{nat})$, given by the expression

$$fix p: \tau is \lambda (x:nat) ifz x \{z \hookrightarrow id \mid s(x') \hookrightarrow s \circ (p(x'))\},$$

where $id = \lambda (y: \mathtt{nat}) y$ is the identity, $e' \circ e = \lambda (x:\tau) e'(e(x))$ is composition, and $s = \lambda (x: \mathtt{nat}) s(x)$ is the successor function. The application $a(\overline{n})$ terminates after three transitions, regardless of the value of n, resulting in a λ -abstraction. When n is positive, the result contains a residual copy of a itself, which is applied to n-1 as a recursive call. The m-bounded version of a, written $a^{(m)}$, is also such that $a^{(m)}(\overline{n})$ terminates in three steps, provided that m>0. But the result is not the same, because the residuals of a appear as $a^{(m-1)}$, rather than as a itself.

Turning now to the proof of compactness, it is helpful to introduce some notation. Suppose that $x: \tau \vdash e_x : \tau$ for some arbitrary abstractor $x.e_x$. Let $f^{(\omega)} = \operatorname{fix} x : \tau \operatorname{is} e_x$, and let $f^{(m)} = \operatorname{fix}^m x : \tau \operatorname{is} e_x$. Observe that $f^{(\omega)} : \tau$ and $f^{(m)} : \tau$ for any $m \ge 0$.

The following technical lemma governing the stack machine allows the bound on occurrences of a recursive expression to be raised without affecting the outcome of evaluation.

Lemma 47.15. For every $m \ge 0$, if $[f^{(m)}/y]k \triangleright [f^{(m)}/y]e \mapsto^* \epsilon \triangleleft \overline{n}$, then $[f^{(m+1)}/y]k \triangleright [f^{(m+1)}/y]e \mapsto^* \epsilon \triangleleft \overline{n}$.

Proof. By induction on $m \ge 0$, and then induction on transition.

Theorem 47.16 (Compactness). Suppose that $y: \tau \vdash e: \mathtt{nat}$ where $y \notin f^{(\omega)}$. If $[f^{(\omega)}/y]e \mapsto^* \overline{n}$, then there exists $m \geq 0$ such that $[f^{(m)}/y]e \mapsto^* \overline{n}$.

Proof. We prove simultaneously the stronger statements that if

$$[f^{(\omega)}/y]k \triangleright [f^{(\omega)}/y]e \longmapsto^* \epsilon \triangleleft \overline{n},$$

then for some $m \ge 0$,

$$[f^{(m)}/y]k \triangleright [f^{(m)}/y]e \longmapsto^* \epsilon \triangleleft \overline{n},$$

47.4 Compactness 459

and if

$$[f^{(\omega)}/y]k \triangleleft [f^{(\omega)}/y]e \longmapsto^* \epsilon \triangleleft \overline{n}$$

then for some $m \ge 0$,

$$[f^{(m)}/y]k \triangleleft [f^{(m)}/y]e \longrightarrow^* \epsilon \triangleleft \overline{n}.$$

(Note that if $[f^{(\omega)}/y]e$ val, then $[f^{(m)}/y]e$ val for all $m \ge 0$.) The result then follows by the correctness of the stack machine (Corollary 28.4).

We proceed by induction on transition. Suppose that the initial state is

$$[f^{(\omega)}/y]k \triangleright f^{(\omega)},$$

which arises when e = y, and the transition sequence is as follows:

$$[f^{(\omega)}/y]k \triangleright f^{(\omega)} \longmapsto [f^{(\omega)}/y]k \triangleright [f^{(\omega)}/x]e_x \longmapsto^* \epsilon \triangleleft \overline{n}.$$

Noting that $[f^{(\omega)}/x]e_x = [f^{(\omega)}/y][y/x]e_x$, we have by induction that there exists $m \ge 0$ such that

$$[f^{(m)}/y]k \triangleright [f^{(m)}/x]e_x \longmapsto^* \epsilon \triangleleft \overline{n}.$$

By Lemma 47.15

$$[f^{(m+1)}/y]k \triangleright [f^{(m)}/x]e_x \longmapsto^* \epsilon \triangleleft \overline{n}$$

and we need only recall that

$$[f^{(m+1)}/y]k \triangleright f^{(m+1)} = [f^{(m+1)}/y]k \triangleright [f^{(m)}/x]e_x$$

to complete the proof. If, on the other hand, the initial step is an unrolling, but $e \neq y$, then we have for some $z \notin f^{(\omega)}$ and $z \neq y$

$$[f^{(\omega)}/y]k \triangleright \mathtt{fix}\,z : \tau \mathtt{is}\,d_\omega \longmapsto [f^{(\omega)}/y]k \triangleright [\mathtt{fix}\,z : \tau \mathtt{is}\,d_\omega/z]d_\omega \longmapsto^* \epsilon \triangleleft \overline{n}.$$

where $d_{\omega} = [f^{(\omega)}/y]d$. By induction there exists $m \geq 0$ such that

$$[f^{(m)}/y]k \triangleright [\operatorname{fix} z : \tau \operatorname{is} d_m/z]d_m \longmapsto^* \epsilon \triangleleft \overline{n},$$

where $d_m = [f^{(m)}/y]d$. But then by Lemma 47.15 we have

$$[f^{(m+1)}/y]k \rhd [\mathtt{fix}\,z : \tau \, \mathtt{is}\, d_{m+1}/z]d_{m+1} \longmapsto^* \epsilon \mathrel{\triangleleft} \overline{n},$$

where $d_{m+1} = [f^{(m+1)}/y]d$, from which the result follows directly.

Corollary 47.17. There exists $m \ge 0$ such that $[f^{(\omega)}/y]e \simeq [f^{(m)}/y]e$.

Proof. If $[f^{(\omega)}/y]e$ diverges, then it suffices to take m to be zero. Otherwise, apply Theorem 47.16 to obtain m, and note that the required Kleene equivalence follows.

47.5 Lazy Natural Numbers

Recall from Chapter 19 that if the successor is evaluated lazily, then the type nat changes its meaning to that of the lazy natural numbers, which we shall write lnat for emphasis. This type contains an "infinite number" ω , which is essentially an endless stack of successors.

To account for the lazy successor the definition of logical equivalence must be reformulated. Rather than being defined *inductively* as the strongest relation closed under specified conditions, it is now defined *coinductively* as the weakest relation consistent with two analogous conditions. We may then show that two expressions are related using the principle of *proof by coinduction*.

The definition of Kleene equivalence must be altered to account for the lazily evaluated successor operation. To account for ω , two computations are compared based solely on the outermost form of their values, if any. We define $e \simeq e'$ to hold iff (a) if $e \longmapsto^* z$, then $e' \longmapsto^* z$, and *vice versa*; and (b) if $e \longmapsto^* s(e_1)$, then $e' \longmapsto^* s(e_1')$, and *vice versa*.

Corollary 47.17 can be proved for the co-natural numbers by essentially the same argument as before.

The definition of logical equivalence at type 1nat is defined to be the *weakest* equivalence relation \mathcal{E} between closed terms of type 1nat satisfying the following *consistency conditions*: if $e \mathcal{E} e'$: 1nat, then

- 1. If $e \mapsto^* z$, then $e' \mapsto^* z$, and *vice versa*.
- 2. If $e \mapsto^* s(e_1)$, then $e' \mapsto^* s(e'_1)$ with $e_1 \mathcal{E} e'_1$: lnat, and *vice versa*.

It is immediate that if $e \sim_{lnat} e'$, then $e \simeq e'$, and so logical equivalence is consistent. It is also strict in that if e and e' are both divergent expressions of type lnat, then $e \sim_{lnat} e'$.

The principle of *proof by coinduction* states that to show $e \sim_{\mathtt{lnat}} e'$, it suffices to exhibit a relation, \mathcal{E} , such that

- 1. $e \mathcal{E} e'$: lnat, and
- 2. \mathcal{E} satisfies the above consistency conditions.

If these requirements hold, then \mathcal{E} is contained in logical equivalence at type lnat, and hence $e \sim_{\mathtt{lnat}} e'$, as required.

As an application of coinduction, let us consider the proof of Theorem 47.8. The overall argument remains as before, but the proof for the type lnat must be altered as follows. Suppose that $\mathcal{A} \sim \mathcal{A}' : \tau_0 \leadsto \mathtt{lnat}$, and let $a = \mathcal{A}\{\mathtt{fix}\,x : \tau_0\,\mathtt{is}\,e\}$ and $a' = \mathcal{A}'\{\mathtt{fix}\,x : \tau_0\,\mathtt{is}\,e'\}$. Writing $a^{(m)} = \mathcal{A}\{\mathtt{fix}^m\,x : \tau_0\,\mathtt{is}\,e\}$ and $a'^{(m)} = \mathcal{A}'\{\mathtt{fix}^m\,x : \tau_0\,\mathtt{is}\,e'\}$, assume that

for every
$$m \ge 0$$
, $a^{(m)} \sim_{lnat} a'^{(m)}$.

We are to show that

$$a \sim_{\text{lnat}} a'$$
.

47.6 Notes 461

Define the functions p_n for $n \ge 0$ on closed terms of type lnat by the following equations:

$$p_0(d) = d$$
 $p_{(n+1)}(d) = \begin{cases} d' & \text{if } p_n(d) \longmapsto^* \mathfrak{s}(d') \\ undefined & \text{otherwise} \end{cases}$

For $n \ge 0$, let $a_n = p_n(a)$ and $a'_n = p_n(a')$. Correspondingly, let $a_n^{(m)} = p_n(a^{(m)})$ and $a'_n^{(m)} = p_n(a'^{(m)})$. Define \mathcal{E} to be the strongest relation such that $a_n \mathcal{E} a'_n$: lnat for all $n \ge 0$. We will show that the relation \mathcal{E} satisfies the consistency conditions, and so it is contained in logical equivalence. Because $a \mathcal{E} a'$: lnat (by construction), the result follows immediately.

To show that \mathcal{E} is consistent, suppose that a_n \mathcal{E} a'_n : lnat for some $n \geq 0$. We have by Corollary 47.17 $a_n \simeq a_n^{(m)}$, for some $m \geq 0$, and hence, by the assumption, $a_n \simeq {a'_n}^{(m)}$, and so by Corollary 47.17 again, $a'_n{}^{(m)} \simeq a'_n$. Now if $a_n \longmapsto^* \mathfrak{s}(b_n)$, then $a_n^{(m)} \longmapsto^* \mathfrak{s}(b_n^{(m)})$ for some $b_n^{(m)}$, and hence there exists $b'_n{}^{(m)}$ such that $a'_n{}^{(m)} \longmapsto^* b'_n{}^{(m)}$, and so there exists $b'_n{}$ such that $a'_n \longmapsto^* \mathfrak{s}(b'_n{})$. But $b_n = p_{n+1}(a)$ and $b'_n = p_{n+1}(a')$, and we have b_n \mathcal{E} $b'_n{}$: lnat by construction, as required.

47.6 Notes

The use of logical relations to characterize observational equivalence for **PCF** is inspired by the treatment of partiality in type theory by Constable and Smith (1987) and by the studies of observational equivalence by Pitts (2000). Although the technical details differ, the proof of compactness here is inspired by Pitts's structurally inductive characterization of termination using an abstract machine. It is critical to restrict attention to transition systems whose states are complete programs (closed expressions of observable type). Structural operational semantics usually does not fulfill this requirement, thereby requiring a considerably more complex argument than given here.

482 47.6 Notes









Appendix B

Background on Finite Sets

We make frequent use of the concepts of a *finite set* of *discrete objects* and of *finite functions* between them. A set X is *discrete* iff equality of its elements is decidable: for every $x,y \in X$, either $x = y \in X$ or $x \neq y \in X$. This condition is to be understood constructively as stating that we may effectively determine whether any two elements of the set X are equal or not. Perhaps the most basic example of a discrete set is the set $\mathbb N$ of natural numbers. A set X is *countable* iff there is a bijection $f: X \cong \mathbb N$ between X and the set of natural numbers, and it is *finite* iff there is a bijection, $f: X \cong \{0, \ldots, n-1\}$, where $n \in \mathbb N$, between it and some initial segment of the natural numbers. This condition is again to be understood constructively in terms of computable mappings, so that countable and finite sets are computably enumerable and, in the finite case, have a computable size.

Given countable sets, U and V, a *finite function* is a computable partial function $\phi: U \to V$ between them. The *domain* $dom(\phi)$ of ϕ is the set $\{u \in U \mid \phi(u) \downarrow \}$, of objects $u \in U$ such that $\phi(u) = v$ for some $v \in V$. Two finite functions, ϕ and ψ , between U and V are *disjoint* iff $dom(\phi) \cap dom(\psi) = \emptyset$. The *empty* finite function, \emptyset , between U and V is the totally undefined partial function between them. If $u \in U$ and $v \in V$, the finite function, $u \hookrightarrow v$, between U and V sends u to v, and is undefined otherwise; its domain is therefore the singleton set $\{u\}$. In some situations we write $u \sim v$ for the finite function $u \hookrightarrow v$.

If ϕ and ψ are two disjoint finite functions from U to V, then $\phi \otimes \psi$ is the finite function from U to V defined by the equation

$$(\phi \otimes \psi)(u) = \begin{cases} \phi(u) & \text{if } u \in dom(\phi) \\ \psi(v) & \text{if } v \in dom(\psi) \\ \text{undefined} & \text{otherwise} \end{cases}$$

If $u_1, \ldots, u_n \in U$ are pairwise distinct, and $v_1, \ldots, v_n \in V$, then we sometimes write $u_1 \hookrightarrow v_1, \ldots, u_n \hookrightarrow v_n$, or $u_1 \sim v_1, \ldots, u_n \sim v_n$, for the finite function $u_1 \hookrightarrow v_1 \otimes \ldots \otimes u_n \hookrightarrow v_n$.



Bibliography

- Martín Abadi and Luca Cardelli. A Theory of Objects. Springer-Verlag, 1996. 183, 248, 254
- Peter Aczel. An introduction to inductive definitions. In Jon Barwise, editor, *Handbook of Mathematical Logic*, chapter C.7, pages 783–818. North-Holland, 1977. 20
- John Allen. Anatomy of LISP. Computer Science Series. McGraw-Hill, 1978. 10, 297
- S.F. Allen, M. Bickford, R.L. Constable, R. Eaton, C. Kreitz, L. Lorigo, and E. Moran. Innovations in computational type theory using Nuprl. *Journal of Applied Logic*, 4(4):428–469, 2006. ISSN 1570-8683. doi: 10.1016/j.jal.2005.10.005. 87
- Stuart Allen. A non-type-theoretic definition of Martin-Löf's types. In *LICS*, pages 215–221, 1987. 452
- Zena M. Ariola and Matthias Felleisen. The call-by-need lambda calculus. *J. Funct. Program.*, 7(3): 265–301, 1997. 339
- Arvind, Rishiyur S. Nikhil, and Keshav Pingali. I-structures: Data structures for parallel computing. In Joseph H. Fasel and Robert M. Keller, editors, *Graph Reduction*, volume 279 of *Lecture Notes in Computer Science*, pages 336–369. Springer, 1986. ISBN 3-540-18420-1. 365
- Arnon Avron. Simple consequence relations. Information and Computation, 92:105–139, 1991. 30
- Henk Barendregt. *The Lambda Calculus, Its Syntax and Semantics*, volume 103 of *Studies in Logic and the Foundations of Mathematics*. North-Holland, 1984. 193
- Henk Barendregt. Lambda calculi with types. In S. Abramsky, D. M. Gabbay, and T. S. E Maibaum, editors, *Handbook of Logic in Computer Science*, volume 2, *Computational Structures*. Oxford University Press, 1992. 39
- Yves Bertot, Gérard Huet, Jean-Jacques Lévy, and Gordon Plotkin, editors. From Semantics to Computer Science: Essays in Honor of Gilles Kahn. Cambridge University Press, 2009. 545
- Guy E. Blelloch. Vector Models for Data-Parallel Computing. MIT Press, 1990. ISBN 0-262-02313-X. 356

542 BIBLIOGRAPHY

Guy E. Blelloch and John Greiner. Parallelism in sequential functional languages. In *FPCA*, pages 226–237, 1995. 356

- Guy E. Blelloch and John Greiner. A provable time and space efficient implementation of NESL. In *ICFP*, pages 213–225, 1996. 61, 356
- Manuel Blum. On the size of machines. Information and Control, 11(3):257-265, September 1967.
- Stephen D. Brookes. The essence of parallel algol. Inf. Comput., 179(1):118-149, 2002. 394
- Samuel R. Buss, editor. Handbook of Proof Theory. Elsevier, Amsterdam, 1998. 542
- Luca Cardelli. Structural subtyping and the notion of power type. In *Proc. ACM Symposium on Principles of Programming Languages*, pages 70–79, 1988. 224
- Luca Cardelli. Program fragments, linking, and modularization. In *Proc. ACM Symposium on Principles of Programming Languages*, pages 266–277, 1997. 408
- Giuseppe Castagna and Benjamin C. Pierce. Decidable bounded quantification. In *Proc. ACM Symposium on Principles of Programming Languages*, pages 151–162, 1994. 429
- Alonzo Church. The Calculi of Lambda-Conversion. Princeton University Press, 1941. 10, 71, 193
- R. L. Constable. *Implementing Mathematics with the Nuprl Proof Development System*. Prentice-Hall, Englewood Cliffs, NJ, 1986. v, 10, 237, 452
- Robert L. Constable. Types in logic, mathematics, and programming. In Buss (1998), chapter X. v.
- Robert L. Constable and Scott F. Smith. Partial objects in constructive type theory. In *LICS*, pages 183–193. IEEE Computer Society, 1987. 461
- William R. Cook. On understanding data abstraction, revisited. In *OOPSLA*, pages 557–572, 2009. 183
- Rowan Davies. *Practical Refinement-Type Checking*. PhD thesis, Carnegie Mellon University School of Computer Science, May 2005. Available as Technical Report CMU–CS–05–110. 237
- Rowan Davies and Frank Pfenning. Intersection types and computational effects. In Martin Odersky and Philip Wadler, editors, *ICFP*, pages 198–208. ACM, 2000. ISBN 1-58113-202-6. 237
- Ewen Denney. Refinement types for specification. In David Gries and Willem P. de Roever, editors, *PROCOMET*, volume 125 of *IFIP Conference Proceedings*, pages 148–166. Chapman & Hall, 1998. ISBN 0-412-83760-9. 237
- Derek Dreyer. *Understanding and Evolving the ML Module System*. PhD thesis, Carnegie Mellon University, Pittsburgh, PA, May 2005. 429
- Joshua Dunfield and Frank Pfenning. Type assignment for intersections and unions in call-by-value languages. In Andrew D. Gordon, editor, FoSSaCS, volume 2620 of Lecture Notes in Computer Science, pages 250–266. Springer, 2003. ISBN 3-540-00897-7. 237

Uffe Engberg and Mogens Nielsen. A calculus of communicating systems with label passing - ten years after. In Gordon D. Plotkin, Colin Stirling, and Mads Tofte, editors, *Proof, Language, and Interaction, Essays in Honour of Robin Milner*, pages 599–622. The MIT Press, 2000. 382

- Matthias Felleisen and Robert Hieb. The revised report on the syntactic theories of sequential control and state. *TCS: Theoretical Computer Science*, 103, 1992. 48, 279
- Tim Freeman and Frank Pfenning. Refinement types for ml. In David S. Wise, editor, *PLDI*, pages 268–277. ACM, 1991. ISBN 0-89791-428-7. 237
- Daniel Friedman and David Wise. The impact of applicative programming on multiprocessing. In *International Conference on Parallel Processing*, 1976. 365
- David Gelernter. Generative communication in Linda. *ACM Trans. Program. Lang. Syst.*, 7(1):80–112, 1985. 394
- Gerhard Gentzen. Investigations into logical deduction. In M. E. Szabo, editor, *The Collected Papers of Gerhard Gentzen*, pages 68–213. North-Holland, Amsterdam, 1969. 39
- J.-Y. Girard. *Interpretation fonctionelle et elimination des coupures de l'arithmetique d'ordre superieur*. These d'etat, Universite Paris VII, 1972. 148, 474
- Jean-Yves Girard. *Proofs and Types*. Cambridge University Press, 1989. Translated by Paul Taylor and Yves Lafont. v, 408
- Kurt Gödel. On a hitherto unexploited extension of the finitary standpoint. *Journal of Philosphical Logic*, 9:133–142, 1980. Translated by Wilfrid Hodges and Bruce Watson. 79
- Michael J. Gordon, Arthur J. Milner, and Christopher P. Wadsworth. *Edinburgh LCF*, volume 78 of *Lecture Notes in Computer Science*. Springer-Verlag, 1979. 10, 271
- John Greiner and Guy E. Blelloch. A provably time-efficient parallel implementation of full speculation. *ACM Trans. Program. Lang. Syst.*, 21(2):240–285, 1999. 365
- Timothy Griffin. A formulae-as-types notion of control. In *Proc. ACM Symposium on Principles of Programming Languages*, pages 47–58, 1990. 119
- Carl Gunter. Semantics of Programming Languages. Foundations of Computing Series. MIT Press, 1992. 183
- Robert H. Halstead, Jr. Multilisp: A language for concurrent symbolic computation. *ACM Trans. Program. Lang. Syst.*, 7(4):501–538, 1985. 365
- Robert Harper. Constructing type systems over an operational semantics. *J. Symb. Comput.*, 14(1): 71–84, 1992. 452
- Robert Harper. A simplified account of polymorphic references. *Inf. Process. Lett.*, 51(4):201–206, 1994. 330

Robert Harper and Mark Lillibridge. A type-theoretic approach to higher-order modules with sharing. In *Proc. ACM Symposium on Principles of Programming Languages*, pages 123–137, 1994. 417, 429, 439

- Robert Harper, John C. Mitchell, and Eugenio Moggi. Higher-order modules and the phase distinction. In *Proc. ACM Symposium on Principles of Programming Languages*, pages 341–354, 1990. 429
- Robert Harper, Furio Honsell, and Gordon Plotkin. A framework for defining logics. *Journal of the Association for Computing Machinery*, 40:194–204, 1993. 10, 30
- Ralf Hinze and Johan Jeuring. Generic haskell: Practice and theory. In Roland Carl Backhouse and Jeremy Gibbons, editors, *Generic Programming*, volume 2793 of *Lecture Notes in Computer Science*, pages 1–56. Springer, 2003. ISBN 3-540-20194-7. 127
- C. A. R. Hoare. Communicating sequential processes. Commun. ACM, 21(8):666–677, 1978. 382
- Tony Hoare. Null references: The billion dollar mistake. Presentation at QCon 2009, August 2009. 95
- S. C. Kleene. Introduction to Metamathematics. van Nostrand, 1952. 10
- Imre Lakatos. *Proofs and Refutations: The Logic of Mathematical Discovery.* Cambridge University Press, 1976. 119
- P. J. Landin. A correspondence between Algol 60 and Church's lambda notation. *CACM*, 8:89–101; 158–165, 1965. 48, 265
- Daniel K. Lee, Karl Crary, and Robert Harper. Towards a mechanized metatheory of standard ml. In *Proc. ACM Symposium on Principles of Programming Languages*, pages 173–184, 2007. 429
- Xavier Leroy. Manifest types, modules, and separate compilation. In *Proc. ACM Symposium on Principles of Programming Languages*, pages 109–122, 1994. 417, 429, 440
- Xavier Leroy. Applicative functors and fully transparent higher-order modules. In *Proc. ACM Symposium on Principles of Programming Languages*, pages 142–153, 1995. 440
- Mark Lillibridge. *Translucent Sums: A Foundation for Higher-Order Module Systems*. PhD thesis, Carnegie Mellon University School of Computer Science, Pittsburgh, PA, May 1997. 429
- Barbara Liskov and Jeannette M. Wing. A behavioral notion of subtyping. *ACM Trans. Program. Lang. Syst.*, 16(6):1811–1841, 1994. 254
- Saunders MacLane. *Categories for the Working Mathematician*. Graduate Texts in Mathematics. Springer-Verlag, second edition, 1998. 127, 136
- David B. MacQueen. Using dependent types to express modular structure. In *Proc. ACM Symposium on Principles of Programming Languages*, pages 277–286, 1986. 429

David B. MacQueen. Kahn networks at the dawn of functional programming. In Bertot et al. (2009), chapter 5. 183

- Yitzhak Mandelbaum, David Walker, and Robert Harper. An effective theory of type refinements. In Runciman and Shivers (2003), pages 213–225. ISBN 1-58113-756-7. 237
- Per Martin-Löf. Constructive mathematics and computer programming. In *Logic, Methodology and Philosophy of Science IV*, pages 153–175. North-Holland, 1980. 39, 54
- Per Martin-Löf. On the meanings of the logical constants and the justifications of the logical laws. Unpublished Lecture Notes, 1983. 20, 30
- Per Martin-Löf. *Intuitionistic Type Theory*. Studies in Proof Theory. Bibliopolis, Naples, Italy, 1984. v, 39, 408
- Per Martin-Löf. Truth of a proposition, evidence of a judgement, validity of a proof. *Synthese*, 73 (3):407–420, 1987. 20, 30
- John McCarthy. LISP 1.5 Programmer's Manual. MIT Press, 1965. 10, 202, 288
- N. P. Mendler. Recursive types and type constraints in second-order lambda calculus. In *LICS*, pages 30–36, 1987. 136
- Robin Milner. A theory of type polymorphism in programming. JCSS, 17:348–375, 1978. 54
- Robin Milner. *Communicating and mobile systems the Pi-calculus*. Cambridge University Press, 1999. ISBN 978-0-521-65869-0. 288, 304, 382, 481
- Robin Milner, Mads Tofte, Robert Harper, and David MacQueen. *The Definition of Standard ML (Revised)*. The MIT Press, 1997. 39, 61, 224, 304, 417, 429, 439
- John C. Mitchell. Coercion and type inference. In *Proc. ACM Symposium on Principles of Programming Languages*, pages 175–185, 1984. 224
- John C. Mitchell. Representation independence and data abstraction. In *Proc. ACM Symposium on Principles of Programming Languages*, pages 263–276, 1986. 157, 474
- John C. Mitchell. Foundations for Programming Languages. MIT Press, 1996. v
- John C. Mitchell and Gordon D. Plotkin. Abstract types have existential type. *ACM Trans. Program. Lang. Syst.*, 10(3):470–502, 1988. 157, 440
- Eugenio Moggi. Computational lambda-calculus and monads. In *LICS*, pages 14–23. IEEE Computer Society, 1989. ISBN 0-8186-1954-6. 318
- Tom Murphy, VII, Karl Crary, Robert Harper, and Frank Pfenning. A symmetric modal lambda calculus for distributed computing. In *LICS*, pages 286–295, 2004. 318, 400
- Chetan R. Murthy. An evaluation semantics for classical proofs. In *LICS*, pages 96–107. IEEE Computer Society, 1991. 119

Aleksandar Nanevski. From dynamic binding to state via modal possibility. In *PPDP*, pages 207–218. ACM, 2003. ISBN 1-58113-705-2. 297

- R. P. Nederpelt, J. H. Geuvers, and R. C. de Vrijer, editors. *Selected Papers on Automath*, volume 133 of *Studies in Logic and the Foundations of Mathematics*. North-Holland, 1994. 10, 30
- B. Nordstrom, K. Petersson, and J. M. Smith. *Programming in Martin-Löf's Type Theory*. Oxford University Press, 1990. URL http://www.cs.chalmers.se/Cs/Research/Logic/book. 10
- OCaml. Ocaml, 2012. URL http://caml.inria.fr/ocaml/. 440
- David Michael Ritchie Park. Concurrency and automata on infinite sequences. In Peter Deussen, editor, *Theoretical Computer Science*, volume 104 of *Lecture Notes in Computer Science*, pages 167–183. Springer, 1981. ISBN 3-540-10576-X. 481
- Frank Pfenning and Rowan Davies. A judgmental reconstruction of modal logic. *Mathematical Structures in Computer Science*, 11(4):511–540, 2001. 318
- Benjamin C. Pierce. Types and Programming Languages. The MIT Press, 2002. v, 87, 224, 248, 254
- Benjamin C. Pierce. Advanced Topics in Types and Programming Languages. The MIT Press, 2004. v
- Andrew M. Pitts. Existential types: Logical relations and operational equivalence. In Kim Guldstrand Larsen, Sven Skyum, and Glynn Winskel, editors, *ICALP*, volume 1443 of *Lecture Notes in Computer Science*, pages 309–326. Springer, 1998. ISBN 3-540-64781-3. 474
- Andrew M. Pitts. Operational semantics and program equivalence. In Gilles Barthe, Peter Dybjer, Luis Pinto, and João Saraiva, editors, *APPSEM*, volume 2395 of *Lecture Notes in Computer Science*, pages 378–412. Springer, 2000. ISBN 3-540-44044-5. 461
- Andrew M. Pitts and Ian D. B. Stark. Observable properties of higher order functions that dynamically create local names, or what's new? In Andrzej M. Borzyszkowski and Stefan Sokolowski, editors, *MFCS*, volume 711 of *Lecture Notes in Computer Science*, pages 122–141. Springer, 1993. ISBN 3-540-57182-5. 10, 288
- G. D. Plotkin. A structural approach to operational semantics. Technical Report DAIMI FN-19, Aarhus University Computer Science Department, 1981. 48, 265
- Gordon D. Plotkin. LCF considered as a programming language. *Theor. Comput. Sci.*, 5(3):223–255, 1977. 175
- Gordon D. Plotkin. The origins of structural operational semantics. *J. of Logic and Algebraic Programming*, 60:3–15, 2004. 48
- John H. Reppy. Concurrent Programming in ML. Cambridge University Press, 1999. 382, 394
- J. C. Reynolds. Types, abstraction, and parametric polymorphism. In *Information Processing '83*, pages 513–523. North-Holland, Amsterdam, 1983. 148, 474

John C. Reynolds. Towards a theory of type structure. In Bernard Robinet, editor, *Symposium on Programming*, volume 19 of *Lecture Notes in Computer Science*, pages 408–423. Springer, 1974. ISBN 3-540-06859-7. 148, 157

- John C. Reynolds. Using category theory to design implicit conversions and generic operators. In Neil D. Jones, editor, *Semantics-Directed Compiler Generation*, volume 94 of *Lecture Notes in Computer Science*, pages 211–258. Springer, 1980. ISBN 3-540-10250-7. 224
- John C. Reynolds. The essence of Algol. In *Proceedings of the 1981 International Symposium on Algorithmic Languages*, pages 345–372. North-Holland, 1981. 317, 329
- John C. Reynolds. The discoveries of continuations. *Lisp and Symbolic Computation*, 6(3-4):233–248, 1993. 279
- John C. Reynolds. *Theories of Programming Languages*. Cambridge University Press, Cambridge, England, 1998. v
- Andreas Rossberg, Claudio V. Russo, and Derek Dreyer. F-ing modules. In Andrew Kennedy and Nick Benton, editors, *TLDI*, pages 89–102. ACM, 2010. ISBN 978-1-60558-891-9. 429
- Colin Runciman and Olin Shivers, editors. *Proceedings of the Eighth ACM SIGPLAN International Conference on Functional Programming, ICFP 2003, Uppsala, Sweden, August 25-29, 2003, 2003.* ACM. ISBN 1-58113-756-7. 545, 548
- Dana Scott. Lambda calculus: Some models, some philosophy. In J. Barwise, H. J. Keisler, and K. Kunen, editors, *The Kleene Symposium*, pages 223–265. North Holland, Amsterdam, 1980a. 193
- Dana S. Scott. Data types as lattices. SIAM J. Comput., 5(3):522-587, 1976. 183
- Dana S Scott. Relating theories of the lambda calculus. *To HB Curry: Essays on combinatory logic, lambda calculus and formalism*, pages 403–450, 1980b. 211
- Dana S. Scott. Domains for denotational semantics. In Mogens Nielsen and Erik Meineche Schmidt, editors, *ICALP*, volume 140 of *Lecture Notes in Computer Science*, pages 577–613. Springer, 1982. ISBN 3-540-11576-5. 183
- Michael B. Smyth and Gordon D. Plotkin. The category-theoretic solution of recursive domain equations. *SIAM J. Comput.*, 11(4):761–783, 1982. 183
- Richard Statman. Logical relations and the typed lambda-calculus. *Information and Control*, 65 (2/3):85–97, 1985. 452
- Guy L. Steele. Common Lisp: The Language. Digital Press, 2nd edition edition, 1990. 271
- Christopher A. Stone and Robert Harper. Extensional equivalence and singleton types. *ACM Trans. Comput. Log.*, 7(4):676–722, 2006. 417, 429

Paul Taylor. *Practical Foundations of Mathematics*. Cambridge Studies in Advanced Mathematics. Cambridge University Press, Cambridge, 1999. 136

- P.W. Trinder, K. Hammond, H.-W. Loidl, and S.L. Peyton Jones. Algorithm + strategy = parallelism. *JOURNAL OF FUNCTIONAL PROGRAMMING*, 8:23–60, 1998. 365
- Jaap van Oosten. Realizability: A historical essay. *Mathematical Structures in Computer Science*, 12 (3):239–263, 2002. 237
- Philip Wadler. Theorems for free! In FPCA, pages 347–359, 1989. 148
- Philip Wadler. Comprehending monads. *Mathematical Structures in Computer Science*, 2(4):461–493, 1992. 318
- Philip Wadler. Call-by-value is dual to call-by-name. In Runciman and Shivers (2003), pages 189–201. ISBN 1-58113-756-7. 119
- Mitchell Wand. Fixed-point constructions in order-enriched categories. *Theor. Comput. Sci.*, 8:13–30, 1979. 183
- Stephen A. Ward and Robert H. Halstead. *Computation structures*. MIT electrical engineering and computer science series. MIT Press, 1990. ISBN 978-0-262-23139-8. 183
- Kevin Watkins, Iliano Cervesato, Frank Pfenning, and David Walker. Specifying properties of concurrent computations in clf. *Electr. Notes Theor. Comput. Sci.*, 199:67–87, 2008. 163
- Andrew K. Wright and Matthias Felleisen. A syntactic approach to type soundness. *Inf. Comput.*, 115(1):38–94, 1994. 54, 330
- Hongwei Xi and Frank Pfenning. Eliminating array bound checking through dependent types. In Jack W. Davidson, Keith D. Cooper, and A. Michael Berman, editors, *PLDI*, pages 249–257. ACM, 1998. ISBN 0-89791-987-4. 237

Index

FPC , see recursive types	abt, see abstract binding tree	
Λ , see untyped λ -calculus	assignables, see Modernized Algol	
F , see universal types	ast, see abstract syntax tree	
MA, see Modernized Algol		
PCF, see Plotkin's PCF	back-patching, see references	
PPCF, see parallelism	benign effects, see references	
T, see Gödel's T	bidirectional typing, 39, 40	
·	boolean type, 92	
abstract binding tree, 3, 6		
abstractor, 7	capabilities, 321	
valence, 7	channel types, see Concurrent Algol	
α -equivalence, 8	combinators	
bound variable, 8	sk basis, 30	
capture, 9	bracket abstraction, 31, 32	
free variable, 8	conversion, 31	
graph representation, 11	substitution, 31	
operator, 7	command types, see Modernized Algol	
arity, 7	compactness, see equality	
parameter, 9	Concurrent Algol, 385	
structural induction, 8	broadcast communication, 388	
substitution, 9	dynamics, 389	
weakening, 11	safety, 389	
abstract binding trees	statics, 389	
closed, 32	class declaration, 394	
abstract syntax tree, 3–5	definability of free assignables, 392	
operator, 4	dynamics, 387	
arity, 4	RS latch, 394	
index, 9	selective communication, 390	
parameter, 9	dynamics, 392	
structural induction, 5	statics, 391	
substitution, 6	statics, 386	
variable, 4	contravariance, see subtyping	
weakening, 10	covariance, see subtyping	
abstract types, see existential types, see also sig-		
natures	definitional equality, see equality	

dynamic types, 196 class dispatch, 200 cons, 199 critique, 201 destructors, 200 dynamics, 196 lists, 202 multiple arguments, 202 multiple results, 202 nultiple results, 202 nil, 199 numeric classes, 199 pairs, 202 predicates, 200 safety, 198 statics, 196 subtyping, 217 dynamics, 35, 41 checked errors, 53 contextual, 44 cost, 60 definitional equality, 46 determinacy, 44 environmental evaluation, 62 equational, 46 equivalence theorem, 45 evaluation, 57 equivalence to transition, 58 evaluation context, 44 induction on transition, 42 inversion principle, 44 structural, 42 transition system, 41 unchecked errors, 53 dynamics types arithmetic, 202	fixed point induction, 455 Kleene equality, 453 Kleene equivalence, 445 logical equivalence, 443, 446, 448, 454 closed, 447 observation, 444 observational equivalence, 443, 445, 448, 453 event types, see Concurrent Algol exceptions, 267, 269 dynamics, 269 evaluation dynamics, 271 exception type, 270, 271 safety, 270, 271 safety, 270, 271 statics, 269 structural dynamics, 271 syntax, 269 existential types, 151 coinductive types, 158 definability from universals, 154 dynamics, 152 modeling data abstraction, 153 parametricity, 158 representation independence, 155, 158 safety, 153 statics, 152 streams, 158 subtyping, 222 failures, see also exceptions, 267 dynamics, 268 safety, 268 statics, 267 finite function, 539 fixed point induction, see equality function types
02	subtyping, 219
enumeration types, 93 equality, 443	future types, 359, 360 future let, 365
coinduction, 446, 460	parallel dynamics, 362
compactness, 455, 457, 458	parallel let, 365
congruence, 445	pipelining, 364
contexts, 444, 453	sequential dynamics, 360
definitional, 46, 145, 171, 187, 443	sparks, 365
equational laws, 450	statics, 360
*	

futures, see future types	weakening, 24
Gödel's T , 73	inductive definition, 13, 14
canonical forms, 79	admissible rule, 25
definability, 76	backward chaining, 16
definitional equality, 76	derivable rule, 23
dynamics, 74	derivation, 15
hereditary termination, 79	forward chaining, 16
iterator, 74	function, 19
recursor, 73	iterated, 18
safety, 75, 79	rule, 14
statics, 74	axiom, 14
termination, 79	conclusion, 14
undefinability, 77	premise, 14
general judgment, 23, 28	rule induction, 15, 16
generic derivability, 28	rule scheme, 14
proliferation, 28	instance, 14
structurality, 28	simultaneous, 18
substitution, 28	inheritance, 251
parametric derivability, 29	class extension, 251
general recursion, 169	class-based, 252
generic inductive definition, 29	method extension, 252
formal generic judgment, 29	method specialization, 255
rule, 29	method-based, 254
rule induction, 29	self-reference, 255
structurality, 29	simple method override, 255
Girard's System F , see universal types	sub-method, 251
31	subclass, 251
hypothetical inductive definition, 26	super-method, 251
formal derivability, 27	superclass, 251
rule, 26	interface, see separate compilation
rule induction, 27	
uniformity of rules, 27	judgment, 13
hypothetical judgment, 23	judgment form, 13
admissibility, 25	predicate, 13
reflexivity, 26	subject, 13
structurality, 26	Vloopo oquality see oquality
transitivity, 26	Kleene equality, see equality
weakening, 26	laziness
derivability, 23	parallel or, 175
reflexivity, 24	linking, see separate compilation
stability, 24	logical equivalence, see equality
structurality, 24	
transitivity, 24	mobile types, 316

mobility condition, 316	implicit parallelism theorem, 346	
rules, 316	multiple fork-join, 349	
Modernized Algol, 309	parallel complexity, 347	
arrays, 318	parallel dynamics, 344	
assignables, 309, 322	parallelizability, 352	
block structure, 312	provably efficient implementation, 351	
classes and objects, 320	sequence types, 349	
command types, 316	cost dynamics, 350	
commands, 309, 315	statics, 350	
control stack, 320	sequential complexity, 347	
data stack, 320	sequential dynamics, 344	
expressions, 309	statics, 343	
free assignables, 324	structural dynamics, 356	
free dynamics, 324	task dynamics, 352, 357	
idioms	work vs. depth, 347	
conditionals, 314	phase distinction, 35, see also signatures	
iteration, 314	Plotkin's PCF, 167	
procedures, 314	Blum size theorem, 174	
sequential composition, 314	bounded recursion, 454	
multiple declaration instances, 318	definability, 172	
own assignables, 319	definitional equality, 171	
passive commands, 318	dynamics, 170	
recursive procedures, 318	eager natural numbers, 173	
scoped dynamics, 311	eagerness and laziness, 173	
scoped safety, 313	halting problem, 175	
separated and consolidated stacks, 320	induction, 173	
stack discipline, 312	mutual recursion, 175	
stack machine, 320	safety, 170	
statics, 310, 316	statics, 169	
modules, see signatures	totality and partiality, 173	
mutual primitive recursion, 86	polarity, 87	
	polymorphic types, see universal types	
null, see option types	primitive recursion, 86	
	product types, 83	
observational equivalence, see equality	dynamics, 84	
option types, 94	finite, 85	
	safety, 84	
parallelism, 343	statics, 83	
binary fork-join, 343	subtyping, 217, 219	
Brent's Theorem, 351		
cost dynamics, 346, 356	recursive types, see also type recursion	
cost dynamics vs. transition dynamics, 348	dynamics, 178	
cost graphs, 346	eager data structures, 179	
exceptions, 356	eager lists, 179	

eager natural numbers, 179	opacity, 420	
lazy data structures, 179	principal signature, 423	
lazy lists, 180	revelation, 420	
lazy natural numbers, 179	sealing, 421	
RS latch, 184	self-recognition, 427	
self-reference, 180	set abstraction, 429	
signals, 184	signature modification, 430	
statics, 177	static part, 421	
subtyping, 220, 225	statics, 425	
reference types, 321	structures, 420	
aliasing, 323	subsignature, 422–424	
free dynamics, 325	syntax, 424	
safety, 323, 326	translucency, 420	
scoped dynamics, 323	transparency, 420	
statics, 323	type abstractions, 419, 420	
references	type classes, 419, 422	
arrays, 330	views, 422	
back-patching, 329	sparks, see future types	
benign effects, 328	speculation types, 361	
mutable data structures, 330	parallel dynamics, 362	
Reynolds's Algol, see Modernized Algol	sequential dynamics, 361	
	statics, 361	
safety	speculations, see speculation types	
evaluation, 59, 60	stack machine, 259	
scoped assignables, see Modernized Algol	correctness, 262	
self types, 180	completeness, 263	
as recursive types, 181	soundness, 263	
deriving general recursion, 182	unraveling, 263	
self-reference, 180	dynamics, 260	
unrolling, 180	frame, 259	
separate compilation, 405	safety, 261	
initialization, 406	stack, 259	
interface, 405	state, 259	
linking, 405	state, 182	
units, 405	from recursion, 182	
signatures, 419	RS latch, 182	
ascription, see sealing	statics, 35	
avoidance problem, 426	canonical forms, 38	
dynamic part, 421	decomposition, 38	
dynamics, 428	induction on typing, 37	
first- vs second-class, 428	introduction and elimination, 38	
graph abstraction, 430	structurality, 37	
graph class, 430	substitution, 37	
instances, 422	type system, 36	

:.ib- 27	20 contrar d 102
unicity, 37	as untyped, 192
weakening, 37	unit
structural subtyping, see subtyping	dynamics, 84
subtyping, 215	statics, 83
bounded quantification, 222	unit type, 83
class types, 217	vs void type, 92
coercion, 225	units, see separate compilation
coherence, 225	universal types, 142
dynamic types, 217	sk combinators, 148
dynamics, 223	Church numerals, 146
function types, 219	definability, 145
numeric types, 216	booleans, 148
product types, 217, 219	inductive types, 149
quantified types, 222	lists, 149
recursive types, 220, 225	natural numbers, 146
safety, 223	products, 145
subsumption, 215	sums, 145
sum types, 217, 219	definitional equality, 145
variance, 218, 224	dynamics, 144
sum types, 89	parametricity, 147, 149
dynamics, 90	safety, 144
finite, 91	statics, 142
statics, 89	subtyping, 222
subtyping, 217, 219	untyped λ -calculus, 187
syntax, 3	Y combinator, 190
abstract, 3	as uni-typed, 192
binding, 3	booleans, 193
chart, 35	bracket abstraction, 194
concrete, 3	Church numerals, 189
structural, 3	definability, 188
surface, 3	definitional equality, 187
System F , see universal types	dynamics, 187
System 1, see universal types	lists, 194
type abstractions, see also existential types	products, 193
type classes, see signatures	Scott's Theorem, 190
type classes, see signatures type recursion, see recursive types	statics, 187
	streams, 194
type safety, 51	sums, 194
canonical forms, 52 checked errors, 54	3um3, 174
· ·	variance, see subtyping
errors, 55	void type, 89
preservation, 51, 55	vs unit type, 92
progress, 52, 55	dynamics, 90
uni typod à calculus 102	statics, 89
uni-typed λ -calculus, 192	outico, o