# 1 *ML*

*By François Pottier and Didier Rémy*

## 1.1  Preliminaries

### Names and renaming

Mathematicians and computer scientists use *names* to refer to arbitrary or unknown objects in the statement of a theorem, to refer to the parameters of a function, *etc.* Names are convenient because they are understandable by humans; nevertheless, they can be tricky. An in-depth treatment of the difficulties associated with names and renaming is beyond the scope of the present chapter: we encourage the reader to study Gabbay and Pitts' excellent series of papers (Gabbay and Pitts, 2002; Pitts, 2002b). Here, we merely recall a few notions that are used throughout this chapter. Consider, for instance, an inductive definition of the abstract syntax of a simple programming language, the pure $\lambda$-calculus:

$$t ::= z \mid \lambda z.t \mid t\ t$$

Here, the *meta-variable* z ranges over an infinite set of *variables*—that is, names—while the meta-variable t ranges over *terms*. As usual in mathematics, we write "the variable z" and "the term t" instead of "the variable denoted by z" and "the term denoted by t". The above definition states that a term may be a variable z, a pair of a variable and a term, written $\lambda z.t$, or a pair of terms, written $t_1\ t_2$. However, this is not quite what we need. Indeed, according to

---

The (currently unfinished) code that accompanies this chapter may be found at `http://pauillac.inria.fr/~remy/mlrow/`. For space reasons, some material, including proofs, exercises, and more, has been left out of this version. In the future, a full version of this chapter that includes the missing material will be made available at the same address. In spite of these omissions, this chapter is still oversize with respect to Benjamin's 100 page barrier: we currently have roughly 135 pages of text and 15 pages of solutions to exercises. We would appreciate comments and suggestions from the proofreaders as to how this chapter could be made shorter, without severely compromising its interest or readability.

this definition, the terms $\lambda z_1.z_1$ and $\lambda z_2.z_2$ are distinct, while we would like them to be a single mathematical object, because we intend $\lambda z.z$ to mean "the function that maps $z$ to $z$"—a meaning that is independent of the name $z$. To achieve this effect, we complete the above definition by stating that the construction $\lambda z.t$ *binds* $z$ within $t$. One may also say that $\lambda z$ is a *binder* whose *scope* is $t$. Then, $\lambda z.t$ is no longer a pair: rather, it is an *abstraction* of the variable $z$ within the term $t$. Abstractions have the property that the identity of the bound variable does not matter; that is, $\lambda z_1.z_1$ and $\lambda z_2.z_2$ are the same term. Informally, we say that terms are considered equal modulo $\alpha$-*conversion*. Once the position and scope of binders are known, several standard notions follow, such as the set of *free variables* of a term $t$, written $fv(t)$, and the *capture-avoiding substitution* of a term $t_1$ for a variable $z$ within a term $t_2$, written $[z \mapsto t_1]t_2$. For conciseness, we write $fv(t_1, t_2)$ for $fv(t_1) \cup fv(t_2)$. A term is said to be *closed* when it has no free variables.

A *renaming* is a total bijective mapping from variables to variables that affects only a finite number of variables. The sole property of a variable is its identity, that is, the fact that it is distinct from other variables. As a result, at a global level, all variables are interchangeable: if a theorem holds in the absence of hypotheses about any particular variable, then any renaming of it holds as well. We often make use of this fact. When proving a theorem $T$, we say that a hypothesis $H$ may be assumed *wihout loss of generality* (*w.l.o.g.*) if the theorem $T$ follows from the theorem $H \Rightarrow T$ via a renaming argument, which is usually left implicit.

If $\bar{z}_1$ and $\bar{z}_2$ are sets of variables, we write $\bar{z}_1 \ \# \ \bar{z}_2$ as a shorthand for $\bar{z}_1 \cap \bar{z}_2 = \varnothing$, and say that $\bar{z}_1$ is *fresh* for $\bar{z}_2$ (or vice-versa). We say that $\bar{z}$ is fresh for $t$ if and only if $\bar{z} \ \# \ fv(t)$ holds.

In this chapter, we work with several distinct varieties of names: program variables, memory locations, and type variables, the latter of which may be further divided into *kinds*. We draw names of different varieties from disjoint sets, each of which is infinite.

## 1.2 What is ML?

The name "ML" appeared during the late seventies. It then referred to a general-purpose programming language that was used as a meta-language (whence its name) within the theorem prover LCF (Gordon, Milner, and Wadsworth, 1979b). Since then, several new programming languages, each of which offers several different implementations, have drawn inspiration from it. So, what does "ML" stand for today?

For a semanticist, "ML" might stand for a programming language featuring first-class functions, data structures built out of products and sums, mutable

memory cells called *references*, exception handling, automatic memory management, and a call-by-value semantics. This view encompasses the Standard ML (Milner, Tofte, and Harper, 1990) and Caml (Leroy, 2000) families of programming languages. We refer to it as *ML-the-programming-language*.

For a type theorist, "ML" might stand for a particular breed of type systems, based on the simply-typed $\lambda$-calculus, but extended with a simple form of polymorphism introduced by `let` declarations. These type systems have decidable type inference; their type inference algorithms crucially rely on first-order unification and can be made efficient in practice. In addition to Standard ML and Caml, this view encompasses programming languages such as Haskell (Hudak, Peyton Jones, Wadler, Boutel, Fairbairn, Fasel, Guzman, Hammond, Hughes, Johnsson, Kieburtz, Nikhil, Partain, and Peterson, 1992) and Clean (Brus, van Eekelen, van Leer, and Plasmeijer, 1987), whose semantics is rather different—indeed, it is pure and lazy—but whose type system fits this description. We refer to it as *ML-the-type-system*. It is also referred to as *Hindley and Milner's type discipline* in the literature.

For us, "ML" might also stand for the particular programming language whose formal definition is given and studied in this chapter. It is a core calculus featuring first-class functions, `let` declarations, and constants. It is equipped with a call-by-value semantics. By customizing constants and their semantics, one may recover data structures, references, and more. We refer to this particular calculus as *ML-the-calculus*.

Why study ML-the-type-system today, such a long time after its initial discovery? One may think of at least two reasons.

First, its treatment in the literature is often cursory, because it is considered either as a simple extension of the simply-typed $\lambda$-calculus (TAPL Chapter 9) or as a subset of Girard and Reynolds' System F (TAPL Chapter 23). The former view is supported by the claim that the `let` construct, which distinguishes ML-the-type-system from the simply-typed $\lambda$-calculus, may be understood as a simple textual expansion facility. However, this view only tells part of the story, because it fails to give an account of the *principal types* property enjoyed by ML-the-type-system, leads to a naïve type inference algorithm whose time complexity is exponential, and breaks down when the language is extended with side effects, such as state or exceptions. The latter view is supported by the fact that every type derivation within ML-the-type-system is also a valid type derivation within an implicity-typed variant of System F. Such a view is correct, but again fails to give an account of type inference for ML-the-type-system, since type inference for System F is undecidable (Wells, 1999).

Second, existing accounts of type inference for ML-the-type-system (Milner, 1978; Damas and Milner, 1982; Tofte, 1988; Leroy, 1992; Lee and Yi, 1998;

Jones, 1999) usually involve heavy manipulations of type substitutions. Such an ubiquitous use of type substitutions is often quite obscure. Furthermore, actual implementations of the type inference algorithm do *not* explicitly manipulate substitutions; instead, they extend a standard first-order unification algorithm, where terms are updated in place as new equations are discovered (Huet, 1976). Thus, it is hard to tell, from these accounts, how to write an efficient type inference algorithm for ML-the-type-system. Yet, in spite of the increasing speed of computers, efficiency remains crucial when ML-the-type-system is extended with expensive features, such as Objective Caml's object types (Rémy and Vouillon, 1998) or polymorphic methods (Garrigue and Rémy, 1999).

For these reasons, we believe it is worth giving an account of ML-the-type-system that focuses on *type inference* and strives to be at once *elegant* and *faithful* to an efficient implementation. To achieve these goals, we forego type substitutions and instead put emphasis on *constraints*, which offer a number of advantanges. First, constraints allow a modular presentation of type inference as the combination of a constraint generator and a constraint solver. Such a decomposition allows reasoning separately about *when* a program is correct, on the one hand, and *how* to check whether it is correct, on the other hand. It has long been standard in the setting of the simply-typed $\lambda$-calculus (TAPL Chapter 22), but, to the best of our knowledge, has never been proposed for ML-the-type-system. Second, it is often natural to define and implement the solver as a constraint rewriting system. Then, the constraint language allows reasoning not only about correctness—is every rewriting step meaning-preserving?—but also about low-level implementation details, since constraints *are* the data structures manipulated throughout the type inference process. For instance, describing unification in terms of *multi-equations* (Jouannaud and Kirchner, 1991) allows reasoning about the sharing of nodes in memory, which a substitution-based approach cannot account for. Last, constraints are more general than type substitutions, and allow describing many extensions of ML-the-type-system, among which extensions with recursive types, rows, subtyping, first-order unification under a mixed prefix, and more.

Before delving into the details of this new presentation of ML-the-type-system, however, it is worth recalling its standard definition. Thus, in what follows, we first define the syntax and operational semantics of the programming language ML-the-calculus, and equip it with a type system, known as *Damas and Milner's type system*.

| x, y | ::= | | *Identifiers:* | | $m$ | | *Memory location* |
| | z | | *Variable* | | $\lambda$z.t | | *Function* |
| | $m$ | | *Memory location* | | c $v_1$ ... $v_k$ | | *Data* |
| | c | | *Constant* | | | | $c \in \mathcal{Q}^+ \wedge k \leq a(c)$ |
| t | ::= | | *Expressions:* | | c $v_1$ ... $v_k$ | | *Partial application* |
| | x | | *Identifier* | | | | $c \in \mathcal{Q}^- \wedge k < a(c)$ |
| | $\lambda$z.t | | *Function* | $\mathcal{E}$ | ::= | | *Evaluation Contexts:* |
| | t t | | *Application* | | $[]$ | | *Empty context* |
| | let z = t in t | | *Local definition* | | $\mathcal{E}$ t | | *Left side of an application* |
| v, w | ::= | | *Values:* | | v $\mathcal{E}$ | | *Right side of an application* |
| | z | | *Variable* | | let z = $\mathcal{E}$ in t | | *Local definition* |

**Figure 1-1: Syntax of ML-the-calculus**

### ML-the-calculus

The syntax of ML-the-calculus is defined in Figure 1-1. It is made up of several syntactic categories.

*Identifiers* group several kinds of names that may be referenced in a program: variables, memory locations, and constants. We let x and y range over identifiers. *Variables*—sometimes called *program variables* to avoid ambiguity—are names that may be bound to values using $\lambda$ or let binding forms; in other words, they are names for function parameters or local definitions. We let z and f range over program variables. We sometimes write _ for a program variable that does not occur free within its scope: for instance, $\lambda$_.t stands for $\lambda$z.t, provided z is fresh for t. *Memory locations* are names that represent memory addresses. By convention, memory locations never appear in source programs, that is, programs that are submitted to a compiler. They only appear during execution, when new memory blocks are allocated. *Constants* are fixed names for primitive values and operations, such as integer literals and integer arithmetic operations. Constants are elements of a finite or infinite set $\mathcal{Q}$. They are never subject to $\alpha$-conversion. Program variables, memory locations, and constants belong to distinct syntactic classes and may never be confused.

The set of constants $\mathcal{Q}$ is kept abstract, so most of our development is independent of its concrete definition. We assume that every constant c has a nonnegative integer *arity* $a(c)$. We further assume that $\mathcal{Q}$ is partitioned into subsets of *constructors* $\mathcal{Q}^+$ and *destructors* $\mathcal{Q}^-$. Constructors and destructors differ in that the former are used to *form* values, while the latter are used to

*operate* on values.

1.2.1  EXAMPLE [INTEGERS]: For every integer $n$, one may introduce a nullary constructor $\hat{n}$. In addition, one may introduce a binary destructor $\hat{+}$, whose applications are written infix, so $\mathsf{t}_1 \hat{+} \mathsf{t}_2$ stands for the double application $\hat{+}\,\mathsf{t}_1\,\mathsf{t}_2$ of the destructor $\hat{+}$ to the expressions $\mathsf{t}_1$ and $\mathsf{t}_2$.  □

*Expressions*—also known as *program terms* or *programs*—are the main syntactic category. Indeed, unlike procedural languages such as C and Java, functional languages, including ML-the-programming-language, suppress the distinction between expressions and statements. Expressions include identifiers, $\lambda$-abstractions, applications, and local definitions. The *$\lambda$-abstraction* $\lambda \mathsf{z}.\mathsf{t}$ represents the function of one parameter named $\mathsf{z}$ whose result is the expression $\mathsf{t}$, or, in other words, the function that maps $\mathsf{z}$ to $\mathsf{t}$. Note that the variable $\mathsf{z}$ is bound within the term $\mathsf{t}$, so (for instance) $\lambda \mathsf{z}_1.\mathsf{z}_1$ and $\lambda \mathsf{z}_2.\mathsf{z}_2$ are the same object. The *application* $\mathsf{t}_1\,\mathsf{t}_2$ represents the result of calling the function $\mathsf{t}_1$ with actual parameter $\mathsf{t}_2$, or, in other words, the result of applying $\mathsf{t}_1$ to $\mathsf{t}_2$. Application is left-associative, that is, $\mathsf{t}_1\,\mathsf{t}_2\,\mathsf{t}_3$ stands for $(\mathsf{t}_1\,\mathsf{t}_2)\,\mathsf{t}_3$. The construct $\mathtt{let}\ \mathsf{z} = \mathsf{t}_1\ \mathtt{in}\ \mathsf{t}_2$ represents the result of evaluating $\mathsf{t}_2$ after binding the variable $\mathsf{z}$ to $\mathsf{t}_1$. Note that the variable $\mathsf{z}$ is bound within $\mathsf{t}_2$, but not within $\mathsf{t}_1$, so for instance $\mathtt{let}\ \mathsf{z}_1 = \mathsf{z}_1\ \mathtt{in}\ \mathsf{z}_1$ and $\mathtt{let}\ \mathsf{z}_2 = \mathsf{z}_1\ \mathtt{in}\ \mathsf{z}_2$ are the same object. The construct $\mathtt{let}\ \mathsf{z} = \mathsf{t}_1\ \mathtt{in}\ \mathsf{t}_2$ has the same meaning as $(\lambda \mathsf{z}.\mathsf{t}_2)\,\mathsf{t}_1$, but is dealt with in a more flexible way by ML-the-type-system. To sum up, the syntax of ML-the-calculus is that of the pure $\lambda$-calculus, extended with memory locations, constants, and the $\mathtt{let}$ construct.

*Values* form a subset of expressions. They are expressions whose evaluation is completed. Values include identifiers, $\lambda$-abstractions, and applications of constants, of the form $\mathsf{c}\,\mathsf{v}_1\,\ldots\,\mathsf{v}_k$, where $k$ does not exceed $\mathsf{c}$'s arity if $\mathsf{c}$ is a constructor, and $k$ is smaller than $\mathsf{c}$'s arity if $\mathsf{c}$ is a destructor. In what follows, we are often interested in closed values, that is, values that do not contain any free program variables. We use the meta-variables $\mathsf{v}$ and $\mathsf{w}$ for values.

1.2.2  EXAMPLE: The integer literals $\ldots, \widehat{-1}, \hat{0}, \hat{1}, \ldots$ are nullary constructors, so they are values. Integer addition $\hat{+}$ is a binary destructor, so it is a value, and so is every partial application $\hat{+}\,\mathsf{v}$. Thus, both $\hat{+}\,\hat{1}$ and $\hat{+}\,\hat{+}$ are values. An application of $\hat{+}$ to two values, such as $\hat{2}\hat{+}\hat{2}$, is not a value.  □

1.2.3  EXAMPLE [PAIRS]: Let $(\cdot, \cdot)$ be a binary constructor. If $\mathsf{t}_1$ are $\mathsf{t}_2$ are expressions, then the double application $(\cdot, \cdot)\,\mathsf{t}_1\,\mathsf{t}_2$ may be called the *pair* of $\mathsf{t}_1$ and $\mathsf{t}_2$, and may be written $(\mathsf{t}_1, \mathsf{t}_2)$. By the definition above, $(\mathsf{t}_1, \mathsf{t}_2)$ is a value if and only if $\mathsf{t}_1$ and $\mathsf{t}_2$ are both values.  □

*Stores* are finite mappings from memory locations to closed values. A store $\mu$ represents what is usually called a heap, that is, a collection of data structures,

each of which is allocated at a particular address in memory and may contain pointers to other elements of the heap. ML-the-programming-language allows overwriting the contents of an existing memory block—an operation sometimes referred to as a *side effect.* In the operational semantics, this effect is achieved by mapping an existing memory location to a new value. We write $\varnothing$ for the empty store. We write $\mu[m \mapsto \mathtt{v}]$ for the store that maps $m$ to $\mathtt{v}$ and otherwise coincides with $\mu$. When $\mu$ and $\mu'$ have disjoint domains, we write $\mu\mu'$ for their union. We write $dom(\mu)$ for the domain of $\mu$ and $range(\mu)$ for the set of memory locations that appear in its codomain.

The operational semantics of a purely functional language, such as the pure $\lambda$-calculus, may be defined as a rewriting system on expressions. Because ML-the-calculus has side effects, however, we define its operational semantics as a rewriting system on *configurations.* A configuration $\mathtt{t}/\mu$ is a pair of an expression $\mathtt{t}$ and a store $\mu$. The memory locations in the domain of $\mu$ are considered bound within $\mathtt{t}/\mu$, so (for instance) $m_1/(m_1 \mapsto \hat{0})$ and $m_2/(m_2 \mapsto \hat{0})$ are the same object. In what follows, we are often interested in *closed* configurations, that is, configurations $\mathtt{t}/\mu$ such that $\mathtt{t}$ has no free program variables and every memory location that appears within $\mathtt{t}$ or within the range of $\mu$ is in the domain of $\mu$. If $\mathtt{t}$ is a source program, its evaluation begins within an empty store—that is, with the configuration $\mathtt{t}/\varnothing$. Because, by convention, source programs do not contain memory locations, this is a closed configuration. Furthermore, we shall see that all reducts of a closed configuration are closed as well. Please note that, instead of separating expressions and stores, it is possible to make store fragments part of the syntax of expressions; this idea, proposed in (Crank and Felleisen, 1991), is reminiscent of the encoding of reference cells in process calculi (Turner, 1995; Fournet and Gonthier, 1996).

A *context* is an expression where a single subexpression has been replaced with a *hole*, written $[]$. *Evaluation contexts* form a strict subset of contexts. In an evaluation context, the hole is meant to highlight a point in the program where it is valid to apply a reduction rule. Thus, the definition of evaluation contexts determines a reduction strategy: it tells where and in what order reduction steps may occur. For instance, the fact that $\lambda\mathtt{z}.[]$ is not an evaluation context means that the body of a function is never evaluated—that is, not until the function is applied, see R-BETA below. The fact that $\mathtt{t}\,\mathcal{E}$ is an evaluation context only if $\mathtt{t}$ is a value means that, to evaluate an application $\mathtt{t}_1\,\mathtt{t}_2$, one should fully evaluate $\mathtt{t}_1$ before attempting to evaluate $\mathtt{t}_2$. More generally, in the case of a multiple application, it means that arguments should be evaluated from left to right. Of course, other choices could be made: for instance, defining $\mathcal{E} ::= \ldots \mid \mathtt{t}\,\mathcal{E} \mid \mathcal{E}\,\mathtt{v} \mid \ldots$ would enforce a right-to-left evaluation order, while defining $\mathcal{E} ::= \ldots \mid \mathtt{t}\,\mathcal{E} \mid \mathcal{E}\,\mathtt{t} \mid \ldots$ would leave the evaluation order unspecified, effectively allowing reduction to alternate between

$$(\lambda \mathtt{z}.\mathtt{t}) \; \mathtt{v} \longrightarrow [\mathtt{z} \mapsto \mathtt{v}]\mathtt{t} \qquad \text{(R-BETA)}$$

$$\mathtt{let} \; \mathtt{z} = \mathtt{v} \; \mathtt{in} \; \mathtt{t} \longrightarrow [\mathtt{z} \mapsto \mathtt{v}]\mathtt{t} \qquad \text{(R-LET)}$$

$$\frac{\mathtt{t}/\mu \xrightarrow{\delta} \mathtt{t}'/\mu'}{\mathtt{t}/\mu \longrightarrow \mathtt{t}'/\mu'} \qquad \text{(R-DELTA)}$$

$$\frac{\begin{array}{c} \mathtt{t}/\mu \longrightarrow \mathtt{t}'/\mu' \\ dom\,(\mu'') \;\#\; dom\,(\mu') \\ range(\mu'') \;\#\; dom\,(\mu' \setminus \mu) \end{array}}{\mathtt{t}/\mu\mu'' \longrightarrow \mathtt{t}'/\mu'\mu''} \qquad \text{(R-EXTEND)}$$

$$\frac{\mathtt{t}/\mu \longrightarrow \mathtt{t}'/\mu'}{\mathcal{E}[\mathtt{t}]/\mu \longrightarrow \mathcal{E}[\mathtt{t}']/\mu'} \qquad \text{(R-CONTEXT)}$$

**Figure 1-2: Semantics of ML-the-calculus**

both subexpressions, and making evaluation nondeterministic. The fact that $\mathtt{let} \; \mathtt{z} = \mathtt{v} \; \mathtt{in} \; \mathcal{E}$ is not an evaluation context means that the body of a local definition is never evaluated—that is, not until the definition itself is reduced, see R-LET below. We write $\mathcal{E}[\mathtt{t}]$ for the expression obtained by replacing the hole in $\mathcal{E}$ with the expression $\mathtt{t}$.

Figure 1-2 defines first a relation $\longrightarrow$ between configurations, then a relation $\longrightarrow\!\!\!\!\!\rightarrow$ between *closed* configurations. If $\mathtt{t}/\mu \longrightarrow \mathtt{t}'/\mu'$ or $\mathtt{t}/\mu \longrightarrow\!\!\!\!\!\rightarrow \mathtt{t}'/\mu'$ holds, then we say that the configuration $\mathtt{t}/\mu$ *reduces* to the configuration $\mathtt{t}'/\mu'$; the ambiguity involved in this definition is benign. If $\mathtt{t}/\mu \longrightarrow \mathtt{t}'/\mu$ holds for every store $\mu$, then we write $\mathtt{t} \longrightarrow \mathtt{t}'$ and say that the reduction is *pure*.

The key reduction rule is R-BETA, which states that a function application $(\lambda \mathtt{z}.\mathtt{t}) \; \mathtt{v}$ reduces to the function body, namely $\mathtt{t}$, where every occurrence of the formal argument $\mathtt{z}$ has been replaced with the actual argument $\mathtt{v}$. The $\lambda$ construct, which prevented the function body $\mathtt{t}$ from being evaluated, disappears, so the new term may (in general) be further reduced. Because ML-the-calculus adopts a *call-by-value* strategy, rule R-BETA is applicable only if the actual argument is a value $\mathtt{v}$. In other words, a function cannot be invoked until its actual argument has been fully evaluated. Rule R-LET is very similar to R-BETA. Indeed, it specifies that $\mathtt{let} \; \mathtt{z} = \mathtt{v} \; \mathtt{in} \; \mathtt{t}$ has the same behavior, with respect to reduction, as $(\lambda \mathtt{z}.\mathtt{t}) \; \mathtt{v}$. We remark that substitution of a value for a program variable throughout a term is expensive, so R-BETA and R-LET are never implemented literally: they are only a simple *specification*. Actual implementations usually employ *runtime environments*, which may be understood as a form of *explicit substitutions* (Abadi, Cardelli, Curien, and Lévy, 1991). Please note that our choice of a call-by-value reduction strategy is fairly arbitrary, and has essentially no impact on the type system; the programming language Haskell (Hudak, Peyton Jones, Wadler, Boutel, Fairbairn, Fasel, Guzman, Hammond, Hughes, Johnsson, Kieburtz,

Nikhil, Partain, and Peterson, 1992), whose reduction strategy is known as *lazy* or *call-by-need*, also relies on Hindley and Milner's type discipline.

Rule R-DELTA describes the semantics of constants. It merely states that a certain relation $\xrightarrow{\delta}$ is a subset of $\longrightarrow$. Of course, since the set of constants is unspecified, the relation $\xrightarrow{\delta}$ must be kept abstract as well. We require that, if $t/\mu \xrightarrow{\delta} t'/\mu'$ holds, then

(i) $t$ is of the form $c\ v_1\ \ldots\ v_n$, where $c$ is a destructor of arity $n$; and

(ii) $dom\,(\mu)$ is a subset of $dom\,(\mu')$.

Condition (i) ensures that $\delta$-reduction concerns full applications of destructors only, and that these are evaluated in accordance with the call-by-value strategy. Condition (ii) ensures that $\delta$-reduction may allocate new memory locations, but not deallocate existing locations. In particular, a "garbage collection" operator, which destroys unreachable memory cells, cannot be made available as a constant. Doing so would not make much sense anyway in the presence of R-EXTEND, which states that any valid reduction is also valid in a larger store. Condition (ii) allows proving that, if $t/\mu$ reduces to $t'/\mu'$, then $dom\,(\mu)$ is a subset of $dom\,(\mu')$; this is left as an exercise to the reader.

1.2.4   EXAMPLE [INTEGERS, CONTINUED]:  The operational semantics of integer addition may be defined as follows:

$$\hat{k}_1 \mathbin{\hat{+}} \hat{k}_2 \xrightarrow{\delta} \widehat{k_1 + k_2} \qquad\qquad \text{(R-ADD)}$$

The left-hand term is the double application $\hat{+}\ \hat{k}_1\ \hat{k}_2$, while the right-hand term is the integer literal $\hat{k}$, where $k$ is the sum of $k_1$ and $k_2$. The distinction between object level and meta level (that is, between $\hat{k}$ and $k$) is needed here to avoid ambiguity. □

1.2.5   EXAMPLE [PAIRS, CONTINUED]:  In addition to the pair constructor defined in Example 1.2.3, we may introduce two destructors $\pi_1$ and $\pi_2$ of arity 1. We may define their operational semantics as follows, for $i \in \{1, 2\}$:

$$\pi_i\ (v_1, v_2) \xrightarrow{\delta} v_i \qquad\qquad \text{(R-PROJ)}$$

Thus, our treatment of constants is general enough to account for pair construction and destruction; we need not build these features explicitly into the language. □

1.2.6   EXERCISE [BOOLEANS, RECOMMENDED, ★★]: Let `true` and `false` be nullary constructors. Let `if` be a ternary destructor. Extend the operational semantics with

$$\text{if true } v_1\ v_2 \xrightarrow{\delta} v_1 \qquad\qquad \text{(R-TRUE)}$$

$$\text{if false } \mathtt{v}_1 \ \mathtt{v}_2 \xrightarrow{\delta} \mathtt{v}_2 \qquad \text{(R-F\textsc{alse})}$$

Let us use the syntactic sugar if $\mathtt{t}_0$ then $\mathtt{t}_1$ else $\mathtt{t}_2$ for the triple application of if $\mathtt{t}_0$ $\mathtt{t}_1$ $\mathtt{t}_2$. Explain why these definitions do not quite provide the expected behavior. Without modifying the semantics of if, suggest a new definition of the syntactic sugar if $\mathtt{t}_0$ then $\mathtt{t}_1$ else $\mathtt{t}_2$ that corrects the problem. □

1.2.7 EXAMPLE [SUMS]: Booleans may in fact be viewed as a special case of the more general concept of *sum*. Let $\mathtt{inj}_1$ and $\mathtt{inj}_2$ be unary constructors, called respectively *left* and *right injections*. Let case be a ternary destructor, whose semantics is defined as follows, for $i \in \{1, 2\}$:

$$\text{case } (\mathtt{inj}_i \ \mathtt{v}) \ \mathtt{v}_1 \ \mathtt{v}_2 \xrightarrow{\delta} \mathtt{v}_i \ \mathtt{v} \qquad \text{(R-C\textsc{ase})}$$

Here, the value $\mathtt{inj}_i$ $\mathtt{v}$ is being scrutinized, while the values $\mathtt{v}_1$ and $\mathtt{v}_2$, which are typically functions, represent the two arms of a standard case construct. The rule selects an appropriate arm (here, $\mathtt{v}_i$) based on whether the value under scrutiny was formed using a left or right injection. The arm $\mathtt{v}_i$ is executed and given access to the data carried by the injection (here, $\mathtt{v}$). □

1.2.8 EXERCISE [★, ↛]: Explain how to encode true, false and the if construct in terms of sums. Check that the behavior of R-T\textsc{rue} and R-F\textsc{alse} is properly emulated. □

1.2.9 EXAMPLE [REFERENCES]: Let ref and ! be unary destructors. Let := be a binary destructor. We write $\mathtt{t}_1$ := $\mathtt{t}_2$ for the double application := $\mathtt{t}_1$ $\mathtt{t}_2$. Define the operational semantics of these three destructors as follows:

$$\text{ref } \mathtt{v}/\varnothing \xrightarrow{\delta} m/(m \mapsto \mathtt{v}) \quad \text{if } m \text{ is fresh for } \mathtt{v} \qquad \text{(R-R\textsc{ef})}$$

$$!m/(m \mapsto \mathtt{v}) \xrightarrow{\delta} \mathtt{v}/(m \mapsto \mathtt{v}) \qquad \text{(R-D\textsc{eref})}$$

$$m := \mathtt{v}/(m \mapsto \mathtt{v}_0) \xrightarrow{\delta} \mathtt{v}/(m \mapsto \mathtt{v}) \qquad \text{(R-A\textsc{ssign})}$$

According to R-R\textsc{ef}, evaluating ref $\mathtt{v}$ allocates a fresh memory location $m$ and binds $\mathtt{v}$ to it. Because configurations are considered equal up to $\alpha$-conversion of memory locations, the choice of the name $m$ is irrelevant, provided it is chosen fresh for $\mathtt{v}$, so as to prevent inadvertent capture of the memory locations that appear free within $\mathtt{v}$. By R-D\textsc{eref}, evaluating $!m$ returns the value bound to the memory location $m$ within the current store. By R-A\textsc{ssign}, evaluating $m := \mathtt{v}$ discards the value $\mathtt{v}_0$ currently bound to $m$ and produces a new store where $m$ is bound to $\mathtt{v}$. Here, the value returned by the assignment $m := \mathtt{v}$ is $\mathtt{v}$ itself; in ML-the-programming-language, it is usually a nullary constructor (), pronounced *unit*. □

1.2.10 EXAMPLE [RECURSION]: Let fix be a binary destructor, whose operational semantics is:

$$\text{fix } v_1 \ v_2 \xrightarrow{\delta} v_1 \ (\text{fix } v_1) \ v_2 \qquad\qquad \text{(R-Fix)}$$

`fix` is a fixpoint combinator: it effectively allows recursive definitions of functions. Indeed, the construct `letrec f` $= \lambda z.t_1$ `in` $t_2$ provided by ML-the-programming-language may be viewed as syntactic sugar for `let f` $=$ `fix` $(\lambda f.\lambda z.t_1)$ `in` $t_2$. □

Rule R-CONTEXT completes the definition of the operational semantics by defining $\longrightarrow$, a relation between closed configurations, in terms of $\longrightarrow$. The rule states that reduction may take place not only at the term's root, but also deep inside it, provided the path from the root to the point where reduction occurs forms an evaluation context. This is how evaluation contexts determine an evaluation strategy. As a purely technical point, because $\longrightarrow$ relates closed configurations only, we do not need to require that the memory locations in $dom(\mu' \setminus \mu)$ be fresh for $\mathcal{E}$: indeed, every memory location that appears within $\mathcal{E}$ must be a member of $dom(\mu)$.

1.2.11 EXERCISE [★★, RECOMMENDED, $\nrightarrow$]: Assuming the availability of Booleans and conditionals, integer literals, subtraction, multiplication, integer comparison, and a fixpoint combinator, most of which were defined in previous examples, define a function that computes the factorial of its integer argument, and apply it to $\hat{3}$. Determine, step by step, how this expression reduces to a value. □

It is straightforward to check that, if $t/\mu$ reduces to $t'/\mu'$, then $t$ is not a value. In other words, values are irreducible: they represent a completed computation. The proof is left as an exercise to the reader. The converse, however, does not hold: if $t/\mu$ is irreducible with respect to $\longrightarrow$, then $t$ is not necessarily a value. In that case, the configuration $t/\mu$ is said to be *stuck*. It represents a *runtime error*, that is, a situation that does not allow computation to proceed, yet is not considered a valid outcome. A closed source program $t$ is said to *go wrong* if and only if the configuration $t/\varnothing$ reduces to a stuck configuration.

1.2.12 EXAMPLE: Runtime errors typically arise when destructors are applied to arguments of an unexpected nature. For instance, the expressions $+ \ 1 \ m$ and $\pi_1 \ 2$ and $!3$ are stuck, regardless of the current store. The program `let z` $=$ $\hat{+} \ \hat{+}$ `in z 1` is not stuck, because $\hat{+} \ \hat{+}$ is a value. However, its reduct through R-LET is $\hat{+} \ \hat{+} \ 1$, which is stuck, so this program goes wrong. The primary purpose of type systems is to prevent such situations from arising. □

1.2.13 EXAMPLE: The configuration $!m/\mu$ is stuck if $m$ is not in the domain of $\mu$. In that case, however, $!m/\mu$ is not closed. Because we consider $\longrightarrow$ as a relation between closed configurations only, this situation cannot arise. In other

words, the semantics of ML-the-calculus never allows the creation of *dangling pointers*. As a result, no particular precautions need be taken to guard against them. Several strongly typed programming languages do nevertheless allow dangling pointers in a controlled fashion (Tofte and Talpin, 1997; Crary, Walker, and Morrisett, 1999b; DeLine and Fähndrich, 2001; Grossman, Morrisett, Jim, Hicks, Wang, and Cheney, 2002a).                                    □

### Damas and Milner's type system

ML-the-type-system was originally defined by Milner (1978). Here, we reproduce the definition given a few years later by Damas and Milner (1982), which is written in a more standard style: typing judgements are defined inductively by a collection of typing rules. We refer to this type system as DM.

To begin, we must define *types*. In DM, like in the simply-typed $\lambda$-calculus, types are first-order terms built out of *type constructors* and *type variables*. We begin with several considerations concerning the specification of type constructors.

First, we do not wish to fix the set of type constructors. Certainly, since ML-the-calculus has functions, we need to be able to form an arrow type $T \to T'$ out of arbitrary types $T$ and $T'$; that is, we need a binary type constructor $\to$. However, because ML-the-calculus includes an unspecified set of constants, we cannot say much else in general. If constants include integer literals and integer operations, as in Example 1.2.1, then a nullary type constructor `int` is needed; if they include pair construction and destruction, as in Examples 1.2.3 and 1.2.5, then a binary type constructor $\times$ is needed; and so on.

Second, it is common to refer to the parameters of a type constructor *by position*, that is, by numeric index. For instance, when one writes $T \to T'$, it is understood that the type constructor $\to$ has arity 2, that $T$ is its *first* parameter, known as its *domain*, and that $T'$ is its *second* parameter, known as its *codomain*. Here, however, we refer to parameters *by names*, known as *directions*. For instance, we define two directions *domain* and *codomain* and let the type constructor $\to$ have arity $\{domain, codomain\}$. The extra generality afforded by directions is exploited in the definition of nonstructural subtyping (Example 1.3.9) and in the definition of rows (Section 1.11).

Last, we allow types to be classified using *kinds*. As a result, every type constructor must come not only with an arity, but with a richer *signature*, which describes the kinds of the types to which it is applicable and the kind of the type that it produces. A distinguished kind $\star$ is associated with "normal" types, that is, types that are directly ascribed to expressions and values. For instance, the signature of the type constructor $\to$ is $\{domain \mapsto \star, codomain \mapsto \star\} \Rightarrow \star$, because it is applicable to two "normal" types and produces a "normal" type.

Introducing kinds other than $\star$ allows viewing some terms as ill-formed types; this is illustrated, for instance, in Section 1.11. In the simplest case, however, $\star$ is really the only kind, so the signature of a type constructor is nothing but its arity (a set of directions), and every term is a well-formed type, provided every application of a type constructor respects its arity.

1.2.14    DEFINITION:  Let $d$ range over a finite or denumerable set of *directions*. Let $\kappa$ range over a finite or denumerable set of *kinds*. Let $\star$ be a distinguished kind. Let $K$ range over partial mappings from directions to kinds. Let $F$ range over a finite or denumerable set of *type constructors*, each of which has a *signature* of the form $K \Rightarrow \kappa$. The domain of $K$ is referred to as the *arity* of $F$, while $\kappa$ is referred to as its *image kind*. We write $\kappa$ instead of $K \Rightarrow \kappa$ when $K$ is empty. Let $\rightarrow$ be a type constructor of signature $\{domain \mapsto \star, codomain \mapsto \star\} \Rightarrow \star$. □

The type constructors and their signatures collectively form a *signature* $\mathcal{S}$. In the following, we assume that a fixed signature $\mathcal{S}$ is given and that every type constructor in it has *finite* arity, so as to ensure that types are machine representable. However, in Section 1.11, we shall explicitly work with several distinct signatures, some of which involve type constructors of denumerable arity.

A *type variable* is a name that is used to stand for a type. For simplicity, we assume that every type variable is branded with a kind, or, in other words, that type variables of distinct kinds are drawn from disjoint sets. Each of these sets of type variables is individually subject to $\alpha$-conversion: that is, renamings must preserve kinds. Attaching kinds to type variables is only a technical convenience: in practice, every operation performed during type inference preserves the property that every type is well-kinded, so it is not necessary to keep track of the kind of every type variable. It is only necessary to check that all types supplied by the user, within type declarations, type annotations, or module interfaces, are well-kinded.

1.2.15    DEFINITION:  For every kind $\kappa$, let $\mathcal{V}_\kappa$ be a disjoint, denumerable set of *type variables*. Let X, Y, and Z range over the set $\mathcal{V}$ of all type variables. Let $\bar{\mathtt{X}}$ and $\bar{\mathtt{Y}}$ range over finite sets of type variables. We write $\bar{\mathtt{X}}\bar{\mathtt{Y}}$ for the set $\bar{\mathtt{X}} \cup \bar{\mathtt{Y}}$ and often write X for the singleton set $\{\mathtt{X}\}$. We write $ftv(o)$ for the set of *free type variables* of an object $o$.                                                                     □

The set of types, ranged over by T, is the free many-kinded term algebra that arises out of the type constructors and type variables.

1.2.16    DEFINITION:  A *type* of kind $\kappa$ is either a member of $\mathcal{V}_\kappa$, or a term of the form $F\,\{d_1 \mapsto \mathtt{T}_1, \ldots, d_n \mapsto \mathtt{T}_n\}$, where $F$ has signature $\{d_1 \mapsto \kappa_1, \ldots, d_n \mapsto \kappa_n\} \Rightarrow \kappa$ and $\mathtt{T}_1, \ldots, \mathtt{T}_n$ are types of kind $\kappa_1, \ldots, \kappa_n$, respectively.                                 □

As a notational convention, we assume that, for every type constructor $F$, the directions that form the arity of $F$ are implicitly ordered, so that we may say that $F$ has signature $\kappa_1 \otimes \ldots \otimes \kappa_n \Rightarrow \kappa$ and employ the syntax $F \, \mathtt{T}_1 \, \ldots \, \mathtt{T}_n$ for applications of $F$. Applications of the type constructor $\rightarrow$ are written infix and associate to the right, so $\mathtt{T} \rightarrow \mathtt{T}' \rightarrow \mathtt{T}''$ stands for $\mathtt{T} \rightarrow (\mathtt{T}' \rightarrow \mathtt{T}'')$.

In order to give meaning to the free type variables of a type, or, more generally, of a typing judgement, traditional presentations of ML-the-type-system, including Damas and Milner's, employ *type substitutions*. Most of our presentation avoids substitutions and uses *constraints* instead. However, we do need substitutions on a few occasions, especially when relating our presentation to Damas and Milner's.

1.2.17   DEFINITION: A *type substitution* $\theta$ is a total, kind-preserving mapping of type variables to types that is the identity everywhere but on a finite subset of $\mathcal{V}$, which we call the *domain* of $\theta$ and write $dom(\theta)$. The *range* of $\theta$, which we write $range(\theta)$, is the set $ftv(\theta(dom(\theta)))$. A type substitution may naturally be viewed as a total, kind-preserving mapping of types to types. In the following, we write $\bar{\mathtt{X}} \mathbin{\#} \theta$ for $\bar{\mathtt{X}} \mathbin{\#} (dom(\theta) \cup range(\theta))$. We write $\theta \setminus \bar{\mathtt{X}}$ for the restriction of $\theta$ outside $\bar{\mathtt{X}}$, that is, the restriction of $\theta$ to $\mathcal{V} \setminus \bar{\mathtt{X}}$. We sometimes let $\varphi$ denote a type substitution.                                                                    □

If $\vec{\mathtt{X}}$ and $\vec{\mathtt{T}}$ are respectively a vector of distinct type variables and a vector of types of the same (finite) length, such that, for every index $i$, $\mathtt{X}_i$ and $\mathtt{T}_i$ have the same kind, then $[\vec{\mathtt{X}} \mapsto \vec{\mathtt{T}}]$ denotes the substitution that maps $\mathtt{X}_i$ to $\mathtt{T}_i$ for every index $i$. The domain of $[\vec{\mathtt{X}} \mapsto \vec{\mathtt{T}}]$ is a subset of $\bar{\mathtt{X}}$, the set underlying the vector $\vec{\mathtt{X}}$. Its range is a subset of $ftv(\bar{\mathtt{T}})$, where $\bar{\mathtt{T}}$ is the set underlying the vector $\vec{\mathtt{T}}$. Every substitution $\theta$ may be written under the form $[\vec{\mathtt{X}} \mapsto \vec{\mathtt{T}}]$, where $\bar{\mathtt{X}} = dom(\theta)$. Then, $\theta$ is *idempotent* if and only if $\bar{\mathtt{X}} \mathbin{\#} ftv(\bar{\mathtt{T}})$ holds.

As pointed out earlier, types are first-order terms; that is, in the grammar of types, none of the productions *binds* a type variable. As a result, every type variable that appears within a type $\mathtt{T}$ appears *free* within $\mathtt{T}$. This situation is identical to that of the simply-typed $\lambda$-calculus. Things become more interesting when we introduce *type schemes*. As its name implies, a type scheme may describe an entire family of types; this effect is achieved via *universal quantification* over a set of type variables.

1.2.18   DEFINITION: A type scheme $\mathtt{S}$ is an object of the form $\forall \bar{\mathtt{X}}.\mathtt{T}$, where $\mathtt{T}$ is a type of kind $\star$ and the type variables $\bar{\mathtt{X}}$ are considered bound within $\mathtt{T}$.           □

One may view the type $\mathtt{T}$ as the trivial type scheme $\forall \varnothing.\mathtt{T}$, where no type variables are universally quantified, so types may be viewed as a subset of type schemes. The type scheme $\forall \bar{\mathtt{X}}.\mathtt{T}$ may be viewed as a finite way of describing the possibly infinite family of types of the form $[\vec{\mathtt{X}} \mapsto \vec{\mathtt{T}}]\mathtt{T}$, where $\vec{\mathtt{T}}$ is arbitrary.

$$\frac{\Gamma(\mathtt{x}) = \mathtt{S}}{\Gamma \vdash \mathtt{x} : \mathtt{S}} \quad (\textsc{dm-Var})$$

$$\frac{\Gamma; \mathtt{z} : \mathtt{T} \vdash \mathtt{t} : \mathtt{T}'}{\Gamma \vdash \lambda \mathtt{z}.\mathtt{t} : \mathtt{T} \to \mathtt{T}'} \quad (\textsc{dm-Abs})$$

$$\frac{\Gamma \vdash \mathtt{t}_1 : \mathtt{T} \to \mathtt{T}' \qquad \Gamma \vdash \mathtt{t}_2 : \mathtt{T}}{\Gamma \vdash \mathtt{t}_1\ \mathtt{t}_2 : \mathtt{T}'} \quad (\textsc{dm-App})$$

$$\frac{\Gamma \vdash \mathtt{t}_1 : \mathtt{S} \qquad \Gamma; \mathtt{z} : \mathtt{S} \vdash \mathtt{t}_2 : \mathtt{T}}{\Gamma \vdash \mathtt{let}\ \mathtt{z} = \mathtt{t}_1\ \mathtt{in}\ \mathtt{t}_2 : \mathtt{T}} \quad (\textsc{dm-Let})$$

$$\frac{\Gamma \vdash \mathtt{t} : \mathtt{T} \qquad \bar{\mathtt{X}}\ \#\ \mathit{ftv}(\Gamma)}{\Gamma \vdash \mathtt{t} : \forall \bar{\mathtt{X}}.\mathtt{T}} \quad (\textsc{dm-Gen})$$

$$\frac{\Gamma \vdash \mathtt{t} : \forall \bar{\mathtt{X}}.\mathtt{T}}{\Gamma \vdash \mathtt{t} : [\vec{\mathtt{X}} \mapsto \vec{\mathtt{T}}]\mathtt{T}} \quad (\textsc{dm-Inst})$$

**Figure 1-3: Typing rules for DM**

Such types are called *instances* of the type scheme $\forall \bar{\mathtt{X}}.\mathtt{T}$. Note that, throughout most of this chapter, we work with *constrained type schemes*, a generalization of DM type schemes (Definition 1.3.2).

*Typing environments*, or environments for short, are used to collect assumptions about an expression's free identifiers.

1.2.19   DEFINITION: An *environment* $\Gamma$ is a finite ordered sequence of pairs of a program identifier and a type scheme. We write $\varnothing$ for the empty environment and ; for the concatenation of environments. An environment may be viewed as a finite mapping from program identifiers to type schemes by letting $\Gamma(\mathtt{x}) = \mathtt{S}$ if and only if $\Gamma$ is of the form $\Gamma_1; \mathtt{x} : \mathtt{S}; \Gamma_2$, where $\Gamma_2$ contains no assumption about $\mathtt{x}$. The set of *defined program identifiers* of an environment $\Gamma$, written $dpi(\Gamma)$, is defined by $dpi(\varnothing) = \varnothing$ and $dpi(\Gamma; \mathtt{x} : \mathtt{S}) = dpi(\Gamma) \cup \{\mathtt{x}\}$. $\qquad\qquad \square$

To complete the definition of Damas and Milner's type system, there remains to define *typing judgements*. A typing judgement takes the form $\Gamma \vdash \mathtt{t} : \mathtt{S}$, where $\mathtt{t}$ is an expression of interest, $\Gamma$ is an environment, which typically contains assumptions about $\mathtt{t}$'s free program identifiers, and $\mathtt{S}$ is a type scheme. Such a judgement may be read: *under assumptions $\Gamma$, the expression $\mathtt{t}$ has the type scheme $\mathtt{S}$*. By abuse of language, it is sometimes said that $\mathtt{t}$ *has type* $\mathtt{S}$. A typing judgement is *valid* (or *holds*) if and only if it may be derived using the rules that appear in Figure 1-3. An expression $\mathtt{t}$ is *well-typed* within the environment $\Gamma$ if and only if some judgement of the form $\Gamma \vdash \mathtt{t} : \mathtt{S}$ holds; it is *ill-typed* within $\Gamma$ otherwise.

Rule DM-VAR allows fetching a type scheme for an identifier $\mathtt{x}$ from the environment. It is equally applicable to program variables, memory locations, and constants. If no assumption concerning $\mathtt{x}$ appears in the environment $\Gamma$, then the rule isn't applicable. In that case, the expression $\mathtt{x}$ is ill-typed within $\Gamma$—can you prove it? Assumptions about constants are usually collected in

a so-called *initial environment* $\Gamma_0$. It is the environment under which closed
programs are typechecked, so every subexpression is typechecked under some
extension $\Gamma$ of $\Gamma_0$. Of course, the type schemes assigned by $\Gamma_0$ to constants
must be consistent with their operational semantics; we say more about this
later (Section 1.7). Rule DM-ABS specifies how to typecheck a $\lambda$-abstraction
$\lambda$z.t. Its premise requires the body of the function, namely t, to be well-typed
under an extra assumption, which causes all free occurrences of z within t to
receive a common type T. Its conclusion forms the arrow type $T \to T'$ out of
the types of the function's formal parameter, namely T, and result, namely
$T'$. It is worth noting that this rule always augments the environment with
a type T—recall that, by convention, types form a subset of type schemes—
but never with a nontrivial type scheme. DM-APP states that the type of a
function application is the codomain of the function's type, provided that the
domain of the function's type is a valid type for the actual argument. DM-
LET closely mirrors the operational semantics: whereas the semantics of the
local definition let z = $t_1$ in $t_2$ is to augment the *runtime* environment
by binding z to the *value* of $t_1$ prior to evaluating $t_2$, the effect of DM-LET
is to augment the *typing* environment by binding z to a *type scheme* for
$t_1$ prior to typechecking $t_2$. DM-GEN turns a type into a type scheme by
universally quantifying over a set of type variables that do not appear free in
the environment; this restriction is discussed in Example 1.2.20 below. DM-
INST, on the contrary, turns a type scheme into one of its instances, which may
be chosen arbitrarily. These two operations are referred to as *generalization*
and *instantiation.* The notion of type scheme and the rules DM-GEN and DM-
INST are characteristic of ML-the-type-system: they distinguish it from the
simply-typed $\lambda$-calculus.

1.2.20    EXAMPLE: It is unsound to allow generalizing type variables that appear free
in the environment. For instance, consider the typing judgement z : X $\vdash$ z :
X **(1)**, which, according to DM-VAR, is valid. Applying an unrestricted version
of DM-GEN to it, we obtain z : X $\vdash$ z : $\forall$X.X **(2)**, whence, by DM-INST, z : X $\vdash$
z : Y **(3)**. By DM-ABS and DM-GEN, we then have $\varnothing \vdash \lambda$z.z : $\forall$XY.X $\to$ Y. In
other words, the identity function has unrelated argument and result types!
Then, the expression $(\lambda$z.z$)\ \hat{0}\ \hat{0}$, which reduces to the stuck expression $\hat{0}\ \hat{0}$,
has type scheme $\forall$Z.Z. So, well-typed programs may cause runtime errors: the
type system is unsound.

     What happened? It is clear that the judgement (1) is correct only because
the type assigned to z is the *same* in its assumption and in its right-hand
side. For the same reason, the judgements (2) and (3)—the former of which
may be written z : X $\vdash$ z : $\forall$Y.Y—are incorrect. Indeed, such judgements de-
feat the very purpose of environments, since they disregard their assumption.

By universally quantifying over X *in the right-hand side only*, we break the connection between occurrences of X in the assumption, which remain free, and occurrences in the right-hand side, which become bound. This is correct only if there are in fact no free occurrences of X in the assumption.                    □

It is a key feature of ML-the-type-system that DM-ABS may only introduce a type T, rather than a type scheme, into the environment. Indeed, this allows the rule's conclusion to form the arrow type $T \to T'$. If instead the rule were to introduce the assumption $z : S$ into the environment, then its conclusion would have to form $S \to T'$, which is not a well-formed type. In other words, this restriction is necessary to preserve the stratification between types and type schemes. If we were to remove this stratification, thus allowing universal quantifiers to appear deep inside types, we would obtain an implicitly-typed version of System F (TAPL Chapter 23). Type inference for System F is undecidable (Wells, 1999), while type inference for ML-the-type-system is decidable, as we show later, so this design choice has a rather drastic impact.

1.2.21     EXERCISE [★, RECOMMENDED]: Build a type derivation for the expression $\lambda z_1.\texttt{let } z_2 = z_1 \texttt{ in } z_2$ within DM.                    □

1.2.22     EXERCISE [★, RECOMMENDED]: Let int be a nullary type constructor of signature $\star$. Let $\Gamma_0$ consist of the bindings $\hat{+} : \texttt{int} \to \texttt{int} \to \texttt{int}$ and $\hat{k} : \texttt{int}$, for every integer $k$. Can you find derivations of the following valid typing judgements? Which of these judgements are valid in the simply-typed $\lambda$-calculus, where $\texttt{let } z = t_1 \texttt{ in } t_2$ is syntactic sugar for $(\lambda z.t_2)\ t_1$?

$$\Gamma_0 \vdash \lambda z.z : \texttt{int} \to \texttt{int}$$
$$\Gamma_0 \vdash \lambda z.z : \forall X.X \to X$$
$$\Gamma_0 \vdash \texttt{let } f = \lambda z.z \hat{+} \hat{1} \texttt{ in } f\ \hat{2} : \texttt{int}$$
$$\Gamma_0 \vdash \texttt{let } f = \lambda z.z \texttt{ in } f\ f\ \hat{2} : \texttt{int}$$

Show that the expressions $\hat{1}\ \hat{2}$ and $\lambda f.(f\ f)$ are ill-typed within $\Gamma_0$. Could these expressions be well-typed in a more powerful type system?                    □

1.2.23     EXERCISE [★★]: In fact, the rules shown in Figure 1-3 are not exactly Damas and Milner's original rules. In (Damas and Milner, 1982), the generalization and instantiation rules are:

$$\frac{\Gamma \vdash t : S \qquad X \notin \mathit{ftv}(\Gamma)}{\Gamma \vdash t : \forall X.S} \tag{DM-GEN'}$$

$$\frac{\Gamma \vdash t : \forall \bar{X}.T \qquad \bar{Y} \mathbin{\#} \mathit{ftv}(\forall \bar{X}.T)}{\Gamma \vdash t : \forall \bar{Y}.[\vec{X} \mapsto \vec{T}]T} \tag{DM-INST'}$$

where $\forall \mathtt{X}.\mathtt{S}$ stands for $\forall \mathtt{X}\bar{\mathtt{X}}.\mathtt{T}$ if $\mathtt{S}$ stands for $\forall \bar{\mathtt{X}}.\mathtt{T}$. Show that the combination of DM-GEN' and DM-INST' is equivalent to the combination of DM-GEN and DM-INST.                                                                                   □

   DM enjoys a number of nice theoretical properties, which have practical implications. First, under suitable hypotheses about the semantics of constants, about the type schemes that they receive in the initial environment, and—in the presence of side effects—under a slight restriction of the syntax of `let` constructs, it is possible to show that the type system is sound: that is, *well-typed (closed) programs do not go wrong.* This essential property ensures that programs that are accepted by the typechecker may be compiled without runtime checks. Furthermore, it is possible to show that there exists an algorithm that, given a (closed) environment $\Gamma$ and a program $\mathtt{t}$, tells whether $\mathtt{t}$ is well-typed with respect to $\Gamma$, and if so, produces a *principal* type scheme $\mathtt{S}$. A principal type scheme is such that (i) it is valid, that is, $\Gamma \vdash \mathtt{t} : \mathtt{S}$ holds, and (ii) it is most general, that is, every judgement of the form $\Gamma \vdash \mathtt{t} : \mathtt{S}'$ follows from $\Gamma \vdash \mathtt{t} : \mathtt{S}$ by DM-INST and DM-GEN. (For the sake of simplicity, we have stated the properties of the type inference algorithm only in the case of a closed environment $\Gamma$; the specification is slightly heavier in the general case.) This implies that *type inference is decidable*: the compiler does not require expressions to be annotated with types. It also implies that, under a fixed environment $\Gamma$, all of the type information associated with an expression $\mathtt{t}$ may be summarized in the form of a single (principal) type scheme, which is very convenient.

## Road map

Before proving the above claims, we first generalize our presentation by moving to a *constraint-based* setting. The necessary tools, namely the constraint language, its interpretation, and a number of constraint equivalence laws, are introduced in Section 1.3. In Section 1.4, we describe the standard constraint-based type system $\mathrm{HM}(X)$ (Odersky, Sulzmann, and Wehr, 1999a; Sulzmann, Müller, and Zenger, 1999; Sulzmann, 2000). We prove that, when constraints are made up of equations between free, finite terms, $\mathrm{HM}(X)$ is a reformulation of DM. In the presence of a more powerful constraint language, $\mathrm{HM}(X)$ is an extension of DM. In Section 1.5, we propose an original reformulation of $\mathrm{HM}(X)$, dubbed $\mathrm{PCB}(X)$, whose distinctive feature is to exploit *type scheme introduction* and *instantiation* constraints. In Section 1.6, we show that, thanks to the extra expressive power afforded by these constraint forms, type inference may be viewed as a combination of constraint generation and constraint solving, as promised earlier. Indeed, we define a constraint generator and relate it with $\mathrm{PCB}(X)$. Then, in Section 1.7, we give a type soundness

theorem. It is stated purely in terms of constraints, but—thanks to the results developed in the previous sections—applies equally to PCB($X$), HM($X$), and DM.

Throughout this core material, the syntax and interpretation of constraints are left partly unspecified. Thus, the development is *parameterized* with respect to them—hence the unknown $X$ in the names HM($X$) and PCB($X$). We really describe a *family* of constraint-based type systems, all of which *share* a common constraint generator and a common type soundness proof. Constraint solving, however, cannot be independent of $X$: on the contrary, the design of an efficient solver is heavily dependent on the syntax and interpretation of constraints. In Section 1.8, we consider constraint solving in the particular case where constraints are made up of equations interpreted in a free tree model, and define a constraint solver on top of a standard first-order unification algorithm.

The remainder of this chapter deals with extensions of the framework. In Section 1.9, we explain how to extend ML-the-calculus with a number of features, including data structures, pattern matching, and type annotations. In Section 1.10, we extend the constraint language with universal quantification and describe a number of extra features that require this extension, including a different flavor of type annotations, polymorphic recursion, and first-class universal and existential types. Last, in Section 1.11, we extend the constraint language with *rows* and describe their applications, which include extensible variants and records.

## 1.3    Constraints

In this section, we define the syntax and logical meaning of constraints. Both are partly unspecified. Indeed, the set of *type constructors* (Definition 1.2.14) must contain at least the binary type constructor $\rightarrow$, but might contain more. Similarly, the syntax of constraints involves a set of so-called *predicates* on types, which we require to contain at least a binary *subtyping* predicate $\leq$, but might contain more. Furthermore, the logical interpretation of type constructors and of predicates is left almost entirely unspecified. This freedom allows reasoning not only about Damas and Milner's type system, but also about a family of constraint-based extensions of it.

Type constructors other than $\rightarrow$ and predicates other than $\leq$ will never explicitly appear in the definition of our constraint-based type systems, precisely because the definition is parametric with respect to them. They can (and usually do) appear in the type schemes assigned to constructors and destructors by the initial environment $\Gamma_0$.

The introduction of subtyping has little impact on the complexity of our

$$
\begin{array}{llr}
\sigma & ::= & \textit{type scheme:}\\
& \forall \bar{\mathtt{x}}[C].\mathtt{T}\\
C, D & ::= & \textit{constraint:}\\
& \mathsf{true} & \textit{truth}\\
& \mathsf{false} & \textit{falsity}\\
& P\,\mathtt{T}_1 \ldots \mathtt{T}_n & \textit{predicate application}\\
& C \wedge C & \textit{conjunction}\\
& \exists \bar{\mathtt{x}}.C & \textit{existential quantification}\\
& \mathsf{def}\;\mathtt{x}:\sigma\;\mathsf{in}\;C & \textit{type scheme introduction}\\
& \mathtt{x} \preceq \mathtt{T} & \textit{type scheme instantiation}\\
\Gamma & ::= & \textit{Typing environments:}
\end{array}
\qquad
\begin{array}{lr}
\varnothing\\
x : \sigma\\
\Gamma; \Gamma\\
C, D \;::= & \textit{Syntactic sugar for constraints:}\\
\ldots & \textit{As before}\\
\sigma \preceq \mathtt{T} & \textit{Definition 1.3.3}\\
\mathsf{let}\;\mathtt{x}:\sigma\;\mathsf{in}\;C & \textit{Definition 1.3.3}\\
\exists \sigma & \textit{Definition 1.3.3}\\
\mathsf{def}\;\Gamma\;\mathsf{in}\;C & \textit{Definition 1.3.4}\\
\mathsf{let}\;\Gamma\;\mathsf{in}\;C & \textit{Definition 1.3.4}\\
\exists \Gamma & \textit{Definition 1.3.4}
\end{array}
$$

**Figure 1-4: Syntax of type schemes and constraints**

proofs, yet increases the framework's expressive power. When subtyping is not desired, we interpret the predicate $\leq$ as equality.

**Syntax**

We now define the syntax of constrained type schemes and of constraints, and introduce some extra constraint forms as syntactic sugar.

1.3.1 DEFINITION: Let $P$ range over a finite or denumerable set of *predicates*, each of which has a *signature* of the form $\kappa_1 \otimes \ldots \otimes \kappa_n \Rightarrow \cdot$, where $n \geq 0$. Let $\leq$ be a distinguished predicate of signature $\star \otimes \star \Rightarrow \cdot$. □

1.3.2 DEFINITION: The syntax of *type schemes* and *constraints* is given in Figure 1-4. It is further restricted by the following requirements. In the type scheme $\forall \bar{\mathtt{x}}[C].\mathtt{T}$ and in the constraint $\mathtt{x} \preceq \mathtt{T}$, the type $\mathtt{T}$ must have kind $\star$. In the constraint $P\,\mathtt{T}_1 \ldots \mathtt{T}_n$, the types $\mathtt{T}_1, \ldots, \mathtt{T}_n$ must have kind $\kappa_1, \ldots, \kappa_n$, respectively, if $P$ has signature $\kappa_1 \otimes \ldots \otimes \kappa_n \Rightarrow \cdot$. We write $\forall \bar{\mathtt{x}}.\mathtt{T}$ for $\forall \bar{\mathtt{x}}[\mathsf{true}].\mathtt{T}$, which allows viewing DM type schemes as a subset of constrained type schemes. □

We write $\mathtt{T}_1 \leq \mathtt{T}_2$ for the binary predicate application $\leq \mathtt{T}_1\,\mathtt{T}_2$, and call it a subtyping constraint. By convention, $\exists$ and $\mathsf{def}$ bind tighter than $\wedge$; that is, $\exists \bar{\mathtt{x}}.C \wedge D$ is $(\exists \bar{\mathtt{x}}.C) \wedge D$ and $\mathsf{def}\;\mathtt{x}:\sigma\;\mathsf{in}\;C \wedge D$ is $(\mathsf{def}\;\mathtt{x}:\sigma\;\mathsf{in}\;C) \wedge D$. In $\forall \bar{\mathtt{x}}[C].\mathtt{T}$, the type variables $\bar{\mathtt{x}}$ are bound within $C$ and $\mathtt{T}$. In $\exists \bar{\mathtt{x}}.C$, the type variables $\bar{\mathtt{x}}$ are bound within $C$. The sets of free type variables of a type scheme $\sigma$ and of a constraint $C$, written $ftv(\sigma)$ and $ftv(C)$, respectively, are defined accordingly. In $\mathsf{def}\;\mathtt{x}:\sigma\;\mathsf{in}\;C$, the identifier $\mathtt{x}$ is bound within $C$. The sets

of free program identifiers of a type scheme $\sigma$ and of a constraint $C$, written $fpi(\sigma)$ and $fpi(C)$, respectively, are defined accordingly. Please note that x occurs free in the constraint $\text{x} \preceq \text{T}$.

We immediately introduce a number of derived constraint forms:

1.3.3    DEFINITION:  Let $\sigma$ be $\forall \bar{\text{X}}[C].\text{T}$. If $\bar{\text{X}} \,\#\, ftv(\text{T}')$ holds, then $\sigma \preceq \text{T}'$ (read: $\text{T}'$ *is an instance of* $\sigma$) stands for the constraint $\exists \bar{\text{X}}.(C \wedge \text{T} \leq \text{T}')$. We write $\exists \sigma$ (read: $\sigma$ *has an instance*) for $\exists \bar{\text{X}}.C$ and $\mathsf{let}\ \text{x} : \sigma\ \mathsf{in}\ C$ for $\exists \sigma \wedge \mathsf{def}\ \text{x} : \sigma\ \mathsf{in}\ C$.    □

Constrained type schemes generalize Damas and Milner's type schemes, while our definition of instantiation constraints generalizes Damas and Milner's instance relation (Definition 1.2.18). Let us draw a comparison. First, Damas and Milner's instance relation yields a "yes/no" answer, and is purely syntactic: for instance, the type $\text{Y} \rightarrow \text{Z}$ is *not* an instance of $\forall \text{X}.\text{X} \rightarrow \text{X}$ in Damas and Milner's sense, because $\text{Y}$ and $\text{Z}$ are distinct type variables. In our presentation, on the other hand, $\forall \text{X}.\text{X} \rightarrow \text{X} \preceq \text{Y} \rightarrow \text{Z}$ is not an assertion; rather, it is a constraint, which by definition is $\exists \text{X}.(\mathsf{true} \wedge \text{X} \rightarrow \text{X} \leq \text{Y} \rightarrow \text{Z})$. We later prove that it is equivalent to $\exists \text{X}.(\text{Y} \leq \text{X} \wedge \text{X} \leq \text{Z})$ and to $\text{Y} \leq \text{Z}$, or, if subtyping is interpreted as equality, to $\text{Y} = \text{Z}$. That is, $\sigma \preceq \text{T}'$ represents a condition on (the types denoted by) the type variables in $ftv(\sigma, \text{T}')$ for $\text{T}'$ to be an instance of $\sigma$, in a logical, rather than purely syntactic, sense. Second, the definition of instantiation constraints involves subtyping, so as to ensure that any supertype of an instance of $\sigma$ is again an instance of $\sigma$ (see rule C-ExTRANS in Figure 1-6 and Lemma 1.3.17). This is consistent with the purpose of subtyping, which is to allow supplying a subtype where a supertype is expected (TAPL Chapter 15). Third and last, every type scheme now carries a constraint. The constraint $C$, whose free type variables may or may not be members of $\bar{\text{X}}$, restricts the instances of the type scheme $\forall \bar{\text{X}}[C].\text{T}$. This is expressed in the instantiation constraint $\exists \bar{\text{X}}.(C \wedge \text{T} \leq \text{T}')$, where the values that $\bar{\text{X}}$ may assume are restricted by the requirement that $C$ be satisfied. This requirement vanishes in the case of DM type schemes, where $C$ is $\mathsf{true}$. Our notions of constrained type scheme and of instantiation constraint are standard: they are exactly those of HM($X$) (Odersky, Sulzmann, and Wehr, 1999a).

The constraint $\mathsf{true}$, which is always satisfied, mainly serves to indicate the absence of a nontrivial constraint, while $\mathsf{false}$, which has no solution, may be understood as the indication of a type error. Composite constraints include conjunction and existential quantification, which have their standard meaning, as well as *type scheme introduction* and *type scheme instantiation* constraints, which similar to Gustavsson and Svenningsson's constraint abstractions (2001b). In short, the construct $\mathsf{def}\ \text{x} : \sigma\ \mathsf{in}\ C$ binds the name x to the type scheme $\sigma$ within the constraint $C$. If $C$ contains a subconstraint of

the form x $\preceq$ T, where this occurrence of x is free in $C$, then this subconstraint
acquires the meaning $\sigma \preceq$ T. Thus, the constraint x $\preceq$ T is indeed an instantia-
tion constraint, where the type scheme that is being instantiated is referred to
by name. The constraint def x : $\sigma$ in $C$ may be viewed as an *explicit substitu-
tion* of the type scheme $\sigma$ for the name x within $C$. Later (Section 1.5), we use
such explicit substitutions to supplant typing environments. That is, where
Damas and Milner's type system augments the current typing environment
(DM-ABS, DM-LET), we introduce a new def binding in the current constraint;
where it looks up the current typing environment (DM-VAR), we employ an in-
stantiation constraint. The point is that it is then up to a constraint solver to
choose a strategy for reducing explicit substitutions—for instance, one might
wish to *simplify* $\sigma$ before substituting it for x within $C$—whereas the use of
environments in standard type systems such as DM and HM($X$) imposes an
eager substitution strategy, which is inefficient and thus never literally imple-
mented. The use of type scheme introduction and instantiation constraints
allows separating constraint generation and constraint solving *without com-
promising efficiency*, or, in other words, without introducing a gap between
the description of the type inference algorithm and its actual implementation.
Although the algorithm that we plan to describe is not new, its description in
terms of constraints is: to the best of our knowledge, the only close relative of
our def constraints is to be found in (Gustavsson and Svenningsson, 2001b).
Fähndrich, Rehof, and Das's instantiation constraints (2000) are also related,
but may be recursive and are meant to be solved using a semi-unification
procedure, as opposed to a unification algorithm extended with facilities for
creating and instantiating type schemes, as in our case.

One consequence of introducing constraints inside type schemes is that some
type schemes have no instances at all, or have instances only if a certain
constraint holds. For instance, the type scheme $\sigma = \forall$X[bool = int].X, where
the nullary type constructors int and bool have distinct interpretations, has
no instances; that is, no constraint of the form $\sigma \preceq$ T$'$ has a solution. The
type scheme $\sigma = \forall$Z[X = Y $\rightarrow$ Z].Z has an instance only if X = Y $\rightarrow$ Z holds
for some Z; in other words, for every T$'$, $\sigma \preceq$ T$'$ entails $\exists$Z.(X = Y $\rightarrow$ Z).
(We define *entailment* on page 29.) We later prove that the constraint $\exists\sigma$
is equivalent to $\exists$Z.$\sigma \preceq$ Z, where Z $\notin$ *ftv*($\sigma$) (Exercise 1.3.23). That is, $\exists\sigma$
expresses the requirement that $\sigma$ have an instance. Type schemes that do not
have an instance indicate a type error, so in many situations, one wishes to
avoid them; for this reason, we often use the constraint form let x : $\sigma$ in $C$,
which requires $\sigma$ to have an instance and at the same time associates it with
the name x. Because the def form is more primitive, it is easier to work with
at a low level, but it is no longer explicitly used after Section 1.3; we always
use let instead.

1.3.4   DEFINITION:  Environments $\Gamma$ remain as in Definition 1.2.19, except DM type schemes S are replaced with constrained type schemes $\sigma$. We write $dfpi(\Gamma)$ for $dpi(\Gamma) \cup fpi(\Gamma)$. We define $\mathsf{def}\ \varnothing\ \mathsf{in}\ C = C$ and $\mathsf{def}\ \Gamma; \mathsf{x} : \sigma\ \mathsf{in}\ C = \mathsf{def}\ \Gamma\ \mathsf{in}\ \mathsf{def}\ \mathsf{x} : \sigma\ \mathsf{in}\ C$. Similarly, we define $\mathsf{let}\ \varnothing\ \mathsf{in}\ C = C$ and $\mathsf{let}\ \Gamma; \mathsf{x} : \sigma\ \mathsf{in}\ C = \mathsf{let}\ \Gamma\ \mathsf{in}\ \mathsf{let}\ \mathsf{x} : \sigma\ \mathsf{in}\ C$. We define $\exists\varnothing = \mathsf{true}$ and $\exists(\Gamma; \mathsf{x} : \sigma) = \exists\Gamma \wedge \mathsf{def}\ \Gamma\ \mathsf{in}\ \exists\sigma$.     □

In order to establish or express certain laws of equivalence between constraints, we need *constraint contexts*. A context is a constraint with zero, one, or several *holes*, written []. The syntax of contexts is as follows:

$$\mathcal{C} ::= [] \mid C \mid \mathcal{C} \wedge \mathcal{C} \mid \exists\bar{\mathsf{x}}.\mathcal{C} \mid \mathsf{def}\ \mathsf{x} : \sigma\ \mathsf{in}\ \mathcal{C} \mid \mathsf{def}\ \mathsf{x} : \forall\bar{\mathsf{x}}[\mathcal{C}].\mathtt{T}\ \mathsf{in}\ C$$

The application of a constraint context $\mathcal{C}$ to a constraint $C$, written $\mathcal{C}[C]$, is defined in the usual way. Because a context may have any number of holes, $C$ may disappear or be duplicated in the process. Because a hole may appear in the scope of a binder, some of $C$'s free type variables and free program identifiers may become bound in $\mathcal{C}[C]$. We write $dtv(\mathcal{C})$ and $dpi(\mathcal{C})$ for the sets of type variables and program identifiers, respectively, that $\mathcal{C}$ may thus capture. We write $\mathsf{let}\ \mathsf{x} : \forall\bar{\mathsf{x}}[\mathcal{C}].\mathtt{T}\ \mathsf{in}\ C$ for $\exists\bar{\mathsf{x}}.\mathcal{C} \wedge \mathsf{def}\ \mathsf{x} : \forall\bar{\mathsf{x}}[\mathcal{C}].\mathtt{T}\ \mathsf{in}\ C$. Being able to state such a definition is why we require multi-hole contexts. We let range over *existential constraint contexts*, defined by $\mathcal{X} ::= [] \mid \exists\bar{\mathsf{x}}.\mathcal{X}$.

### Meaning

We have defined the syntax of constraints and given an informal description of their meaning. We now give a formal definition of the interpretation of constraints. We begin with the definition of a *model*:

1.3.5   DEFINITION:  For every kind $\kappa$, let $\mathcal{M}_\kappa$ be a nonempty set, whose elements are the *ground types* of kind $\kappa$. In the following, $t$ ranges over $\mathcal{M}_\kappa$, for some $\kappa$ that may be determined from the context. For every type constructor $F$ of signature $K \Rightarrow \kappa$, let $F$ denote a total function from $\mathcal{M}_K$ into $\mathcal{M}_\kappa$, where the indexed product $\mathcal{M}_K$ is the set of all mappings of domain $dom(K)$ that map every $d \in dom(K)$ to an element of $\mathcal{M}_{K(d)}$. For every predicate $P$ of signature $\kappa_1 \otimes \ldots \otimes \kappa_n \Rightarrow \cdot$, let $P$ denote a predicate on $\mathcal{M}_{\kappa_1} \times \ldots \times \mathcal{M}_{\kappa_n}$. We require the predicate $\leq$ on $\mathcal{M}_\star \times \mathcal{M}_\star$ to be a partial order.     □

For the sake of convenience, we abuse notation and write $F$ for both the type constructor and its interpretation; similarly for predicates. We freely assume that a binary equality predicate, whose interpretation is equality on $\mathcal{M}_\kappa$, is available at every kind $\kappa$, so $\mathtt{T}_1 = \mathtt{T}_2$, where $\mathtt{T}_1$ and $\mathtt{T}_2$ have kind $\kappa$, is a well-formed constraint.

By varying the set of type constructors, the set of predicates, the set of ground types, and the interpretation of type constructors and predicates, one may define an entire family of related type systems. We informally refer to the collection of these choices as $X$. Thus, the type systems $\text{HM}(X)$ and $\text{PCB}(X)$, described in Sections 1.4 and 1.5, are *parameterized* by $X$.

The following examples give standard ways of defining the set of ground types and the interpretation of type constructors.

1.3.6    EXAMPLE [SYNTACTIC MODELS]: For every kind $\kappa$, let $\mathcal{M}_\kappa$ consist of the *closed* types of kind $\kappa$. Then, ground types are types that do not have any free type variables, and form the so-called *Herbrand universe*. Let every type constructor $F$ be interpreted as itself. Models that define ground types and interpret type constructors in this manner are referred to as *syntactic*.     □

1.3.7    EXAMPLE [TREE MODELS]: Let a *path* $\pi$ be a finite sequence of directions. The empty path is written $\epsilon$ and the concatenation of the paths $\pi$ and $\pi'$ is written $\pi \cdot \pi'$. Let a *tree* be a partial function $t$ from paths to type constructors whose domain is nonempty and prefix-closed and such that, for every path $\pi$ in the domain of $t$, if the type constructor $t(\pi)$ has signature $K \Rightarrow \kappa$, then $\pi \cdot d \in dom(t)$ is equivalent to $d \in dom(K)$ and, furthermore, for every $d \in dom(K)$, the type constructor $t(\pi \cdot d)$ has image kind $K(d)$. If $\pi$ is in the domain of $t$, then the *subtree* of $t$ rooted at $\pi$, written $t/\pi$, is the partial function $\pi' \mapsto t(\pi \cdot \pi')$. A tree is *finite* if and only if it has finite domain. A tree is *regular* if and only if it has a finite number of distinct subtrees. Every finite tree is thus regular. Let $\mathcal{M}_\kappa$ consist of the *finite* (resp. *regular*) trees $t$ such that $t(\epsilon)$ has image kind $\kappa$: then, we have a *finite* (resp. *regular*) *tree model*.

If $F$ has signature $K \Rightarrow \kappa$, one may interpret $F$ as the function that maps $T \in \mathcal{M}_K$ to the ground type $t \in \mathcal{M}_\kappa$ defined by $t(\epsilon) = F$ and $t/d = T(d)$ for $d \in dom(T)$, that is, the unique ground type whose head symbol is $F$ and whose subtree rooted at $d$ is $T(d)$. Then, we have a *free* tree model. Please note that free finite tree models coincide with syntactic models, as defined in the previous example.     □

Rows (Section 1.11) are interpreted in a tree model, albeit not a free one. The following examples suggest different ways of interpreting the subtyping predicate.

1.3.8    EXAMPLE [EQUALITY MODELS]: The simplest way of interpreting the subtyping predicate is to let $\leq$ denote equality on every $\mathcal{M}_\kappa$. Models that do so are referred to as *equality models*. When no predicate other than equality is available, we say that the model is *equality-only*.     □

1.3.9   EXAMPLE [STRUCTURAL, NONSTRUCTURAL SUBTYPING]: Let a *variance* $\nu$ be a nonempty subset of $\{-, +\}$, written $-$ (*contravariant*), $+$ (*covariant*), or $\pm$ (*invariant*) for short. Define the *composition* of two variances as an associative commutative operation with $+$ as neutral element and such that $-- = +$ and $\pm- = \pm\pm = \pm$. Now, consider a free (finite or regular) tree model, where every direction $d$ comes with a fixed variance $\nu(d)$. Define the variance $\nu(\pi)$ of a path $\pi$ as the composition of the variances of its elements. Let $\leqslant$ be a partial order on type constructors such that (i) if $F_1 \leqslant F_2$ holds and $F_1$ and $F_2$ have signature $K_1 \Rightarrow \kappa_1$ and $K_2 \Rightarrow \kappa_2$, respectively, then $K_1$ and $K_2$ agree on the intersection of their domains and $\kappa_1$ and $\kappa_2$ coincide; and (ii) $F_0 \leqslant F_1 \leqslant F_2$ implies $dom(F_0) \cap dom(F_2) \subseteq dom(F_1)$. Let $\leqslant^+$, $\leqslant^-$, and $\leqslant^\pm$ stand for $\leqslant$, $\geqslant$, and $=$, respectively. Then, define the interpretation of subtyping as follows: if $t_1, t_2 \in \mathcal{M}_\kappa$, let $t_1 \leq t_2$ hold if and only if, for every path $\pi \in dom(t_1) \cap dom(t_2)$, $t_1(\pi) \leqslant^{\nu(\pi)} t_2(\pi)$ holds. It is not difficult to check that $\leq$ is a partial order on every $\mathcal{M}_\kappa$. The reader is referred to (Kozen, Palsberg, and Schwartzbach., 1995) for more details about this construction. Models that define subtyping in this manner are referred to as *nonstructural subtyping models.*

A simple nonstructural subtyping model is obtained by letting the directions *domain* and *codomain* be contra- and covariant, respectively, and introducing, in addition to the type constructor $\rightarrow$, two type constructors $\bot$ and $\top$ of signature $\star$. This gives rise to a model where $\bot$ is the least ground type, $\top$ is the greatest ground type, and the arrow type constructor is, as usual, contravariant in its domain and covariant in its codomain.

A typical use of nonstructural subtyping is in type systems for records. One may, for instance, introduce a covariant direction *content* of kind $\star$, a kind $\diamond$, a type constructor abs of signature $\diamond$, a type constructor pre of signature $\{content \mapsto \star\} \Rightarrow \diamond$, and let pre $\leqslant$ abs. This gives rise to a model where pre $t \leq$ abs holds for every $t \in \mathcal{M}_\star$. This form of subtyping is called *nonstructural* because comparable ground types may have different shapes, such as $\bot$ and $\bot \rightarrow \top$, or pre $\top$ and abs. Nonstructural subtyping has been studied, for example, in (Kozen, Palsberg, and Schwartzbach., 1995; Palsberg, Wand, and O'Keefe, 1997; Pottier, 2001b; Niehren and Priesnitz, 2003). Section 1.11 says more about typechecking operations on records.

An important particular case arises when *any two type constructors related by $\leqslant$ have the same arity.* In that case, it is not difficult to show that *any two ground types related by subtyping must have the same shape*, that is, if $t_1 \leq t_2$ holds, then $dom(t_1)$ and $dom(t_2)$ coincide. For this reason, such an interpretation of subtyping is usually referred to as *atomic* or *structural* subtyping. It has been studied in the finite (Mitchell, 1984, 1991b; Frey, 1997; Rehof, 1997; Kuncak and Rinard, 2003; Simonet, 2003) and regular (Tiuryn and Wand,

1993) cases. Structural subtyping is often used in automated program analyses that enrich standard types with atomic annotations without altering their shape. □

Our last example suggests a predicate other than equality and subtyping.

1.3.10 EXAMPLE [CONDITIONAL CONSTRAINTS]: Consider a nonstructural subtyping model. For every type constructor $F$ of image kind $\kappa$ and for every kind $\kappa'$, let $(F \leqslant \cdot \Rightarrow \cdot \leq \cdot)$ be a predicate of signature $\kappa \otimes \kappa' \otimes \kappa' \Rightarrow \cdot$. Thus, if $\mathtt{T}_0$ has kind $\kappa$ and $\mathtt{T}_1, \mathtt{T}_2$ have the same kind, then $F \leqslant \mathtt{T}_0 \Rightarrow \mathtt{T}_1 \leq \mathtt{T}_2$ is a well-formed constraint, called a *conditional subtyping constraint*. Its interpretation is defined as follows: if $t_0 \in \mathcal{M}_\kappa$ and $t_1, t_2 \in \mathcal{M}_{\kappa'}$, then $F \leqslant t_0 \Rightarrow t_1 \leq t_2$ holds if and only if $F \leqslant t_0(\epsilon)$ implies $t_1 \leq t_2$. In other words, if $t_0$'s head symbol exceeds $F$ according to the ordering on type constructors, then the subtyping constraint $t_1 \leq t_2$ must hold; otherwise, the conditional constraint holds vacuously. Conditional constraints have been studied *e.g.* in (Reynolds, 1969a; Heintze, 1993; Aiken, Wimmers, and Lakshman, 1994; Pottier, 2000; Su and Aiken, 2001). □

Many other kinds of constraints exist; see *e.g.* (Comon, 1993).

Throughout this chapter, we assume (unless stated otherwise) that the set of type constructors, the set of predicates, and the model—which, together, form the parameter $X$—are arbitrary and fixed.

As usual, the meaning of a constraint is a function of the meaning of its free type variables, which is given by a *ground assignment*. The meaning of free program identifiers may be defined as part of the constraint, if desired, using a def prefix, so it need not be given by a separate assignment.

1.3.11 DEFINITION: A *ground assignment* $\phi$ is a total, kind-preserving mapping from $\mathcal{V}$ into $\mathcal{M}$. Ground assignments are extended to types by $\phi(F\,\mathtt{T}_1\,\ldots\,\mathtt{T}_n) = F(\phi(\mathtt{T}_1), \ldots, \phi(\mathtt{T}_n))$. Then, for every type $\mathtt{T}$ of kind $\kappa$, $\phi(\mathtt{T})$ is a ground type of kind $\kappa$. Whether a constraint $C$ holds under a ground assignment $\phi$, written $\phi \vdash C$ (read: $\phi$ *satisfies* $C$), is defined by the rules in Figure 1-5. A constraint $C$ is *satisfiable* if and only if $\phi \vdash C$ holds for some $\phi$. It is *false* if and only if $\phi \vdash \mathsf{def}\ \Gamma\ \mathsf{in}\ C$ holds for *no* ground assignment $\phi$ and environment $\Gamma$. □

Let us now explain the rules that define constraint satisfaction (Figure 1-5). They are syntax-directed: that is, to a given constraint, at most one rule applies. It is determined by the nature of the first construct that appears under a maximal def prefix. CM-TRUE states that a constraint of the form def $\Gamma$ in true is a tautology, that is, holds under every ground assignment. No rule matches constraints of the form def $\Gamma$ in false, which means that such constraints do not have a solution. CM-PREDICATE states that the meaning

$$\phi \vdash \mathsf{def}\ \Gamma\ \mathsf{in}\ \mathsf{true} \qquad (\text{CM-T{\scriptsize RUE}})$$

$$\frac{P(\phi(\mathtt{T}_1),\ldots,\phi(\mathtt{T}_n))}{\phi \vdash \mathsf{def}\ \Gamma\ \mathsf{in}\ P\ \mathtt{T}_1\ \ldots\ \mathtt{T}_n}\ (\text{CM-P{\scriptsize REDICATE}})$$

$$\frac{\phi \vdash \mathsf{def}\ \Gamma\ \mathsf{in}\ C_1 \qquad \phi \vdash \mathsf{def}\ \Gamma\ \mathsf{in}\ C_2}{\phi \vdash \mathsf{def}\ \Gamma\ \mathsf{in}\ (C_1 \wedge C_2)}\ (\text{CM-A{\scriptsize ND}})$$

$$\frac{\begin{array}{c}\phi[\vec{\mathtt{X}} \mapsto \vec{t}] \vdash \mathsf{def}\ \Gamma\ \mathsf{in}\ C \\ \bar{\mathtt{X}}\ \#\ ftv(\Gamma)\end{array}}{\phi \vdash \mathsf{def}\ \Gamma\ \mathsf{in}\ \exists\bar{\mathtt{X}}.C}\ (\text{CM-E{\scriptsize XISTS}})$$

$$\frac{\begin{array}{c}\phi \vdash \mathsf{def}\ \Gamma_1\ \mathsf{in}\ \sigma \preceq \mathtt{T}' \\ \mathtt{x} \notin dpi(\Gamma_2)\end{array}}{\phi \vdash \mathsf{def}\ \Gamma_1; \mathtt{x}:\sigma; \Gamma_2\ \mathsf{in}\ \mathtt{x} \preceq \mathtt{T}'}\ (\text{CM-I{\scriptsize NSTANCE}})$$

**Figure 1-5: Meaning of constraints**

of a predicate application is given by the predicate's interpretation within the model. More specifically, if $P$'s signature is $\kappa_1 \otimes \ldots \otimes \kappa_n \Rightarrow \cdot$, then, by well-formedness of the constraint, every $\mathtt{T}_i$ is of kind $\kappa_i$, so $\phi(\mathtt{T}_i)$ is a ground type in $\mathcal{M}_{\kappa_i}$. By Definition 1.3.5, $P$ denotes a predicate on $\mathcal{M}_{\kappa_1} \times \ldots \times \mathcal{M}_{\kappa_n}$, so the rule's premise is mathematically well-formed. It is independent of $\Gamma$, which is natural, since a predicate application has no free program identifiers. CM-A{\scriptsize ND} requires each of the conjuncts to be valid in isolation. The information in $\Gamma$ is made available to each branch. CM-E{\scriptsize XISTS} allows the type variables $\vec{\mathtt{X}}$ to denote arbitrary ground types $\vec{t}$ within $C$, independently of their image through $\phi$. We implicitly require $\vec{\mathtt{X}}$ and $\vec{t}$ to have matching kinds, so that $\phi[\vec{\mathtt{X}} \mapsto \vec{t}]$ remains a kind-preserving ground assignment. The side condition $\bar{\mathtt{X}}\ \#\ ftv(\Gamma)$—which may always be satisfied by suitable $\alpha$-conversion of the constraint $\exists\bar{\mathtt{X}}.C$—prevents free occurrences of the type variables $\bar{\mathtt{X}}$ within $\Gamma$ from being unduly affected. CM-I{\scriptsize NSTANCE} concerns constraints of the form $\mathsf{def}\ \Gamma\ \mathsf{in}\ \mathtt{x} \preceq \mathtt{T}'$. The constraint $\mathtt{x} \preceq \mathtt{T}'$ is turned into $\sigma \preceq \mathtt{T}'$, where, according to the second premise, $\sigma$ is $\Gamma(\mathtt{x})$. Please recall that constraints of such a form were introduced in Definition 1.3.3. The environment $\Gamma$ is replaced with a suitable prefix of itself, namely $\Gamma_1$, so that the free program identifiers of $\sigma$ retain their meaning.

It is intuitively clear that the constraints $\mathsf{def}\ \mathtt{x}:\sigma\ \mathsf{in}\ C$ and $[\mathtt{x} \mapsto \sigma]C$ have the same meaning, where the latter denotes the capture-avoiding substitution of $\sigma$ for $\mathtt{x}$ throughout $C$. As a matter of fact, it would have been possible to use this equivalence as a *definition* of the meaning of $\mathsf{def}$ constraints, but the present style is pleasant as well. This confirms our (informal) claim that the $\mathsf{def}$ form is an explicit substitution form.

It is possible for a constraint to be neither satisfiable nor false. Consider, for instance, the constraint $\exists \mathtt{Z}.\mathtt{x} \preceq \mathtt{Z}$. Because the identifier $\mathtt{x}$ is free, CM-I{\scriptsize NSTANCE} is not applicable, so the constraint is not satisfiable. Furthermore,

placing it within the context let x : ∀X.X in [] makes it satisfied by every ground assignment, so it is not false. Here, the assertions "$C$ is satisfiable" and "$C$ is false" are opposite when $fpi(C) = \varnothing$ holds, whereas in a standard first-order logic, they always are.

In a judgement of the form $\phi \vdash C$, the ground assignment $\phi$ applies to the free type variables of $C$. This is made precise by the following statements. In the second one, $\circ$ is composition and $\theta(C)$ is the capture-avoiding application of the type substitution $\theta$ to $C$.

1.3.12 LEMMA: If $\bar{\mathbf{x}} \mathbin{\#} ftv(C)$ holds, then $\phi \vdash C$ and $\phi[\vec{\mathbf{x}} \mapsto \vec{t}] \vdash C$ are equivalent. □

1.3.13 LEMMA: $\phi \circ \theta \vdash C$ and $\phi \vdash \theta(C)$ are equivalent. □

### Reasoning with constraints

Because constraints lie at the heart of our treatment of ML-the-type-system, most of our proofs involve establishing logical properties of constraints, that is, *entailment* or *equivalence* assertions. Let us first define these notions.

1.3.14 DEFINITION: We write $C_1 \Vdash C_2$, and say that $C_1$ *entails* $C_2$, if and only if, for every ground assignment $\phi$ and for every environment $\Gamma$, $\phi \vdash$ def $\Gamma$ in $C_1$ implies $\phi \vdash$ def $\Gamma$ in $C_2$. We write $C_1 \equiv C_2$, and say that $C_1$ and $C_2$ are *equivalent*, if and only if $C_1 \Vdash C_2$ and $C_2 \Vdash C_1$ hold. □

This definition measures the strength of a constraint by the set of pairs $(\phi, \Gamma)$ that satisfy it, and considers a constraint stronger if fewer such pairs satisfy it. In other words, $C_1$ entails $C_2$ when $C_1$ imposes stricter requirements on its free type variables and program identifiers than $C_2$ does. We remark that $C$ is false if and only if $C \equiv$ false holds. It is straightforward to check that entailment is reflexive and transitive and that $\equiv$ is indeed an equivalence relation.

We immediately exploit the notion of constraint equivalence to define what it means for a type constructor to be covariant, contravariant, or invariant with respect to one of its parameters. Let $F$ be a type constructor of signature $\kappa_1 \otimes \ldots \otimes \kappa_n \Rightarrow \kappa$. Let $i \in \{1, \ldots, n\}$. $F$ is *covariant* (resp. *contravariant*, *invariant*) with respect to its $i^{\text{th}}$ parameter if and only if, for all types $\mathtt{T}_1, \ldots, \mathtt{T}_n$ and $\mathtt{T}_i'$ of appropriate kinds, the constraint $F\,\mathtt{T}_1 \ldots \mathtt{T}_i \ldots \mathtt{T}_n \leq F\,\mathtt{T}_1 \ldots \mathtt{T}_i' \ldots \mathtt{T}_n$ is equivalent to $\mathtt{T}_i \leq \mathtt{T}_i'$ (resp. $\mathtt{T}_i' \leq \mathtt{T}_i$, $\mathtt{T}_i = \mathtt{T}_i'$). We let the reader check the following facts: (i) in an equality model, these three notions coincide; (ii) in an equality free tree model, every type constructor is invariant with respect to each of its parameters; and (iii) in a nonstructural subtyping model, if the direction $d$ has been declared covariant (resp. contravariant, invariant), then every type constructor whose arity includes $d$ is covariant (resp. contravariant,

invariant) with respect to $d$. In the following, *we require the type constructor*
$\rightarrow$ *to be contravariant with respect to its domain and covariant with respect to*
*its codomain*—a standard requirement in type systems with subtyping (TAPL
Chapter 15). These properties are summed up by the following equivalence
law:

$$\text{T}_1 \rightarrow \text{T}_2 \leq \text{T}'_1 \rightarrow \text{T}'_2 \equiv \text{T}'_1 \leq \text{T}_1 \wedge \text{T}_2 \leq \text{T}'_2 \qquad\qquad (\text{C-Arrow})$$

Please note that this is a high-level requirement about the interpretation of
types and of the subtyping predicate. In an equality free tree model, for in-
stance, it is always satisfied. In a nonstructural subtyping model, it boils
down to requiring that the directions *domain* and *codomain* be declared con-
travariant and covariant, respectively. In the general case, we do not have any
knowledge of the model, and cannot formulate a more precise requirement.
Thus, it is up to the designer of the model to ensure that C-Arrow holds.

   We also exploit the notion of constraint equivalence to define what it means
for two type constructors to be incompatible. Two type constructors $F_1$ and $F_2$
with the same image kind are *incompatible* if and only if all constraints of the
form $F_1 \vec{\text{T}}_1 \leq F_2 \vec{\text{T}}_2$ and $F_2 \vec{\text{T}}_2 \leq F_1 \vec{\text{T}}_1$ are false; then, we write $F_1 \bowtie F_2$. Please
note that in an equality free tree model, any two distinct type constructors are
incompatible. In the following, we often indicate that a newly introduced type
constructor must be *isolated*. We implicitly require that, whenever each of $F_1$
and $F_2$ is isolated, $F_1$ and $F_2$ be incompatible. Thus, the notion of "isolation"
provides a concise and modular way of stating a collection of incompatibility
requirements. We consider the type constructor $\rightarrow$ isolated.

   Entailment is preserved by arbitrary constraint contexts, as stated by
the following theorem. As a result, constraint equivalence is a congruence.
Throughout this chapter, these facts are often used implicitly.

1.3.15   Theorem [Congruence]:  $C_1 \Vdash C_2$ implies $\mathcal{C}[C_1] \Vdash \mathcal{C}[C_2]$.                    □

   We now give a series of lemmas that provide useful entailment laws.
   The following is a standard property of existential quantification.

1.3.16   Lemma:  $C \Vdash \exists \bar{\text{X}}.C$.                    □

   The following lemma states that any supertype of an instance of $\sigma$ is also
an instance of $\sigma$.

1.3.17   Lemma:  $\sigma \preceq \text{T} \wedge \text{T} \leq \text{T}' \Vdash \sigma \preceq \text{T}'$.                    □

   The next lemma gives another interesting simplification law.

1.3.18   Lemma:  $\text{X} \notin \mathit{ftv}(\text{T})$ implies $\exists \text{X}.(\text{X} = \text{T}) \equiv \mathsf{true}$.                    □

   The following lemma states that, provided $D$ is satisfied, the type T is an
instance of the constrained type scheme $\forall \bar{\text{X}}[D].\text{T}$.

1.3.19    LEMMA:  $D \Vdash \forall \bar{\mathtt{x}}[D].\mathtt{T} \preceq \mathtt{T}$.                                                               □

This technical lemma helps justify Definition 1.3.21 below.

1.3.20    LEMMA:  Let $\mathtt{Z} \notin ftv(C, \sigma, \mathtt{T})$. Then, $C \Vdash \sigma \preceq \mathtt{T}$ holds if and only if $C \wedge \mathtt{T} \leq$ $\mathtt{Z} \Vdash \sigma \preceq \mathtt{Z}$ holds.                                                                   □

It is useful to define what it means for a type scheme $\sigma_1$ to be *more general* than a type scheme $\sigma_2$. Our informal intent is for $\sigma_1 \preceq \sigma_2$ to mean: *every instance of $\sigma_2$ is an instance of $\sigma_1$*. In Definition 1.3.3, we have introduced the constraint form $\sigma \preceq \mathtt{T}$ as syntactic sugar. Similarly, one might wish to make $\sigma_1 \preceq \sigma_2$ a derived constraint form; however, this is impossible, because neither universal quantification nor implication are available in the constraint language. We can, however, exploit the fact that these logical connectives are implicit in entailment assertions by defining a judgement of the form $C \Vdash \sigma_1 \preceq \sigma_2$, whose meaning is: *under the constraint $C$, $\sigma_1$ is more general than $\sigma_2$*.

1.3.21    DEFINITION:  We write $C \Vdash \sigma_1 \preceq \sigma_2$ if and only if $\mathtt{Z} \notin ftv(C, \sigma_1, \sigma_2)$ implies $C \wedge \sigma_2 \preceq \mathtt{Z} \Vdash \sigma_1 \preceq \mathtt{Z}$. We write $C \Vdash \sigma_1 \equiv \sigma_2$ when both $C \Vdash \sigma_1 \preceq \sigma_2$ and $C \Vdash \sigma_2 \preceq \sigma_1$ hold.                                                                   □

This notation is not ambiguous because the assertion $C \Vdash \sigma \preceq \mathtt{T}$, whose meaning was initially given by Definitions 1.3.3 and 1.3.14, retains the same meaning under the new definition—this is shown by Lemma 1.3.20 above.

The next lemma provides a way of exploiting the ordering between type schemes introduced by Definition 1.3.21. It states that a type scheme occurs in *contravariant* position when it is within a def prefix. In other words, the more general the type scheme, the weaker the entire constraint.

1.3.22    LEMMA:  $C \Vdash \sigma_1 \preceq \sigma_2$ implies $C \wedge \mathsf{def}\ \mathtt{x} : \sigma_2\ \mathsf{in}\ D \Vdash \mathsf{def}\ \mathtt{x} : \sigma_1\ \mathsf{in}\ D$.              □

The following exercise generalizes this result to let forms.

1.3.23    EXERCISE [★★, ↛]:  Prove that $\mathtt{Z} \notin ftv(\sigma)$ implies $\exists \sigma \equiv \exists \mathtt{Z}.\sigma \preceq \mathtt{Z}$. Explain why, as a result, $C \Vdash \sigma_1 \preceq \sigma_2$ implies $C \wedge \exists \sigma_2 \Vdash \exists \sigma_1$. Use this fact to prove that $C \Vdash \sigma_1 \preceq \sigma_2$ implies $C \wedge \mathsf{let}\ \mathtt{x} : \sigma_2\ \mathsf{in}\ D \Vdash \mathsf{let}\ \mathtt{x} : \sigma_1\ \mathsf{in}\ D$.              □

The next lemma states that, modulo equivalence, the only constraint that constrains x without explicitly referring to it is false.

1.3.24    LEMMA:  $C \Vdash \mathtt{x} \preceq \mathtt{T}$ and $\mathtt{x} \notin fpi(C)$ imply $C \equiv \mathsf{false}$.                           □

The following lemma states that the more universal quantifiers are present, the more general the type scheme.

1.3.25   LEMMA:  let $\mathtt{x} : \forall \bar{\mathtt{X}}[C_1].\mathtt{T}$ in $C_2 \Vdash$ let $\mathtt{x} : \forall \bar{\mathtt{X}}\bar{\mathtt{Y}}[C_1].\mathtt{T}$ in $C_2$.                          □

Conversely, and perhaps surprisingly, it is sometimes possible to *remove* some type variables from the universal quantifier prefix of a type scheme without compromising its generality. This is the case when the value of these type variables is *determined* in a unique way. In short, $C$ determines $\bar{\mathtt{Y}}$ if and only if, given the values of $ftv(C) \setminus \bar{\mathtt{Y}}$ and given that $C$ holds, it is possible to reconstruct, in a unique way, the values of $\bar{\mathtt{Y}}$.

1.3.26   DEFINITION:  $C$ *determines* $\bar{\mathtt{Y}}$ if and only if, for every environment $\Gamma$, two ground assignments that satisfy def $\Gamma$ in $C$ and that coincide outside $\bar{\mathtt{Y}}$ must coincide on $\bar{\mathtt{Y}}$ as well.                          □

Two concrete instances of determinacy, one of which is valid only in free tree models, are given by Lemma 1.8.7 on page 82. Determinacy is exploited by the equivalence law C-LETALL in Figure 1-6.

We now give a toolbox of constraint equivalence laws. It is worth noting that they do not form a complete axiomatization of constraint equivalence—in fact, they cannot, since the syntax and meaning of constraints is partly unspecified.

1.3.27   THEOREM:  All equivalence laws in Figure 1-6 hold.                          □

Let us explain. C-AND and C-ANDAND state that conjunction is commutative and associative. C-DUP states that redundant conjuncts may be freely added or removed, where a conjunct is redundant if and only if it is entailed by another conjunct. Throughout this chapter, these three laws are often used implicitly. C-EXEX and C-EX* allow grouping consecutive existential quantifiers and suppressing redundant ones, where a quantifier is redundant if and only if it does not occur free within its scope. C-EXAND allows conjunction and existential quantification to commute, provided no capture occurs; it is known as a *scope extrusion* law. When the rule is oriented from left to right, its side-condition may always be satisfied by suitable $\alpha$-conversion. C-EXTRANS states that it is equivalent for a type $\mathtt{T}'$ to be an instance of $\sigma$ or to be a supertype of some instance of $\sigma$. We remark that the instances of a monotype are its supertypes, that is, by Definition 1.3.3, $\mathtt{T} \preceq \mathtt{T}'$ and $\mathtt{T} \leq \mathtt{T}'$ are equivalent. As a result, specializing C-EXTRANS to the case where $\sigma$ is a monotype, we find that $\mathtt{T} \leq \mathtt{T}'$ is equivalent to $\exists \mathtt{Z}.(\mathtt{T} \leq \mathtt{Z} \wedge \mathtt{Z} \leq \mathtt{T}')$, for fresh $\mathtt{Z}$, a standard equivalence law. When oriented from left to right, it becomes an interesting *simplification* law: in a chain of subtyping constraints, an intermediate variable such as $\mathtt{Z}$ may be suppressed, provided it is *local*, as witnessed by the existential quantifier $\exists \mathtt{Z}$. C-INID states that, within the scope of the binding $\mathtt{x} : \sigma$, every *free* occurrence of $\mathtt{x}$ may be safely replaced with $\sigma$. The restriction to free occurrences stems from the side-condition $\mathtt{x} \notin dpi(\mathcal{C})$. When the

$$C_1 \wedge C_2 \quad\equiv\quad C_2 \wedge C_1 \qquad\qquad\qquad\qquad\qquad\qquad\qquad\text{(C-AND)}$$

$$(C_1 \wedge C_2) \wedge C_3 \quad\equiv\quad C_1 \wedge (C_2 \wedge C_3) \qquad\qquad\qquad\qquad\qquad\text{(C-ANDAND)}$$

$$C_1 \wedge C_2 \quad\equiv\quad C_1 \qquad\qquad\qquad \text{if } C_1 \Vdash C_2 \qquad\text{(C-DUP)}$$

$$\exists \bar{\mathtt{X}}.\exists \bar{\mathtt{Y}}.C \quad\equiv\quad \exists \bar{\mathtt{X}}\bar{\mathtt{Y}}.C \qquad\qquad\qquad\qquad\qquad\qquad\text{(C-EXEX)}$$

$$\exists \bar{\mathtt{X}}.C \quad\equiv\quad C \qquad\qquad\qquad \text{if } \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(C) \qquad\text{(C-EX*)}$$

$$(\exists \bar{\mathtt{X}}.C_1) \wedge C_2 \quad\equiv\quad \exists \bar{\mathtt{X}}.(C_1 \wedge C_2) \qquad\qquad \text{if } \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(C_2) \qquad\text{(C-EXAND)}$$

$$\exists \mathtt{Z}.(\sigma \preceq \mathtt{Z} \wedge \mathtt{Z} \leq \mathtt{T}') \quad\equiv\quad \sigma \preceq \mathtt{T}' \qquad\qquad \text{if } \mathtt{Z} \notin \mathit{ftv}(\sigma, \mathtt{T}') \qquad\text{(C-EXTRANS)}$$

$$\text{let } \mathtt{x} : \sigma \text{ in } \mathcal{C}[\mathtt{x} \preceq \mathtt{T}'] \quad\equiv\quad \text{let } \mathtt{x} : \sigma \text{ in } \mathcal{C}[\sigma \preceq \mathtt{T}'] \qquad\qquad\qquad\qquad\qquad\text{(C-INID)}$$
$$\text{if } \mathtt{x} \notin \mathit{dpi}(\mathcal{C}) \text{ and } \mathit{dtv}(\mathcal{C}) \mathbin{\#} \mathit{ftv}(\sigma) \text{ and } \{\mathtt{x}\} \cup \mathit{dpi}(\mathcal{C}) \mathbin{\#} \mathit{fpi}(\sigma)$$

$$\text{let } \Gamma \text{ in } C \quad\equiv\quad \exists \Gamma \wedge C \qquad\qquad \text{if } \mathit{dpi}(\Gamma) \mathbin{\#} \mathit{fpi}(C) \qquad\text{(C-IN*)}$$

$$\text{let } \Gamma \text{ in } (C_1 \wedge C_2) \quad\equiv\quad (\text{let } \Gamma \text{ in } C_1) \wedge (\text{let } \Gamma \text{ in } C_2) \qquad\qquad\qquad\qquad\text{(C-INAND)}$$

$$\text{let } \Gamma \text{ in } (C_1 \wedge C_2) \quad\equiv\quad (\text{let } \Gamma \text{ in } C_1) \wedge C_2 \qquad\qquad \text{if } \mathit{dpi}(\Gamma) \mathbin{\#} \mathit{fpi}(C_2) \qquad\text{(C-INAND*)}$$

$$\text{let } \Gamma \text{ in } \exists \bar{\mathtt{X}}.C \quad\equiv\quad \exists \bar{\mathtt{X}}.\text{let } \Gamma \text{ in } C \qquad\qquad \text{if } \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(\Gamma) \qquad\text{(C-INEX)}$$

$$\text{let } \Gamma_1; \Gamma_2 \text{ in } C \quad\equiv\quad \text{let } \Gamma_2; \Gamma_1 \text{ in } C \qquad\qquad\qquad\qquad\qquad\qquad\text{(C-LETLET)}$$
$$\text{if } \mathit{dpi}(\Gamma_1) \mathbin{\#} \mathit{dpi}(\Gamma_2) \text{ and } \mathit{dpi}(\Gamma_2) \mathbin{\#} \mathit{fpi}(\Gamma_1) \text{ and } \mathit{dpi}(\Gamma_1) \mathbin{\#} \mathit{fpi}(\Gamma_2)$$

$$\text{let } \mathtt{x} : \forall \bar{\mathtt{X}}[C_1 \wedge C_2].\mathtt{T} \text{ in } C_3 \quad\equiv\quad C_1 \wedge \text{let } \mathtt{x} : \forall \bar{\mathtt{X}}[C_2].\mathtt{T} \text{ in } C_3 \qquad \text{if } \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(C_1) \qquad\text{(C-LETAND)}$$

$$\text{let } \Gamma; \mathtt{x} : \forall \bar{\mathtt{X}}[C_1].\mathtt{T} \text{ in } C_2 \quad\equiv\quad \text{let } \Gamma; \mathtt{x} : \forall \bar{\mathtt{X}}[\text{let } \Gamma \text{ in } C_1].\mathtt{T} \text{ in } C_2 \qquad\qquad\text{(C-LETDUP)}$$
$$\text{if } \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(\Gamma) \text{ and } \mathit{dpi}(\Gamma) \mathbin{\#} \mathit{fpi}(\Gamma)$$

$$\text{let } \mathtt{x} : \forall \bar{\mathtt{X}}[\exists \bar{\mathtt{Y}}.C_1].\mathtt{T} \text{ in } C_2 \quad\equiv\quad \text{let } \mathtt{x} : \forall \bar{\mathtt{X}}\bar{\mathtt{Y}}[C_1].\mathtt{T} \text{ in } C_2 \qquad \text{if } \bar{\mathtt{Y}} \mathbin{\#} \mathit{ftv}(\mathtt{T}) \qquad\text{(C-LETEX)}$$

$$\text{let } \mathtt{x} : \forall \bar{\mathtt{X}}\bar{\mathtt{Y}}[C_1].\mathtt{T} \text{ in } C_2 \quad\equiv\quad \exists \bar{\mathtt{Y}}.\text{let } \mathtt{x} : \forall \bar{\mathtt{X}}[C_1].\mathtt{T} \text{ in } C_2 \qquad\qquad\text{(C-LETALL)}$$
$$\text{if } \bar{\mathtt{Y}} \mathbin{\#} \mathit{ftv}(C_2) \text{ and } \exists \bar{\mathtt{X}}.C_1 \text{ determines } \bar{\mathtt{Y}}$$

$$\exists \mathtt{X}.(\mathtt{T} \leq \mathtt{X} \wedge \text{let } \mathtt{x} : \mathtt{X} \text{ in } C) \quad\equiv\quad \text{let } \mathtt{x} : \mathtt{T} \text{ in } C \qquad \text{if } \mathtt{X} \notin \mathit{ftv}(\mathtt{T}, C) \qquad\text{(C-LETSUB)}$$

$$\vec{\mathtt{X}} = \vec{\mathtt{T}} \wedge [\vec{\mathtt{X}} \mapsto \vec{\mathtt{T}}]C \quad\equiv\quad \vec{\mathtt{X}} = \vec{\mathtt{T}} \wedge C \qqua\qquad\qquad\qquad\qquad\text{(C-EQ)}$$

$$\text{true} \quad\equiv\quad \exists \bar{\mathtt{X}}.(\vec{\mathtt{X}} = \vec{\mathtt{T}}) \qquad\qquad \text{if } \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(\bar{\mathtt{T}}) \qquad\text{(C-NAME)}$$

$$[\vec{\mathtt{X}} \mapsto \vec{\mathtt{T}}]C \quad\equiv\quad \exists \bar{\mathtt{X}}.(\vec{\mathtt{X}} = \vec{\mathtt{T}} \wedge C) \qquad\qquad \text{if } \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(\bar{\mathtt{T}}) \qquad\text{(C-NAMEEQ)}$$

**Figure 1-6: Constraint equivalence laws**

rule is oriented from left to right, its other side-conditions, which require the context $\mathsf{let}\ x : \sigma\ \mathsf{in}\ \mathcal{C}$ not to capture $\sigma$'s free type variables or free program identifiers, may always be satisfied by suitable $\alpha$-conversion. C-In* complements the previous rule by allowing redundant $\mathsf{let}$ bindings to be simplified. We remark that C-InId and C-In* provide a simple procedure for eliminating $\mathsf{let}$ forms. C-InAnd states that the $\mathsf{let}$ form commutes with conjunction; C-InAnd* spells out a common particular case. C-InEx states that it commutes with existential quantification. When the rule is oriented from left to right, its side-condition may always be satisfied by suitable $\alpha$-conversion. C-LetLet states that $\mathsf{let}$ forms may commute, provided they bind distinct program identifiers and provided no free program identifiers are captured in the process. C-LetAnd allows the conjunct $C_1$ to be moved outside of the constrained type scheme $\forall \bar{\mathsf{x}}[C_1 \wedge C_2].\mathsf{T}$, provided it does not involve any of the universally quantified type variables $\bar{\mathsf{x}}$. When oriented from left to right, the rule yields an important simplification law: indeed, taking an instance of $\forall \bar{\mathsf{x}}[C_2].\mathsf{T}$ is less expensive than taking an instance of $\forall \bar{\mathsf{x}}[C_1 \wedge C_2].\mathsf{T}$, since the latter involves creating a copy of $C_1$, while the former does not. C-LetDup allows pushing a series of $\mathsf{let}$ bindings into a constrained type scheme, provided no capture occurs in the process. It is not used as a simplification law but as a tool in some proofs. C-LetEx states that it does not make any difference for a set of type variables $\bar{\mathsf{Y}}$ to be existentially quantified inside a constrained type scheme or part of the type scheme's universal quantifiers. Indeed, in either case, taking an instance of the type scheme means producing a constraint where $\bar{\mathsf{Y}}$ is existentially quantified. C-LetAll provides a restricted converse of Lemma 1.3.25. Together, C-LetEx and C-LetAll allow—in some situations only—to hoist existential quantifiers out of the *left*-hand side of a $\mathsf{let}$ form.

1.3.28    Example: C-LetAll would be invalid without the condition that $\exists \bar{\mathsf{x}}.C_1$ determines $\bar{\mathsf{Y}}$. Consider, for instance, the constraint $\mathsf{let}\ x : \forall \mathsf{Y}.\mathsf{Y} \to \mathsf{Y}\ \mathsf{in}\ (\mathsf{x} \preceq \mathsf{int} \to \mathsf{int} \wedge \mathsf{x} \preceq \mathsf{bool} \to \mathsf{bool})$ **(1)**, where $\mathsf{int}$ and $\mathsf{bool}$ are incompatible nullary type constructors. By C-InId and C-In*, it is equivalent to $\exists \mathsf{Y}.(\mathsf{Y} \to \mathsf{Y} \leq \mathsf{int} \to \mathsf{int}) \wedge \exists \mathsf{Y}.(\mathsf{Y} \to \mathsf{Y} \leq \mathsf{bool} \to \mathsf{bool})$, that is, $\mathsf{true}$. Now, if C-LetAll was valid without its side-condition, then (1) would also be equivalent to $\exists \mathsf{Y}.\mathsf{let}\ x : \mathsf{Y} \to \mathsf{Y}\ \mathsf{in}\ (\mathsf{x} \preceq \mathsf{int} \to \mathsf{int} \wedge \mathsf{x} \preceq \mathsf{bool} \to \mathsf{bool})$, which by C-InId and C-In* is $\exists \mathsf{Y}.(\mathsf{Y} \to \mathsf{Y} \leq \mathsf{int} \to \mathsf{int} \wedge \mathsf{Y} \to \mathsf{Y} \leq \mathsf{bool} \to \mathsf{bool})$. By C-Arrow and C-ExTrans, this is $\mathsf{int} = \mathsf{bool}$, that is, $\mathsf{false}$. Thus, the law is invalid in this case. It is easy to see why: when the type scheme $\sigma$ contains a $\forall \mathsf{Y}$ quantifier, every instance of $\sigma$ receives its own $\exists \mathsf{Y}$ quantifier, making $\mathsf{Y}$ a distinct (local) type variable; when $\mathsf{Y}$ is not universally quantified, however, all instances of $\sigma$ *share* references to a single (global) type variable $\mathsf{Y}$. This corresponds to the

intuition that, in the former case, $\sigma$ is *polymorphic* in Y, while in the latter case, it is *monomorphic* in Y. Lemma 1.3.25 states that, when deprived of its side-condition, C-LetAll is only an entailment law, as opposed to an equivalence law. Similarly, it is in general invalid to hoist an existential quantifier out of the left-hand side of a let form. To see this, one may study the (equivalent) constraint let $x : \forall X[\exists Y.X = Y \to Y].X$ in $(x \preceq \text{int} \to \text{int} \wedge x \preceq \text{bool} \to \text{bool})$.

Naturally, in the above examples, the side-condition "true determines Y" does *not* hold: by Definition 1.3.26, it is equivalent to "two ground assignments that coincide outside Y must coincide on Y as well", which is false as soon as $\mathcal{M}_\star$ contains two distinct elements, such as int and bool here. There are cases, however, where the side-condition does hold. For instance, we later prove that $\exists X.Y = \text{int}$ determines Y; see Lemma 1.8.7. As a result, C-LetAll states that let $x : \forall XY[Y = \text{int}].Y \to X$ in $C$ **(1)** is equivalent to $\exists Y.$let $x : \forall X[Y = \text{int}].Y \to X$ in $C$ **(2)**, provided $Y \notin \mathit{ftv}(C)$. The intuition is simple: because Y is forced to assume the value int by the equation $Y = \text{int}$, it makes no difference whether Y is or isn't universally quantified. We remark that, by C-LetAnd, (2) is equivalent to $\exists Y.(Y = \text{int} \wedge \text{let } x : \forall X.Y \to X \text{ in } C)$ **(3)**. In an efficient constraint solver, simplifying (1) into (3) *before* using C-InId to eliminate the let form is worthwhile, since doing so obviates the need for copying the type variable Y and the equation $Y = \text{int}$ at every free occurrence of $x$ inside $C$. $\square$

C-LetSub is the analogue of an environment strengthening lemma: roughly speaking, it states that, if a constraint holds under the assumption that $x$ has type X, where X is some supertype of T, then it also holds under the assumption that $x$ has type T. The last three rules deal with the equality predicate. C-Eq states that it is valid to replace equals with equals; note the absence of a side-condition. When oriented from left to right, C-Name allows introducing fresh names $\vec{X}$ for the types $\vec{T}$. As always, $\vec{X}$ stands for a vector of *distinct* type variables. Of course, this makes sense only if the definition is not circular, that is, if the type variables $\bar{X}$ do not occur free within the terms $\bar{T}$. When oriented from right to left, C-Name may be viewed as a simplification law: it allows eliminating type variables whose value has been determined. C-NameEq is a combination of C-Eq and C-Name. It shows that applying an idempotent substitution to a constraint $C$ amounts to placing $C$ within a certain context. This immediately yields a proof of the following fact:

1.3.29   LEMMA: $C \Vdash D$ implies $\theta(C) \Vdash \theta(D)$. $\square$

It is important to stress that, because the effect of a type substitution may be emulated using equations, conjunction, and existential quantification, there is no need ever to employ type substitutions in the definition of a constraint-based type system—it is possible, instead, to express every concept in terms

of constraints. In this chapter, we follow this route, and use type substitutions only when dealing with the type system DM, whose historical formulation is substitution-based.

So far, we have considered def a primitive constraint form and defined the let form in terms of def, conjunction, and existential quantification. The motivation for this approach was to simplify the proof of several constraint equivalence laws. However, in the remainder of this chapter, we work with let forms exclusively and never employ the def construct. As a result, *it is possible, from here on, to discard* def *and pretend that* let *is primitive*. This change in perspective offers us a few extra properties, stated in the next two lemmas. First, every constraint that contains a false subconstraint must be false. Second, no satisfiable constraint has a free program identifier.

1.3.30    LEMMA:  $\mathcal{C}[\mathsf{false}] \equiv \mathsf{false}$.                                                  □

1.3.31    LEMMA:  If $C$ is satisfiable, then $fpi(C) = \varnothing$.                                □

### Reasoning with constraints in an equality-only syntactic model

We have given a number of equivalence laws that are valid with respect to any interpretation of constraints, that is, within any model. However, an important special case is that of *equality-only syntactic models*. Indeed, in that specific setting, our constraint-based type systems are in close correspondence with DM. In short, we aim to prove that every satisfiable constraint admits a *canonical solved form*, to show that this notion corresponds to the standard concept of a *most general unifier*, and to establish a few technical properties of most general unifiers.

Thus, let us now assume that constraints are interpreted in an equality-only syntactic model. Let us further assume that, for every kind $\kappa$, (i) there are at least *two* type constructors of image kind $\kappa$ and (ii) for every type constructor $F$ of image kind $\kappa$, there exists $t \in \mathcal{M}_\kappa$ such that $t(\epsilon) = F$. We refer to models that violate (i) or (ii) as *degenerate*; one may argue that such models are of little interest. The assumption that the model is nondegenerate is used in the proof of Lemmas 1.3.32 and 1.3.39.

Under these new assumptions, the interpretation of equality coincides with its syntax: every equation that holds in the model is in fact a syntactic truism. The converse, of course, holds in every model.

1.3.32    LEMMA:  If $\mathsf{true} \Vdash \mathsf{T} = \mathsf{T}'$ holds, then $\mathsf{T}$ and $\mathsf{T}'$ coincide.                  □

In a syntactic model, ground types are finite trees. As a result, cyclic equations, such as $\mathsf{X} = \mathsf{int} \to \mathsf{X}$, are false.

1.3.33    LEMMA:  $\mathsf{X} \in ftv(\mathsf{T})$ and $\mathsf{T} \not\in \mathcal{V}$ imply $(\mathsf{X} = \mathsf{T}) \equiv \mathsf{false}$.                  □

A *solved form* is a conjunction of equations, where the left-hand sides are *distinct* type variables that do not appear in the right-hand sides, possibly surrounded by a number of existential quantifiers. Our definition is identical to Lassez, Maher and Marriott's solved forms (1988) and to Jouannaud and Kirchner's *tree* solved forms (1991), except we allow for prenex existential quantifiers, which are made necessary by our richer constraint language. Jouannaud and Kirchner also define *dag* solved forms, which may be exponentially smaller. Because we define solved forms only for proof purposes, we need not take performance into account at this point. The efficient constraint solver presented in Section 1.8 does manipulate graphs, rather than trees. Type scheme introduction and instantiation constructs cannot appear within solved forms; indeed, provided the constraint at hand has no free program identifiers, they can be expanded away. For this reason, their presence in the constraint language has no impact on the results contained in this section.

1.3.34  DEFINITION: A *solved form* is of the form $\exists \vec{Y}.(\vec{X} = \vec{T})$, where $\vec{X} \mathrel{\#} ftv(\vec{T})$.  □

Solved forms offer a convenient way of reasoning about constraints because every satisfiable constraint is equivalent to one. In other words, every constraint is equivalent to either a solved form or false. This property is established by the following lemma, whose proof provides a simple but effective procedure to rewrite a constraint to either a solved form or false.

1.3.35  LEMMA: Let $fpi(C) = \varnothing$. Then, $C$ is equivalent to either a solved form or false.  □

*Proof:* We first establish that every conjunction of equations is equivalent to either a solved form or false. To do so, we present Robinson's unification algorithm (1971) as a rewriting system. The system's invariant is to operate on constraints of the form either $\vec{X} = \vec{T}; C$, where $\vec{X} \mathrel{\#} ftv(\vec{T}, C)$ and the semicolon is interpreted as a distinguished conjunction, or false. We identify equations in $C$ up to commutativity. The system is defined as follows:

$$
\begin{aligned}
\vec{X} = \vec{T}; & & \mathtt{X} = \mathtt{X} \wedge C & \;\to\; & \vec{X} = \vec{T}; C \\
\vec{X} = \vec{T}; & & F\,\vec{T}_1 = F\,\vec{T}_2 \wedge C & \;\to\; & \vec{X} = \vec{T}; \vec{T}_1 = \vec{T}_2 \wedge C \\
\vec{X} = \vec{T}; & & F_1\,\vec{T}_1 = F_2\,\vec{T}_2 \wedge C & \;\to\; & \mathsf{false} \\
& & & & \text{if } F_1 \neq F_2 \\
\vec{X} = \vec{T}; & & \mathtt{X} = \mathtt{T} \wedge C & \;\to\; & \vec{X} = [\mathtt{X} \mapsto \mathtt{T}]\vec{T} \wedge \mathtt{X} = \mathtt{T}; [\mathtt{X} \mapsto \mathtt{T}]C \\
& & & & \text{if } \mathtt{X} \notin ftv(\mathtt{T}) \\
\vec{X} = \vec{T}; & & \mathtt{X} = \mathtt{T} \wedge C & \;\to\; & \mathsf{false} \\
& & & & \text{if } \mathtt{X} \in ftv(\mathtt{T}) \text{ and } \mathtt{T} \notin \mathcal{V}
\end{aligned}
$$

It is straightforward to check that the above invariant is indeed preserved by the rewriting system. Let us check that constraint equivalence is also preserved. For the first rule, this is immediate. For the second and third rules, it

follows from the fact that we have assumed a free tree model; for the fourth rule, a consequence of C-Eq; for the last rule, a consequence of Lemma 1.3.33. Furthermore, the system is terminating; this is witnessed by an ordering where false is the least element and where constraints of the form $\vec{X} = \vec{T}; C$ are ordered lexicographically, first by the number of type variables that appear free within $C$, second by the size of $C$. Last, a normal form for this rewriting system must be of the form either $\vec{X} = \vec{T};$ true, where (by the invariant) $\bar{X}\ \#\ ftv(\bar{T})$—that is, a solved form, or false.

Next, we show that the present lemma holds when $C$ is built out of equations, conjunction, and existential quantification. Orienting C-ExAnd from left to right yields a terminating rewriting system that preserves constraint equivalence. The normal form of $C$ must be $\exists\bar{Y}.C'$, where $C'$ is a conjunction of equations. By the previous result, $C'$ is equivalent to either a solved form or false. Because solved forms are preserved by existential quantification and because $\exists\bar{Y}.$false is false, the same holds of $C$.

Last, we establish the result in the general case. We assume $fpi(C) = \varnothing$ **(1)**. Orienting C-InId and C-In* from left to right yields a terminating rewriting system that preserves constraint equivalence. The normal form $C'$ of $C$ cannot contain any type scheme introduction forms; given (1), it cannot contain any instantiation forms either. Thus, $C'$ is built out of equations, conjunction, and existential quantification only. By the previous result, it is equivalent to either a solved form or false, which implies that the same holds of $C$.          □

It is possible to impose further restrictions on solved forms. A solved form $\exists\bar{Y}.(\vec{X} = \vec{T})$ is *canonical* if and only if its free type variables are exactly $\bar{X}$. This is stated, in an equivalent way, by the following definition.

1.3.36   Definition:   A *canonical solved form* is a constraint of the form $\exists\bar{Y}.(\vec{X} = \vec{T})$, where $ftv(\bar{T}) \subseteq \bar{Y}$ and $\bar{X}\ \#\ \bar{Y}$.          □

1.3.37   Lemma:   Every solved form is equivalent to a canonical solved form.          □

It is easy to describe the solutions of a canonical solved form: they are the ground refinements of the substitution $[\vec{X} \mapsto \vec{T}]$.

1.3.38   Lemma:   A ground assignment $\phi$ satisfies a canonical solved form $\exists\bar{Y}.(\vec{X} = \vec{T})$ if and only if there exists a ground assignment $\phi'$ such that $\phi(\vec{X}) = \phi'(\vec{T})$. As a result, every canonical solved form is satisfiable.          □

*Proof:*   Let $\exists\bar{Y}.(\vec{X} = \vec{T})$ be a canonical solved form. By CM-Exists and CM-Predicate, $\phi$ satisfies $\exists\bar{Y}.(\vec{X} = \vec{T})$ if and only if there exists $\vec{t}$ such that $\phi[\vec{Y} \mapsto \vec{t}](\vec{X}) = \phi[\vec{Y} \mapsto \vec{t}](\vec{T})$. Thanks to the hypotheses $\bar{X}\ \#\ \bar{Y}$ and $ftv(\bar{T}) \subseteq \bar{Y}$, this is equivalent to the existence of a ground assignment $\phi'$ such that $\phi(\vec{X}) = \phi'(\vec{T})$.

Thus, for every ground assignment $\phi'$, $\phi'[\vec{X} \mapsto \phi'(\vec{T})]$ satisfies $\exists \bar{Y}.(\vec{X} = \vec{T})$, which proves that this constraint is satisfiable. □

Together, Lemmas 1.3.37 and 1.3.38 imply that every solved form is satisfiable. Our interest in canonical solved forms stems from the following fundamental property, which provides a *syntactic* characterization of entailment between canonical solved forms: if $\exists \bar{Y}_1.(\vec{X} = \vec{T}_1)$ is more specific than $\exists \bar{Y}_2.(\vec{X} = \vec{T}_2)$, in a logical sense, then $\vec{T}_1$ refines $\vec{T}_2$, in a syntactic sense. The converse also holds (can you prove it?), but is not needed here.

1.3.39 LEMMA: If $\exists \bar{Y}_1.(\vec{X} = \vec{T}_1) \Vdash \exists \bar{Y}_2.(\vec{X} = \vec{T}_2)$, where both sides are canonical solved forms, then there exists a type substitution $\varphi$ such that $\vec{T}_1 = \varphi(\vec{T}_2)$. □

As a corollary, we find that canonical solved forms are *unique* up to $\alpha$-conversion and up to C-Ex*, *provided* the set $\bar{X}$ of their free type variables is fixed.

1.3.40 LEMMA: If the canonical solved forms $\exists \bar{Y}_1.(\vec{X} = \vec{T}_1)$ and $\exists \bar{Y}_2.(\vec{X} = \vec{T}_2)$ are equivalent, then there exists a renaming $\rho$ such that $\vec{T}_1 = \rho(\vec{T}_2)$. □

Please note that the fact that the canonical solved forms $\exists \bar{Y}_1.(\vec{X}_1 = \vec{T}_1)$ and $\exists \bar{Y}_2.(\vec{X}_2 = \vec{T}_2)$ are equivalent does *not* imply that $\bar{X}_1$ and $\bar{X}_2$ coincide. Consider, for example, the canonical solved forms $\mathtt{true}$ and $\exists Y.(X = Y)$, which by C-NameEq are equivalent. One might wish to further restrict canonical solved forms by requiring $\bar{X}$ to be the set of *essential* type variables of the constraint $\exists \bar{Y}.(\vec{X} = \vec{T})$, that is, the set of the type variables that appear free in *all* equivalent constraints. However, as far our technical development is concerned, it seems more convenient not to do so. Instead, we show that it is possible to explicitly restrict or extend $\bar{X}$ when needed (Lemma 1.3.43).

The following definition allows entertaining a dual view of canonical solved forms, either as constraints or as idempotent type substitutions. The latter view is commonly found in standard treatments of unification (Lassez, Maher, and Marriott, 1988; Jouannaud and Kirchner, 1991) and in classic presentations of ML-the-type-system.

1.3.41 DEFINITION: If $[\vec{X} \mapsto \vec{T}]$ is an idempotent substitution of domain $\bar{X}$, let $\exists[\vec{X} \mapsto \vec{T}]$ denote the canonical solved form $\exists \bar{Y}.(\vec{X} = \vec{T})$, where $\bar{Y} = ftv(\vec{T})$. An idempotent substitution $\theta$ is a *most general unifier* of the constraint $C$ if and only if $\exists \theta$ and $C$ are equivalent. □

By definition, equivalent constraints admit the same most general unifiers. Many properties of canonical solved forms may be reformulated in terms of most general unifiers. By Lemmas 1.3.31, 1.3.35, and 1.3.37, every satisfiable constraint admits a most general unifier. By Lemma 1.3.40, if $[\vec{X} \mapsto \vec{T}_1]$ and

$[\vec{\mathtt{X}} \mapsto \vec{\mathtt{T}}_2]$ are most general unifiers of $C$, then $\vec{\mathtt{T}}_1$ and $\vec{\mathtt{T}}_2$ coincide up to a renaming. Conversely, if $[\vec{\mathtt{X}} \mapsto \vec{\mathtt{T}}]$ is a most general unifier of $C$ and if $\bar{\mathtt{X}} \mathrel{\#} \rho$ holds, then $[\vec{\mathtt{X}} \mapsto \rho\vec{\mathtt{T}}]$ is also a most general unifier of $C$; indeed, these two substitutions correspond to $\alpha$-equivalent canonical solved forms.

The following result relates the substitution $\theta$ to the canonical solved form $\exists\theta$, stating that every ground refinement of the former satisfies the latter.

1.3.42     LEMMA: $\theta(\exists\theta) \equiv \mathsf{true}$.                                      □

The following lemma offers two technical results: the domain of a most general unifier of $C$ may be restricted so as to become a subset of $ftv(C)$; it may also be extended to include arbitrary fresh variables. The next lemma is a simple corollary.

1.3.43     LEMMA: Let $\theta$ be a most general unifier of $C$. If $\bar{\mathtt{Z}} \mathrel{\#} ftv(C)$, then $\theta \setminus \bar{\mathtt{Z}}$ is also a most general unifier of $C$. If $\bar{\mathtt{Z}} \mathrel{\#} \theta$, then there exists a most general unifier of $C$ that extends $\theta$ and whose domain is $dom(\theta) \cup \bar{\mathtt{Z}}$.        □

1.3.44     LEMMA: Let $\theta_1$ and $\theta_2$ be most general unifiers of $C$. Let $\bar{\mathtt{X}} = dom(\theta_1) \cap dom(\theta_2)$. Then, $\theta_1(\bar{\mathtt{X}})$ and $\theta_2(\bar{\mathtt{X}})$ coincide up to a renaming.        □

Our last technical result relates the most general unifiers of $C$ with the most general unifiers of $\exists\mathtt{X}.C$. It states that the former are extensions of the latter. Furthermore, under a few freshness conditions, *every* most general unifier of $\exists\mathtt{X}.C$ may be extended to yield a most general unifier of $C$.

1.3.45     LEMMA: If $\theta$ is a most general unifier of $C$, then $\theta \setminus \mathtt{X}$ is a most general unifier of $\exists\mathtt{X}.C$. Conversely, if $\theta$ is a most general unifier of $\exists\mathtt{X}.C$ and $\mathtt{X} \mathrel{\#} \theta$ and $ftv(\exists\mathtt{X}.C) \subseteq dom(\theta)$, then there exists a type substitution $\theta'$ such that $\theta'$ extends $\theta$, $\theta'$ is a most general unifier of $C$, and $dom(\theta') = dom(\theta) \cup \mathtt{X}$.     □

## 1.4    HM($X$)

Constraint-based type systems appeared during the 1980s (Mitchell, 1984; Fuh and Mishra, 1988) and were widely studied during the following decade (Curtis, 1990; Aiken and Wimmers, 1993; Jones, 1994a; Smith, 1994; Palsberg, 1995; Trifonov and Smith, 1996; Fähndrich, 1999; Pottier, 2001b). We now present one such system, baptized HM($X$) because it is a *parameterized* extension of Hindley and Milner's type discipline; the meaning of the parameter $X$ was explained on page 24. Its original description is due to Odersky, Sulzmann, and Wehr (1999a). Since then, it has been completed in a number of works (Sulzmann, Müller, and Zenger, 1999; Sulzmann, 2000; Pottier, 2001a;

Skalka and Pottier, 2002). Each of these presentations introduces minor variations. Here, we follow (Pottier, 2001a), which is itself inspired by (Sulzmann, Müller, and Zenger, 1999).

### Definition

Our presentation of HM($X$) relies on the constraint language introduced in section 1.3. Technically, our approach of constraints is more direct than that of (Odersky, Sulzmann, and Wehr, 1999a). We interpret constraints within a model, give conjunction and existential quantification their standard meaning, and derive a number of equivalence laws (Section 1.3). Odersky *et al.*, on the other hand, do not explicitly rely on a logical interpretation; instead, they axiomatize constraint equivalence, that is, they consider a number of equivalence laws as axioms. Thus, they ensure that their high-level proofs, such as type soundness and correctness and completeness of type inference, are independent of the low-level details of the logical interpretation of constraints. Their approach is also more general, since it allows dealing with other logical interpretations—such as "open-world" interpretations, where constraints are interpreted not within a fixed model, but within a *family* of extensions of a "current" model. In this chapter, we have avoided this extra layer of abstraction, for the sake of definiteness; however, the changes required to adopt Odersky *et al.*'s approach would not be extensive, since the forthcoming proofs do indeed rely mostly on constraint equivalence laws, rather than on low-level details of the logical interpretation of constraints.

Another slight departure from Odersky *et al.*'s work lies in the fact that we have enriched the constraint language with type scheme introduction and instantiation forms, which were absent in the original presentation of HM($X$). To prevent this addition from affecting HM($X$), we require the constraints that appear in HM($X$) typing judgements to *have no free program identifiers*. Please note that this does not prevent them from containing let forms; we shall in fact exploit this feature when establishing an equivalence between HM($X$) and the type system presented in section 1.5, where the new constraint forms are effectively used.

The type system HM($X$) consists of a four-place *judgement* whose parameters are a constraint $C$, an environment $\Gamma$, an expression t, and a type scheme $\sigma$. A judgement is written $C, \Gamma \vdash t : \sigma$ and is read: *under the assumptions $C$ and $\Gamma$, the expression t has type $\sigma$*. One may view $C$ as an assumption about the judgement's free type variables and $\Gamma$ as an assumption about t's free program identifiers. Please recall that $\Gamma$ now contains *constrained* type schemes, and that $\sigma$ is a *constrained* type scheme.

We would like the validity of a typing judgement to depend not on the

$$\frac{\Gamma(\mathtt{x}) = \sigma \qquad C \Vdash \exists \sigma}{C, \Gamma \vdash \mathtt{x} : \sigma} \quad (\textsc{hmx-Var})$$

$$\frac{C, (\Gamma; \mathtt{z} : \mathtt{T}) \vdash \mathtt{t} : \mathtt{T}'}{C, \Gamma \vdash \lambda \mathtt{z}.\mathtt{t} : \mathtt{T} \to \mathtt{T}'} \quad (\textsc{hmx-Abs})$$

$$\frac{C, \Gamma \vdash \mathtt{t_1} : \mathtt{T} \to \mathtt{T}' \qquad C, \Gamma \vdash \mathtt{t_2} : \mathtt{T}}{C, \Gamma \vdash \mathtt{t_1}\,\mathtt{t_2} : \mathtt{T}'} \quad (\textsc{hmx-App})$$

$$\frac{C, \Gamma \vdash \mathtt{t_1} : \sigma \qquad C, (\Gamma; \mathtt{z} : \sigma) \vdash \mathtt{t_2} : \mathtt{T}}{C, \Gamma \vdash \mathtt{let}\ \mathtt{z} = \mathtt{t_1}\ \mathtt{in}\ \mathtt{t_2} : \mathtt{T}} \quad (\textsc{hmx-Let})$$

$$\frac{C \wedge D, \Gamma \vdash \mathtt{t} : \mathtt{T} \qquad \bar{\mathtt{x}} \mathbin{\#} \mathit{ftv}(C, \Gamma)}{C \wedge \exists \bar{\mathtt{x}}.D, \Gamma \vdash \mathtt{t} : \forall \bar{\mathtt{x}}[D].\mathtt{T}} \quad (\textsc{hmx-Gen})$$

$$\frac{C, \Gamma \vdash \mathtt{t} : \forall \bar{\mathtt{x}}[D].\mathtt{T}}{C \wedge D, \Gamma \vdash \mathtt{t} : \mathtt{T}} \quad (\textsc{hmx-Inst})$$

$$\frac{C, \Gamma \vdash \mathtt{t} : \mathtt{T} \qquad C \Vdash \mathtt{T} \le \mathtt{T}'}{C, \Gamma \vdash \mathtt{t} : \mathtt{T}'} \quad (\textsc{hmx-Sub})$$

$$\frac{C, \Gamma \vdash \mathtt{t} : \sigma \qquad \bar{\mathtt{x}} \mathbin{\#} \mathit{ftv}(\Gamma, \sigma)}{\exists \bar{\mathtt{x}}.C, \Gamma \vdash \mathtt{t} : \sigma} \quad (\textsc{hmx-Exists})$$

**Figure 1-7: Typing rules for HM**($X$)

*syntax*, but only on the *meaning* of its constraint assumption. We enforce this point of view by considering judgements equal modulo equivalence of their constraint assumptions. In other words, the typing judgements $C, \Gamma \vdash \mathtt{t} : \sigma$ and $D, \Gamma \vdash \mathtt{t} : \sigma$ are considered identical when $C \equiv D$ holds. As a result, it does not make sense to analyze the syntax of a judgement's constraint assumption. A judgement is *valid*, or *holds*, if and only if it is derivable via the rules given in Figure 1-7. Please note that a valid judgement may involve an unsatisfiable constraint. A program $\mathtt{t}$ is *well-typed* within the environment $\Gamma$ if and only if a judgement of the form $C, \Gamma \vdash \mathtt{t} : \sigma$ holds for some *satisfiable* constraint $C$.

Let us now explain the rules. Like DM-VAR, HMX-VAR looks up the environment to determine the type scheme associated with the program identifier $\mathtt{x}$. The constraint $C$ that appears in the conclusion must be strong enough to guarantee that $\sigma$ has an instance; this is expressed by the second premise. This technical requirement is used in the proof of Lemma 1.4.1. HMX-ABS, HMX-APP, and HMX-LET are identical to DM-ABS, DM-APP, and DM-LET, respectively, except that the assumption $C$ is made available to every sub-derivation. We recall that the type $\mathtt{T}$ may be viewed as the type scheme $\forall \varnothing[\mathsf{true}].\mathtt{T}$ (Definitions 1.2.18 and 1.3.2). As a result, types form a subset of type schemes, which implies that $\Gamma; \mathtt{z} : \mathtt{T}$ is a well-formed environment and $C, \Gamma \vdash \mathtt{t} : \mathtt{T}$ a well-formed typing judgement. To understand HMX-GEN, it is best to first consider the particular case where $C$ is $\mathsf{true}$. This yields the following, simpler rule:

$$\frac{D, \Gamma \vdash \mathtt{t} : \mathtt{T} \qquad \bar{\mathtt{x}} \mathbin{\#} \mathit{ftv}(\Gamma)}{\exists \bar{\mathtt{x}}.D, \Gamma \vdash \mathtt{t} : \forall \bar{\mathtt{x}}[D].\mathtt{T}} \quad (\textsc{hmx-Gen'})$$

The second premise is identical to that of DM-GEN: the type variables that are generalized must not occur free within the environment. The conclusion forms the type scheme $\forall \bar{\mathtt{X}}[D].\mathtt{T}$, where the type variables $\bar{\mathtt{X}}$ have become universally quantified, but are still subject to the constraint $D$. Please note that the type variables that occur free in $D$ may include not only $\bar{\mathtt{X}}$, but also other type variables, typically free in $\Gamma$. The rule's conclusion carries the constraint $\exists \bar{\mathtt{X}}.D$, thus recording the requirement that the newly formed type scheme should have an instance; again, this is used in the proof of Lemma 1.4.1. HMX-GEN may be viewed as a more liberal version of HMX-GEN', whereby part of the current constraint, namely $C$, need not be copied if it does not concern the type variables that are being generalized, namely $\bar{\mathtt{X}}$. This optimization is important in practice, because $C$ may be very large. An intuitive explanation for its correctness is given by the constraint equivalence law C-LETAND, which expresses the same optimization in terms of let constraints. Because HM($X$) does not use let constraints, the optimization is hard-wired into the typing rule. HMX-INST allows taking an instance of a type scheme. The reader may be surprised to find that, contrary to DM-INST, it does not involve a type substitution. Instead, the rule merely drops the universal quantifier, which amounts to applying the identity substitution $\vec{\mathtt{X}} \mapsto \vec{\mathtt{X}}$. One should recall, however, that type schemes are considered equal modulo $\alpha$-conversion, so it is possible to *rename* the type scheme's universal quantifiers prior to using HMX-INST. The reason why this provides sufficient expressive power appears in the proof of Theorem 1.4.7 below. The constraint $D$ carried by the type scheme is recorded as part of the current constraint in HMX-INST's conclusion. The *subsumption* rule HMX-SUB allows a type $\mathtt{T}$ to be replaced at any time with an arbitrary supertype $\mathtt{T}'$. Because both $\mathtt{T}$ and $\mathtt{T}'$ may have free type variables, whether $\mathtt{T} \leq \mathtt{T}'$ holds depends on the current assumption $C$, which is why the rule's second premise is an *entailment* assertion. An operational explanation of HMX-SUB is that it requires all uses of subsumption to be explicitly recorded in the current constraint. Please note that HMX-SUB remains a useful and necessary rule even when subtyping is interpreted as equality: then, it allows exploiting the type *equations* found in $C$. Last, HMX-EXISTS allows the type variables that occur only within the current constraint to become existentially quantified. As a result, these type variables no longer occur free in the rule's conclusion; in other words, they have become *local* to the subderivation rooted at the premise. One may prove that the presence of HMX-EXISTS in the type system does not augment the set of well-typed programs, but does augment the set of valid typing judgements; it is a pleasant technical convenience. Indeed, because judgements are considered equal modulo constraint equivalence, constraints may be transparently *simplified* at any time. (By *simplifying* a constraint, we mean replacing it with an equiva-

lent constraint whose syntactic representation is considered simpler.) Bearing this fact in mind, one finds that an effect of rule HMX-EXISTS is to enable *more* simplifications: because constraint equivalence is a congruence, $C \equiv D$ implies $\exists \bar{\mathtt{x}}.C \equiv \exists \bar{\mathtt{x}}.D$, but the converse does not hold in general. For instance, there is in general no way of simplifying the judgement $\mathtt{X} \leq \mathtt{Y} \leq \mathtt{Z}, \Gamma \vdash \mathtt{t} : \sigma$, but if it is known that $\mathtt{Y}$ does not appear free in $\Gamma$ or $\sigma$, then HMX-EXISTS allows deriving $\exists \mathtt{Y}.(\mathtt{X} \leq \mathtt{Y} \leq \mathtt{Z}), \Gamma \vdash \mathtt{t} : \sigma$, which is the same judgement as $\mathtt{X} \leq \mathtt{Z}, \Gamma \vdash \mathtt{t} : \sigma$. Thus, an interesting simplification has been enabled. Please note that $\mathtt{X} \leq \mathtt{Y} \leq \mathtt{Z} \equiv \mathtt{X} \leq \mathtt{Z}$ does *not* hold, while, according to C-EXTRANS, $\exists \mathtt{Y}.(\mathtt{X} \leq \mathtt{Y} \leq \mathtt{Z}) \equiv \mathtt{X} \leq \mathtt{Z}$ does.

We now establish a few simple properties of the type system $\mathrm{HM}(X)$. Our first lemma is a minor technical property.

1.4.1    LEMMA: $C, \Gamma \vdash \mathtt{t} : \sigma$ implies $C \Vdash \exists \sigma$.           □

The next lemma states that *strengthening* a judgement's constraint assumption preserves its validity. In other words, *weakening* a judgement preserves its validity. It is worth noting that in traditional presentations, which rely more heavily on type substitutions, the analogue of this result is a *type substitution* lemma; see for instance (Tofte, 1988, Lemma 2.7), (Leroy, 1992, Proposition 1.2), (Skalka and Pottier, 2002, Lemma 3.4). Here, the lemma further states that weakening a judgement does not alter the shape of its derivation, a useful property when reasoning by induction on type derivations.

1.4.2    LEMMA [WEAKENING]: If $C' \Vdash C$, then every derivation of $C, \Gamma \vdash \mathtt{t} : \sigma$ may be turned into a derivation of $C', \Gamma \vdash \mathtt{t} : \sigma$ with the same shape.     □

*Proof:*   The proof is by structural induction on a derivation of $C, \Gamma \vdash \mathtt{t} : \sigma$. In each proof case, we adopt the notations of Figure 1-7.

∘ *Case* HMX-VAR. The rule's conclusion is $C, \Gamma \vdash \mathtt{x} : \sigma$. Its premises are $\Gamma(\mathtt{x}) = \sigma$ **(1)** and $C \Vdash \exists \sigma$ **(2)**. By hypothesis, we have $C' \Vdash C$ **(3)**. By transitivity of entailment, (3) and (2) imply $C' \Vdash \exists \sigma$ **(4)**. By HMX-VAR, (1) and (4) yield $C', \Gamma \vdash \mathtt{x} : \sigma$.

∘ *Cases* HMX-ABS, HMX-APP, HMX-LET. By the induction hypothesis and by HMX-ABS, HMX-APP, or HMX-LET, respectively.

∘ *Case* HMX-GEN. The rule's conclusion is $C \wedge \exists \bar{\mathtt{x}}.D, \Gamma \vdash \mathtt{t} : \forall \bar{\mathtt{x}}[D].\mathtt{T}$. Its premises are $C \wedge D, \Gamma \vdash \mathtt{t} : \mathtt{T}$ **(1)** and $\bar{\mathtt{x}} \mathbin{\#} ftv(C, \Gamma)$ **(2)**. By hypothesis, we have $C' \Vdash C \wedge \exists \bar{\mathtt{x}}.D$ **(3)**. We may assume, *w.l.o.g.*, $\bar{\mathtt{x}} \mathbin{\#} ftv(C')$ **(4)**. Applying the induction hypothesis to (1) and to the entailment assertion $C' \wedge C \wedge D \Vdash C \wedge D$, we obtain $C' \wedge C \wedge D, \Gamma \vdash \mathtt{t} : \mathtt{T}$ **(5)**. By HMX-GEN, applied to (5), (2) and (4), we get $C' \wedge C \wedge \exists \bar{\mathtt{x}}.D, \Gamma \vdash \mathtt{t} : \forall \bar{\mathtt{x}}[D].\mathtt{T}$ **(6)**. By (3) and C-DUP, the constraints $C' \wedge C \wedge \exists \bar{\mathtt{x}}.D$ and $C'$ are equivalent, so (6) is the goal $C', \Gamma \vdash \mathtt{t} : \forall \bar{\mathtt{x}}[D].\mathtt{T}$.

∘ *Case* HMX-INST. The rule's conclusion is $C \wedge D, \Gamma \vdash \mathtt{t} : \mathtt{T}$. Its premise is $C, \Gamma \vdash \mathtt{t} : \forall \bar{\mathtt{X}}[D].\mathtt{T}$ **(1)**. By hypothesis, $C'$ entails $C \wedge D$ **(2)**. Because (2) implies $C' \Vdash C$, the induction hypothesis may be applied to (1), yielding $C', \Gamma \vdash \mathtt{t} : \forall \bar{\mathtt{X}}[D].\mathtt{T}$ **(3)**. By HMX-INST, we obtain $C' \wedge D, \Gamma \vdash \mathtt{t} : \mathtt{T}$ **(4)**. Because (2) implies $C' \equiv C' \wedge D$, (4) is the goal $C', \Gamma \vdash \mathtt{t} : \mathtt{T}$.

∘ *Case* HMX-SUB. The rule's conclusion is $C, \Gamma \vdash \mathtt{t} : \mathtt{T}'$. Its premises are $C, \Gamma \vdash \mathtt{t} : \mathtt{T}$ **(1)** and $C \Vdash \mathtt{T} \leq \mathtt{T}'$ **(2)**. By hypothesis, we have $C' \Vdash C$ **(3)**. Applying the induction hypothesis to (1) and (3) yields $C', \Gamma \vdash \mathtt{t} : \mathtt{T}$ **(4)**. By transitivity of entailment, (3) and (2) imply $C' \Vdash \mathtt{T} \leq \mathtt{T}'$ **(5)**. By HMX-SUB, (4) and (5) yield $C', \Gamma \vdash \mathtt{t} : \mathtt{T}'$.

∘ *Case* HMX-EXISTS. The rule's conclusion is $\exists \bar{\mathtt{X}}.C, \Gamma \vdash \mathtt{t} : \sigma$. Its premises are $C, \Gamma \vdash \mathtt{t} : \sigma$ **(1)** and $\bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(\Gamma, \sigma)$ **(2)**. By hypothesis, we have $C' \Vdash \exists \bar{\mathtt{X}}.C$ **(3)**. We may assume, *w.l.o.g.*, $\bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(C')$ **(4)**. Applying the induction hypothesis to (1) and to the entailment assertion $C' \wedge C \Vdash C$, we obtain $C' \wedge C, \Gamma \vdash \mathtt{t} : \sigma$ **(5)**. By HMX-EXISTS, (5) and (2) yield $\exists \bar{\mathtt{X}}.(C' \wedge C), \Gamma \vdash \mathtt{t} : \sigma$ **(6)**. By (4) and C-EXAND, the constraint $\exists \bar{\mathtt{X}}.(C' \wedge C)$ is equivalent to $C' \wedge \exists \bar{\mathtt{X}}.C$, which, by (3) and C-DUP, is equivalent to $C'$. Thus, (6) is the goal $C', \Gamma \vdash \mathtt{t} : \sigma$. □

We do not give a direct type soundness proof for $\mathrm{HM}(X)$. Instead, in section 1.5, we prove that it is equivalent to another type system, which later is itself proven sound. A direct type soundness result, based on a denotational semantics, may be found in (Odersky, Sulzmann, and Wehr, 1999a). Another type soundness proof, which follows Wright and Felleisen's syntactic approach (1994b), appears in (Skalka and Pottier, 2002). Last, a hybrid approach, which combines some of the advantages of the previous two, is given in (Pottier, 2001a).

## An alternate presentation of HM($X$)

The presentation of $\mathrm{HM}(X)$ given in Figure 1-7 has only four syntax-directed rules out of eight. It is a good specification of the type system, but it is far from an algorithmic description. As a first step towards such a description, we provide an alternate presentation of $\mathrm{HM}(X)$, where generalization is performed only at $\mathtt{let}$ expressions and instantiation takes place only at references to program identifiers (Figure 1-8). It has the property that all judgements are of the form $C, \Gamma \vdash \mathtt{t} : \mathtt{T}$, rather than $C, \Gamma \vdash \mathtt{t} : \sigma$. The following theorem states that the two presentations are indeed equivalent.

1.4.3 THEOREM: $C, \Gamma \vdash \mathtt{t} : \mathtt{T}$ is derivable via the rules of Figure 1-8 if and only if it is a valid $\mathrm{HM}(X)$ judgement. □

$$\frac{\Gamma(\mathtt{x}) = \forall \bar{\mathtt{x}}[D].\mathtt{T}}{C \wedge D, \Gamma \vdash \mathtt{x} : \mathtt{T}} \quad (\textsc{hmd-VarInst})$$

$$\frac{C, (\Gamma; \mathtt{z} : \mathtt{T}) \vdash \mathtt{t} : \mathtt{T}'}{C, \Gamma \vdash \lambda \mathtt{z}.\mathtt{t} : \mathtt{T} \to \mathtt{T}'} \quad (\textsc{hmd-Abs})$$

$$\frac{C, \Gamma \vdash \mathtt{t}_1 : \mathtt{T} \to \mathtt{T}' \qquad C, \Gamma \vdash \mathtt{t}_2 : \mathtt{T}}{C, \Gamma \vdash \mathtt{t}_1 \ \mathtt{t}_2 : \mathtt{T}'} \quad (\textsc{hmd-App})$$

$$\frac{\begin{array}{c} C \wedge D, \Gamma \vdash \mathtt{t}_1 : \mathtt{T}_1 \qquad \bar{\mathtt{x}} \mathbin{\#} \mathit{ftv}(C, \Gamma) \\ C \wedge \exists \bar{\mathtt{x}}.D, (\Gamma; \mathtt{z} : \forall \bar{\mathtt{x}}[D].\mathtt{T}_1) \vdash \mathtt{t}_2 : \mathtt{T}_2 \end{array}}{C \wedge \exists \bar{\mathtt{x}}.D, \Gamma \vdash \mathtt{let}\ \mathtt{z} = \mathtt{t}_1\ \mathtt{in}\ \mathtt{t}_2 : \mathtt{T}_2} \quad (\textsc{hmd-LetGen})$$

$$\frac{C, \Gamma \vdash \mathtt{t} : \mathtt{T} \qquad C \Vdash \mathtt{T} \le \mathtt{T}'}{C, \Gamma \vdash \mathtt{t} : \mathtt{T}'} \quad (\textsc{hmd-Sub})$$

$$\frac{C, \Gamma \vdash \mathtt{t} : \mathtt{T} \qquad \bar{\mathtt{x}} \mathbin{\#} \mathit{ftv}(\Gamma, \mathtt{T})}{\exists \bar{\mathtt{x}}.C, \Gamma \vdash \mathtt{t} : \mathtt{T}} \quad (\textsc{hmd-Exists})$$

**Figure 1-8: An alternate presentation of HM($X$)**

This theorem shows that the rule sets of Figures 1-7 and 1-8 derive the same monomorphic judgements, that is, the same judgements of the form $C, \Gamma \vdash \mathtt{t} : \mathtt{T}$. The fact that judgements of the form $C, \Gamma \vdash \mathtt{t} : \sigma$, where $\sigma$ is a not a monotype, cannot be derived using the new rule set is a technical simplification, without deep significance; the first two exercises below shed some light on this issue.

1.4.4    EXERCISE [★★]: Show that both rule sets lead to the same set of *well-typed* programs. □

1.4.5    EXERCISE [★★]: Show that, if hmx-Gen is added to the rule set of Figure 1-8, then both rule sets derive exactly the same judgements. □

1.4.6    EXERCISE [★★★, ↛]: Show that it is possible to simplify the presentation of Damas and Milner's type system in an analogous manner. That is, define an alternate set of typing rules for DM, which allows deriving judgements of the form $\Gamma \vdash \mathtt{t} : \mathtt{T}$; then, show that this new rule set is equivalent to the previous one, in the same sense as above. Which auxiliary properties of DM does your proof require? A solution is given in (Clément, Despeyroux, Despeyroux, and Kahn, 1986). □

### Relating HM($X$) with Damas and Milner's type system

In order to explain our interest in HM($X$), we wish to show that it is more general than Damas and Milner's type system. Since HM($X$) really is a *family* of type systems, we must make this statement more precise. First, every member of the HM($X$) family contains DM. Conversely, DM contains HM(=), the

constraint-based type system obtained by specializing $\mathrm{HM}(X)$ to the setting of an equality-only syntactic model.

The first of these assertions is easy to prove, because the mapping from DM judgements to $\mathrm{HM}(X)$ judgements is essentially the identity: every valid DM judgement may be viewed as a valid $\mathrm{HM}(X)$ judgement under the trivial assumption true. This statement relies on the fact that the DM type scheme $\forall \bar{\mathsf{X}}.\mathsf{T}$ is identified with the constrained type scheme $\forall \bar{\mathsf{X}}[\mathsf{true}].\mathsf{T}$, so DM type schemes (resp. environments) form a subset of $\mathrm{HM}(X)$ type schemes (resp. environments). Its proof is routine, except perhaps in the case of DM-INST, where it is shown how the effect of applying a substitution in DM is emulated by strengthening the current constraint in $\mathrm{HM}(X)$.

1.4.7   THEOREM: If $\Gamma \vdash \mathsf{t} : \mathsf{S}$ holds in DM, then $\mathsf{true}, \Gamma \vdash \mathsf{t} : \mathsf{S}$ holds in $\mathrm{HM}(X)$.   □

*Proof:*   The proof is by structural induction on a derivation of $\Gamma \vdash \mathsf{t} : \mathsf{S}$. In each proof case, we adopt the notations of Figure 1-3.

○ *Case* DM-VAR. The rule's conclusion is $\Gamma \vdash \mathsf{x} : \mathsf{S}$. Its premise is $\Gamma(\mathsf{x}) = \mathsf{S}$ **(1)**. By definition and by C-EX*, the constraint $\exists \mathsf{S}$ is equivalent to true. By applying HMX-VAR to (1) and to the assertion $\mathsf{true} \Vdash \mathsf{true}$, we obtain $\mathsf{true}, \Gamma \vdash \mathsf{x} : \mathsf{S}$.

○ *Cases* DM-ABS, DM-APP, DM-LET. By the induction hypothesis and by HMX-ABS, HMX-APP or HMX-LET, respectively.

○ *Case* DM-GEN. The rule's conclusion is $\Gamma \vdash \mathsf{t} : \forall \bar{\mathsf{X}}.\mathsf{T}$. Its premises are $\Gamma \vdash \mathsf{t} : \mathsf{T}$ **(1)** and $\bar{\mathsf{X}} \mathbin{\#} \mathit{ftv}(\Gamma)$ **(2)**. Applying the induction hypothesis to (1) yields $\mathsf{true}, \Gamma \vdash \mathsf{t} : \mathsf{T}$ **(3)**. Furthermore, (2) implies $\bar{\mathsf{X}} \mathbin{\#} \mathit{ftv}(\mathsf{true}, \Gamma)$ **(4)**. By HMX-GEN, (3) and (4) yield $\mathsf{true}, \Gamma \vdash \mathsf{t} : \forall \bar{\mathsf{X}}[\mathsf{true}].\mathsf{T}$.

○ *Case* DM-INST. The rule's conclusion is $\Gamma \vdash \mathsf{t} : [\vec{\mathsf{X}} \mapsto \vec{\mathsf{T}}]\mathsf{T}$. Its premise is $\Gamma \vdash \mathsf{t} : \forall \bar{\mathsf{X}}.\mathsf{T}$ **(1)**. We may assume, *w.l.o.g.*, $\bar{\mathsf{X}} \mathbin{\#} \mathit{ftv}(\Gamma, \bar{\mathsf{T}})$ **(2)**. Applying the induction hypothesis to (1) yields $\mathsf{true}, \Gamma \vdash \mathsf{t} : \forall \bar{\mathsf{X}}[\mathsf{true}].\mathsf{T}$ **(3)**. By HMX-INST, (3) implies $\mathsf{true}, \Gamma \vdash \mathsf{t} : \mathsf{T}$ **(4)**. By Lemma 1.4.2, we may weaken this judgement so as to obtain $\vec{\mathsf{X}} = \vec{\mathsf{T}}, \Gamma \vdash \mathsf{t} : \mathsf{T}$ **(5)**. Using C-EQ, C-EXTRANS, and C-EXAND, it is possible to establish $\vec{\mathsf{X}} = \vec{\mathsf{T}} \Vdash \mathsf{T} = [\vec{\mathsf{X}} \mapsto \vec{\mathsf{T}}]\mathsf{T}$ **(6)**. Applying HMX-SUB to (5) and (6), we find $\vec{\mathsf{X}} = \vec{\mathsf{T}}, \Gamma \vdash \mathsf{t} : [\vec{\mathsf{X}} \mapsto \vec{\mathsf{T}}]\mathsf{T}$ **(7)**. Last, (2) implies $\bar{\mathsf{X}} \mathbin{\#} \mathit{ftv}(\Gamma, [\vec{\mathsf{X}} \mapsto \vec{\mathsf{T}}]\mathsf{T})$ **(8)**. Applying HMX-EXISTS to (7) and (8), we obtain $\exists \bar{\mathsf{X}}.(\vec{\mathsf{X}} = \vec{\mathsf{T}}), \Gamma \vdash \mathsf{t} : [\vec{\mathsf{X}} \mapsto \vec{\mathsf{T}}]\mathsf{T}$ **(9)**. By (2) and C-NAME, the constraint $\exists \bar{\mathsf{X}}.(\vec{\mathsf{X}} = \vec{\mathsf{T}})$ is equivalent to true, so (9) is the goal.   □

We are now interested in proving that $\mathrm{HM}(=)$, as defined above, is contained within DM. To this end, we must translate every $\mathrm{HM}(=)$ judgement to a DM judgement. It quickly turns out that this is possible if the original judgement's constraint assumption is *satisfiable*.

We begin by explaining how an HM($=$) is translated into a DM type scheme. Such a translation is made possible by the fact that the definition of HM($=$) assumes an equality-only syntactic model. Indeed, in that setting, every satisfiable constraint admits a most general unifier (Definition 1.3.41), whose properties we make essential use of.

In fact, we must not only translate a type scheme, but also apply a type substitution to it. Instead of separating these steps, we perform both at once, and parameterize the translation by a type substitution $\theta$. (It does not appear that separating them would help.) The definition of $[\![\sigma]\!]_\theta$ is somewhat involved: it is given in the statement of the following lemma, whose proof establishes that the definition is indeed well-formed.

1.4.8    LEMMA: Consider a type scheme $\sigma$ and an idempotent type substitution $\theta$ such that $ftv(\sigma) \subseteq dom(\theta)$ **(1)** and $\exists\theta \Vdash \exists\sigma$ **(2)**. Write $\sigma = \forall\bar{\mathrm{X}}[D].\mathrm{T}$, where $\bar{\mathrm{X}} \mathbin{\#} \theta$ **(3)**. Then, there exists a type substitution $\theta'$ such that $\theta'$ extends $\theta$, $dom(\theta')$ is $dom(\theta) \cup \bar{\mathrm{X}}$, and $\theta'$ is a most general unifier of $\exists\theta \wedge D$. Let $\bar{\mathrm{Y}} = ftv(\theta'(\bar{\mathrm{X}})) \setminus range(\theta)$. Then, the *translation* of $\sigma$ under $\theta$, written $[\![\sigma]\!]_\theta$, is the DM type scheme $\forall\bar{\mathrm{Y}}.\theta'(\mathrm{T})$. This is a well-formed definition. Furthermore, $ftv([\![\sigma]\!]_\theta) \subseteq range(\theta)$ holds.                          $\square$

*Proof:*   By (2), $\exists\theta$ is equivalent to $\exists\theta \wedge \exists\sigma$, which may be written $\exists\theta \wedge \exists\bar{\mathrm{X}}.D$. By (3) and C-EXAND, this is $\exists\bar{\mathrm{X}}.(\exists\theta \wedge D)$. Thus, because $\theta$ is a most general unifier of $\exists\theta$, $\theta$ is also a most general unifier of $\exists\bar{\mathrm{X}}.(\exists\theta \wedge D)$ **(4)**. Furthermore, $ftv(\exists\bar{\mathrm{X}}.(\exists\theta \wedge D))$ is $ftv(\exists\theta \wedge \exists\sigma)$, which by definition of $\exists\theta$ and by (1) is a subset of $dom(\theta)$ **(5)**. By (4), (3), (5), and Lemma 1.3.45, there exists a type substitution $\theta'$ such that $\theta'$ extends $\theta$ **(6)** and $\theta'$ is a most general unifier of $\exists\theta \wedge D$ **(7)** and $dom(\theta') = dom(\theta) \cup \bar{\mathrm{X}}$ **(8)**.

Let us now define $\bar{\mathrm{Y}} = ftv(\theta'(\bar{\mathrm{X}})) \setminus range(\theta)$ and $[\![\sigma]\!]_\theta = \forall\bar{\mathrm{Y}}.\theta'(\mathrm{T})$. By (1), we have $ftv(\mathrm{T}) \subseteq \bar{\mathrm{X}} \cup dom(\theta)$. Applying $\theta'$ and exploiting (6), we find $ftv(\theta'(\mathrm{T})) \subseteq ftv(\theta'(\bar{\mathrm{X}})) \cup range(\theta)$, which by definition of $\bar{\mathrm{Y}}$ may be written $ftv(\theta'(\mathrm{T})) \subseteq \bar{\mathrm{Y}} \cup range(\theta)$. Subtracting $\bar{\mathrm{Y}}$ on each side, we find $ftv([\![\sigma]\!]_\theta) \subseteq range(\theta)$ **(9)**.

To show that the definition of $[\![\sigma]\!]_\theta$ is valid, there remains to show that it does not depend on the choice of $\bar{\mathrm{X}}$ or $\theta'$. To prove the former, it suffices to establish $\bar{\mathrm{X}} \mathbin{\#} ftv([\![\sigma]\!]_\theta)$, which indeed follows from (3) and (9). As for the latter, because of the constraints imposed by (6), (7), and (8), and by Lemma 1.3.44, distinct choices of $\theta'$ may differ only by a renaming of $ftv(\theta'(\bar{\mathrm{X}})) \setminus range(\theta)$, that is, $\bar{\mathrm{Y}}$. So, we must check $\bar{\mathrm{Y}} \mathbin{\#} ftv([\![\sigma]\!]_\theta)$, which holds by definition.        $\square$

Please note that if $\sigma$ is in fact a type $\mathrm{T}$, where $ftv(\mathrm{T}) \subseteq dom(\theta)$, then $\bar{\mathrm{X}}$ is empty, so $\theta'$ is $\theta$, $\bar{\mathrm{Y}}$ is empty, and $[\![\mathrm{T}]\!]_\theta = \theta(\mathrm{T})$. In other words, the translation of a type under $\theta$ is its image through $\theta$. More generally, the translation of an unconstrained type scheme (that is, a type scheme whose constraint is `true`) is its image through $\theta$, as stated by the following exercise.

1.4.9    EXERCISE [★★, ↛]:  Prove that $[\![ \forall \bar{\mathtt{X}}.\mathtt{T} ]\!]_\theta$, when defined, is $\theta(\forall \bar{\mathtt{X}}.\mathtt{T})$. □

The translation becomes more than a mere type substitution when applied to a nontrivial constrained type scheme. Some examples of this situation are given below.

1.4.10    EXAMPLE: Let $\sigma = \forall \mathtt{XY}[\mathtt{X} = \mathtt{Y} \to \mathtt{Y}].\mathtt{X}$. Let $\theta$ be the identity substitution. The type scheme $\sigma$ is closed and the constraint $\exists \sigma$ is equivalent to true, so $[\![ \sigma ]\!]_\theta$ is defined. We must find a type substitution $\theta'$ whose domain is $\mathtt{XY}$ and that is a most general unifier of $\mathtt{X} = \mathtt{Y} \to \mathtt{Y}$. All such substitutions are of the form $[\mathtt{X} \mapsto (\mathtt{Z} \to \mathtt{Z}), \mathtt{Y} \mapsto \mathtt{Z}]$, where $\mathtt{Z}$ is fresh. We have $ftv(\theta'(\mathtt{XY})) = \mathtt{Z}$, whence $[\![ \sigma ]\!]_\theta = \forall \mathtt{Z}.\mathtt{Z} \to \mathtt{Z}$. Note that the choice of $\mathtt{Z}$ does not matter, since it is bound in $[\![ \sigma ]\!]_\theta$. Roughly speaking, the effect of the translation was to replace the body $\mathtt{X}$ of the constrained type scheme with its most general solution under the constraint $\mathtt{X} = \mathtt{Y} \to \mathtt{Y}$.

Let $\sigma = \forall \mathtt{XY}_1[\mathtt{X} = \mathtt{Y}_1 \to \mathtt{Y}_2].\mathtt{X}$. Let $\theta = [\mathtt{Y}_2 \mapsto \mathtt{Z}_2]$. We have $ftv(\sigma) = \mathtt{Y}_2 \subseteq dom(\theta)$. The constraint $\exists \sigma$ is equivalent to true, so $[\![ \sigma ]\!]_\theta$ is defined. We must find a type substitution $\theta'$ whose domain is $\mathtt{XY}_1\mathtt{Y}_2$ that extends $\theta$ and that is a most general unifier of $\mathtt{X} = \mathtt{Y}_1 \to \mathtt{Y}_2$. All such substitutions are of the form $[\mathtt{X} \mapsto (\mathtt{Z}_1 \to \mathtt{Z}_2), \mathtt{Y}_1 \mapsto \mathtt{Z}_1, \mathtt{Y}_2 \mapsto \mathtt{Z}_2]$, where $\mathtt{Z}_1$ is fresh. We have $ftv(\theta'(\mathtt{XY}_1)) \setminus range(\theta) = \mathtt{Z}_1\mathtt{Z}_2 \setminus \mathtt{Z}_2 = \mathtt{Z}_1$, whence $[\![ \sigma ]\!]_\theta = \forall \mathtt{Z}_1.\mathtt{Z}_1 \to \mathtt{Z}_2$. The type variable $\mathtt{Z}_2$ is *not* universally quantified—even though it appears in the image of $\mathtt{X}$, which *was* universally quantified in $\sigma$—because $\mathtt{Z}_2$ is the image of $\mathtt{Y}_2$, which was free in $\sigma$. □

Before attacking the main theorem, let us establish a couple of technical properties of the translation. First, $[\![ \sigma ]\!]_\theta$ is insensitive to the behavior of $\theta$ outside $ftv(\sigma)$, a natural property, since our informal intent is for $\theta$ to be applied to $\sigma$.

1.4.11    LEMMA:  If $\theta_1$ and $\theta_2$ coincide on $ftv(\sigma)$, then $[\![ \sigma ]\!]_{\theta_1}$ and $[\![ \sigma ]\!]_{\theta_2}$ are either both undefined, or both defined and identical. □

Second, if $C \Vdash \sigma \preceq \mathtt{T}'$ holds, then the translations of $\sigma$ and $\mathtt{T}'$ under a most general unifier of $C$ are in Damas and Milner's instance relation. One might say, roughly speaking, that the instance relation is preserved by the translation.

1.4.12    LEMMA:  Let $ftv(\sigma, \mathtt{T}') \subseteq dom(\theta)$ **(1)** and $\exists \theta \Vdash \exists \sigma$ **(2)**. Let $\exists \theta \Vdash \sigma \preceq \mathtt{T}'$ **(3)**. Then, $\theta(\mathtt{T}')$ is an instance of the DM type scheme $[\![ \sigma ]\!]_\theta$. □

*Proof:*  Write $\sigma = \forall \bar{\mathtt{X}}[D].\mathtt{T}$, where $\bar{\mathtt{X}} \mathbin{\#} \theta$ **(4)** and $\bar{\mathtt{X}} \mathbin{\#} ftv(\mathtt{T}')$ **(5)**. By **(1)**, **(2)**, and **(4)**, one may define $\theta'$, $\bar{\mathtt{Y}}$, and $[\![ \sigma ]\!]_\theta$ exactly as in the statement of Lemma 1.4.8. By **(5)** and Definition 1.3.3, **(3)** is synonymous with $\exists \theta \Vdash \exists \bar{\mathtt{X}}.(D \wedge$

$\mathtt{T} = \mathtt{T}'$). Reasoning in the same manner as in the first paragraph of the proof of Lemma 1.4.8, we find that there exists a type substitution $\theta''$ such that $\theta''$ extends $\theta$, $dom(\theta'')$ is $dom(\theta) \cup \bar{\mathtt{x}}$, and $\theta''$ is a most general unifier of $\exists \theta \wedge D \wedge \mathtt{T} = \mathtt{T}'$.

We have $dom(\theta') = dom(\theta'')$ **(6)**. Furthermore, $\theta'$ is a most general unifier of $\exists \theta \wedge D$, while $\theta''$ is a most general unifier of $\exists \theta \wedge D \wedge \mathtt{T} = \mathtt{T}'$, which implies $\exists \theta'' \Vdash \exists \theta'$ **(7)**. By Lemma 1.3.39, $\theta''$ refines $\theta'$. That is, there exists a type substitution $\varphi$ such that $\theta''$ is the restriction of $\varphi \circ \theta'$ to $dom(\theta) \cup \bar{\mathtt{x}}$ **(8)**. We may require $dom(\varphi) \subseteq range(\theta) \cup ftv(\theta'(\bar{\mathtt{x}}))$ **(9)** without compromising (8).

Consider $\mathtt{X} \in dom(\theta)$. Because $\theta''$ extends $\theta$, we have $\theta''(\mathtt{X}) = \theta(\mathtt{X})$ **(10)**. Furthermore, by (8), we have $\theta''(\mathtt{X}) = (\varphi \circ \theta')(\mathtt{X}) = (\varphi \circ \theta)(\mathtt{X})$ **(11)**. Using (10) and (11), we find $\theta(\mathtt{X}) = \varphi(\theta(\mathtt{X}))$. Because this holds for every $\mathtt{X} \in dom(\theta)$, $\varphi$ must be the identity over $range(\theta)$; that is, $dom(\varphi) \# range(\theta)$ **(12)** holds. Combining (9) and (12), we find $dom(\varphi) \subseteq ftv(\theta'(\bar{\mathtt{x}})) \setminus range(\theta)$, that is, $dom(\varphi) \subseteq \bar{\mathtt{Y}}$ **(13)**.

By construction of $\theta''$, we have $\exists \theta'' \Vdash \mathtt{T} = \mathtt{T}'$. By Lemma 1.3.29, this implies $\theta''(\exists \theta'') \Vdash \theta''(\mathtt{T}) = \theta''(\mathtt{T}')$, which by Lemma 1.3.42 may be read $\mathsf{true} \Vdash \theta''(\mathtt{T}) = \theta''(\mathtt{T}')$. By Lemma 1.3.32, $\theta''(\mathtt{T})$ and $\theta''(\mathtt{T}')$ coincide. Because by (1) $ftv(\mathtt{T})$ is a subset of $dom(\theta) \cup \bar{\mathtt{x}}$ and by (8), the former may be written $\varphi(\theta'(\mathtt{T}))$. By (1) and because $\theta''$ extends $\theta$, the latter is $\theta(\mathtt{T}')$. Thus, we have $\varphi(\theta'(\mathtt{T})) = \theta(\mathtt{T}')$. Together with (13), this establishes that $\theta(\mathtt{T}')$ is an instance of $\forall \bar{\mathtt{Y}}.\theta'(\mathtt{T})$, that is, $[\![\sigma]\!]_\theta$. $\qquad\square$

We extend the translation to environments as follows. $[\![\varnothing]\!]_\theta$ is $\varnothing$. If $\exists \theta \Vdash \exists \sigma$ holds, then $[\![\Gamma; \mathtt{x} : \sigma]\!]_\theta$ is $[\![\Gamma]\!]_\theta; \mathtt{x} : [\![\sigma]\!]_\theta$, otherwise it is $[\![\Gamma]\!]_\theta$. Notice that $[\![\Gamma]\!]_\theta$ contains fewer bindings than $\Gamma$, which ensures that bindings $x : \sigma$ for which $\exists \theta \Vdash \exists \sigma$ does not hold will not be used in the translation. Please note that $[\![\Gamma]\!]_\theta$ is defined when $ftv(\Gamma) \subseteq dom(\theta)$ holds.

We are now ready to prove the main theorem. Please note that, by requiring $\theta$ to be a most general unifier of $C$, we also require $C$ to be satisfiable. Judgements that carry an unsatisfiable constraint cannot be translated.

1.4.13   THEOREM: Let $C, \Gamma \vdash \mathtt{t} : \sigma$ hold in HM(=). Let $\theta$ be a most general unifier of $C$ such that $ftv(\Gamma, \sigma) \subseteq dom(\theta)$. Then, $[\![\Gamma]\!]_\theta \vdash \mathtt{t} : [\![\sigma]\!]_\theta$ holds in DM. $\qquad\square$

*Proof:*  Let us first remark that, by Lemma 1.4.1, we have $C \Vdash \exists \sigma$. This may be written $\exists \theta \Vdash \exists \sigma$, which guarantees that $[\![\sigma]\!]_\theta$ is defined. The proof is by structural induction on an HM(=) typing derivation. We assume that the derivation is expressed in terms of the rules of Figure 1-8, but split HMD-LETGEN into HMX-LET and HMX-GEN for the sake of readability.

∘ *Case* HMD-VARINST. The rule's conclusion is $C \wedge D, \Gamma \vdash \mathtt{x} : \mathtt{T}$. By hypothesis, $\theta$ is a most general unifier of $C \wedge D$ **(1)**, and $ftv(\mathtt{T}) \subseteq dom(\theta)$ **(2)**

holds. The rule's premise is $\Gamma(\mathtt{x}) = \sigma$ **(3)**, where $\sigma$ stands for $\forall \bar{\mathtt{X}}[D].\mathtt{T}$. By (1), we have $\exists \theta \equiv C \wedge D \Vdash D \Vdash \exists \bar{\mathtt{X}}.D \equiv \exists \sigma$ **(4)**. Furthermore, we have $ftv(\sigma) \subseteq ftv(\Gamma) \subseteq dom(\theta)$ **(5)**. These facts show that $[\![\sigma]\!]_\theta$ is defined. To-gether with (3), this implies $[\![\Gamma]\!]_\theta(\mathtt{x}) = [\![\sigma]\!]_\theta$. By DM-VAR, $[\![\Gamma]\!]_\theta \vdash \mathtt{x} : [\![\sigma]\!]_\theta$ **(6)** follows. Now, by Lemma 1.3.19, we have $D \Vdash \sigma \preceq \mathtt{T}$, which, combined with $\exists \theta \Vdash D$, yields $\exists \theta \Vdash \sigma \preceq \mathtt{T}$ **(7)**. By (7), (4), (5), (2), and Lemma 1.4.12, we find that $\theta(\mathtt{T})$ is an instance of $[\![\sigma]\!]_\theta$. Thus, applying DM-INST to (6) yields $[\![\Gamma]\!]_\theta \vdash \mathtt{t} : \theta(\mathtt{T})$.

○ *Case* HMD-ABS. The rule's conclusion is $C, \Gamma \vdash \lambda \mathtt{z}.\mathtt{t} : \mathtt{T} \to \mathtt{T}'$. Its premise is $C, (\Gamma; \mathtt{z} : \mathtt{T}) \vdash \mathtt{t} : \mathtt{T}'$. Applying the induction hypothesis to it yields $[\![\Gamma]\!]_\theta; \mathtt{z} : \theta(\mathtt{T}) \vdash \mathtt{t} : \theta(\mathtt{T}')$. By DM-ABS, this implies $[\![\Gamma]\!]_\theta \vdash \lambda \mathtt{z}.\mathtt{t} : \theta(\mathtt{T}) \to \theta(\mathtt{T}')$, that is, $[\![\Gamma]\!]_\theta \vdash \lambda \mathtt{z}.\mathtt{t} : \theta(\mathtt{T} \to \mathtt{T}')$.

○ *Case* HMD-APP. By an extension of $dom(\theta)$ to include $ftv(\mathtt{T})$, by the induction hypothesis, and by DM-APP.

○ *Case* HMX-LET. By an extension of $dom(\theta)$ to include $ftv(\sigma)$, by the induction hypothesis, and by DM-LET.

○ *Case* HMX-GEN. The rule's conclusion is $C \wedge \exists \sigma, \Gamma \vdash \mathtt{t} : \sigma$, where $\sigma$ stands for $\forall \bar{\mathtt{X}}[D].\mathtt{T}$. By hypothesis, $\theta$ is a most general unifier of $C \wedge \exists \sigma$ **(1)**, and $ftv(\Gamma, \sigma) \subseteq dom(\theta)$ **(2)** holds. The rule's premises are $C \wedge D, \Gamma \vdash \mathtt{t} : \mathtt{T}$ **(3)** and $\bar{\mathtt{X}} \# ftv(C, \Gamma)$ **(4)**. We may further assume, *w.l.o.g.*, $\bar{\mathtt{X}} \# \theta$ **(5)**. Given (1), (2), and (5), we may define $\theta'$ and $\bar{\mathtt{Y}}$ exactly as in Lemma 1.4.8. Then, $\theta'$ is a most general unifier of $\exists \theta \wedge D$, that is, $C \wedge D$. Furthermore, $dom(\theta')$ is $dom(\theta) \cup \bar{\mathtt{X}}$, which by (2) is a superset of $ftv(\Gamma, \mathtt{T})$. Thus, the induction hypothesis applies to $\theta'$ and to (3), yielding $[\![\Gamma]\!]_{\theta'} \vdash \mathtt{t} : \theta'(\mathtt{T})$. Because $\theta'$ extends $\theta$, by (2) and by Lemma 1.4.11, this may be read $[\![\Gamma]\!]_\theta \vdash \mathtt{t} : \theta'(\mathtt{T})$ **(6)**. According to Lemma 1.4.8, we have $ftv([\![\Gamma]\!]_\theta) \subseteq range(\theta)$, which by construction of $\bar{\mathtt{Y}}$ implies $\bar{\mathtt{Y}} \# ftv([\![\Gamma]\!]_\theta)$ **(7)**. By DM-GEN, (6) and (7) yield $[\![\Gamma]\!]_\theta \vdash \mathtt{t} : \forall \bar{\mathtt{Y}}.\theta'(\mathtt{T})$, that is, $[\![\Gamma]\!]_\theta \vdash \mathtt{t} : [\![\sigma]\!]_\theta$.

○ *Case* HMD-SUB. The rule's conclusion is $C, \Gamma \vdash \mathtt{t} : \mathtt{T}'$. By hypothesis, $\theta$ is a most general unifier of $C$ **(1)**, and $ftv(\Gamma, \mathtt{T}') \subseteq dom(\theta)$ **(2)** holds. The goal is $[\![\Gamma]\!]_\theta \vdash \mathtt{t} : \theta(\mathtt{T}')$ **(3)**. The rule's premises are $C, \Gamma \vdash \mathtt{t} : \mathtt{T}$ **(4)** and $C \Vdash \mathtt{T} = \mathtt{T}'$ **(5)**. We may assume, *w.l.o.g.*, $ftv(\mathtt{T}) \# range(\theta)$ **(6)**. Then, by (6) and Lemma 1.3.43, we may extend the domain of $\theta$, so as to achieve $ftv(\mathtt{T}) \subseteq dom(\theta)$ **(7)**, without compromising (1) or (2) or affecting the goal (3). By (1), (2), and (7), the induction hypothesis applies to (4), yielding $[\![\Gamma]\!]_\theta \vdash \mathtt{t} : \theta(\mathtt{T})$ **(8)**. Now, thanks to (1), (5) may be read $\exists \theta \Vdash \mathtt{T} = \mathtt{T}'$, which by Lemmas 1.3.29 and 1.3.42 implies $\mathtt{true} \Vdash \theta(\mathtt{T}) = \theta(\mathtt{T}')$. Then, Lemma 1.3.32 shows that $\theta(\mathtt{T})$ and $\theta(\mathtt{T}')$ coincide. As a result, (8) is the goal (3).

○ *Case* HMD-EXISTS. The rule's conclusion is $\exists \bar{\mathtt{X}}.C, \Gamma \vdash \mathtt{t} : \mathtt{T}$. By hypothesis, $\theta$ is a most general unifier of $\exists \bar{\mathtt{X}}.C$ **(1)**, and $ftv(\Gamma, \mathtt{T}) \subseteq dom(\theta)$ **(2)** holds. The

rule's premises are $C, \Gamma \vdash \mathtt{t} : \mathtt{T}$ **(3)** and $\bar{\mathtt{X}} \mathrel{\#} \mathit{ftv}(\Gamma, \mathtt{T})$. We may assume, *w.l.o.g.*, $\bar{\mathtt{X}} \mathrel{\#} \theta$ **(4)**. As in the previous case, we may extend the domain of $\theta$ to guarantee $\mathit{ftv}(\exists\bar{\mathtt{X}}.C) \subseteq \mathit{dom}(\theta)$ **(5)**. By (1), (4), (5), and Lemma 1.3.45, there exists a type substitution $\theta'$ such that $\theta'$ extends $\theta$ **(6)** and $\theta'$ is a most general unifier of $C$. Applying the induction hypothesis to $\theta'$ and to (3) yields $[\![\Gamma]\!]_{\theta'} \vdash \mathtt{t} : \theta'(\mathtt{T})$. By (2), (6), and Lemma 1.4.11, this may be read $[\![\Gamma]\!]_\theta \vdash \mathtt{t} : \theta(\mathtt{T})$. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

Together, Theorems 1.4.7 and 1.4.13 yield a precise correspondence between DM and HM(=): there exists a compositional translation from each to the other. In other words, they may be viewed as two equivalent formulations of a single type system. One might also say that HM(=) is a constraint-based formulation of DM. Furthermore, Theorem 1.4.7 states that every member of the HM($X$) family is an extension of DM. This explains our double interest in HM($X$), as an alternate formulation of DM, which we believe is more pleasant, for reasons already discussed, and as a more expressive framework.

## 1.5   A purely constraint-based type system: PCB($X$)

In the previous section, we have presented HM($X$), an elegant constraint-based extension of Damas and Milner's type system. However, HM($X$), as described there, suffers from a drawback. A typing judgement involves both a constraint, which represents an assumption about its free type variables, and an environment, which represents an assumption about its free program identifiers. At a let node, HMD-LETGEN turns a part of the current constraint, namely $D$, into a type scheme, namely $\forall \bar{\mathtt{X}}[D].\mathtt{T}$, and stores it into the environment. Then, at every occurrence of the let-bound variable, HMD-VARINST retrieves this type scheme from the environment and adds a copy of $D$ back to the current constraint. In practice, it is important to *simplify* the type scheme $\forall \bar{\mathtt{X}}[D].\mathtt{T}$ *before* it is stored in the environment, because it would be inefficient to copy an unsimplified constraint. In other words, it appears that, in order to preserve efficiency, constraint generation and constraint simplification cannot be separated.

Of course, in practice, it is not difficult to intermix these phases, so the problem is not technical, but pedagogical. Indeed, we argued earlier that it is natural and desirable to separate them. *Type scheme introduction and elimination constraints*, which we introduced in Section 1.3 but did not use in the specification of HM($X$), are intended as a means of solving this problem. In the present section, we exploit them to give a novel formulation of HM($X$), which no longer requires copying constraints back and forth between the environment and the constraint assumption. In fact, the environment is

$$\frac{C \Vdash \mathtt{x} \preceq \mathtt{T}}{C \vdash \mathtt{x} : \mathtt{T}} \quad (\textsc{Var})$$

$$\frac{C \vdash \mathtt{t} : \mathtt{T}'}{\mathsf{let}\ \mathtt{z} : \mathtt{T}\ \mathsf{in}\ C \vdash \lambda \mathtt{z}.\mathtt{t} : \mathtt{T} \to \mathtt{T}'} \quad (\textsc{Abs})$$

$$\frac{C_1 \vdash \mathtt{t}_1 : \mathtt{T} \to \mathtt{T}' \qquad C_2 \vdash \mathtt{t}_2 : \mathtt{T}}{C_1 \wedge C_2 \vdash \mathtt{t}_1\ \mathtt{t}_2 : \mathtt{T}'} \quad (\textsc{App})$$

$$\frac{C_1 \vdash \mathtt{t}_1 : \mathtt{T}_1 \qquad C_2 \vdash \mathtt{t}_2 : \mathtt{T}_2}{\substack{\mathsf{let}\ \mathtt{z} : \forall \mathcal{V}[C_1].\mathtt{T}_1\ \mathsf{in}\ C_2 \\ \vdash\ \mathsf{let}\ \mathtt{z} = \mathtt{t}_1\ \mathsf{in}\ \mathtt{t}_2 : \mathtt{T}_2}} \quad (\textsc{Let})$$

$$\frac{C \vdash \mathtt{t} : \mathtt{T}}{C \wedge \mathtt{T} \leq \mathtt{T}' \vdash \mathtt{t} : \mathtt{T}'} \quad (\textsc{Sub})$$

$$\frac{C \vdash \mathtt{t} : \mathtt{T} \qquad \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(\mathtt{T})}{\exists \bar{\mathtt{X}}.C \vdash \mathtt{t} : \mathtt{T}} \quad (\textsc{Exists})$$

**Figure 1-9: Typing rules for** $\mathbf{PCB}(X)$

suppressed altogether: taking advantage of the new constraint forms, we encode information about program identifiers within the constraint assumption.

### Presentation

We now employ the full constraint language (Section 1.3). Typing judgements take the form $C \vdash \mathtt{t} : \mathtt{T}$, where $C$ may have free type variables *and* free program identifiers. The rules that allow deriving such judgements appear in Figure 1-9. As before, we identify judgements up to constraint equivalence.

Let us review the rules. VAR states that $\mathtt{x}$ has type $\mathtt{T}$ under any constraint that entails $\mathtt{x} \preceq \mathtt{T}$. Note that we no longer consult the type scheme associated with $\mathtt{x}$ in the environment—indeed, there is no environment. Instead, we let the constraint assumption record the fact that the type scheme should admit $\mathtt{T}$ as one of its instances. Thus, in a judgement $C \vdash \mathtt{t} : \mathtt{T}$, any program identifier that occurs free within $\mathtt{t}$ typically also occurs free within $C$. ABS requires the body $\mathtt{t}$ of a $\lambda$-abstraction to have type $\mathtt{T}'$ under assumption $C$. Although no explicit assumption about $\mathtt{z}$ appears in the premise, $C$ typically contains a number of instantiation constraints bearing on $\mathtt{z}$, of the form $\mathtt{z} \preceq \mathtt{T}_i$. In the rule's conclusion, $C$ is wrapped within the prefix $\mathsf{let}\ \mathtt{z} : \mathtt{T}\ \mathsf{in}\ []$, where $\mathtt{T}$ is the type assigned to $\mathtt{z}$. This effectively requires every $\mathtt{T}_i$ to denote a super-type of $\mathtt{T}$, as desired. Please note that $\mathtt{z}$ does not occur free in the constraint $\mathsf{let}\ \mathtt{z} : \mathtt{T}\ \mathsf{in}\ C$, which is natural, since it does not occur free in $\lambda \mathtt{z}.\mathtt{t}$. APP exhibits a minor stylistic difference with respect to HMX-APP: its constraint assumption is split between its premises. It is not difficult to prove that, when weakening holds (see Lemma 1.5.2 below), this choice does not affect the set of valid judgements. This new presentation encourages reading the rules in Figure 1-9 as the specification of an algorithm, which, given $\mathtt{t}$ and $\mathtt{T}$, pro-

duces $C$ such that $C \vdash \mathtt{t} : \mathtt{T}$ holds. In the case of App, the algorithm invokes itself recursively for each of the two subexpressions, yielding the constraints $C_1$ and $C_2$, then constructs their conjunction. Let is analogous to Abs: by wrapping $C_2$ within a let prefix, it gives meaning to the instantiation constraints bearing on $\mathtt{z}$ within $C_2$. The difference is that $\mathtt{z}$ may now be assigned a type scheme, as opposed to a monotype. An appropriate type scheme is built in the most straightforward manner from the constraint $C_1$ and the type $\mathtt{T}_1$ that describe $\mathtt{t}_1$. All of the type variables that appear free in the left-hand premise are generalized, hence the notation $\forall \mathcal{V}[C_1].\mathtt{T}_1$, which is a convenient shorthand for $\forall ftv(C_1, \mathtt{T}_1)[C_1].\mathtt{T}_1$. The side-condition that "type variables that occur free in the environment must not be generalized", which was present in DM and $\mathrm{HM}(X)$, naturally disappears, since judgements no longer involve an environment. Sub again exhibits a minor stylistic difference with respect to hmx-Sub: the comments made about App above apply here as well. Exists is essentially identical to hmx-Exists.

In the standard specification of $\mathrm{HM}(X)$, hmd-Abs and hmd-LetGen accumulate information in the environment. Through the environment, this information is made available to hmd-VarInst, which retrieves and copies it. Here, instead, no information is explicitly transmitted. Where a program identifier is bound, a type scheme introduction constraint is built; where a program identifier is used, a type scheme instantiation constraint is produced. The two are related only by our definition of the meaning of constraints.

The reader may be puzzled by the fact that Let allows *all* type variables that occur free in its left-hand premise to be generalized. The following exercise sheds some light on this issue.

1.5.1   Exercise [★, Recommended]: Build a type derivation for the expression $\lambda \mathtt{z}_1.\mathtt{let}\ \mathtt{z}_2 = \mathtt{z}_1\ \mathtt{in}\ \mathtt{z}_2$ within $\mathrm{PCB}(X)$. Draw a comparison with the solution of Exercise 1.2.21.                                                          □

The following lemma is an analogue of Lemma 1.4.2.

1.5.2   Lemma [Weakening]: If $C' \Vdash C$, then every derivation of $C \vdash \mathtt{t} : \mathtt{T}$ may be turned into a derivation of $C' \vdash \mathtt{t} : \mathtt{T}$ with the same shape.                     □

*Proof:*   The proof is by structural induction on a derivation of $C \vdash \mathtt{t} : \mathtt{T}$. In each proof case, we adopt the notations of Figure 1-9.

   ∘ *Case* Var. By transitivity of entailment.

   ∘ *Case* Abs. The rule's conclusion is $\mathtt{let}\ \mathtt{z} : \mathtt{T}\ \mathtt{in}\ C \vdash \lambda \mathtt{z}.\mathtt{t} : \mathtt{T} \to \mathtt{T}'\ \textbf{(1)}$. By hypothesis, we have $C' \Vdash \mathtt{let}\ \mathtt{z} : \mathtt{T}\ \mathtt{in}\ C\ \textbf{(2)}$. We may assume, *w.l.o.g.*, $\mathtt{z} \notin fpi(C')\ \textbf{(3)}$. The rule's premise is $C \vdash \mathtt{t} : \mathtt{T}'\ \textbf{(4)}$. Applying the induction hypothesis to (4) yields $C \wedge C' \vdash \mathtt{t} : \mathtt{T}'$, which by Abs implies $\mathtt{let}\ \mathtt{z} : \mathtt{T}\ \mathtt{in}\ (C\ \wedge$

$C'$) ⊢ $\lambda$z.t : T → T' **(5)**. By (3) and C-INAND*, let z : T in ($C \wedge C'$) is equivalent to (let z : T in $C$) $\wedge\ C'$, which by (2) and C-DUP is equivalent to $C'$. Thus, (5) is the goal $C'$ ⊢ $\lambda$z.t : T → T'.

○ *Case* APP. By applying the induction hypothesis to each premise, using the fact that $C'$ ⊩ $C_1 \wedge C_2$ implies $C'$ ⊩ $C_1$ and $C'$ ⊩ $C_2$.

○ *Case* LET. Analogous to the case of ABS. The induction hypothesis is applied to the second premise only.

○ *Case* SUB. Analogous to the case of APP.

○ *Case* EXISTS. See the corresponding case in the proof of Lemma 1.4.2. ☐

### Relating PCB($X$) with HM($X$)

Let us now provide evidence for our claim that PCB($X$) is an alternate presentation of HM($X$). The next two theorems define an effective translation from HM($X$) to PCB($X$) and back.

The first theorem states that if, within HM($X$), t has type T under assumptions $C$ and $\Gamma$, then, within PCB($X$), t also has type T, under some assumption $C'$. The relationship $C$ ⊩ let $\Gamma$ in $C'$ states that $C$ entails the residual constraint obtained by confronting $\Gamma$, which provides information about the free program identifiers in t, with $C'$, which contains instantiation constraints bearing on these program identifiers. The statement requires $C$ and $\Gamma$ to have no free program identifiers, which is natural, since they are part of an HM($X$) judgement. The hypothesis $C$ ⊩ $\exists\Gamma$ excludes the somewhat pathological situation where $\Gamma$ contains constraints not apparent in $C$. This hypothesis vanishes when $\Gamma$ is the initial environment; see Definition 1.7.3.

1.5.3 THEOREM: Let $C$ ⊩ $\exists\Gamma$. Assume $fpi(C, \Gamma) = \varnothing$. If $C, \Gamma$ ⊢ t : T holds in HM($X$), then there exists a constraint $C'$ such that $C'$ ⊢ t : T holds in PCB($X$) and $C$ entails let $\Gamma$ in $C'$. ☐

*Proof:* The proof is by structural induction on a derivation of $C, \Gamma$ ⊢ t : T. In each proof case, we adopt the notations of Figure 1-8.

○ *Case* HMD-VARINST. The rule's conclusion is $C \wedge D, \Gamma$ ⊢ x : T. By hypothesis, we have $C \wedge D$ ⊩ $\exists\Gamma$ **(1)** and $fpi(C, D, \Gamma) = \varnothing$ **(2)**. The rule's premise is $\Gamma(x) = \forall\bar{x}[D].T$ **(3)**. By VAR, we have x $\preceq$ T ⊢ x : T, so there remains to establish $C \wedge D$ ⊩ let $\Gamma$ in x $\preceq$ T **(4)**. By (3), (2), and C-INID, the constraint let $\Gamma$ in x $\preceq$ T is equivalent to let $\Gamma$ in $\forall\bar{x}[D].T \preceq$ T, which, by (2) and C-IN*, is itself equivalent to $\exists\Gamma \wedge \forall\bar{x}[D].T \preceq$ T **(5)**. By (1) and Lemma 1.3.19, $C \wedge D$ entails (5). We have established (4).

○ *Case* HMD-ABS. The rule's conclusion is $C, \Gamma$ ⊢ $\lambda$z.t : T → T'. Its premise is $C, (\Gamma; z : T)$ ⊢ t : T' **(1)**. The constraints $\exists\Gamma$ and $\exists(\Gamma; z : T)$ are equivalent,

so the induction hypothesis applies to (1) and yields a constraint $C'$ such that $C' \vdash \mathsf{t} : \mathsf{T}'$ **(2)** and $C \Vdash \mathsf{let}\ \Gamma; \mathsf{z} : \mathsf{T}\ \mathsf{in}\ C'$ **(3)**. Applying ABS to (2) yields $\mathsf{let}\ \mathsf{z} : \mathsf{T}\ \mathsf{in}\ C' \vdash \lambda \mathsf{z}.\mathsf{t} : \mathsf{T} \to \mathsf{T}'$. There remains to check that $C$ entails $\mathsf{let}\ \Gamma\ \mathsf{in}\ \mathsf{let}\ \mathsf{z} : \mathsf{T}\ \mathsf{in}\ C'$—but that is precisely (3).

∘ *Case* HMD-APP. The rule's conclusion is $C, \Gamma \vdash \mathsf{t}_1\ \mathsf{t}_2 : \mathsf{T}'$. Its premises are $C, \Gamma \vdash \mathsf{t}_1 : \mathsf{T} \to \mathsf{T}'$ **(1)** and $C, \Gamma \vdash \mathsf{t}_2 : \mathsf{T}$ **(2)**. Applying the induction hypothesis to (1) and (2), we obtain constraints $C'_1$ and $C'_2$ such that $C'_1 \vdash \mathsf{t}_1 : \mathsf{T} \to \mathsf{T}'$ **(3)** and $C'_2 \vdash \mathsf{t}_2 : \mathsf{T}$ **(4)** and $C \Vdash \mathsf{let}\ \Gamma\ \mathsf{in}\ C'_1$ **(5)** and $C \Vdash \mathsf{let}\ \Gamma\ \mathsf{in}\ C'_2$ **(6)**. By APP, (3) and (4) imply $C'_1 \wedge C'_2 \vdash \mathsf{t}_1\ \mathsf{t}_2 : \mathsf{T}'$. Furthermore, by C-INAND, (5) and (6) yield $C \Vdash \mathsf{let}\ \Gamma\ \mathsf{in}\ C'_1 \wedge C'_2$.

∘ *Case* HMD-LETGEN. The rule's conclusion is $C \wedge \exists \bar{\mathsf{x}}.D, \Gamma \vdash \mathsf{let}\ \mathsf{z} = \mathsf{t}_1\ \mathsf{in}\ \mathsf{t}_2 : \mathsf{T}_2$. By hypothesis, we have $C \wedge \exists \bar{\mathsf{x}}.D \Vdash \exists \Gamma$ **(1)** and $fpi(C, D, \Gamma) = \varnothing$ **(2)**. The rule's premises are $C \wedge D, \Gamma \vdash \mathsf{t}_1 : \mathsf{T}_1$ **(3)** and $\bar{\mathsf{x}}\ \#\ ftv(C, \Gamma)$ **(4)** and $C \wedge \exists \bar{\mathsf{x}}.D, \Gamma' \vdash \mathsf{t}_2 : \mathsf{T}_2$ **(5)**, where $\Gamma'$ is $\Gamma; \mathsf{z} : \forall \bar{\mathsf{x}}[D].\mathsf{T}_1$. Applying the induction hypothesis to (3) yields a constraint $C'_1$ such that $C'_1 \vdash \mathsf{t}_1 : \mathsf{T}_1$ **(6)** and $C \wedge D \Vdash \mathsf{let}\ \Gamma\ \mathsf{in}\ C'_1$ **(7)**. By (1), (2), and C-IN*, we have $C \wedge \exists \bar{\mathsf{x}}.D \Vdash \exists \Gamma'$. Thus, the induction hypothesis applies to (5) and yields a constraint $C'_2$ such that $C'_2 \vdash \mathsf{t}_2 : \mathsf{T}_2$ **(8)** and $C \wedge \exists \bar{\mathsf{x}}.D \Vdash \mathsf{let}\ \Gamma'\ \mathsf{in}\ C'_2$ **(9)**. By LET, (6) and (8) imply $\mathsf{let}\ \mathsf{z} : \forall \mathcal{V}[C'_1].\mathsf{T}_1\ \mathsf{in}\ C'_2 \vdash \mathsf{let}\ \mathsf{z} = \mathsf{t}_1\ \mathsf{in}\ \mathsf{t}_2 : \mathsf{T}_2$ **(10)**. By Lemmas 1.3.25 and 1.5.2, (10) yields $\mathsf{let}\ \mathsf{z} : \forall \bar{\mathsf{x}}[C'_1].\mathsf{T}_1\ \mathsf{in}\ C'_2 \vdash \mathsf{let}\ \mathsf{z} = \mathsf{t}_1\ \mathsf{in}\ \mathsf{t}_2 : \mathsf{T}_2$ **(11)**, where the universal quantification is over $\bar{\mathsf{x}}$ only. There remains to establish that $C \wedge \exists \bar{\mathsf{x}}.D$ entails $\mathsf{let}\ \Gamma; \mathsf{z} : \forall \bar{\mathsf{x}}[C'_1].\mathsf{T}_1\ \mathsf{in}\ C'_2$ **(12)**. By (4), (2), and C-LETDUP, the constraint (12) is equivalent to $\mathsf{let}\ \Gamma; \mathsf{z} : \forall \bar{\mathsf{x}}[\mathsf{let}\ \Gamma\ \mathsf{in}\ C'_1].\mathsf{T}_1\ \mathsf{in}\ C'_2$. By (7), this constraint is entailed by $\mathsf{let}\ \Gamma; \mathsf{z} : \forall \bar{\mathsf{x}}[C \wedge D].\mathsf{T}_1\ \mathsf{in}\ C'_2$, which by (4) and C-LETAND, is equivalent to $C \wedge \mathsf{let}\ \Gamma; \mathsf{z} : \forall \bar{\mathsf{x}}[D].\mathsf{T}_1\ \mathsf{in}\ C'_2$, that is, $C \wedge \mathsf{let}\ \Gamma'\ \mathsf{in}\ C'_2$. By (9), this constraint is entailed by $C \wedge \exists \bar{\mathsf{x}}.D$.

∘ *Case* HMD-SUB. The rule's conclusion is $C, \Gamma \vdash \mathsf{t} : \mathsf{T}'$. Its premises are $C, \Gamma \vdash \mathsf{t} : \mathsf{T}$ **(1)** and $C \Vdash \mathsf{T} \leq \mathsf{T}'$ **(2)**. Applying the induction hypothesis to (1) yields a constraint $C'$ such that $C' \vdash \mathsf{t} : \mathsf{T}$ **(3)** and $C \Vdash \mathsf{let}\ \Gamma\ \mathsf{in}\ C'$ **(4)**. By SUB, (3) implies $C' \wedge \mathsf{T} \leq \mathsf{T}' \vdash \mathsf{t} : \mathsf{T}'$. There remains to establish $C \Vdash \mathsf{let}\ \Gamma\ \mathsf{in}\ (C' \wedge \mathsf{T} \leq \mathsf{T}')$, which follows from (4) and (2) by C-INAND*.

∘ *Case* HMD-EXISTS. The rule's conclusion is $\exists \bar{\mathsf{x}}.C, \Gamma \vdash \mathsf{t} : \mathsf{T}$. Its premises are $C, \Gamma \vdash \mathsf{t} : \mathsf{T}$ **(1)** and $\bar{\mathsf{x}}\ \#\ ftv(\Gamma, \mathsf{T})$ **(2)**. By hypothesis, we have $\exists \bar{\mathsf{x}}.C \Vdash \exists \Gamma$, which by Lemma 1.3.16 implies $C \Vdash \exists \Gamma$. Thus, the induction hypothesis applies to (1) and yields a constraint $C'$ such that $C' \vdash \mathsf{t} : \mathsf{T}$ **(3)** and $C \Vdash \mathsf{let}\ \Gamma\ \mathsf{in}\ C'$ **(4)**. By EXISTS, (3) and (2) imply $\exists \bar{\mathsf{x}}.C' \vdash \mathsf{t} : \mathsf{T}$. There remains to establish $\exists \bar{\mathsf{x}}.C \Vdash \mathsf{let}\ \Gamma\ \mathsf{in}\ \exists \bar{\mathsf{x}}.C'$. By congruence of entailment, (4) implies $\exists \bar{\mathsf{x}}.C \Vdash \exists \bar{\mathsf{x}}.\mathsf{let}\ \Gamma\ \mathsf{in}\ C'$. The result follows by (2) and C-INEX. □

The second theorem states that if, within $\mathrm{PCB}(X)$, $\mathsf{t}$ has type $\mathsf{T}$ under assumption $C$, then, within $\mathrm{HM}(X)$, $\mathsf{t}$ also has type $\mathsf{T}$, under assumptions

let $\Gamma$ in $C$ and $\Gamma$. The idea is simple: the constraint $C$ represents a combined assumption about the initial judgement's free type variables and free program identifiers. In HM($X$), these two kinds of assumptions must be maintained separately. So, we split them into a pair of an environment $\Gamma$, which may be chosen arbitrarily, provided it satisfies $fpi(C) \subseteq dpi(\Gamma)$—that is, provided it defines all program variables of interest, and the residual constraint let $\Gamma$ in $C$, which has no free program identifiers, thus represents an assumption about the new judgement's type variables only. Distinct choices of $\Gamma$ give rise to distinct HM($X$) judgements, which may be incomparable; this is related to the fact that ML-the-type-system does not have principal typings (Jim, 1995). Again, the hypothesis $fpi(\Gamma) = fpi(\text{let } \Gamma \text{ in } C) = \varnothing$ is natural, since we wish $\Gamma$ and let $\Gamma$ in $C$ to appear in an HM($X$) judgement.

1.5.4    THEOREM: Assume $fpi(\Gamma) = fpi(\text{let } \Gamma \text{ in } C) = \varnothing$ and $C \not\equiv \mathsf{false}$. If $C \vdash \mathsf{t} : \mathsf{T}$ holds in PCB($X$), then let $\Gamma$ in $C, \Gamma \vdash \mathsf{t} : \mathsf{T}$ holds in HM($X$).                □

*Proof:*   The proof is by structural induction on a derivation of $C \vdash \mathsf{t} : \mathsf{T}$. In each proof case, we adopt the notations of Figure 1-9.

By Lemma 1.3.30, the hypothesis $C \not\equiv \mathsf{false}$ is preserved whenever the induction hypothesis is invoked. It is explicitly used only in case VAR, where it guarantees that the identifier at hand is bound in $\Gamma$.

○ *Case* VAR. The rule's conclusion is $C \vdash \mathsf{x} : \mathsf{T}$. Its premise is $C \Vdash \mathsf{x} \preceq \mathsf{T}$ **(1)**. By Lemma 1.3.24, (1) and the hypothesis $C \not\equiv \mathsf{false}$ imply $\mathsf{x} \in fpi(C)$. Because let $\Gamma$ in $C$ has no free program identifiers, this implies $\mathsf{x} \in dpi(\Gamma)$, that is, the environment $\Gamma$ must define $\mathsf{x}$. Let $\Gamma(\mathsf{x}) = \forall\bar{\mathsf{x}}[D].\mathsf{T}'$ **(2)**, where $\bar{\mathsf{x}} \mathrel{\#} ftv(\Gamma, \mathsf{T})$ **(3)**. By (2), HMD-VARINST, and HMD-SUB, we have $D \wedge \mathsf{T}' \leq \mathsf{T}, \Gamma \vdash \mathsf{x} : \mathsf{T}$. By (3) and HMD-EXISTS, this implies $\exists\bar{\mathsf{x}}.(D \wedge \mathsf{T}' \leq \mathsf{T}), \Gamma \vdash \mathsf{x} : \mathsf{T}$ **(4)**. Now, by (3), the constraint $\exists\bar{\mathsf{x}}.(D \wedge \mathsf{T}' \leq \mathsf{T})$ may be written $\forall\bar{\mathsf{x}}[D].\mathsf{T}' \preceq \mathsf{T}$ **(5)**. The hypothesis $fpi(\Gamma) = \varnothing$ implies $fpi(D) = \varnothing$ **(6)**. By (6), C-INID and C-IN*, (5) is equivalent to let $\Gamma$ in $\mathsf{x} \preceq \mathsf{T}$. Thus, (4) may be written let $\Gamma$ in $\mathsf{x} \preceq \mathsf{T}, \Gamma \vdash \mathsf{x} : \mathsf{T}$. By (1), by congruence of entailment, and by Lemma 1.4.2, this implies let $\Gamma$ in $C, \Gamma \vdash \mathsf{x} : \mathsf{T}$.

○ *Case* ABS. The rule's conclusion is let $\mathsf{z} : \mathsf{T}$ in $C \vdash \lambda\mathsf{z}.\mathsf{t} : \mathsf{T} \rightarrow \mathsf{T}'$. Its premise is $C \vdash \mathsf{t} : \mathsf{T}'$ **(1)**. Let $\Gamma'$ stand for $\Gamma; \mathsf{z} : \mathsf{T}$. Applying the induction hypothesis to (1) yields let $\Gamma'$ in $C, \Gamma' \vdash \mathsf{t} : \mathsf{T}'$. By HMD-ABS, this implies let $\Gamma'$ in $C, \Gamma \vdash \lambda\mathsf{z}.\mathsf{t} : \mathsf{T} \rightarrow \mathsf{T}'$.

○ *Case* APP. The rule's conclusion is $C_1 \wedge C_2 \vdash \mathsf{t}_1\, \mathsf{t}_2 : \mathsf{T}'$. Its premises are $C_1 \vdash \mathsf{t}_1 : \mathsf{T} \rightarrow \mathsf{T}'$ and $C_2 \vdash \mathsf{t}_2 : \mathsf{T}$. Applying the induction hypothesis yields respectively let $\Gamma$ in $C_1, \Gamma \vdash \mathsf{t}_1 : \mathsf{T} \rightarrow \mathsf{T}'$ and let $\Gamma$ in $C_2, \Gamma \vdash \mathsf{t}_2 : \mathsf{T}$, which by Lemma 1.4.2 and HMD-APP imply let $\Gamma$ in $(C_1 \wedge C_2), \Gamma \vdash \mathsf{t}_1\, \mathsf{t}_2 : \mathsf{T}'$.

○ *Case* LET. The rule's conclusion is let $\mathsf{z} : \forall\mathcal{V}[C_1].\mathsf{T}_1$ in $C_2 \vdash \mathtt{let\ z =}$

$\mathtt{t_1}$ in $\mathtt{t_2}$ : $\mathtt{T_2}$. Its premises are $C_1 \vdash \mathtt{t_1} : \mathtt{T_1}$ **(1)** and $C_2 \vdash \mathtt{t_2} : \mathtt{T_2}$ **(2)**. Let $\bar{\mathtt{x}}$ stand for $ftv(C_1, \mathtt{T_1})$. We may require, *w.l.o.g.*, $\bar{\mathtt{x}} \mathrel{\#} ftv(\Gamma, C_2)$ **(3)**. By hypothesis, we have $fpi(\Gamma) = \varnothing$ **(4)**. We also have $fpi(\text{let } \Gamma; \mathtt{z} : \forall \mathcal{V}[C_1].\mathtt{T_1} \text{ in } C_2) = \varnothing$, which implies $fpi(\text{let } \Gamma \text{ in } C_1) = \varnothing$. Thus, the induction hypothesis applies to (1) and yields $\text{let } \Gamma \text{ in } C_1, \Gamma \vdash \mathtt{t_1} : \mathtt{T_1}$ **(5)**. Now, let $\sigma$ stand for $\forall \bar{\mathtt{x}}[\text{let } \Gamma \text{ in } C_1].\mathtt{T_1}$ and $\Gamma'$ stand for $\Gamma; \mathtt{z} : \sigma$. We have $fpi(\Gamma') = fpi(\text{let } \Gamma' \text{ in } C_2) = \varnothing$. Thus, the induction hypothesis applies to (2) and yields $\text{let } \Gamma' \text{ in } C_2, \Gamma' \vdash \mathtt{t_2} : \mathtt{T_2}$ **(6)**. Let us now weaken (5) and (6) so as to make them suitable premises for HMD-LETGEN. Applying Lemma 1.4.2 to (5) yields $(\text{let } \Gamma' \text{ in } C_2) \wedge (\text{let } \Gamma \text{ in } C_1), \Gamma \vdash \mathtt{t_1} : \mathtt{T_1}$ **(7)**. Applying Lemma 1.4.2 to (6) yields $(\text{let } \Gamma' \text{ in } C_2) \wedge \exists \bar{\mathtt{x}}.(\text{let } \Gamma \text{ in } C_1), \Gamma' \vdash \mathtt{t_2} : \mathtt{T_2}$ **(8)**. Last, (3) implies $\bar{\mathtt{x}} \mathrel{\#} ftv(\Gamma, \text{let } \Gamma' \text{ in } C_2)$ **(9)**. Applying HMD-LETGEN to (7), (9) and (8), we obtain $(\text{let } \Gamma' \text{ in } C_2) \wedge \exists \bar{\mathtt{x}}.(\text{let } \Gamma \text{ in } C_1), \Gamma \vdash \text{let } \mathtt{z} = \mathtt{t_1} \text{ in } \mathtt{t_2} : \mathtt{T_2}$ **(10)**. Now, by (4), (3), and C-LETDUP, $\text{let } \Gamma' \text{ in } C_2$ is equivalent to $\text{let } \Gamma; \mathtt{z} : \forall \bar{\mathtt{x}}[C_1].\mathtt{T_1} \text{ in } C_2$. Using this fact, as well as (3), C-INEX, and C-INAND, we find that the constraint $(\text{let } \Gamma' \text{ in } C_2) \wedge \exists \bar{\mathtt{x}}.(\text{let } \Gamma \text{ in } C_1)$ is equivalent to $\text{let } \Gamma \text{ in } (\text{let } \mathtt{z} : \forall \bar{\mathtt{x}}[C_1].\mathtt{T_1} \text{ in } C_2 \wedge \exists \bar{\mathtt{x}}.C_1)$, which by definition of the let form, is itself equivalent to $\text{let } \Gamma; \mathtt{z} : \forall \bar{\mathtt{x}}[C_1].\mathtt{T_1} \text{ in } C_2$. Last, by definition of $\bar{\mathtt{x}}$, this constraint is $\text{let } \Gamma; \mathtt{z} : \forall \mathcal{V}[C_1].\mathtt{T_1} \text{ in } C_2$. Thus, (10) is the goal.

   ◦ *Case* SUB. The rule's conclusion is $C \wedge \mathtt{T} \le \mathtt{T'} \vdash \mathtt{t} : \mathtt{T'}$. Its premise is $C \vdash \mathtt{t} : \mathtt{T}$ **(1)**. Applying the induction hypothesis to (1) yields $\text{let } \Gamma \text{ in } C, \Gamma \vdash \mathtt{t} : \mathtt{T}$ **(2)**. By Lemma 1.4.2 and HMD-SUB, (2) implies $(\text{let } \Gamma \text{ in } C) \wedge \mathtt{T} \le \mathtt{T'}, \Gamma \vdash \mathtt{t} : \mathtt{T'}$, which by C-INAND* may be written $\text{let } \Gamma \text{ in } (C \wedge \mathtt{T} \le \mathtt{T'}), \Gamma \vdash \mathtt{t} : \mathtt{T'}$.

   ◦ *Case* EXISTS. The rule's conclusion is $\exists \bar{\mathtt{x}}.C \vdash \mathtt{t} : \mathtt{T}$. Its premises are $C \vdash \mathtt{t} : \mathtt{T}$ **(1)** and $\bar{\mathtt{x}} \mathrel{\#} ftv(\mathtt{T})$ **(2)**. We may further require, *w.l.o.g.*, $\bar{\mathtt{x}} \mathrel{\#} ftv(\Gamma)$ **(3)**. Applying the induction hypothesis to (1) yields $\text{let } \Gamma \text{ in } C, \Gamma \vdash \mathtt{t} : \mathtt{T}$ **(4)**. Applying HMD-EXISTS to (2), (3), and (4), we find $\exists \bar{\mathtt{x}}.\text{let } \Gamma \text{ in } C, \Gamma \vdash \mathtt{t} : \mathtt{T}$, which, by (3) and C-INEX, may be written $\text{let } \Gamma \text{ in } \exists \bar{\mathtt{x}}.C, \Gamma \vdash \mathtt{t} : \mathtt{T}$.      □

As a corollary, we find that, for closed programs, the type systems $\mathrm{HM}(X)$ and $\mathrm{PCB}(X)$ coincide. In particular, a program is well-typed with respect to one if and only if it is well-typed with respect to the other. This supports the view that $\mathrm{PCB}(X)$ is an alternate formulation of $\mathrm{HM}(X)$.

1.5.5    THEOREM: Assume $fpi(C) = \varnothing$ and $C \not\equiv \mathsf{false}$. Then, $C, \varnothing \vdash \mathtt{t} : \mathtt{T}$ holds in $\mathrm{HM}(X)$ if and only if $C \vdash \mathtt{t} : \mathtt{T}$ holds in $\mathrm{PCB}(X)$.      □

## 1.6   Constraint generation

We now explain how to reduce type inference problems for $\mathrm{PCB}(X)$ to constraint solving problems. A type inference problem consists of an expression

$$
\begin{aligned}
[\![x : T]\!] &= x \preceq T \\
[\![\lambda z.t : T]\!] &= \exists X_1 X_2.(\text{let } z : X_1 \text{ in } [\![t : X_2]\!] \wedge X_1 \rightarrow X_2 \leq T) \\
[\![t_1\ t_2 : T]\!] &= \exists X_2.([\![t_1 : X_2 \rightarrow T]\!] \wedge [\![t_2 : X_2]\!]) \\
[\![\text{let } z = t_1 \text{ in } t_2 : T]\!] &= \text{let } z : \forall X[[\![t_1 : X]\!]].X \text{ in } [\![t_2 : T]\!]
\end{aligned}
$$

**Figure 1-10: Constraint generation**

t and a type T of kind $\star$. The problem is to determine whether t is well-typed with type T, that is, whether there exists a satisfiable constraint $C$ such that $C \vdash t : T$ holds. This formulation of the problem may seem to require an appropriate type T to be known in advance; this is not really the case, since T may be a type variable. A constraint solving problem consists of a constraint $C$. The problem is to determine whether $C$ is satisfiable. To reduce a type inference problem $(t, T)$ to a constraint solving problem, we must produce a constraint $C$ that is both *sufficient* and *necessary* for $C \vdash t : T$ to hold. Below, we explain how to compute such a constraint, which we write $[\![t : T]\!]$. We check that it is indeed *sufficient* by proving $[\![t : T]\!] \vdash t : T$. That is, the constraint $[\![t : T]\!]$ is specific enough to guarantee that t has type T. We say that constraint generation is *sound*. We check that it is indeed *necessary* by proving that, for every constraint $C$, $C \vdash t : T$ implies $C \Vdash [\![t : T]\!]$. That is, every constraint that guarantees that t has type T is at least as specific as $[\![t : T]\!]$. We say that constraint generation is *complete*. Together, these properties mean that $[\![t : T]\!]$ is the *least specific* constraint that guarantees that t has type T.

We now see how to reduce a type inference problem to a constraint solving problem. Indeed, if there exists a satisfiable constraint $C$ such that $C \vdash t : T$ holds, then, by the completeness property, $C \Vdash [\![t : T]\!]$ holds, so $[\![t : T]\!]$ is satisfiable. Conversely, the soundness property states that $[\![t : T]\!] \vdash t : T$ holds, so, if $[\![t : T]\!]$ is satisfiable, then there exists a satisfiable constraint $C$ such that $C \vdash t : T$ holds. In other words, t is well-typed with type T if and only if $[\![t : T]\!]$ is satisfiable.

The existence of such a constraint is the analogue of the existence of *principal type schemes* in classic presentations of ML-the-type-system (Damas and Milner, 1982). Indeed, a principal type scheme is least specific in the sense that all valid types are substitution instances of it. Here, the constraint $[\![t : T]\!]$ is least specific in the sense that all valid constraints entail it. Earlier, we have established a connection between constraint entailment and refinement of type substitutions, in the specific case of equality constraints interpreted over a free algebra of finite types; see Lemma 1.3.39.

The constraint $[\![t : T]\!]$ is defined in Figure 1-10 by induction on the structure of the expression $t$. We refer to these defining equations as the *constraint generation rules*. The definition is quite terse. It is perhaps even simpler than the declarative specification of $\mathrm{PCB}(X)$ given in Figure 1-9; yet, we prove below that the two are equivalent.

Before explaining the definition, we state the requirements that bear on the type variables $X_1$, $X_2$, and $X$, which appear bound in the right-hand sides of the second, third, and fourth equations. These type variables must have kind $\star$. They must be chosen distinct (that is, $X_1 \neq X_2$ in the second equation) and *fresh* in the following sense: *type variables that appear bound in an equation's right-hand side must not appear free in the equation's left-hand side*. Provided this restriction is obeyed, different choices of $X_1$, $X_2$, and $X$ lead to $\alpha$-equivalent constraints—that is, to the same constraint, since we identify objects up to $\alpha$-conversion—which guarantees that the above equations make sense. We remark that, since expressions do not have free type variables, the freshness requirement may be simplified to: type variables that appear bound in an equation's right-hand side must not appear free in $T$. However, this simplification is rendered invalid by the introduction of type annotations within expressions (page 102). Please note that we are able to state a *formal* freshness requirement. This is made possible by the fact that $[\![t : T]\!]$ has no free type variables other than those of $T$, which in turn depends on our explicit use of existential quantification.

Let us now review the four equations. The first one simply mirrors VAR. The second one requires $t$ to have type $X_2$ under the hypothesis that $z$ has type $X_1$, and forms the arrow type $X_1 \to X_2$; this corresponds to ABS. Here, $X_1$ and $X_2$ must be fresh type variables, because we cannot in general guess the expected types of $z$ and $t$. The expected type $T$ is required to be a supertype of $X_1 \to X_2$; this corresponds to SUB. We must bind the fresh type variables $X_1$ and $X_2$, so as to guarantee that the generated constraint is unique up to $\alpha$-conversion. Furthermore, we must bind them *existentially*, because we intend the constraint solver to choose some appropriate value for them. This is justified by EXISTS. The third equation uses the fresh type variable $X_2$ to stand for the unknown type of $t_2$. The subexpression $t_1$ is expected to have type $X_2 \to T$. This corresponds to APP. The fourth equation, which corresponds to LET, is most interesting. It summons a fresh type variable $X$ and produces $[\![t_1 : X]\!]$. This constraint, whose sole free type variable is $X$, is the *least specific* constraint that must be imposed on $X$ so as to make it a valid type for $t_1$. As a result, the type scheme $\forall X[\![\![t_1 : X]\!]\!].X$, abbreviated $\sigma$ in the following, is a *principal* type scheme for $t_1$. There remains to place $[\![t_2 : T]\!]$ inside the context let $z : \sigma$ in $[\,]$. Indeed, when placed inside this context, an instantiation constraint of the form $z \preceq T'$ acquires the meaning

$\sigma \preceq T'$, which by definition of $\sigma$ and by Lemma 1.6.4 (see below) is equivalent to $[\![t_1 : T']\!]$. Thus, the constraint produced by the fourth equation simulates a textual expansion of the `let` construct, whereby every occurrence of `z` would be replaced with $t_1$. Thanks to type scheme introduction and instantiation constraints, however, this effect is achieved without duplication of source code or constraints. In other words, constraint generation has linear time and space complexity; duplication may take place during constraint solving only.

1.6.1   EXERCISE [★, ↛]: Define the *size* of an expression, of a type, and of a constraint, viewed as abstract syntax trees. Check that the size of $[\![t : T]\!]$ is linear in the sum of the sizes of `t` and `T`.   □

We now establish several properties of constraint generation. We begin with soundness, whose proof is straightforward.

1.6.2   THEOREM [SOUNDNESS]: $[\![t : T]\!] \vdash t : T$.   □

*Proof:*   By induction on the structure of `t`.

∘ *Case* x. The goal $x \preceq T \vdash x : T$ follows from VAR.

∘ *Case* λz.t. By the induction hypothesis, we have $[\![t : X_2]\!] \vdash t : X_2$. By ABS, this implies $\text{let } z : X_1 \text{ in } [\![t : X_2]\!] \vdash \lambda z.t : X_1 \to X_2$. By SUB, this implies $\text{let } z : X_1 \text{ in } [\![t : X_2]\!] \wedge X_1 \to X_2 \leq T \vdash \lambda z.t : T$. Lastly, because $X_1 X_2 \mathbin{\#} ftv(T)$ holds, EXISTS applies and yields $[\![\lambda z.t : T]\!] \vdash \lambda z.t : T$.

∘ *Case* $t_1\ t_2$. By the induction hypothesis, we have $[\![t_1 : X_2 \to T]\!] \vdash t_1 : X_2 \to T$ and $[\![t_2 : X_2]\!] \vdash t_2 : X_2$. By APP, this implies $[\![t_1 : X_2 \to T]\!] \wedge [\![t_2 : X_2]\!] \vdash t_1\ t_2 : T$. Because $X_2 \notin ftv(T)$ holds, EXISTS applies and yields $[\![t_1\ t_2 : T]\!] \vdash t_1\ t_2 : T$.

∘ *Case* $\text{let } z = t_1 \text{ in } t_2$. By the induction hypothesis, we have $[\![t_1 : X]\!] \vdash t_1 : X$ and $[\![t_2 : T]\!] \vdash t_2 : T$. By LET, these imply $\text{let } z : \forall \mathcal{V}[\![t_1 : X]\!].X \text{ in } [\![t_2 : T]\!] \vdash \text{let } z = t_1 \text{ in } t_2 : T$. Because $ftv([\![t_1 : X]\!])$ is X, the universal quantification on $\mathcal{V}$ really bears on X alone. We have proved $[\![\text{let } z = t_1 \text{ in } t_2 : T]\!] \vdash \text{let } z = t_1 \text{ in } t_2 : T$.   □

The following lemmas are used in the proof of the completeness property and in a number of other occasions. The first two state that $[\![t : T]\!]$ is *covariant* with respect to `T`. Roughly speaking, this means that enough subtyping constraints are generated to achieve completeness with respect to SUB.

1.6.3   LEMMA: $[\![t : T]\!] \wedge T \leq T'$ entails $[\![t : T']\!]$.   □

1.6.4   LEMMA: $X \notin ftv(T)$ implies $\exists X.([\![t : X]\!] \wedge X \leq T) \equiv [\![t : T]\!]$.   □

The next lemma gives a simplified version of the second constraint genera-tion rule, in the specific case where the expected type is an arrow type. Then, fresh type variables need not be generated; one may directly use the arrow's domain and codomain instead.

1.6.5    Lemma:  $[\![\lambda z.t : T_1 \to T_2]\!]$ is equivalent to $\mathsf{let}\ z : T_1\ \mathsf{in}\ [\![t : T_2]\!]$.                    □

We conclude with the completeness property.

1.6.6    Theorem [Completeness]:  if $C \vdash t : T$, then $C \Vdash [\![t : T]\!]$.                    □

*Proof:*   By induction on the derivation of $C \vdash t : T$.

∘ *Case* Var. The rule's conclusion is $C \vdash x : T$. Its premise is $C \Vdash x \preceq T$, which is also the goal.

∘ *Case* Abs. The rule's conclusion is $\mathsf{let}\ z : T\ \mathsf{in}\ C \vdash \lambda z.t : T \to T'$. Its premise is $C \vdash t : T'$. By the induction hypothesis, we have $C \Vdash [\![t : T']\!]$. By congruence of entailment, this implies $\mathsf{let}\ z : T\ \mathsf{in}\ C \Vdash \mathsf{let}\ z : T\ \mathsf{in}\ [\![t : T']\!]$, which, by Lemma 1.6.5, may be written $\mathsf{let}\ z : T\ \mathsf{in}\ C \Vdash [\![\lambda z.t : T \to T']\!]$.

∘ *Case* App. The rule's conclusion is $C_1 \wedge C_2 \vdash t_1\ t_2 : T'$. Its premises are $C_1 \vdash t_1 : T \to T'$ and $C_2 \vdash t_2 : T$. By the induction hypothesis, we have $C_1 \Vdash [\![t_1 : T \to T']\!]$ and $C_2 \Vdash [\![t_2 : T]\!]$. Thus, $C_1 \wedge C_2$ entails $[\![t_1 : T \to T']\!] \wedge [\![t_2 : T]\!]$, which, by C-NameEq, may be written $\exists X_2.(X_2 = T \wedge [\![t_1 : X_2 \to T']\!] \wedge [\![t_2 : X_2]\!])$, where $X_2 \notin ftv(T, T')$. Forgetting about the equation $X_2 = T$, we find that $C_1 \wedge C_2$ entails $\exists X_2.([\![t_1 : X_2 \to T']\!] \wedge [\![t_2 : X_2]\!])$, which is precisely $[\![t_1\ t_2 : T']\!]$.

∘ *Case* Let. The rule's conclusion is $\mathsf{let}\ z : \forall \mathcal{V}[C_1].T_1\ \mathsf{in}\ C_2 \vdash \mathtt{let}\ z = t_1\ \mathtt{in}\ t_2 : T_2$. Its premises are $C_1 \vdash t_1 : T_1$ and $C_2 \vdash t_2 : T_2$. By the induction hypothesis, we have $C_1 \Vdash [\![t_1 : T_1]\!]$ and $C_2 \Vdash [\![t_2 : T_2]\!]$, which implies $\mathsf{let}\ z : \forall \mathcal{V}[C_1].T_1\ \mathsf{in}\ C_2 \Vdash \mathsf{let}\ z : \forall \mathcal{V}[[\![t_1 : T_1]\!]].T_1\ \mathsf{in}\ [\![t_2 : T_2]\!]$ **(1)**.

Now, let us establish $\mathsf{true} \Vdash \forall X[[\![t_1 : X]\!]].X \preceq \forall \mathcal{V}[[\![t_1 : T_1]\!]].T_1$ **(2)**. By definition, this requires proving $\exists \bar{X}_1.([\![t_1 : T_1]\!] \wedge T_1 \leq Z) \Vdash \exists X.([\![t_1 : X]\!] \wedge X \leq Z)$ **(3)**, where $\bar{X}_1 = ftv(T_1)$ and $Z \notin X\bar{X}_1$ **(4)**. By Lemma 1.6.3, (4), and C-Ex*, the left-hand side of (3) entails $[\![t_1 : Z]\!]$. By (4) and Lemma 1.6.4, the right-hand side of (3) is $[\![t_1 : Z]\!]$. Thus, (3) holds, and so does (2).

By (2) and Lemma 1.3.22, we have $\mathsf{let}\ z : \forall \mathcal{V}[[\![t_1 : T_1]\!]].T_1\ \mathsf{in}\ [\![t_2 : T_2]\!] \Vdash \mathsf{let}\ z : \forall X[[\![t_1 : X]\!]].X\ \mathsf{in}\ [\![t_2 : T_2]\!]$ **(5)**. By transitivity of entailment, (1) and (5) yield $\mathsf{let}\ z : \forall \mathcal{V}[C_1].T_1\ \mathsf{in}\ C_2 \Vdash [\![\mathtt{let}\ z = t_1\ \mathtt{in}\ t_2 : T_2]\!]$.

∘ *Case* Sub. The rule's conclusion is $C \wedge T \leq T' \vdash t : T'$. Its premise is $C \vdash t : T$. By the induction hypothesis, we have $C \Vdash [\![t : T]\!]$, which implies $C \wedge T \leq T' \Vdash [\![t : T]\!] \wedge T \leq T'$. By lemma 1.6.3 and by transitivity of entailment, we obtain $C \wedge T \leq T' \Vdash [\![t : T']\!]$.

∘ *Case* Exists. The rule's conclusion is $\exists \bar{X}.C \vdash t : T$. Its premises are $C \vdash t : T$ and $\bar{X}\ \#\ ftv(T)$ **(1)**. By the induction hypothesis, we have $C \Vdash [\![t : T]\!]$.

By congruence of entailment, this implies $\exists \bar{x}.C \Vdash \exists \bar{x}.[\![\texttt{t} : \texttt{T}]\!]$ **(2)**. Furthermore, (1) implies $\bar{x} \mathbin{\#} \mathit{ftv}([\![\texttt{t} : \texttt{T}]\!])$ **(3)**. By (3) and C-Ex*, (2) may be written $\exists \bar{x}.C \Vdash [\![\texttt{t} : \texttt{T}]\!]$. $\qquad\qquad\square$

## 1.7  Type soundness

We are now ready to establish type soundness for our type system. The statement that we wish to prove is sometimes known as *Milner's slogan*: *well-typed programs do not go wrong* (Milner, 1978). Below, we define well-typedness in terms of our constraint generation rules, for the sake of convenience, and establish type soundness with respect to that particular definition. Theorems 1.4.7, 1.5.4, and 1.6.6 imply that type soundness also holds when well-typedness is defined with respect to the typing judgements of DM, HM($X$), or PCB($X$). We establish type soundness by following Wright and Felleisen's so-called *syntactic approach* (1994b). The approach consists in isolating two independent properties. *Subject reduction*, whose exact statement will be given below, implies that well-typedness is preserved by reduction. *Progress* states that no stuck configuration is well-typed. It is immediate to check that, if both properties hold, then no well-typed program can reduce to a stuck configuration. Subject reduction itself depends on a key lemma, usually known as a (term) *substitution lemma*. We immediately give two versions of this lemma: the former is stated in terms of PCB($X$) judgements, while the latter is stated in terms of the constraint generation rules.

1.7.1  LEMMA [SUBSTITUTION]: $C \vdash \texttt{t} : \texttt{T}$ and $C_0 \vdash \texttt{t}_0 : \texttt{T}_0$ imply $\mathsf{let}\ z_0 : \forall \bar{x}_0[C_0].\texttt{T}_0\ \mathsf{in}\ C \vdash [z_0 \mapsto \texttt{t}_0]\texttt{t} : \texttt{T}$. $\qquad\qquad\square$

*Proof:*  The proof is by structural induction on the derivation of $C \vdash \texttt{t} : \texttt{T}$. In each proof case, we adopt the notations of Figure 1-9. We write $\sigma_0$ for $\forall \bar{x}_0[C_0].\texttt{T}_0$. We refer to the hypothesis $C_0 \vdash \texttt{t}_0 : \texttt{T}_0$ as **(1)**. We assume, *w.l.o.g.*, $\bar{x}_0 \mathbin{\#} \mathit{ftv}(C, \texttt{T})$ **(2)** and $z_0 \notin \mathit{fpi}(\sigma_0)$ **(3)**.

○ *Case* VAR. The rule's conclusion is $C \vdash x : \texttt{T}$ **(4)**. Its premise is $C \Vdash x \preceq \texttt{T}$ **(5)**. Two subcases arise.

*Subcase* $x$ is $z_0$. Applying SUB to (1) yields $C_0 \wedge \texttt{T}_0 \leq \texttt{T} \vdash \texttt{t}_0 : \texttt{T}$. By (2) and EXISTS, this implies $\exists \bar{x}_0.(C_0 \wedge \texttt{T}_0 \leq \texttt{T}) \vdash \texttt{t}_0 : \texttt{T}$ **(6)**. Furthermore, by (2) again, the constraint $\exists \bar{x}_0.(C_0 \wedge \texttt{T}_0 \leq \texttt{T})$ is $\sigma_0 \preceq \texttt{T}$, which is equivalent to $\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ z_0 \preceq \texttt{T}$. As a result, (6) may be written $\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ x \preceq \texttt{T} \vdash [z_0 \mapsto \texttt{t}_0]x : \texttt{T}$ **(7)**.

*Subcase* $x$ isn't $z_0$. Then, $[z_0 \mapsto \texttt{t}_0]x$ is $x$. Thus, VAR yields $\exists \sigma_0 \wedge x \preceq \texttt{T} \vdash [z_0 \mapsto \texttt{t}_0]x : \texttt{T}$. By C-IN*, this may be read $\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ x \preceq \texttt{T} \vdash [z_0 \mapsto \texttt{t}_0]x : \texttt{T}$, that is, again (7).

In either subcase, by (5), by congruence of entailment, and by Lemma 1.5.2, (7) implies $\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ C \vdash [z_0 \mapsto t_0]t : T$.

◦ *Case* ABS. The rule's conclusion is $\mathsf{let}\ z : T\ \mathsf{in}\ C \vdash \lambda z.t : T \to T'$. Its premise is $C \vdash t : T'$ **(8)**. We may assume, *w.l.o.g.*, that $z$ is distinct from $z_0$ and does not occur free within $t_0$ or $\sigma_0$ **(9)**. Applying the induction hypothesis to (8) yields $\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ C \vdash [z_0 \mapsto t_0]t : T'$, which, by ABS, implies $\mathsf{let}\ z : T\ \mathsf{in}\ (\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ C) \vdash \lambda z.[z_0 \mapsto t_0]t : T \to T'$. By (9) and C-LETLET, this may be written $\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ (\mathsf{let}\ z : T\ \mathsf{in}\ C) \vdash [z_0 \mapsto t_0](\lambda z.t) : T \to T'$.

◦ *Case* APP. By the induction hypothesis, by APP, and by C-INAND.

◦ *Case* LET. The rule's conclusion is $\mathsf{let}\ z : \forall \bar{X}_1[C_1].T_1\ \mathsf{in}\ C_2 \vdash \mathtt{let\ z} = \mathtt{t_1\ in\ t_2} : T_2$, where $\bar{X}_1 = ftv(C_1, T_1)$. Its premises are $C_1 \vdash t_1 : T_1$ **(10)** and $C_2 \vdash t_2 : T_2$ **(11)**. We may assume, *w.l.o.g.*, that $z$ is distinct from $z_0$ and does not occur free within $t_0$ or $\sigma_0$ **(12)**. We may also assume, *w.l.o.g.*, $\bar{X}_1 \mathrel{\#} ftv(\sigma_0)$ **(13)**. Applying the induction hypothesis to (10) and (11) respectively yields $\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ C_1 \vdash [z_0 \mapsto t_0]t_1 : T_1$ **(14)** and $\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ C_2 \vdash [z_0 \mapsto t_0]t_2 : T_2$ **(15)**. Applying LET to (14) and (15) produces $\mathsf{let}\ z : \forall \mathcal{V}[\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ C_1].T_1\ \mathsf{in}\ \mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ C_2 \vdash [z_0 \mapsto t_0](\mathtt{let\ z} = \mathtt{t_1\ in\ t_2}) : T_2$ **(16)**. Now, we have

$$
\begin{aligned}
& \mathsf{let}\ z_0 : \sigma_0; z : \forall \bar{X}_1[C_1].T_1\ \mathsf{in}\ C_2 \\
\equiv\ & \mathsf{let}\ z_0 : \sigma_0; z : \forall \bar{X}_1[\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ C_1].T_1\ \mathsf{in}\ C_2 \quad \textbf{(17)} \\
\equiv\ & \mathsf{let}\ z : \forall \bar{X}_1[\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ C_1].T_1; z_0 : \sigma_0\ \mathsf{in}\ C_2 \quad \textbf{(18)} \\
\Vdash\ & \mathsf{let}\ z : \forall \mathcal{V}\ [\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ C_1].T_1; z_0 : \sigma_0\ \mathsf{in}\ C_2 \quad \textbf{(19)}
\end{aligned}
$$

where (17) follows from (13), (3), and C-LETDUP; (18) follows from (12) and C-LETLET; and (19) is by Lemma 1.3.25. Thus, applying Lemma 1.5.2 to (16) yields $\mathsf{let}\ z_0 : \sigma_0; z : \forall \bar{X}_1[C_1].T_1\ \mathsf{in}\ C_2 \vdash [z_0 \mapsto t_0](\mathtt{let\ z} = \mathtt{t_1\ in\ t_2}) : T_2$.

◦ *Case* SUB. By the induction hypothesis, by SUB, and by C-INAND*.

◦ *Case* EXISTS. The rule's conclusion is $\exists \bar{X}.C \vdash t : T$. Its premises are $C \vdash t : T$ **(20)** and $\bar{X} \mathrel{\#} ftv(T)$ **(21)**. We may assume, *w.l.o.g.*, $\bar{X} \mathrel{\#} ftv(\sigma_0)$ **(22)**. Applying the induction hypothesis to (20) yields $\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ C \vdash [z_0 \mapsto t_0]t : T$, which, by (21) and EXISTS, implies $\exists \bar{X}.\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ C \vdash [z_0 \mapsto t_0]t : T$ **(23)**. By (22) and C-INEX, (23) is $\mathsf{let}\ z_0 : \sigma_0\ \mathsf{in}\ \exists \bar{X}.C \vdash [z_0 \mapsto t_0]t : T$. □

1.7.2    LEMMA: $\mathsf{let}\ z : \forall \bar{X}[\llbracket t_2 : T_2 \rrbracket].T_2\ \mathsf{in}\ \llbracket t_1 : T_1 \rrbracket$ entails $\llbracket [z \mapsto t_2]t_1 : T_1 \rrbracket$. □

Before going on, let us give a few definitions and formulate several requirements. First, we must define an *initial environment* $\Gamma_0$, which assigns a type scheme to every constant. A couple of requirements must be made to ensure that $\Gamma_0$ is consistent with the semantics of constants, as specified by $\xrightarrow{\delta}$. Second, we must extend constraint generation and well-typedness to *configurations*, as opposed to programs, since reduction operates on configurations.

Last, we must formulate a *restriction* to tame the interaction between side effects and `let`-polymorphism, which is unsound if unrestricted.

1.7.3   DEFINITION: Let $\Gamma_0$ be an environment whose domain is the set of constants $\mathcal{Q}$. We require $ftv(\Gamma_0) = \varnothing$, $fpi(\Gamma_0) = \varnothing$, and $\exists \Gamma_0 \equiv \mathsf{true}$. We refer to $\Gamma_0$ as the *initial* typing environment. □

1.7.4   DEFINITION: Let $\mathsf{ref}$ be an isolated, invariant type constructor of signature $\star \Rightarrow \star$. A *store type* $M$ is a finite mapping from memory locations to types. We write $\mathsf{ref}\, M$ for the environment that maps $m$ to $\mathsf{ref}\, M\,(m)$ when $m$ is in the domain of $M$. Assuming $dom\,(\mu)$ and $dom\,(M)$ coincide, the constraint $[\![\mu : M]\!]$ is defined as the conjunction of the constraints $[\![\mu(m) : M(m)]\!]$, where $m$ ranges over $dom\,(\mu)$. Under the same assumption, the constraint $[\![\mathsf{t}/\mu : \mathsf{T}/M]\!]$ is defined as $[\![\mathsf{t} : \mathsf{T}]\!] \wedge [\![\mu : M]\!]$. A configuration $\mathsf{t}/\mu$ is *well-typed* if and only if there exist a type $\mathsf{T}$ and a store type $M$ such that $dom\,(\mu) = dom\,(M)$ and the constraint $\mathsf{let}\ \Gamma_0; \mathsf{ref}\, M\ \mathsf{in}\ [\![\mathsf{t}/\mu : \mathsf{T}/M]\!]$ is satisfiable. □

The type $\mathsf{ref}\,\mathsf{T}$ is the type of references (that is, memory locations) that store data of type $\mathsf{T}$. It must be *invariant* in its parameter, reflecting the fact that references may be *read* and *written*.

A store is a complex object: it may contain values that indirectly refer to each other via memory locations. In fact, it is a representation of the graph formed by objects and pointers in memory, which may contain cycles. We rely on store types to deal with such cycles. In the definition of well-typedness, the store type $M$ imposes a constraint on the contents of the store—the value $\mu(m)$ must have type $M(m)$—but also plays the role of a hypothesis: by placing the constraint $[\![\mathsf{t}/\mu : \mathsf{T}/M]\!]$ within the context $\mathsf{let}\ \mathsf{ref}\, M\ \mathsf{in}\ [\,]$, we give meaning to free occurrences of memory locations within $[\![\mathsf{t}/\mu : \mathsf{T}/M]\!]$, and stipulate that it is valid to assume that $m$ has type $M(m)$. In other words, we essentially view the store as a large, mutually recursive binding of locations to values. Since no satisfiable constraint may have a free program identifier (Lemma 1.3.31), every well-typed configuration must be closed. The context $\mathsf{let}\ \Gamma_0\ \mathsf{in}\ [\,]$ gives meaning to occurrences of constants within $[\![\mathsf{t}/\mu : \mathsf{T}/M]\!]$.

We now define a relation between configurations that plays a key role in the statement of the subject reduction property. The point of subject reduction is to guarantee that well-typedness is preserved by reduction. However, such a simple statement is too weak to be amenable to inductive proof. Thus, for the purposes of the proof, we must be more specific. To begin, let us consider the simpler case of a pure semantics, that is, a semantics without stores. Then, we must state that if an expression $\mathsf{t}$ has type $\mathsf{T}$ under a certain constraint, then its reduct $\mathsf{t}'$ has type $\mathsf{T}$ under the same constraint. In terms of generated constraints, this statement becomes: $\mathsf{let}\ \Gamma_0\ \mathsf{in}\ [\![\mathsf{t} : \mathsf{T}]\!]$ entails $\mathsf{let}\ \Gamma_0\ \mathsf{in}\ [\![\mathsf{t}' : \mathsf{T}]\!]$.

Let us now return to the general case, where a store is present. Then, the statement of well-typedness for a configuration $t/\mu$ involves a store type $M$ whose domain is that of $\mu$. So, the statement of well-typedness for its reduct $t'/\mu'$ must involve a store type $M'$ whose domain is that of $\mu'$—which is larger if allocation occurred. The types of existing memory locations must not change: we must request that $M$ and $M'$ agree on $dom(M)$, that is, $M'$ must extend $M$. Furthermore, the types assigned to new memory locations in $dom(M') \setminus dom(M)$ might involve new type variables, that is, variables that do not appear free in $M$ or $T$. We must allow these variables to be hidden—that is, existentially quantified—otherwise the entailment assertion cannot hold. These considerations lead us to the following definition:

1.7.5    DEFINITION: $t/\mu \sqsubseteq t'/\mu'$ holds if and only if, for every type $T$ and for every store type $M$ such that $dom(\mu) = dom(M)$, there exist a set of type variables $\bar{Y}$ and a store type $M'$ such that $\bar{Y} \mathbin{\#} ftv(T, M)$ and $ftv(M') \subseteq \bar{Y} \cup ftv(M)$ and $dom(M') = dom(\mu')$ and $M'$ extends $M$ and

$$
\begin{aligned}
&\mathsf{let}\ \Gamma_0\,;\mathsf{ref}\,M\ \mathsf{in}\ [\![ t/\mu : T/M ]\!] \\
\Vdash\ &\exists \bar{Y}.\mathsf{let}\ \Gamma_0\,;\mathsf{ref}\,M'\ \mathsf{in}\ [\![ t'/\mu' : T/M' ]\!]
\end{aligned}
$$

The relation $\sqsubseteq$ is intended to express a connection between a configuration and its reduct. Thus, subject reduction may be stated as: $(\longrightarrow) \subseteq (\sqsubseteq)$, that is, $\sqsubseteq$ is indeed a conservative description of reduction.    $\square$

We have introduced an initial environment $\Gamma_0$ and used it in the definition of well-typedness, but we haven't yet ensured that the type schemes assigned to constants are an adequate description of their semantics. We now formulate two requirements that relate $\Gamma_0$ with $\xrightarrow{\delta}$. They are specializations of the subject reduction and progress properties to configurations that involve an application of a constant. They represent proof obligations that must be discharged when concrete definitions of $\mathcal{Q}$, $\xrightarrow{\delta}$, and $\Gamma_0$ are given.

1.7.6    DEFINITION: We require (i) $(\xrightarrow{\delta}) \subseteq (\sqsubseteq)$; and (ii) if the configuration $c\ v_1\ \ldots\ v_k/\mu$ (where $k \geq 0$) is well-typed, then either it is reducible, or $c\ v_1\ \ldots\ v_k$ is a value.    $\square$

The last point that remains to be settled before proving type soundness is the interaction between side effects and $\mathsf{let}$-polymorphism. The following example illustrates the problem:

$$
\mathsf{let}\ r = \mathsf{ref}\ \lambda z.z\ \mathsf{in}\ \mathsf{let}\ \_ = (r := \lambda z.(z \mathbin{\hat{+}} \hat{1}))\ \mathsf{in}\ !r\ \mathsf{true}
$$

This expression reduces to $\mathsf{true} \mathbin{\hat{+}} \hat{1}$, so it must not be well-typed. Yet, if natural type schemes are assigned to $\mathsf{ref}$, $!$, and $:=$ (see Example 1.9.5), then

it *is* well-typed with respect to the rules given so far, because r receives the polymorphic type scheme $\forall \mathtt{X}.\mathtt{ref}\,(\mathtt{X} \to \mathtt{X})$, which allows writing a function of type int $\to$ int into r and reading it back with type bool $\to$ bool. The problem is that let-polymorphism simulates a textual duplication of the let-bound expression ref $\lambda \mathtt{z}.\mathtt{z}$, while the semantics first reduces it to a value $m$, causing a new binding $m \mapsto \lambda \mathtt{z}.\mathtt{z}$ to appear in the store, then duplicates the address $m$. The new store binding is not duplicated: both copies of $m$ refer to the same memory cell. For this reason, generalization is unsound in this case, and must be restricted. Many authors have attempted to come up with a sound type system that accepts *all* pure programs and remains flexible enough in the presence of side effects (Tofte, 1988; Leroy, 1992). These proposals are often complex, which is why they have been abandoned in favor of an extremely simple *syntactic* restriction, known as the *value restriction* (Wright, 1995).

1.7.7   DEFINITION: A program satisfies the *value restriction* if and only if all subexpressions of the form let z = $\mathtt{t}_1$ in $\mathtt{t}_2$ are in fact of the form let z = $\mathtt{v}_1$ in $\mathtt{t}_2$. In the following, we assume that either all constants have pure semantics, or all programs satisfy the value restriction.                    □

Put slightly differently, the value restriction states that only values may be generalized. This eliminates the problem altogether, since duplicating values does not affect a program's semantics. Note that any program that does not satisfy the value restriction can be turned into one that does and has the same semantics: it suffices to change let z = $\mathtt{t}_1$ in $\mathtt{t}_2$ into $(\lambda \mathtt{z}.\mathtt{t}_2)\,\mathtt{t}_1$ when $\mathtt{t}_1$ is not a value. Of course, such a transformation may cause the program to become ill-typed. In other words, the value restriction causes some perfectly safe programs to be rejected. In particular, as stated above, it prevents generalizing applications of the form c $\mathtt{v}_1 \ldots \mathtt{v}_k$, where c is a destructor of arity $k$. This is excessive, because many destructors have pure semantics; only a few, such as ref, allocate new mutable storage. Furthermore, we use pure destructors to encode numerous language features (Section 1.9). Fortunately, it is easy to relax the restriction to allow generalizing not only values, but also a more general class of *nonexpansive* expressions, whose syntax guarantees that such expressions cannot allocate new mutable storage (that is, *expand* the domain of the store). The term *nonexpansive* was coined by Tofte (1988). Nonexpansive expressions may include applications of the form c $\mathtt{t}_1 \ldots \mathtt{t}_k$, where c is a pure destructor of arity $k$ and $\mathtt{t}_1, \ldots, \mathtt{t}_k$ are nonexpansive. Experience shows that this slightly relaxed restriction is acceptable in practice. Some other improvements to the value restriction exist; see *e.g.* Exercise (Garrigue, 2002). Another frequent limitation of the value restriction are constructor functions, that is, functions that only build values, which are treated as ordinary functions and not as constructors, and their applications are not considered to be

values. For instance, in the expression `let f = c v in let z = f w in t` where `c` is a constructor of arity 2, the partial application `c v` bound to `f` is a constructor function (of arity 1), but `f w` is treated as a regular application and cannot be generalized. Technically, the effect of the (strict) value restriction is summarized by the following result.

1.7.8   LEMMA: Under the value restriction, the production $\mathcal{E} ::= $ `let z = ` $\mathcal{E}$ ` in t` may be suppressed from the grammar of evaluation contexts (Figure 1-1) without altering the operational semantics.                                              □

We are done with definitions and requirements. We now come to the bulk of the type soundness proof.

1.7.9   THEOREM [SUBJECT REDUCTION]: $(\longrightarrow) \subseteq (\sqsubseteq)$.                                      □

*Proof:*   Because $\longrightarrow$ and $\longrightarrow\!\!\!\!\!\longrightarrow$ are the smallest relations that satisfy the rules of Figure 1-2, it suffices to prove that $\sqsubseteq$ satisfies these rules as well. We remark that if, for every type T, $[\![t : T]\!] \Vdash [\![t' : T]\!]$ holds, then $t/\mu \sqsubseteq t'/\mu$ holds. (Take $\bar{Y} = \varnothing$ and $M' = M$ and use the fact that entailment is a congruence to check that the conditions of Definition 1.7.5 are met.) We make use of this fact in cases R-BETA and R-LET below.

○ *Case* R-BETA. We have

$$
\begin{array}{ll}
& [\![(\lambda z.t)\ v : T]\!] \\
\equiv & \exists X.([\![\lambda z.t : X \to T]\!] \wedge [\![v : X]\!]) & \textbf{(1)} \\
\equiv & \exists X.(\text{let } z : X \text{ in } [\![t : T]\!] \wedge [\![v : X]\!]) & \textbf{(2)} \\
\equiv & \exists X.\text{let } z : \forall\varnothing[[\![v : X]\!]].X \text{ in } [\![t : T]\!] & \textbf{(3)} \\
\Vdash & [\![[z \mapsto v]t : T]\!] & \textbf{(4)}
\end{array}
$$

where (1) is by definition of constraint generation; (2) is by Lemma 1.6.5; (3) is by C-LETAND; (4) is by Lemma 1.7.2 and C-EX*.

○ *Case* R-LET. We have

$$
\begin{array}{ll}
& [\![\text{let } z = v \text{ in } t : T]\!] \\
= & \text{let } z : \forall X[[\![v : X]\!]].X \text{ in } [\![t : T]\!] & \textbf{(1)} \\
\Vdash & [\![[z \mapsto v]t : T]\!] & \textbf{(2)}
\end{array}
$$

where (1) is by definition of constraint generation and (2) is by Lemma 1.7.2.

○ *Case* R-DELTA. This case is exactly requirement (i) in Definition 1.7.6.

○ *Case* R-EXTEND. Our hypotheses are $t/\mu \sqsubseteq t'/\mu'$ **(1)** and $dom(\mu'')$ # $dom(\mu')$ **(2)** and $range(\mu'')$ # $dom(\mu' \setminus \mu)$ **(3)**. Because $dom(\mu)$ must be a subset of $dom(\mu')$, it is also disjoint with $dom(\mu'')$. Our goal is $t/\mu\mu'' \sqsubseteq t'/\mu'\mu''$ **(4)**. Thus, let us introduce a type T and a store type of domain

$dom(\mu\mu'')$, or (equivalently) two store types $M$ and $M''$ whose domains are respectively $dom(\mu)$ and $dom(\mu'')$. By (1), there exist type variables $\bar{\text{Y}}$ and a store type $M'$ such that $\bar{\text{Y}} \mathrel{\#} ftv(\text{T}, M)$ **(5)** and $ftv(M') \subseteq \bar{\text{Y}} \cup ftv(M)$ and $dom(M') = dom(\mu')$ and $M'$ extends $M$ **(6)** and let $\Gamma_0; \mathsf{ref}\, M$ in $[\![\mathsf{t}/\mu :$ $\text{T}/M]\!] \Vdash \exists \bar{\text{Y}}.\mathsf{let}\, \Gamma_0; \mathsf{ref}\, M'$ in $[\![\mathsf{t}'/\mu' : \text{T}/M']\!]$. We may further require, *w.l.o.g.*, $\bar{\text{Y}} \mathrel{\#} ftv(M'')$ **(7)**. Let us now add the conjunct $\mathsf{let}\, \Gamma_0; \mathsf{ref}\, M$ in $[\![\mu'' : M'']\!]$ to each side of this entailment assertion. On the left-hand side, by C-INAND and by Definition 1.7.4, we obtain $\mathsf{let}\, \Gamma_0; \mathsf{ref}\, M$ in $[\![\mathsf{t}/\mu\mu'' : \text{T}/MM'']\!]$ **(8)**. On the right-hand side, by (5), (7), C-EXAND, and C-INAND, we obtain $\exists \bar{\text{Y}}.\mathsf{let}\, \Gamma_0$ in $(\mathsf{let}\, \mathsf{ref}\, M'$ in $[\![\mathsf{t}'/\mu' : \text{T}/M']\!] \wedge \mathsf{let}\, \mathsf{ref}\, M$ in $[\![\mu' : M'']\!])$ **(9)**. Now, recall that $M'$ extends $M$ (6) and, furthermore, (3) implies $fpi([\![\mu'' :$ $M'']\!]) \mathrel{\#} dpi(M' \setminus M)$ **(10)**. By (10), C-INAND\*, and C-INAND, (9) is equivalent to $\exists \bar{\text{Y}}.\mathsf{let}\, \Gamma_0; \mathsf{ref}\, M'$ in $([\![\mathsf{t}'/\mu' : \text{T}/M']\!] \wedge [\![\mu'' : M'']\!])$, that is, $\exists \bar{\text{Y}}.\mathsf{let}\, \Gamma_0; \mathsf{ref}\, M'$ in $[\![\mathsf{t}'/\mu'\mu'' : \text{T}/M'M'']\!]$ **(11)**. Thus, we have established that (8) entails (11). Let us now place this entailment assertion within the constraint context $\mathsf{let}\,\, \mathsf{ref}\, M''$ in $[\![\,]\!]$. On the left-hand side, because $fpi(\Gamma_0, M, M'') = \varnothing$ and $dpi(M'') \cap dpi(\Gamma_0, M) \subseteq dom(\mu'') \cap (\mathcal{Q} \cup dom(\mu)) = \varnothing$, C-LETLET applies, yielding $\mathsf{let}\, \Gamma_0; \mathsf{ref}\, MM''$ in $[\![\mathsf{t}/\mu\mu'' : \text{T}/MM'']\!]$ **(12)**. On the right-hand side, by (7), C-INEX, and by analogous reasoning, we obtain $\exists \bar{\text{Y}}.\mathsf{let}\, \Gamma_0; \mathsf{ref}\, M'M''$ in $[\![\mathsf{t}'/\mu'\mu'' : \text{T}/M'M'']\!]$ **(13)**. Thus, (12) entails (13). Given (5), (7), given $ftv(M'M'') \subseteq \bar{\text{Y}} \cup ftv(MM'')$, and given that $M'M''$ extends $MM''$, this establishes the goal (4).

$\circ$ *Case* R-CONTEXT. The hypothesis is $\mathsf{t}/\mu \sqsubseteq \mathsf{t}'/\mu'$. The goal is $\mathcal{E}[\mathsf{t}]/\mu \sqsubseteq \mathcal{E}[\mathsf{t}']/\mu'$. Because $\longrightarrow$ relates closed configurations only, we may assume that the configuration $\mathcal{E}[\mathsf{t}]/\mu$ is closed, so the memory locations that appear free within $\mathcal{E}$ are members of $dom(\mu)$. Let us now reason by induction on the structure of $\mathcal{E}$.

*Subcase* $\mathcal{E} = [\,]$. The hypothesis and the goal coincide.

*Subcase* $\mathcal{E} = \mathcal{E}_1\, \mathsf{t}_1$. The induction hypothesis is $\mathcal{E}_1[\mathsf{t}]/\mu \sqsubseteq \mathcal{E}_1[\mathsf{t}']/\mu'$ **(1)**. Let us introduce a type $\text{T}$ and a store type $M$ such that $dom(M) = dom(\mu)$. Consider the constraint $\mathsf{let}\, \Gamma_0; \mathsf{ref}\, M$ in $[\![\mathcal{E}[\mathsf{t}]/\mu : \text{T}/M]\!]$ **(2)**. By definition of constraint generation, C-EXAND, C-INEX, and C-INAND, it is equivalent to

$$\exists \text{X}.(\mathsf{let}\, \Gamma_0; \mathsf{ref}\, M \text{ in } [\![\mathcal{E}_1[\mathsf{t}]/\mu : \text{X} \to \text{T}/M]\!] \wedge \mathsf{let}\, \Gamma_0; \mathsf{ref}\, M \text{ in } [\![\mathsf{t}_1 : \text{X}]\!]) \; \textbf{(3)}$$

where $\text{X} \notin ftv(\text{T}, M)$ **(4)**. By (1), there exist type variables $\bar{\text{Y}}$ and a store type $M'$ such that $\bar{\text{Y}} \mathrel{\#} ftv(\text{X}, \text{T}, M)$ **(5)** and $ftv(M') \subseteq \bar{\text{Y}} \cup ftv(M)$ **(6)** and $dom(M') = dom(\mu')$ and $M'$ extends $M$ and (3) entails

$$\exists \text{X}.(\exists \bar{\text{Y}}.\mathsf{let}\, \Gamma_0; \mathsf{ref}\, M' \text{ in } [\![\mathcal{E}_1[\mathsf{t}']/\mu' : \text{X} \to \text{T}/M']\!] \wedge \mathsf{let}\, \Gamma_0; \mathsf{ref}\, M \text{ in } [\![\mathsf{t}_1 : \text{X}]\!]) \; \textbf{(7)}.$$

We pointed out earlier that the memory locations that appear free in $\mathsf{t}_1$ are members of $dom(M)$, which implies $\mathsf{let}\, \mathsf{ref}\, M$ in $[\![\mathsf{t}_1 : \text{X}]\!] \equiv \mathsf{let}\, \mathsf{ref}\, M'$ in $[\![\mathsf{t}_1 :$

X⟧ **(8)**. By (5), C-ExAnd, (8), C-InAnd, and by definition of constraint generation, we find that (7) is equivalent to

$$\exists X \bar{Y}.\mathsf{let}\ \Gamma_0;\mathsf{ref}\ M'\ \mathsf{in}\ (\llbracket \mathcal{E}_1[\mathsf{t}'] : X \to T \rrbracket \wedge \llbracket \mathsf{t}_1 : X \rrbracket \wedge \llbracket \mu' : M' \rrbracket)\ (\mathbf{9}).$$

(4), (5) and (6) imply $X \notin ftv(M')$. Thus, by C-InEx and C-ExAnd, (9) may be written

$$\exists \bar{Y}.\mathsf{let}\ \Gamma_0;\mathsf{ref}\ M'\ \mathsf{in}\ (\exists X.(\llbracket \mathcal{E}_1[\mathsf{t}'] : X \to T \rrbracket \wedge \llbracket \mathsf{t}_1 : X \rrbracket) \wedge \llbracket \mu' : M' \rrbracket),$$

which, by definition of constraint generation, is

$$\exists \bar{Y}.\mathsf{let}\ \Gamma_0;\mathsf{ref}\ M'\ \mathsf{in}\ \llbracket \mathcal{E}[\mathsf{t}']/\mu' : T/M' \rrbracket\ (\mathbf{10}).$$

Thus, we have proved that (2) entails (10). By Definition 1.7.5, this establishes $\mathcal{E}[\mathsf{t}]/\mu \sqsubseteq \mathcal{E}[\mathsf{t}']/\mu'$.

*Subcase* $\mathcal{E} = \mathsf{v}\ \mathcal{E}_1$. Analogous to the previous subcase.

*Subcase* $\mathcal{E} = \mathtt{let}\ \mathtt{z} = \mathcal{E}_1\ \mathtt{in}\ \mathtt{t}_1$. The induction hypothesis is $\mathcal{E}_1[\mathsf{t}]/\mu \sqsubseteq \mathcal{E}_1[\mathsf{t}']/\mu'$ **(1)**. This subcase is particularly interesting, because it is where `let`-polymorphism and side effects interact. In the previous two subcases, we relied on the fact that the $\exists \bar{Y}$ quantifier, which hides the types of the memory cells created by the reduction step, *commutes* with the connectives $\exists$ and $\wedge$ introduced by application contexts. However, it does not in general (left-)commute with the `let` connective (Example 1.3.28). Fortunately, under the value restriction, this subcase *never arises* (Lemma 1.7.8). By Definition 1.7.7, this subcase may arise only if all constants have pure semantics, which implies $\mu = \mu' = \varnothing$. Then, we have

$$
\begin{aligned}
&\phantom{\equiv}\ \ \mathsf{let}\ \Gamma_0\ \mathsf{in}\ \llbracket \mathcal{E}[\mathsf{t}] : T \rrbracket \\
&= \ \ \mathsf{let}\ \Gamma_0; \mathsf{z} : \forall X[\llbracket \mathcal{E}_1[\mathsf{t}] : X \rrbracket].X\ \mathsf{in}\ \llbracket \mathsf{t}_1 : T \rrbracket \quad\quad\quad (\mathbf{2})\\
&\equiv \ \ \mathsf{let}\ \Gamma_0; \mathsf{z} : \forall X[\mathsf{let}\ \Gamma_0\ \mathsf{in}\ \llbracket \mathcal{E}_1[\mathsf{t}] : X \rrbracket].X\ \mathsf{in}\ \llbracket \mathsf{t}_1 : T \rrbracket \quad (\mathbf{3})\\
&\Vdash \ \ \mathsf{let}\ \Gamma_0; \mathsf{z} : \forall X[\mathsf{let}\ \Gamma_0\ \mathsf{in}\ \llbracket \mathcal{E}_1[\mathsf{t}'] : X \rrbracket].X\ \mathsf{in}\ \llbracket \mathsf{t}_1 : T \rrbracket \quad (\mathbf{4})\\
&\equiv \ \ \mathsf{let}\ \Gamma_0\ \mathsf{in}\ \llbracket \mathcal{E}[\mathsf{t}'] : T \rrbracket \quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad (\mathbf{5})
\end{aligned}
$$

where (2) is by definition of constraint generation; (3) follows from $ftv(\Gamma_0) = fpi(\Gamma_0) = \varnothing$ and C-LetDup; (4) follows from (1), specialized to the case of a pure semantics; and (5) is obtained by performing these steps in reverse.    □

1.7.10    Exercise [Recommended, ★★★]:  Try to carry out the last subcase of the above proof in the case of an impure semantics and in the absence of the value restriction. Find out why it fails. Show that it succeeds if $\bar{Y}$ is assumed to be empty. Use this fact to prove that generalization is still safe when restricted to *nonexpansive* expressions, provided (i) evaluating a nonexpansive expression

cannot cause new memory cells to be allocated, (ii) nonexpansive expressions are stable by substitution of values for variables, and (iii) nonexpansive expressions are preserved by reduction. □

Subject reduction ensures that well-typedness is preserved by reduction.

1.7.11   LEMMA: Let $\mathtt{t}/\mu \longrightarrow \mathtt{t}'/\mu'$. If $\mathtt{t}/\mu$ is well-typed, then so is $\mathtt{t}'/\mu'$. □

*Proof:* Assume $\mathtt{t}/\mu \longrightarrow \mathtt{t}'/\mu'$ **(1)** and $\mathtt{t}/\mu$ is well-typed **(2)**. By (2) and Definition 1.7.4, there exist a type $\mathtt{T}$ and a store type $M$ such that $dom(\mu) = dom(M)$ and the constraint $\mathsf{let}\ \Gamma_0; \mathsf{ref}\,M\ \mathsf{in}\ [\![\mathtt{t}/\mu : \mathtt{T}/M]\!]$ **(3)** is satisfiable. By Theorem 1.7.9 and Definition 1.7.5, (1) implies that there exist a set of type variables $\bar{\mathtt{Y}}$ and a store type $M'$ such that $dom(M') = dom(\mu')$ **(4)** and the constraint (3) entails $\exists\bar{\mathtt{Y}}.\mathsf{let}\ \Gamma_0; \mathsf{ref}\,M'\ \mathsf{in}\ [\![\mathtt{t}'/\mu' : \mathtt{T}/M']\!]$ **(5)**. Because (3) is satisfiable, so is (5), which implies that $\mathsf{let}\ \Gamma_0; \mathsf{ref}\,M'\ \mathsf{in}\ [\![\mathtt{t}'/\mu' : \mathtt{T}/M']\!]$ is satisfiable **(6)**. By (4) and (6) and Definition 1.7.4, $\mathtt{t}'/\mu'$ is well-typed. □

Let us now establish the progress property.

1.7.12   LEMMA: If $\mathtt{t}_1\ \mathtt{t}_2$ is well-typed, then $\mathtt{t}_1/\mu$ and $\mathtt{t}_2/\mu$ are well-typed. If $\mathsf{let}\ \mathtt{z} = \mathtt{t}_1\ \mathsf{in}\ \mathtt{t}_2/\mu$ is well-typed, then $\mathtt{t}_1/\mu$ is well-typed. □

1.7.13   THEOREM [PROGRESS]: If $\mathtt{t}/\mu$ is well-typed, then either it is reducible, or $\mathtt{t}$ is a value. □

*Proof:* The proof is by induction on the structure of $\mathtt{t}$.

  ∘ *Case* $\mathtt{t} = \mathtt{z}$. Well-typed configurations are closed: this case cannot occur.

  ∘ *Case* $\mathtt{t} = m$. $\mathtt{t}$ is a value.

  ∘ *Case* $\mathtt{t} = \mathtt{c}$. By requirement (ii) of Definition 1.7.6.

  ∘ *Case* $\mathtt{t} = \lambda\mathtt{z}.\mathtt{t}_1$. $\mathtt{t}$ is a value.

  ∘ *Case* $\mathtt{t} = \mathtt{t}_1\ \mathtt{t}_2$. By Lemma 1.7.12, $\mathtt{t}_1/\mu$ is well-typed. By the induction hypothesis, either it is reducible, or $\mathtt{t}_1$ is a value. If the former, by R-CONTEXT and because every context of the form $\mathcal{E}\ \mathtt{t}_2$ is an evaluation context, the configuration $\mathtt{t}/\mu$ is reducible as well. Thus, let us assume $\mathtt{t}_1$ is a value. By Lemma 1.7.12, $\mathtt{t}_2/\mu$ is well-typed. By the induction hypothesis, either it is reducible, or $\mathtt{t}_2$ is a value. If the former, by R-CONTEXT and because every context of the form $\mathtt{t}_1\ \mathcal{E}$—where $\mathtt{t}_1$ is a value—is an evaluation context, the configuration $\mathtt{t}/\mu$ is reducible as well. Thus, let us assume $\mathtt{t}_2$ is a value. Let us now reason by cases on the structure of $\mathtt{t}_1$.

  *Subcase* $\mathtt{t}_1 = \mathtt{z}$. Again, this subcase cannot occur.

  *Subcase* $\mathtt{t}_1 = m$. Because $\mathtt{t}/\mu$ is well-typed, a constraint of the form $\mathsf{let}\ \Gamma_0; \mathsf{ref}\,M\ \mathsf{in}\ (\exists\mathtt{X}.(m \preceq \mathtt{X} \to \mathtt{T} \wedge [\![\mathtt{t}_2 : \mathtt{X}]\!]) \wedge [\![\mu : M]\!])$ must be satisfiable. This implies that $m$ is a member of $dom(M)$ and that the constraint

ref $M(m) \leq \mathtt{X} \to \mathtt{T}$ is satisfiable. Because the type constructors ref and $\to$ are incompatible, this is a contradiction. So, this subcase cannot occur.

*Subcase* $\mathtt{t}_1 = \lambda\mathtt{z}.\mathtt{t}_1'$. By R-BETA, $\mathtt{t}/\mu$ is reducible.

*Subcase* $\mathtt{t}_1 = \mathtt{c}\ \mathtt{v}_1\ \ldots\ \mathtt{v}_k$. Then, $\mathtt{t}$ is of the form $\mathtt{c}\ \mathtt{v}_1\ \ldots\ \mathtt{v}_{k+1}$. The result follows by requirement (ii) of Definition 1.7.6.

$\circ$ *Case* $\mathtt{t} = \mathtt{let}\ \mathtt{z} = \mathtt{t}_1\ \mathtt{in}\ \mathtt{t}_2$. By Lemma 1.7.12, $\mathtt{t}_1/\mu$ is well-typed. By the induction hypothesis, either $\mathtt{t}_1/\mu$ is reducible, or $\mathtt{t}_1$ is a value. If the former, by R-CONTEXT and because every context of the form $\mathtt{let}\ \mathtt{z} = \mathcal{E}\ \mathtt{in}\ \mathtt{t}_2$ is an evaluation context, the configuration $\mathtt{t}/\mu$ is reducible as well. If the latter, then $\mathtt{t}/\mu$ is reducible by R-LET. $\qquad\qquad\square$

We may now conclude:

1.7.14   THEOREM [TYPE SOUNDNESS]: Well-typed   source   programs   do   not   go wrong. $\qquad\qquad\square$

*Proof:*   We say that a source program $\mathtt{t}$ is well-typed if and only if the configuration $\mathtt{t}/\varnothing$ is well-typed, that is, if and only if $\exists \mathtt{X}.\mathtt{let}\ \Gamma_0\ \mathtt{in}\ [\![\mathtt{t} : \mathtt{X}]\!] \equiv \mathtt{true}$ holds. By Lemma 1.7.11, all reducts of $\mathtt{t}/\varnothing$ are well-typed. By Theorem 1.7.13, none is stuck. $\qquad\qquad\square$

Let us recall that this result holds only if the requirements of Definition 1.7.6 are met. In other words, some proof obligations remain to be discharged when concrete definitions of $\mathcal{Q}$, $\xrightarrow{\delta}$, and $\Gamma_0$ are given. This is illustrated by several examples in the next section.

## 1.8   Constraint solving

We have introduced a parameterized constraint language, given equivalence laws that describe the interaction between its logical connectives, and exploited them to prove theorems about type inference and type soundness, which are valid independently of the nature of primitive constraints—the so-called predicate applications. However, there would be little point in proposing a parameterized constraint solver, because much of the difficulty of designing an efficient constraint solver precisely lies in the treatment of primitive constraints and in its interaction with let-polymorphism. For this reason, in this section, we focus on constraint solving in the setting of an *equality-only free tree model*. Thus, the constraint solver developed here allows performing type inference for HM(=) (that is, for Damas and Milner's type system) and for its extension with recursive types. Of course, some of its mechanisms may be useful in other settings. Other constraint solvers used in program analysis or type inference are described *e.g.* in (Aiken and Wimmers, 1992; Niehren,

Müller, and Podelski, 1997; Fähndrich, 1999; Melski and Reps, 2000; Müller, Niehren, and Treinen, 2001; Pottier, 2001b; Nielson, Nielson, and Seidl, 2002; McAllester, 2002, 2003).

We begin with a rule-based presentation of a standard, efficient first-order unification algorithm. This yields a constraint solver for a subset of the constraint language, deprived of type scheme introduction and instantiation forms. On top of it, we build a full constraint solver, which corresponds to the code that accompanies this chapter.

### Unification

Unification is the process of solving equations between terms. We now present a unification algorithm due to Huet (1976) as a (nondeterministic) system of constraint rewriting rules. The specification is almost the same in the case of *finite* and *regular* tree models: only one rule, which implements the *occurs check*, must be removed in the latter case. In other words, the algorithm works with *possibly cyclic* terms, and does not rely in an essential way on the occurs check. In order to accurately reflect the behavior of the actual algorithm, which relies on a *union-find* data structure (Tarjan, 1975), we modify the syntax of constraints by replacing equations with *multi-equations*. A multi-equation is an equation that involves an arbitrary number of types, as opposed to exactly two.

1.8.1   DEFINITION: Let there be, for every kind $\kappa$ and for every $n \geq 1$, a predicate $=_\kappa^n$, of signature $\kappa^n \Rightarrow \cdot$, whose interpretation is ($n$-ary) equality. The predicate constraint $=_\kappa^n \; \mathtt{T}_1 \ldots \mathtt{T}_n$ is written $\mathtt{T}_1 = \ldots = \mathtt{T}_n$, and called a *multi-equation*. We consider the constraint $\mathsf{true}$ as a multi-equation of length 0. In the following, we identify multi-equations up to permutations of their members, so a multi-equation $\epsilon$ of kind $\kappa$ may be viewed as a finite *multiset* of types of kind $\kappa$. We write $\epsilon = \epsilon'$ for the multi-equation obtained by concatenating $\epsilon$ and $\epsilon'$.                                                                                   □

Thus, we are interested in the following subset of the constraint language:

$$U ::= \mathsf{true} \mid \mathsf{false} \mid \epsilon \mid U \wedge U \mid \exists \bar{\mathsf{x}}.U$$

Equations are replaced with multi-equations; no other predicates are available. Type scheme introduction and instantiation forms are absent.

1.8.2   DEFINITION: A multi-equation is *standard* if and only if its variable members are distinct and it has at most one nonvariable member. A constraint $U$ is *standard* if and only if every multi-equation inside $U$ is standard and every variable that occurs (free or bound) in $U$ is a member of at most one multi-equation inside $U$.                                                                   □

A union-find algorithm maintains equivalence classes (that is, disjoint sets) of variables, and associates, with each class, a *descriptor*, which in our case is either absent or a nonvariable term. Thus, a *standard* constraint represents a state of the union-find algorithm. A constraint that is *not* standard may be viewed as a superposition of a state of the union-find algorithm, on the one hand, and of control information, on the other hand. For instance, a multi-equation of the form $\epsilon = \mathtt{T}_1 = \mathtt{T}_2$, where $\mathtt{T}_1$ and $\mathtt{T}_2$ are nonvariable terms, may be viewed, roughly speaking, as the equivalence class $\epsilon = \mathtt{T}_1$, together with a pending request to solve $\mathtt{T}_1 = \mathtt{T}_2$ and to update the class's descriptor accordingly. Because multi-equations encode both state and control, our specification of unification is rather high-level. It would be possible to give a lower-level description, where state (standard conjunctions of multi-equations) and control (pending binary equations) are distinguished.

1.8.3    DEFINITION:  Let $U$ be a conjunction of multi-equations. Y is *dominated* by X with respect to $U$ (written: $\mathtt{Y} \prec_U \mathtt{X}$) if and only if $U$ contains a conjunct of the form $\mathtt{X} = F\,\vec{\mathtt{T}} = \epsilon$, where $\mathtt{Y} \in \mathit{ftv}(\vec{\mathtt{T}})$. $U$ is *cyclic* if and only if the graph of $\prec_U$ exhibits a cycle.                                                                   □

The specification of the unification algorithm consists of a set of constraint rewriting rules, given in Figure 1-11. Rewriting is performed modulo $\alpha$-conversion, modulo permutations of the members of a multi-equation, modulo commutativity and associativity of conjunction, and under an arbitrary context. The specification is nondeterministic: several rule instances may be simultaneously applicable.

S-EXAND is a directed version of C-EXAND, whose effect is to float up all existential quantifiers. In the process, all multi-equations become part of a single conjunction, possibly causing rules whose left-hand side is a conjunction of multi-equations, namely S-FUSE and S-CYCLE, to become applicable. S-FUSE identifies two multi-equations that share a common variable X, and fuses them. The new multi-equation is not necessarily standard, even if the two original multi-equations were. Indeed, it may have repeated variables or contain two nonvariable terms. The purpose of the next few rules, whose left-hand side consists of a single multi-equation, is to deal with these situations. S-STUTTER eliminates redundant variables. It only deals with variables, as opposed to terms of arbitrary size, so as to have constant time cost. The comparison of nonvariable terms is implemented by S-DECOMPOSE and S-CLASH. S-DECOMPOSE decomposes an equation between two terms whose head symbols match. It produces a conjunction of equations between their subterms, namely $\vec{\mathtt{X}} = \vec{\mathtt{T}}$. Only one of the two terms remains in the original multi-equation, which may thus become standard. The terms $\vec{\mathtt{X}}$ are copied— there are two occurrences of $\vec{\mathtt{X}}$ on the right-hand side. For this reason, we

$$(\exists \bar{\mathtt{X}}.U_1) \wedge U_2 \quad \rightarrow \quad \exists \bar{\mathtt{X}}.(U_1 \wedge U_2) \qquad \text{(S-ExAnd)}$$
$$\text{if } \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(U_2)$$

$$\mathtt{X} = \epsilon \wedge \mathtt{X} = \epsilon' \quad \rightarrow \quad \mathtt{X} = \epsilon = \epsilon' \qquad \text{(S-Fuse)}$$

$$\mathtt{X} = \mathtt{X} = \epsilon \quad \rightarrow \quad \mathtt{X} = \epsilon \qquad \text{(S-Stutter)}$$

$$F\,\vec{\mathtt{X}} = F\,\vec{\mathtt{T}} = \epsilon \quad \rightarrow \quad \vec{\mathtt{X}} = \vec{\mathtt{T}} \wedge F\,\vec{\mathtt{X}} = \epsilon \qquad \text{(S-Decompose)}$$

$$F\,\mathtt{T}_1 \ldots \mathtt{T}_i \ldots \mathtt{T}_n = \epsilon \quad \rightarrow \quad \exists \mathtt{X}.(\mathtt{X} = \mathtt{T}_i \wedge F\,\mathtt{T}_1 \ldots \mathtt{X} \ldots \mathtt{T}_n = \epsilon) \qquad \text{(S-Name-1)}$$
$$\text{if } \mathtt{T}_i \notin \mathcal{V} \wedge \mathtt{X} \notin \mathit{ftv}(\mathtt{T}_1, \ldots, \mathtt{T}_n, \epsilon)$$

$$F\,\vec{\mathtt{T}} = F'\,\vec{\mathtt{T}}' = \epsilon \quad \rightarrow \quad \mathsf{false} \qquad \text{(S-Clash)}$$
$$\text{if } F \neq F'$$

$$\mathtt{T} \quad \rightarrow \quad \mathsf{true} \qquad \text{(S-Single)}$$
$$\text{if } \mathtt{T} \notin \mathcal{V}$$

$$U \wedge \mathsf{true} \quad \rightarrow \quad U \qquad \text{(S-True)}$$

$$U \quad \rightarrow \quad \mathsf{false} \qquad \text{(S-Cycle)}$$
$$\text{if the model is syntactic and } U \text{ is cyclic}$$

$$\mathcal{U}[\mathsf{false}] \quad \rightarrow \quad \mathsf{false} \qquad \text{(S-Fail)}$$
$$\text{if } \mathcal{U} \neq []$$

**Figure 1-11: Unification**

require them to be type variables, as opposed to terms of arbitrary size. (We slightly abuse notation by using $\vec{\mathtt{X}}$ to denote a vector of type variables whose elements are *not* necessarily distinct.) By doing so, we allow explicitly reasoning about *sharing*: since a variable represents a pointer to an equivalence class, we explicitly specify that only *pointers*, not whole terms, are copied. As a result of this decision, S-Decompose is not applicable when both terms at hand have a nonvariable subterm. S-Name-1 remedies this problem by introducing a fresh variable that stands for one such subterm. When repeatedly applied, S-Name-1 yields a unification problem composed of so-called *small terms* only—that is, where sharing has been made fully explicit. S-Clash complements S-Decompose by dealing with the case where two terms with different head symbols are equated; in a free tree model, such an equation is false, so failure is signaled. S-Single and S-True suppress multi-equations of size 1 and 0, respectively, which are tautologies. S-Single is restricted to non-variable terms so as not to break the property that every variable is a member

of exactly one multi-equation (Definition 1.8.2). S-Cycle is the occurs check: that is, it signals failure if the constraint is cyclic. It is applicable only in the case of syntactic unification, that is, when ground types are finite trees. It is a global check: its left-hand side is an entire conjunction of multi-equations. S-Fail propagates failure; $\mathcal{U}$ ranges over unification constraint contexts.

The constraint rewriting system in Figure 1-11 enjoys the following properties. First, rewriting is strongly normalizing, so the rules define a (nondeterministic) algorithm. Second, rewriting is meaning-preserving. Third, every normal form is either false or of the form $\exists \bar{x}.U$, where $U$ is satisfiable. The latter two properties indicate that the algorithm is indeed a constraint solver.

1.8.4    Lemma:  The rewriting system $\rightarrow$ is strongly normalizing.                                   □

1.8.5    Lemma:  $U_1 \rightarrow U_2$ implies $U_1 \equiv U_2$.                                   □

1.8.6    Lemma:  Every normal form is either false or of the form $\mathcal{X}[U]$, where $\mathcal{X}$ is an existential constraint context, $U$ is a standard conjunction of multi-equations and, if the model is syntactic, $U$ is acyclic. These conditions imply that $U$ is satisfiable.                                   □

### A constraint solver

On top of the unification algorithm, we now define a constraint solver. Its specification is independent of the rules and strategy employed by the unification algorithm. However, the structure of the unification algorithm's normal forms, as well as the logical properties of multi-equations, are exploited when performing generalization, that is, when creating and simplifying type schemes. Like the unification algorithm, the constraint solver is specified in terms of a *reduction system*. However, the objects that are subject to rewriting are not just constraints: they have more complex structure. Working with such richer *states* allows distinguishing the solver's external language—namely, the full constraint language, which is used to express the problem that one wishes to solve—and an internal language, introduced below, which is used to describe the solver's private data structures. In the following, $C$ and $D$ range over *external* constraints, that is, constraints that were part of the solver's input. External constraints are to be viewed as abstract syntax trees, subject to no implicit laws other than $\alpha$-conversion. As a simplifying assumption, we require external constraints not to contain any occurrence of false—otherwise the problem at hand is clearly false. *Internal* data structures include unification constraints $U$, as previously studied, and *stacks*. Stacks form a subset of constraint contexts, defined on page 24. Their syntax is as follows:

$$S ::= [] \mid S[[] \wedge C] \mid S[\exists \bar{x}.[]] \mid S[\text{let } x : \forall \bar{x}[[]].\mathsf{T} \text{ in } C] \mid S[\text{let } x : \sigma \text{ in } []]$$

In the second and fourth productions, $C$ is an external constraint. In the last production, we require $\sigma$ to be of the form $\forall \bar{\mathtt{x}}[U].\mathtt{X}$, and we demand $\exists \sigma \equiv \mathsf{true}$. A stack may be viewed as a list of *frames*. Frames may be added and deleted at the inner end of a stack, that is, near the hole of the constraint context that it represents. We refer to the four kinds of frames as *conjunction*, *existential*, $\mathsf{let}$, and *environment* frames, respectively. A *state* of the constraint solver is a triple $S; U; C$, where $S$ is a stack, $U$ is a unification constraint, and $C$ is an external constraint. The state $S; U; C$ is to be understood as a representation of the constraint $S[U \wedge C]$. The notion of $\alpha$-equivalence between states is defined accordingly. In particular, one may rename type variables in $dtv(S)$, provided $U$ and $C$ are renamed as well. In short, the three components of a state play the following roles. $C$ is an external constraint that the solver intends to examine next. $U$ is the internal state of the underlying unification algorithm: one might think of it as the knowledge that has been obtained so far. $S$ tells where the type variables that occur free in $U$ and $C$ are bound, associates type schemes with the program variables that occur free in $C$, and records what should be done after $C$ is solved. The solver's initial state is usually of the form $[]; \mathsf{true}; C$, where $C$ is the external constraint that one wishes to solve—that is, whose satisfiability one wishes to determine. For simplicity, we make the (unessential) assumption that states have no free type variables.

The solver consists of a (nondeterministic) state rewriting system, given in Figure 1-12. Rewriting is performed modulo $\alpha$-conversion. S-Unify makes the unification algorithm a component of the constraint solver, and allows the current unification problem $U$ to be solved at any time. Rules S-Ex-1 to S-Ex-4 float existential quantifiers out of the unification problem into the stack, and through the stack up to the nearest enclosing $\mathsf{let}$ frame, if there is any, or to the outermost level, otherwise. Their side-conditions prevent capture of type variables, and may always be satisfied by suitable $\alpha$-conversion of the left-hand state. If $S; U; C$ is a normal form with respect to the above five rules, then every type variable in $dtv(S)$ is either universally quantified at a $\mathsf{let}$ frame, or existentially bound at the outermost level. (Recall that, by assumption, states have no free type variables.) In other words, provided these rules are applied in an eager fashion, *there is no need for existential frames to appear in the machine representation of stacks.* Instead, it suffices to maintain, at every $\mathsf{let}$ frame and at the outermost level, a list of the type variables that are bound at this point; and, conversely, to annotate every type variable in $dtv(S)$ with an integer *rank*, which allows telling, in constant time, where the variable is bound: type variables of rank 0 are bound at the outermost level, and type variables of rank $k \geq 1$ are bound at the $k^{\mathrm{th}}$ $\mathsf{let}$ frame down in the stack $S$. The code that accompanies this chapter adopts this convention. Ranks were

$$S; U; C \quad\rightarrow\quad S; U'; C \qquad\qquad\qquad\text{(S-\textsc{Unify})}$$
$$\text{if } U \rightarrow U'$$

$$S; \exists \bar{\mathtt{X}}.U; C \quad\rightarrow\quad S[\exists \bar{\mathtt{X}}.[]]; U; C \qquad\qquad\qquad\text{(S-\textsc{Ex-1})}$$
$$\text{if } \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(C)$$

$$S[(\exists \bar{\mathtt{X}}.[]) \wedge C] \quad\rightarrow\quad S[\exists \bar{\mathtt{X}}.([] \wedge C)] \qquad\qquad\qquad\text{(S-\textsc{Ex-2})}$$
$$\text{if } \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(C)$$

$$S[\mathsf{let}\ \mathtt{x} : \forall\bar{\mathtt{X}}[\exists \bar{\mathtt{Y}}.[]].\mathtt{T}\ \mathsf{in}\ C] \quad\rightarrow\quad S[\mathsf{let}\ \mathtt{x} : \forall\bar{\mathtt{X}}\bar{\mathtt{Y}}[[]].\mathtt{T}\ \mathsf{in}\ C] \qquad\text{(S-\textsc{Ex-3})}$$
$$\text{if } \bar{\mathtt{Y}} \mathbin{\#} \mathit{ftv}(\mathtt{T})$$

$$S[\mathsf{let}\ \mathtt{x} : \sigma\ \mathsf{in}\ \exists \bar{\mathtt{X}}.[]] \quad\rightarrow\quad S[\exists \bar{\mathtt{X}}.\mathsf{let}\ \mathtt{x} : \sigma\ \mathsf{in}\ []] \qquad\qquad\text{(S-\textsc{Ex-4})}$$
$$\text{if } \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(\sigma)$$

$$S; U; \mathtt{T}_1 = \mathtt{T}_2 \quad\rightarrow\quad S; U \wedge \mathtt{T}_1 = \mathtt{T}_2; \mathsf{true} \qquad\qquad\text{(S-\textsc{Solve-Eq})}$$

$$S; U; \mathtt{x} \preceq \mathtt{T} \quad\rightarrow\quad S; U; S(\mathtt{x}) \preceq \mathtt{T} \qquad\qquad\qquad\text{(S-\textsc{Solve-Id})}$$

$$S; U; C_1 \wedge C_2 \quad\rightarrow\quad S[[] \wedge C_2]; U; C_1 \qquad\qquad\text{(S-\textsc{Solve-And})}$$

$$S; U; \exists \bar{\mathtt{X}}.C \quad\rightarrow\quad S[\exists \bar{\mathtt{X}}.[]]; U; C \qquad\qquad\qquad\text{(S-\textsc{Solve-Ex})}$$
$$\text{if } \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(U)$$

$$S; U; \mathsf{let}\ \mathtt{x} : \forall\bar{\mathtt{X}}[D].\mathtt{T}\ \mathsf{in}\ C \quad\rightarrow\quad S[\mathsf{let}\ \mathtt{x} : \forall\bar{\mathtt{X}}[[]].\mathtt{T}\ \mathsf{in}\ C]; U; D \qquad\text{(S-\textsc{Solve-Let})}$$
$$\text{if } \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(U)$$

$$S[[] \wedge C]; U; \mathsf{true} \quad\rightarrow\quad S; U; C \qquad\qquad\qquad\text{(S-\textsc{Pop-And})}$$

$$S[\mathsf{let}\ \mathtt{x} : \forall\bar{\mathtt{X}}[[]].\mathtt{T}\ \mathsf{in}\ C]; U; \mathsf{true} \quad\rightarrow\quad S[\mathsf{let}\ \mathtt{x} : \forall\bar{\mathtt{X}}\mathtt{X}[[]].\mathtt{X}\ \mathsf{in}\ C];$$
$$U \wedge \mathtt{X} = \mathtt{T}; \mathsf{true} \qquad\qquad\text{(S-\textsc{Name-2})}$$
$$\text{if } \mathtt{X} \notin \mathit{ftv}(U, \mathtt{T}) \wedge \mathtt{T} \notin \mathcal{V}$$

$$S[\mathsf{let}\ \mathtt{x} : \forall\bar{\mathtt{X}}\mathtt{Y}[[]].\mathtt{X}\ \mathsf{in}\ C]; \mathtt{Y} = \mathtt{Z} = \epsilon \wedge U; \mathsf{true} \quad\rightarrow\quad S[\mathsf{let}\ \mathtt{x} : \forall\bar{\mathtt{X}}\mathtt{Y}[[]].\theta(\mathtt{X})\ \mathsf{in}\ C];$$
$$\mathtt{Y} \wedge \mathtt{Z} = \theta(\epsilon) \wedge \theta(U); \mathsf{true} \qquad\text{(S-\textsc{Compress})}$$
$$\text{if } \mathtt{Y} \neq \mathtt{Z} \wedge \theta = [\mathtt{Y} \mapsto \mathtt{Z}]$$

$$S[\mathsf{let}\ \mathtt{x} : \forall\bar{\mathtt{X}}\mathtt{Y}[[]].\mathtt{X}\ \mathsf{in}\ C]; \mathtt{Y} = \epsilon \wedge U; \mathsf{true} \quad\rightarrow\quad S[\mathsf{let}\ \mathtt{x} : \forall\bar{\mathtt{X}}[[]].\mathtt{X}\ \mathsf{in}\ C]; \epsilon \wedge U; \mathsf{true} \qquad\text{(S-\textsc{UnName})}$$
$$\text{if } \mathtt{Y} \notin \mathtt{X} \cup \mathit{ftv}(\epsilon, U)$$

$$S[\mathsf{let}\ \mathtt{x} : \forall\bar{\mathtt{X}}\bar{\mathtt{Y}}[[]].\mathtt{X}\ \mathsf{in}\ C]; U; \mathsf{true} \quad\rightarrow\quad S[\exists \bar{\mathtt{Y}}.\mathsf{let}\ \mathtt{x} : \forall\bar{\mathtt{X}}[[]].\mathtt{X}\ \mathsf{in}\ C]; U; \mathsf{true} \qquad\text{(S-\textsc{LetAll})}$$
$$\text{if } \bar{\mathtt{Y}} \mathbin{\#} \mathit{ftv}(C) \wedge \exists \bar{\mathtt{X}}.U \text{ determines } \bar{\mathtt{Y}}$$

$$S[\mathsf{let}\ \mathtt{x} : \forall\bar{\mathtt{X}}[[]].\mathtt{X}\ \mathsf{in}\ C]; U_1 \wedge U_2; \mathsf{true} \quad\rightarrow\quad S[\mathsf{let}\ \mathtt{x} : \forall\bar{\mathtt{X}}[U_2].\mathtt{X}\ \mathsf{in}\ []]; U_1; C \qquad\text{(S-\textsc{Pop-Let})}$$
$$\text{if } \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(U_1) \wedge \exists \bar{\mathtt{X}}.U_2 \equiv \mathsf{true}$$

$$S[\mathsf{let}\ \mathtt{x} : \sigma\ \mathsf{in}\ []]; U; \mathsf{true} \quad\rightarrow\quad S; U; \mathsf{true} \qquad\qquad\qquad\text{(S-\textsc{Pop-Env})}$$

**Figure 1-12: A constraint solver**

initially described in (Rémy, 1992a), and also appear in (McAllester, 2003).

Rules S-Solve-Eq to S-Solve-Let encode an analysis of the structure of the third component of the current state. There is one rule for each possible case, except false, which by assumption cannot arise, and true, which is dealt with further on. S-Solve-Eq discovers an equation and makes it available to the unification algorithm. S-Solve-Id discovers an instantiation constraint $\text{x} \preceq \text{T}$ and replaces it with $\sigma \preceq \text{T}$, where the type scheme $\sigma = S(\text{x})$ is the type scheme carried by the nearest environment frame that defines x in the stack $S$. It is defined as follows:

$$
\begin{array}{rcl}
S[[] \wedge C](\text{x}) & = & S(\text{x}) \\
S[\exists \bar{\text{x}}.[]](\text{x}) & = & S(\text{x}) \quad \text{if } \bar{\text{x}} \mathbin{\#} \textit{ftv}(S(\text{x})) \\
S[\text{let y} : \forall \bar{\text{x}}[[]].\text{T in } C](\text{x}) & = & S(\text{x}) \quad \text{if } \bar{\text{x}} \mathbin{\#} \textit{ftv}(S(\text{x})) \\
S[\text{let y} : \sigma \text{ in } []](\text{x}) & = & S(\text{x}) \quad \text{if } \text{x} \neq \text{y} \\
S[\text{let x} : \sigma \text{ in } []](\text{x}) & = & \sigma
\end{array}
$$

If $\text{x} \in \textit{dpi}(S)$ does not hold, then $S(\text{x})$ is undefined and the rule is not applicable. If it does hold, then the rule may always be made applicable by suitable $\alpha$-conversion of the left-hand state. Please recall that, if $\sigma$ is of the form $\forall \bar{\text{x}}[U].\text{X}$, where $\bar{\text{x}} \mathbin{\#} \textit{ftv}(\text{T})$, then $\sigma \preceq \text{T}$ stands for $\exists \bar{\text{x}}.(U \wedge \text{X} = \text{T})$. The process of constructing this constraint is informally referred to as "taking an instance of $\sigma$". It involves taking fresh copies of the type variables $\bar{\text{x}}$, of the unification constraint $U$, and of the body X. In the worst case, this process is just as inefficient as textually expanding the corresponding let construct in the program's source code, and leads to exponential time complexity (Mairson, Kanellakis, and Mitchell, 1991). In practice, however, the unification constraint $U$ is often compact, because it was simplified before the environment frame let $\text{x} : \sigma$ in $[]$ was created. which is why the solver usually performs well. (The creation of environment frames, performed by S-Pop-Let, is discussed below.) S-Solve-And discovers a conjunction. It arbitrarily chooses to explore the left branch first, and pushes a conjunction frame onto the stack, so as to record that the right branch should be explored afterwards. S-Solve-Ex discovers an existential quantifier and enters it, creating a new existential frame to record its existence. Similarly, S-Solve-Let discovers a let form and enters its left-hand side, creating a new let frame to record its existence. The choice of examining the left-hand side first is *not* arbitrary. Indeed, examining the right-hand side first would require creating an environment frame—but environment frames must contain *simplified* type schemes of the form $\forall \bar{\text{x}}[U].\text{X}$, whereas the type scheme $\forall \bar{\text{x}}[D].\text{T}$ is arbitrary. In other words, our strategy is to simplify type schemes prior to allowing them to be copied by S-Solve-Id, so as to avoid any duplication of effort. The side-conditions of S-Solve-Ex and S-Solve-Let may always be satisfied by suitable $\alpha$-conversion of the left-hand state.

Rules S-SOLVE-EQ to S-SOLVE-LET may be referred to as *forward* rules, because they "move down into" the external constraint, causing the stack to grow. This process stops when the external constraint at hand becomes true. Then, part of the work has been finished, and the solver must examine the stack in order to determine what to do next. This task is performed by the last series of rules, which may be referred to as *backward* rules, because they "move back out", causing the stack to shrink, and possibly scheduling new external constraints for examination. These rules encode an analysis of the structure of the innermost stack frame. There are three cases, corresponding to conjunction, let, and environment frames. The case of existential stack frames need not be considered, because rules S-EX-2 to S-EX-4 allow either fusing them with let frames or floating them up to the outermost level, where they shall remain inert. S-POP-AND deals with conjunction frames. The frame is popped, and the external constraint that it carries is scheduled for examination. S-POP-ENV deals with environment frames. Because the right-hand side of the let construct at hand has been solved—that is, turned into a unification constraint $U$—it cannot contain an occurrence of x. Furthermore, by assumption, $\exists\sigma$ is true. Thus, this environment frame is no longer useful: it is destroyed. The remaining rules deal with let frames. Roughly speaking, their purpose is to change the state $S[\text{let } \mathtt{x} : \forall\bar{\mathtt{x}}[[]].\mathtt{T} \text{ in } C]; U; \text{true}$ into $S[\text{let } \mathtt{x} : \forall\bar{\mathtt{x}}[U].\mathtt{T} \text{ in } []]; \text{true}; C$, that is, to turn the current unification constraint $U$ into a type scheme, turn the let frame into an environment frame, and schedule the right-hand side of the let construct (that is, the external constraint $C$) for examination. In fact, the process is more complex, because the type scheme $\forall\bar{\mathtt{x}}[U].\mathtt{T}$ must be *simplified* before becoming part of an environment frame. The simplification process is described by rules S-NAME-2 to S-POP-LET. In the following, we refer to type variables in $\bar{\mathtt{x}}$ as *young* and to type variables in $dtv(S) \setminus \bar{\mathtt{x}}$ as *old*. The former are the universal quantifiers of the type scheme that is being created; the latter are its free type variables.

S-NAME-2 ensures that the body T of the type scheme that is being created is a type variable, as opposed to an arbitrary term. If it isn't, then it is replaced with a fresh variable X, and the equation X = T is added so as to recall that X stands for T. Thus, the rule moves the term T into the current unification problem, where it potentially becomes subject to S-NAME-1. This ensures that sharing is made explicit everywhere. S-COMPRESS determines that the (young) type variable Y is an alias for the type variable Z. Then, every free occurrence of Y other than its defining occurrence is replaced with Z. In an actual implementation, this occurs transparently when the union-find algorithm performs *path compression* (Tarjan, 1975, 1979), provided we are careful never to create a link from a variable to a variable of higher rank. This requires making the unification algorithm aware of ranks, but is otherwise

easily achieved. S-UnName determines that the (young) type variable Y has
no occurrences other than its defining occurrence in the current type scheme.
(This occurs, in particular, when S-Compress has just been applied.) Then,
Y is suppressed altogether. In the particular case where the remaining multi-
equation $\epsilon$ has cardinal 1, it may then be suppressed by S-Single. In other
words, the combination of S-UnName and S-Single is able to suppress young
unused type variables as well as the term that they stand for. This may,
in turn, cause new type variables to become eligible for elimination by S-
UnName. In fact, assuming the current unification constraint is acyclic, an
inductive argument shows that every young type variable may be suppressed
unless it is dominated either by X or by an old type variable. (In the setting
of a regular tree model, it is possible to extend the rule so that young cycles
that are not dominated either by X or by an old type variable are suppressed
as well.) S-LetAll is a directed version of C-LetAll. It turns the young
type variables $\bar{\text{Y}}$ into old variables. How to tell whether $\exists \bar{\text{X}}.U$ determines $\bar{\text{Y}}$
is discussed later (see Lemma 1.8.7). Why S-LetAll is an interesting and
important rule will be explained shortly. S-Pop-Let is meant to be applied
when the current state has become a normal form with respect to S-Unify, S-
Name-2, S-Compress, S-UnName, and S-LetAll, that is, when the type
scheme that is about to be created is fully simplified. It splits the current
unification constraint into two components $U_1$ and $U_2$, where $U_1$ is made up
entirely of *old* variables—as expressed by the side-condition $\bar{\text{X}} \mathrel{\#} \mathit{ftv}(U_1)$—
and $U_2$ constrains *young* variables only—as expressed by the side-condition
$\exists \bar{\text{X}}.U_2 \equiv \mathsf{true}$. Please note that $U_2$ may still contain free occurrences of old type
variables, so the type scheme $\forall \bar{\text{X}}[U_2].\text{X}$ that appears on the right-hand side is
not necessarily closed. It is not obvious why such a decomposition must exist;
the proof of Lemma 1.8.11 sheds more light on this issue. Let us say, for now,
that S-LetAll plays a role in guaranteeing its existence, whence part of its
importance. Once the decomposition $U_1 \wedge U_2$ is obtained, the behavior of S-
Pop-Let is simple. The unification constraint $U_1$ concerns old variables only,
that is, variables that are not quantified in the current let frame; thus, it need
not become part of the new type scheme, and may instead remain part of the
current unification constraint. This is justified by C-LetAnd and C-InAnd*
(see the proof of Lemma 1.8.10) and corresponds to the difference between
hmx-Gen' and hmx-Gen discussed in Section 1.4. The unification constraint
$U_2$, on the other hand, becomes part of the newly built type scheme $\forall \bar{\text{X}}[U_2].\text{X}$.
The property $\exists \bar{\text{X}}.U_2 \equiv \mathsf{true}$ guarantees that the newly created environment
frame meets the requirements imposed on such frames. Please note that, the
more type variables are considered old, the larger $U_1$ may become, and the
smaller $U_2$. This is another reason why S-LetAll is interesting: by allowing
more variables to be considered old, it decreases the size of the type scheme

$\forall \bar{\mathtt{X}}[U_2].\mathtt{X}$, making it cheaper to take instances of.

To complete our description of the constraint solver, there remains to explain how to decide when $\exists \bar{\mathtt{X}}.U$ determines $\bar{\mathtt{Y}}$, since this predicate occurs in the side-condition of S-LETALL. The following lemma describes two important situations where, by examining the structure of an equation, it is possible to discover that a constraint $C$ *determines* some of its free type variables $\bar{\mathtt{Y}}$ (Definition 1.3.26). In the first situation, the type variables $\bar{\mathtt{Y}}$ are *equated* with or *dominated* by a distinct type variable $\mathtt{X}$ that occurs *free* in $C$. In that case, because the model is a free tree model, the values of the type variables $\bar{\mathtt{Y}}$ are determined by the value of $\mathtt{X}$—they are subtrees of it at specific positions. For instance, $\mathtt{X} = \mathtt{Y}_1 \to \mathtt{Y}_2$ determines $\mathtt{Y}_1\mathtt{Y}_2$, while $\exists \mathtt{Y}_1.(\mathtt{X} = \mathtt{Y}_1 \to \mathtt{Y}_2)$ determines $\mathtt{Y}_2$. In the second situation, the type variables $\bar{\mathtt{Y}}$ are equated with a term $\mathtt{T}$, *all* of whose free type variables are *free* in $C$. Again, the value of the type variables $\bar{\mathtt{Y}}$ is then determined by the values of the type variables $ftv(\mathtt{T})$—indeed, the term $\mathtt{T}$ itself defines a function that maps the latter to the former. For instance, $\mathtt{X} = \mathtt{Y}_1 \to \mathtt{Y}_2$ determines $\mathtt{X}$, while $\exists \mathtt{Y}_1.(\mathtt{X} = \mathtt{Y}_1 \to \mathtt{Y}_2)$ does not. In the second situation, no assumption is in fact made about the model. Please note that $\mathtt{X} = \mathtt{Y}_1 \to \mathtt{Y}_2$ determines $\mathtt{Y}_1\mathtt{Y}_2$ and determines $\mathtt{X}$, but does *not* simultaneously determine $\mathtt{X}\mathtt{Y}_1\mathtt{Y}_2$.

1.8.7   LEMMA: Let $\bar{\mathtt{X}} \mathbin{\#} \bar{\mathtt{Y}}$. Assume either $\epsilon$ is $\mathtt{X} = \epsilon'$, where $\mathtt{X} \notin \bar{\mathtt{X}}\bar{\mathtt{Y}}$ and $\bar{\mathtt{Y}} \subseteq ftv(\epsilon')$, or $\epsilon$ is $\bar{\mathtt{Y}} = \mathtt{T} = \epsilon'$, where $ftv(\mathtt{T}) \mathbin{\#} \bar{\mathtt{X}}\bar{\mathtt{Y}}$. Then, $\exists \bar{\mathtt{X}}.(C \wedge \epsilon)$ determines $\bar{\mathtt{Y}}$.            □

*Proof:* Let $\bar{\mathtt{X}} \mathbin{\#} \bar{\mathtt{Y}}$ **(1)**. Let $\phi \vdash \mathsf{def}\ \Gamma\ \mathsf{in}\ \exists \bar{\mathtt{X}}.(C \wedge \epsilon)$ **(2)** and $\phi' \vdash \mathsf{def}\ \Gamma\ \mathsf{in}\ \exists \bar{\mathtt{X}}.(C \wedge \epsilon)$ **(3)**, where $\phi$ and $\phi'$ coincide outside of $\bar{\mathtt{Y}}$. We may assume, *w.l.o.g.*, $\bar{\mathtt{X}} \mathbin{\#} ftv(\Gamma)$ **(4)**. By (2), (4), CM-EXISTS, and CM-AND, we obtain $\phi_1 \vdash \mathsf{def}\ \Gamma\ \mathsf{in}\ \epsilon$ **(5)**, where $\phi$ and $\phi_1$ coincide outside $\bar{\mathtt{X}}$. By CM-PREDICATE, (5) implies that all members of $\epsilon$ have the same image through $\phi_1$. Similarly, exploiting (3) and (4), we find that all members of $\epsilon$ have the same image through $\phi'_1$, where $\phi'$ and $\phi'_1$ coincide outside $\bar{\mathtt{X}}$. Now, we claim that $\phi_1$ and $\phi'_1$ coincide on $\bar{\mathtt{Y}}$. Once the claim is established, by (1), there follows that $\phi$ and $\phi'$ must coincide on $\bar{\mathtt{Y}}$ as well, which is the goal. So, there only remains to establish the claim; we distinguish two subcases.

*Subcase* $\epsilon$ is $\mathtt{X} = \epsilon'$ and $\mathtt{X} \notin \bar{\mathtt{X}}\bar{\mathtt{Y}}$ **(6)** and $\bar{\mathtt{Y}} \subseteq ftv(\epsilon')$ **(7)**. Because $\phi_1$ and $\phi'_1$ coincide outside $\bar{\mathtt{X}}\bar{\mathtt{Y}}$ and by (6), we have $\phi_1(\mathtt{X}) = \phi'_1(\mathtt{X})$. As a result, all members of $\epsilon'$ have the same image through $\phi_1$ and $\phi'_1$. In a free tree model, where decomposition is valid, a simple inductive argument shows that $\phi_1$ and $\phi'_1$ must coincide on $ftv(\epsilon')$, hence—by (7)—also on $\bar{\mathtt{Y}}$.

*Subcase* $\epsilon$ is $\bar{\mathtt{Y}} = \mathtt{T} = \epsilon'$ and $ftv(\mathtt{T}) \mathbin{\#} \bar{\mathtt{X}}\bar{\mathtt{Y}}$ **(8)**. Because $\phi_1$ and $\phi'_1$ coincide outside $\bar{\mathtt{X}}\bar{\mathtt{Y}}$ and by (8), we have $\phi_1(\mathtt{T}) = \phi'_1(\mathtt{T})$. Thus, for every $\mathtt{Y} \in \bar{\mathtt{Y}}$, we have $\phi_1(\mathtt{Y}) = \phi_1(\mathtt{T}) = \phi'_1(\mathtt{T}) = \phi'_1(\mathtt{Y})$. That is, $\phi_1$ and $\phi'_1$ coincide on $\bar{\mathtt{Y}}$.            □

Thanks to Lemma 1.8.7, a straightforward implementation of S-LETALL

comes to mind. The problem is, given a constraint $\exists \bar{X}.U$, where $U$ is a standard conjunction of multi-equations, to determine the greatest subset $\bar{Y}$ of $\bar{X}$ such that $\exists (\bar{X} \setminus \bar{Y}).U$ determines $\bar{Y}$. By the first part of the lemma, it is safe for $\bar{Y}$ to include all members of $\bar{X}$ that are directly or indirectly dominated (with respect to $U$) by some free variable of $\exists \bar{X}.U$. Those can be found, in time linear in the size of $U$, by a top-down traversal of the graph of $\prec_U$. By the second part of the lemma, it is safe to close $\bar{Y}$ under the closure law $X \in \bar{X} \wedge (\forall Y \quad Y \prec_U X \Rightarrow Y \in \bar{Y}) \Rightarrow X \in \bar{Y}$. That is, it is safe to also include all members of $\bar{X}$ whose descendants (with respect to $U$) have already been found to be members of $\bar{Y}$. This closure computation may be performed, again in linear time, by a bottom-up traversal of the graph of $\prec_U$. When $U$ is acyclic, it is possible to show that this procedure is complete, that is, does compute the greatest subset $\bar{Y}$ that meets our requirement. This is the topic of the following exercise.

1.8.8   EXERCISE [★★★, ↛]: Assuming $U$ is acyclic, prove that the above procedure computes the greatest subset $\bar{Y}$ of $\bar{X}$ such that $\exists (\bar{X} \setminus \bar{Y}).U$ determines $\bar{Y}$. In the setting of a regular tree model, exhibit a satisfiable constraint $U$ such that the above procedure is incomplete. Can you define a complete procedure in that setting?                                                               □

The above discussion has shown that when Y and Z are equated, if Y is young and Z is old, then S-LETALL allows making Y old as well. If binding information is encoded in terms of integer ranks, as suggested earlier, then this remark may be formulated as follows: when Y and Z are equated, if the rank of Y exceeds that of Z, then it may be decreased so that both ranks match. As a result, it is possible to attach ranks with *multi-equations*, rather than with variables. When two multi-equations are fused, the smaller rank is kept.

S-SOLVE-LET and S-NAME-2 to S-POP-LET are unnecessarily complex when x is assigned a *monotype* T, rather than an arbitrary type scheme $\forall \bar{X}[D].$T. In that case, the combined effect of these rules may be obtained directly via the following two new rules, which may be implemented in a more efficient way:

$$S; U; \mathsf{let}\ x : \mathtt{T}\ \mathsf{in}\ C \quad \rightarrow \quad S[\exists \mathtt{X}.[]]; U \wedge \mathtt{X} = \mathtt{T}; \mathsf{let}\ x : \mathtt{X}\ \mathsf{in}\ C$$
$$\text{(S-NAME-2-MONO)}$$
$$\text{if}\ \mathtt{X} \notin \mathit{ftv}(U, \mathtt{T}, C) \wedge \mathtt{T} \notin \mathcal{V}$$

$$S; U; \mathsf{let}\ x : \mathtt{X}\ \mathsf{in}\ C \quad \rightarrow \quad S[\mathsf{let}\ x : \mathtt{X}\ \mathsf{in}\ []]; U; C \qquad \text{(S-SOLVE-LET-MONO)}$$

If T isn't a variable, it is replaced with a fresh variable X, together with the equation X = T. This corresponds to the effect of S-NAME-2. Then, we directly

create an environment frame for x, without bothering to create and discard a
let frame, since there is no way the type scheme X may be further simplified.

Let us now state and establish the properties of the constraint solver. First,
the reduction system is terminating, so it defines an algorithm.

1.8.9    LEMMA: The reduction system $\to$ is strongly normalizing.                    □

Second, every rewriting step preserves the meaning of the constraint that
the current state represents. We recall that the state $S; U; C$ is meant to
represent the constraint $S[U \wedge C]$.

1.8.10    LEMMA: $S; U; C \to S'; U'; C'$ implies $S[U \wedge C] \equiv S'[U' \wedge C']$.                    □

*Proof:*   By examination of every rule.

∘ *Case* S-UNIFY. By Lemma 1.8.5.

∘ *Case* S-EX-1, S-EX-2, S-SOLVE-EX. By C-EXAND.

∘ *Case* S-EX-3. By C-LETEX.

∘ *Case* S-EX-4. By C-INEX.

∘ *Case* S-SOLVE-EQ, S-POP-AND. By C-DUP.

∘ *Case* S-SOLVE-ID. Because $\sigma$ is of the form $\forall \bar{\mathtt{X}}[U].\mathtt{X}$, we have $fpi(\sigma) = \varnothing$.
The result follows by C-INID.

∘ *Case* S-SOLVE-AND. By C-ANDAND.

∘ *Case* S-SOLVE-LET. By C-LETAND.

∘ *Case* S-NAME-2. By Definition 1.3.21 and C-NAMEEQ, $\mathtt{X} \notin ftv(U, \mathtt{T})$ implies $\mathtt{true} \Vdash \forall \bar{\mathtt{X}}[U].\mathtt{T} \equiv \forall \bar{\mathtt{X}}\mathtt{X}[U \wedge \mathtt{X} = \mathtt{T}].\mathtt{X}$. The result follows by Lemma 1.3.22.

∘ *Case* S-COMPRESS. Let $\theta = [\mathtt{Y} \mapsto \mathtt{Z}]$. By Definition 1.3.21 and C-NAMEEQ, $\mathtt{Y} \neq \mathtt{Z}$ implies $\mathtt{true} \Vdash \forall \bar{\mathtt{X}}\mathtt{Y}[\mathtt{Y} = \mathtt{Z} = \epsilon \wedge U].\mathtt{X} \equiv \forall \bar{\mathtt{X}}\mathtt{Y}[\mathtt{Y} \wedge \mathtt{Z} = \theta(\epsilon) \wedge \theta(U)].\theta(\mathtt{X})$. The result follows by Lemma 1.3.22.

∘ *Case* S-UNNAME. Using Lemma 1.3.18, it is straightforward to check that $\mathtt{Y} \notin ftv(\epsilon)$ implies $\exists \mathtt{Y}.(\mathtt{Y} = \epsilon) \equiv \epsilon$. The result follows by C-EXAND and C-LETEX.

∘ *Case* S-LETALL. By C-LETALL.

∘ *Case* S-POP-LET. By C-LETAND and C-INAND*.

∘ *Case* S-POP-ENV. By C-IN*, recalling that $\exists \sigma$ must be true.                    □

Last, we classify the normal forms of the reduction system:

1.8.11    LEMMA: A normal form for the reduction system $\to$ is one of (i) $S; U; \mathtt{x} \preceq \mathtt{T}$,
where $\mathtt{x} \notin dpi(S)$; (ii) $S; \mathtt{false}; \mathtt{true}$; or (iii) $\mathcal{X}; U; \mathtt{true}$, where $\mathcal{X}$ is an existential
constraint context and $U$ a satisfiable conjunction of multi-equations.                    □

*Proof:*   Because, by definition, $S; U;\mathsf{false}$ is not a valid state, a normal form for S-Solve-Eq, S-Solve-Id, S-Solve-And, S-Solve-Ex, and S-Solve-Let must be either an instance of the left-hand side of S-Solve-Id, with $\mathsf{x} \notin \mathit{dpi}(S)$, which corresponds to case (i), or of the form $S; U;\mathsf{true}$. Let us consider the latter case. Because $S; U;\mathsf{true}$ is a normal form with respect to S-Unify, by Lemma 1.8.6, $U$ must be either $\mathsf{false}$ of the form $\mathcal{X}[U']$, where $U'$ is a standard conjunction of multi-equations and, if the model is syntactic, $U'$ is acyclic. The former case corresponds to (ii); thus, let us consider the latter case. Because $S; \mathcal{X}[U'];\mathsf{true}$ is a normal form with respect to S-Ex-1, the context $\mathcal{X}$ must in fact be empty, and $U'$ is $U$. If $S$ is an existential constraint context, then we are in situation (iii). Otherwise, because $S; U;\mathsf{true}$ is a normal form with respect to S-Ex-2, S-Ex-3, and S-Ex-4, the stack $S$ does not end with an existential frame. Because $S; U;\mathsf{true}$ is a normal form with respect to S-Pop-And and S-Pop-Env, $S$ must then be of the form $S'[\mathsf{let}\ \mathsf{x} : \forall\bar{\mathsf{x}}[[]].\mathsf{T}\ \mathsf{in}\ C]$. Because $S; U;\mathsf{true}$ is a normal form with respect to S-Name-2, $\mathsf{T}$ must be a type variable $\mathsf{X}$. Let us write $U$ as $U_1 \wedge U_2$, where $\bar{\mathsf{x}}\ \#\ \mathit{ftv}(U_1)$, and where $U_1$ is maximal for this criterion. Then, consider a multi-equation $\epsilon \in U$. By the first part of Lemma 1.8.7, if one variable member of $\epsilon$ is free (that is, outside $\bar{\mathsf{x}}$), then $\exists\bar{\mathsf{x}}.U$ determines all other variables in $\mathit{ftv}(\epsilon)$. Because $S; U;\mathsf{true}$ is a normal form with respect to S-LetAll, all variables in $\mathit{ftv}(\epsilon)$ must then be free (that is, outside $\bar{\mathsf{x}}$). By definition of $U_1$, this implies $\epsilon \in U_1$. By contraposition, for every multi-equation $\epsilon \in U_2$, *all variable members of $\epsilon$ are in $\bar{\mathsf{x}}$*. Furthermore, let us recall that $U_2$ is a standard conjunction of multi-equations and, if the model is syntactic, $U_2$ is acyclic. We let the reader check that this implies $\exists\bar{\mathsf{x}}.U_2 \equiv \mathsf{true}$; the proof is a slight generalization of the last part of that of Lemma 1.8.6. Then, $S; U;\mathsf{true}$ is reducible via S-Pop-Let. This is a contradiction, so this last case cannot arise.                                                                                                      □

In case (i), the constraint $S[U \wedge C]$ has a free program identifier $\mathsf{x}$, so it is not satisfiable. In other words, the source program contains an unbound program identifier. Such an error could of course be detected prior to constraint solving, if desired. In case (ii), the unification algorithm failed. By Lemma 1.3.30, the constraint $S[U \wedge C]$ is then false. In case (iii), the constraint $S[U \wedge C]$ is equivalent to $\mathcal{X}[U]$, where $U$ is satisfiable, so it is satisfiable as well. Thus, each of the three classes of normal forms may be immediately identified as denoting success or failure. Thus, Lemmas 1.8.10 and 1.8.11 indeed prove that the algorithm is a constraint solver.

## 1.9   From ML-the-calculus to ML-the-programming-language

In this section, we explain how to extend the framework developed so far to accommodate operations on values of base type (such as integers), pairs, sums, references, and recursive function definitions. Then, we describe more complex extensions, namely algebraic data type definitions, pattern matching, and type annotations. Last, the issues associated with recursive types are briefly discussed. Exceptions are not discussed; the reader is referred to (TAPL Chapter 14).

### Simple extensions

Many features of ML-the-programming-language may be introduced into ML-the-calculus by introducing new constants and extending $\overset{\delta}{\longrightarrow}$ and $\Gamma_0$ appropriately. In each case, it is necessary to check that the requirements of Definition 1.7.6 are met, that is, the new initial environment faithfully reflects the nature of the new constants as well as the behavior of the new reduction rules. Below, we describe several such extensions in isolation.

1.9.1   EXERCISE [INTEGERS, RECOMMENDED, ★★]: Integer literals and integer addition have been introduced and given an operational semantics in Examples 1.2.1, 1.2.2 and 1.2.4. Let us now introduce an isolated type constructor int of signature ⋆ and extend the initial environment $\Gamma_0$ with the bindings $\hat{n}$ : int, for every integer $n$, and $\hat{+}$ : int → int → int. Check that these definitions meet the requirements of Definition 1.7.6.                    □

1.9.2   EXERCISE [BOOLEANS, RECOMMENDED, ★★, ↛]: Booleans and conditionals have been introduced and given an operational semantics in Exercise 1.2.6. Introduce an isolated type constructor bool to represent Boolean values and explain how to extend the initial environment. Check that your definitions meet the requirements of Definition 1.7.6. What is the constraint generation rule for the syntactic sugar if $t_0$ then $t_1$ else $t_2$?                    □

1.9.3   EXERCISE [PAIRS, ★★, ↛]: Pairs and pair projections have been introduced and given an operational semantics in Examples 1.2.3 and 1.2.5. Let us now introduce an isolated type constructor × of signature ⋆ ⊗ ⋆ ⇒ ⋆, covariant in both of its parameters, and extend the initial environment $\Gamma_0$ with the following bindings:

$$(\cdot,\cdot) : \quad \forall \text{XY}.\text{X} \to \text{Y} \to \text{X} \times \text{Y}$$
$$\pi_1 \ : \quad \forall \text{XY}.\text{X} \times \text{Y} \to \text{X}$$
$$\pi_2 \ : \quad \forall \text{XY}.\text{X} \times \text{Y} \to \text{Y}$$

Check that these definitions meet the requirements of Definition 1.7.6.     □

1.9.4    EXERCISE [SUMS, ★★, ⇸]: Sums have been introduced and given an oper-
ational semantics in Example 1.2.7. Let us now introduce an isolated type
constructor + of signature $\star \otimes \star \Rightarrow \star$, covariant in both of its parameters, and
extend the initial environment $\Gamma_0$ with the following bindings:

$$\begin{aligned}
\texttt{inj}_1 &: \quad \forall \texttt{XY}.\texttt{X} \to \texttt{X} + \texttt{Y} \\
\texttt{inj}_2 &: \quad \forall \texttt{XY}.\texttt{Y} \to \texttt{X} + \texttt{Y} \\
\texttt{case} &: \quad \forall \texttt{XYZ}.(\texttt{X} + \texttt{Y}) \to (\texttt{X} \to \texttt{Z}) \to (\texttt{Y} \to \texttt{Z}) \to \texttt{Z}
\end{aligned}$$

Check that these definitions meet the requirements of Definition 1.7.6.    □

1.9.5    EXERCISE [REFERENCES, ★★★]: References have been introduced and
given an operational semantics in Example 1.2.9. The type constructor ref has
been introduced in Definition 1.7.4. Let us now extend the initial environment
$\Gamma_0$ with the following bindings:

$$\begin{aligned}
\texttt{ref} &: \quad \forall \texttt{X}.\texttt{X} \to \textsf{ref}\,\texttt{X} \\
\texttt{!} &: \quad \forall \texttt{X}.\textsf{ref}\,\texttt{X} \to \texttt{X} \\
\texttt{:=} &: \quad \forall \texttt{X}.\textsf{ref}\,\texttt{X} \to \texttt{X} \to \texttt{X}
\end{aligned}$$

Check that these definitions meet the requirements of Definition 1.7.6.    □

1.9.6    EXERCISE [RECURSION, RECOMMENDED, ★★★]: The fixpoint combinator
`fix` has been introduced and given an operational semantics in Exam-
ple 1.2.10. Let us now extend the initial environment $\Gamma_0$ with the following
binding:

$$\texttt{fix} : \quad \forall \texttt{XY}.((\texttt{X} \to \texttt{Y}) \to (\texttt{X} \to \texttt{Y})) \to \texttt{X} \to \texttt{Y}$$

Check that these definitions meet the requirements of Definition 1.7.6. Recall
how the `letrec` syntactic sugar was defined in Example 1.2.10, and check that
this gives rise to the following constraint generation rule:

$$\begin{aligned}
&\textsf{let } \Gamma_0 \textsf{ in } [\![\texttt{letrec f} = \lambda\texttt{z}.\texttt{t}_1 \textsf{ in } \texttt{t}_2 : \texttt{T}]\!] \\
\equiv \quad &\textsf{let } \Gamma_0 \textsf{ in let f} : \forall\texttt{XY}[\textsf{let f} : \texttt{X} \to \texttt{Y}; \texttt{z} : \texttt{X} \textsf{ in } [\![\texttt{t}_1 : \texttt{Y}]\!]].\texttt{X} \to \texttt{Y} \textsf{ in } [\![\texttt{t}_2 : \texttt{T}]\!]
\end{aligned}$$

Note the somewhat peculiar structure of this constraint: the program variable
`f` is bound twice in it, with different type schemes. The constraint requires
all occurrences of `f` within $\texttt{t}_1$ to be assigned the *monomorphic* type $\texttt{X} \to \texttt{Y}$.
This type is generalized and turned into a type scheme before inspecting $\texttt{t}_2$,
however, so every occurrence of `f` within $\texttt{t}_2$ may receive a different type, as
usual with `let`-polymorphism. A more powerful way of typechecking recursive
function definitions is discussed in Section 1.10 (page 113).    □

### Algebraic data types

Exercises 1.9.3 and 1.9.4 have shown how to extend the language with binary, anonymous products and sums. These constructs are quite general, but still have several shortcomings. First, they are only binary, while we would like to have $k$-ary products and sums, for arbitrary $k \geq 0$. Such a generalization is of course straightforward. Second, more interestingly, their components must be referred to by numeric index (as in "please extract the *second* component of the pair"), rather than by name ("extract the component named y"). In practice, it is crucial to use names, because they make programs more readable and more robust in the face of changes. One could introduce a mechanism that allows defining names as syntactic sugar for numeric indices. That would help a little, but not much, because these names would not appear in *types*, which would still be made of anonymous products and sums. Third, in the absence of recursive types, products and sums do not have sufficient expressiveness to allow defining unbounded data structures, such as lists. Indeed, it is easy to see that every value whose type `T` is composed of base types (`int`, `bool`, *etc.*), products, and sums must have bounded size, where the bound $|\,T\,|$ is a function of `T`. More precisely, up to a constant factor, we have $|\,\texttt{int}\,| = |\,\texttt{bool}\,| = 1$, $|\,T_1 \times T_2\,| = 1 + |\,T_1\,| + |\,T_2\,|$, and $|\,T_1 + T_2\,| = 1 + \max(|\,T_1\,|, |\,T_2\,|)$. The following example describes another facet of the same problem.

1.9.7   EXAMPLE:  A list is either empty, or a pair of an element and another list. So, it seems natural to try and encode the type of lists as a sum of some arbitrary type (say, `unit`), on the one hand, and of a product of some element type and of the type of lists itself, on the other hand. With this encoding in mind, we can go ahead and write code—for instance, a function that computes the length of a list:

$$\texttt{letrec length} = \lambda \texttt{l.case l } (\lambda\_.\hat{0})\ (\lambda \texttt{z.}\hat{1} \mathbin{\hat{+}} \texttt{length } (\pi_2\ \texttt{z}))$$

We have used integers, pairs, sums, and the `letrec` construct introduced in the previous section. The code analyzes the list `l` using a `case` construct. If the left branch is taken, the list is empty, so 0 is returned. If the right branch is taken, then `z` becomes bound to a pair of some element and the tail of the list. The latter is obtained using the projection operator $\pi_2$. Its length is computed using a recursive call to `length` and incremented by 1. This code makes perfect sense. However, applying the constraint generation and constraint solving algorithms eventually leads to an equation of the form $\texttt{X} = \texttt{Y} + (\texttt{Z} \times \texttt{X})$, where `X` stands for the type of `l`. This equation accurately reflects our encoding of the type of lists. However, in a syntactic model, it has no solution, so our definition of `length` is ill-typed. It is possible to adopt a free

regular tree model,thus introducing *equirecursive* types into the system (TAPL Chapter 20); however, there are good reasons not to do so (page 106).      □

To work around this problem, ML-the-programming-language offers *algebraic data type* definitions, whose elegance lies in the fact that, while representing only a modest theoretical extension, they do solve the three problems mentioned above. An algebraic data type may be viewed as an *abstract type* that is declared to be *isomorphic* to a ($k$-ary) product or sum type with named components. The type of each component is declared as well, and may refer to the algebraic data type that is being defined: thus, algebraic data types are *isorecursive* (TAPL Chapter 20). In order to allow sufficient flexibility when declaring the type of each component, algebraic data type definitions may be *parameterized* by a number of type variables. Last, in order to allow the description of complex data structures, it is necessary to allow several algebraic data types to be defined at once; the definitions may then be *mutually recursive*. In fact, in order to simplify this formal presentation, we assume that *all* algebraic data types are defined at once at the beginning of the program. This decision is of course at odds with modular programming, but will not otherwise be a problem.

In the following, $D$ ranges over a set of *data types*. We assume that data types form a subset of type constructors. We require each of them to be isolated and to have a signature of the form $\vec{\kappa} \Rightarrow \star$. Furthermore, $\ell$ ranges over a set $\mathcal{L}$ of *labels*, which we use indifferently as *data constructors* and as *record labels*. An *algebraic data type definition* is either a *variant type* definition or a *record type* definition, whose respective forms are

$$D\,\vec{X} \approx \sum_{i=1}^{k} \ell_i : T_i \qquad \text{and} \qquad D\,\vec{X} \approx \prod_{i=1}^{k} \ell_i : T_i.$$

In either case, $k$ must be nonnegative. If $D$ has signature $\vec{\kappa} \Rightarrow \star$, then the type variables $\vec{X}$ must have kind $\vec{\kappa}$. Every $T_i$ must have kind $\star$. We refer to $\bar{X}$ as the *parameters* and to $\vec{T}$ (the vector formed by $T_1, \dots, T_k$) as the *components* of the definition. The parameters are bound within the components, and the definition must be closed, that is, $ftv(\vec{T}) \subseteq \bar{X}$ must hold. Last, for an algebraic data type definition to be valid, the behavior of the type constructor $D$ with respect to subtyping must match its definition. This requirement is clarified below.

1.9.8   DEFINITION: Consider an algebraic data type definition whose parameters and components are respectively $\vec{X}$ and $\vec{T}$. Let $\vec{X}'$ and $\vec{T}'$ be their images under an arbitrary renaming. Then, $D\,\vec{X} \leq D\,\vec{X}' \Vdash \vec{T} \leq \vec{T}'$ must hold.      □

The above requirement bears on the definition of subtyping in the model. The idea is, since $D\,\vec{X}$ is declared to be isomorphic to (a sum or a product of)

$\vec{\text{T}}$, whenever two types built with D are comparable, their unfoldings should be comparable as well. The reverse entailment assertion is not required for type soundness, and it is sometimes useful to declare algebraic data types that do not validate it—so-called *phantom types* (Fluet and Pucella, 2002). Note that the requirement may always be satisfied by making the type constructor D *invariant* in all of its parameters. Indeed, in that case, $\text{D}\,\vec{\text{X}} \leq \text{D}\,\vec{\text{X}}'$ entails $\vec{\text{X}} = \vec{\text{X}}'$, which must entail $\vec{\text{T}} = \vec{\text{T}}'$ since $\vec{\text{T}}'$ is precisely $[\vec{\text{X}} \mapsto \vec{\text{X}}']\vec{\text{T}}$. In an equality free tree model, every type constructor is naturally invariant, so the requirement is trivially satisfied. In other settings, however, it is often possible to satisfy the requirement of Definition 1.9.8 while assigning D a less restrictive variance. The following example illustrates such a case.

1.9.9 EXAMPLE: Let list be a data type of signature $\star \Rightarrow \star$. Let Nil and Cons be data constructors. Then, the following is a definition of list as a variant type:

$$\text{list X} \approx \Sigma\,(\texttt{Nil} : \texttt{unit}; \texttt{Cons} : \text{X} \times \text{list X})$$

Because data types form a subset of type constructors, it is valid to form the type list X in the right-hand side of the definition, even though we are still in the process of defining the meaning of list. In other words, data type definitions may be recursive. However, because $\approx$ is not interpreted as equality, the type list X is *not* a recursive type: it is nothing but an application of the unary type constructor list to the type variable X. To check that the definition of list satisfies the requirement of Definition 1.9.8, we must ensure that

$$\text{list X} \leq \text{list X}' \Vdash \texttt{unit} \leq \texttt{unit} \wedge \text{X} \times \text{list X} \leq \text{X}' \times \text{list X}'$$

holds. This assertion is equivalent to $\text{list X} \leq \text{list X}' \Vdash \text{X} \leq \text{X}'$. To satisfy the requirement, it is sufficient to make list a *covariant* type constructor, that is, to define subtyping in the model so that $\text{list X} \leq \text{list X}' \equiv \text{X} \leq \text{X}'$ holds.

Let tree be a data type of signature $\star \Rightarrow \star$. Let root and sons be record labels. Then, the following is a definition of tree as a record type:

$$\text{tree X} \approx \Pi\,(\texttt{root} : \text{X}; \texttt{sons} : \text{list}\,(\text{tree X}))$$

This definition is again recursive, and relies on the previous definition. Because list is covariant, it is straightforward to check that the definition of tree is valid if tree is made a covariant type constructor as well. □

1.9.10 EXERCISE [★★, ↛]: Consider a nonrecursive algebraic data type definition, where the variance of every type constructor that appears on the right-hand side is known. Can you systematically determine, for each of the parameters, the least restrictive variance that makes the definition valid? Generalize this procedure to the case of recursive and mutually recursive algebraic data type definitions. □

A *prologue* is a set of algebraic data type definitions, where each data type is defined at most once and where each data constructor or record label appears at most once. A *program* is a pair of a prologue and an expression. The effect of a prologue is to enrich the programming language with new constants. That is, a variant type definition extends the operational semantics with several injections and a `case` construct, as in Example 1.2.7. A record type definition extends it with a record formation construct and several projections, as in Examples 1.2.3 and 1.2.5. In either case, the initial typing environment $\Gamma_0$ is extended with information about these new constants. Thus, algebraic data type definitions might be viewed as a simple configuration language that allows specifying in which instance of ML-the-calculus the expression that follows the prologue should be typechecked and interpreted. Let us now give a precise account of this phenomenon.

To begin, suppose the prologue contains the definition $\mathsf{D}\,\vec{\mathsf{X}} \approx \sum_{i=1}^{k} \ell_i : \mathsf{T}_i$. Then, for each $i \in \{1, \ldots, k\}$, a constructor of arity 1, named $\ell_i$, is introduced. Furthermore, a destructor of arity $k+1$, named $\mathsf{case_D}$, is introduced. When $k > 0$, it is common to write $\mathsf{case}\ \mathsf{t}\ [\ell_i : \mathsf{t}_i]_{i=1}^{k}$ for the application $\mathsf{case_D}\ \mathsf{t}\ \mathsf{t}_1\ \ldots\ \mathsf{t}_n$. The operational semantics is extended with the following reduction rules, for $i \in \{1, \ldots, k\}$:

$$\mathsf{case}\ (\ell_i\ \mathsf{v})\ [\ell_j : \mathsf{v}_j]_{j=1}^{k} \xrightarrow{\delta} \mathsf{v}_i\ \mathsf{v} \qquad\qquad \text{(R-ALG-CASE)}$$

For each $i \in \{1, \ldots, k\}$, the initial environment is extended with the binding $\ell_i : \forall \vec{\mathsf{X}}.\mathsf{T}_i \rightarrow \mathsf{D}\,\vec{\mathsf{X}}$. It is further extended with the binding $\mathsf{case_D} : \forall \vec{\mathsf{X}}\mathsf{Z}.\mathsf{D}\,\vec{\mathsf{X}} \rightarrow (\mathsf{T}_1 \rightarrow \mathsf{Z}) \rightarrow \ldots (\mathsf{T}_k \rightarrow \mathsf{Z}) \rightarrow \mathsf{Z}$.

Now, suppose the prologue contains the definition $\mathsf{D}\,\vec{\mathsf{X}} \approx \prod_{i=1}^{k} \ell_i : \mathsf{T}_i$. Then, for each $i \in \{1, \ldots, k\}$, a destructor of arity 1, named $\ell_i$, is introduced. Furthermore, a constructor of arity $k$, named $\mathsf{make_D}$, is introduced. It is common to write $\mathsf{t}.\ell$ for the application $\ell\ \mathsf{t}$ and, when $k > 0$, to write $\{\ell_i = \mathsf{t}_i\}_{i=1}^{k}$ for the application $\mathsf{make_D}\ \mathsf{t}_1\ \ldots\ \mathsf{t}_k$. The operational semantics is extended with the following reduction rules, for $i \in \{1, \ldots, k\}$:

$$(\{\ell_j = \mathsf{v}_j\}_{j=1}^{k}).\ell_i \xrightarrow{\delta} \mathsf{v}_i \qquad\qquad \text{(R-ALG-PROJ)}$$

For each $i \in \{1, \ldots, k\}$, the initial environment is extended with the binding $\ell_i : \forall \vec{\mathsf{X}}.\mathsf{D}\,\vec{\mathsf{X}} \rightarrow \mathsf{T}_i$. It is further extended with the binding $\mathsf{make_D} : \forall \vec{\mathsf{X}}.\mathsf{T}_1 \rightarrow \ldots \rightarrow \mathsf{T}_k \rightarrow \mathsf{D}\,\vec{\mathsf{X}}$.

1.9.11   EXAMPLE:  The effect of defining $\mathsf{list}$ (Example 1.9.9) is to make `Nil` and `Cons` data constructors of arity 1 and to introduce a binary destructor $\mathsf{case_{list}}$. The definition also extends the initial environment as follows:

$$
\begin{aligned}
\mathtt{Nil} : &\quad \forall \mathsf{X}.\mathtt{unit} \rightarrow \mathsf{list}\ \mathsf{X} \\
\mathtt{Cons} : &\quad \forall \mathsf{X}.\mathsf{X} \times \mathsf{list}\ \mathsf{X} \rightarrow \mathsf{list}\ \mathsf{X} \\
\mathsf{case_{list}} : &\quad \forall \mathsf{XZ}.\mathsf{list}\ \mathsf{X} \rightarrow (\mathtt{unit} \rightarrow \mathsf{Z}) \rightarrow (\mathsf{X} \times \mathsf{list}\ \mathsf{X} \rightarrow \mathsf{Z}) \rightarrow \mathsf{Z}
\end{aligned}
$$

Thus, the value $\text{Cons}(\hat{0}, \text{Nil}())$, an integer list of length 1, has type $\text{list int}$. A function that computes the length of a list may now be written as follows:

$$\text{letrec length} = \lambda\text{l.case l}\,[\,\text{Nil} : \lambda\_.\hat{0}\mid\text{Cons} : \lambda\text{z}.\hat{1}\mathbin{\hat{+}}\text{length}\,(\pi_2\,\text{z})\,]$$

Recall that this notation is syntactic sugar for

$$\text{letrec length} = \lambda\text{l.case}_{\text{list}}\,\text{l}\,(\lambda\_.\hat{0})\,(\lambda\text{z}.\hat{1}\mathbin{\hat{+}}\text{length}\,(\pi_2\,\text{z}))$$

The difference with the code in Example 1.9.7 appears minimal: the $\text{case}$ construct is now annotated with the data type $\text{list}$. As a result, the type inference algorithm employs the type scheme assigned to $\text{case}_{\text{list}}$, which is derived from the definition of $\text{list}$, instead of the type scheme assigned to the anonymous $\text{case}$ construct, given in Exercise 1.9.4. This is good for a couple of reasons. First, the former is more informative than the latter, because it contains the type $\text{T}_i$ associated with the data constructor $\ell_i$. Here, for instance, the generated constraint requires the type of $\text{z}$ to be $\text{X}\times\text{list X}$ for some $\text{X}$, so a good error message would be given if a mistake was made in the second branch, such as omitting the use of $\pi_2$. Second, and more fundamentally, *the code is now well-typed*, even in the absence of recursive types. In Example 1.9.7, a cyclic equation was produced because $\text{case}$ required the type of $\text{l}$ to be a sum type and because a sum type carries the types of its left and right branches as subterms. Here, instead, $\text{case}_{\text{list}}$ requires $\text{l}$ to have type $\text{list X}$ for some $\text{X}$. This is an abstract type: it does not explicitly contain the types of the branches. As a result, the generated constraint no longer involves a cyclic equation. It is, in fact, satisfiable; the reader may check that $\text{length}$ has type $\forall\text{X.list X}\to\text{int}$, as expected.                                                                                   □

Example 1.9.11 stresses the importance of using *declared, abstract* types, as opposed to *anonymous, concrete* sum or product types, in order to obviate the need for recursive types. The essence of the trick lies in the fact that the type schemes associated with operations on algebraic data types implicitly *fold* and *unfold* the data type's definition. More precisely, let us recall the type scheme assigned to the $i^{\text{th}}$ injection in the setting of ($k$-ary) anonymous sums: it is $\forall\text{X}_1\ldots\text{X}_k.\text{X}_i\to\text{X}_1+\ldots+\text{X}_k$, or, more concisely, $\forall\text{X}_1\ldots\text{X}_k.\text{X}_i\to\sum_{i=1}^k\text{X}_i$. By instantiating each $\text{X}_i$ with $\text{T}_i$ and generalizing again, we find that a more specific type scheme is $\forall\bar{\text{X}}.\text{T}_i\to\sum_{i=1}^k\text{T}_i$. Perhaps this could have been the type scheme assigned to $\ell_i$? Instead, however, it is $\forall\bar{\text{X}}.\text{T}_i\to\text{D}\,\vec{\text{X}}$. We now realize that this type scheme not only reflects the operational behavior of the $i^{\text{th}}$ injection, but also *folds* the definition of the algebraic data type $\text{D}$ by turning the anonymous sum $\sum_{i=1}^k\text{T}_i$—which forms the definition's right-hand side—into the parameterized abstract type $\text{D}\,\vec{\text{X}}$—which is the definition's left-hand side. Conversely, the type scheme assigned to $\text{case}_{\text{D}}$ *unfolds* the definition. The

situation is identical in the case of record types: in either case, *constructors fold, destructors unfold.* In other words, occurrences of data constructors and record labels in the code may be viewed as explicit instructions for the type-checker to fold or unfold an algebraic data type definition. This mechanism is characteristic of *isorecursive* types.

1.9.12   EXERCISE [★, ↛]: For a fixed $k$, check that all of the machinery associated with $k$-ary anonymous products—that is, constructors, destructors, reduction rules, and extensions to the initial typing environment—may be viewed as the result of a single algebraic data type definition. Conduct a similar check in the case of $k$-ary anonymous sums.                                                    □

1.9.13   EXERCISE [★★★, ↛]: Check that the above definitions meet the requirements of Definition 1.7.6.                                                    □

1.9.14   EXERCISE [★★★, ↛]: For sake of simplicity, we have assumed that data constructors are always of arity one. It is indeed possible to allow data constructors of any arity and define variants as $D\,\vec{X} \approx \sum_{i=1}^{k} \ell_i : \vec{T}_i$. For instance, the definition of $\mathsf{list}$ could then be $\mathsf{list\,X} \approx \Sigma\,(\mathtt{Nil}; \mathtt{Cons} : \mathtt{X} \times \mathsf{list\,X})$ and for instance $\mathtt{Cons}(\hat{0}, \mathtt{Nil})$ would be a list value. Make the necessary changes in the definitions above and check that they still meet the requirements of Definition 1.7.6.                                                    □

In this formal presentation of algebraic data types, we have assumed that all algebraic data type definitions are known before the program is typechecked. This simplifying assumption is forced on us by the fact that we interpret constraints in a fixed model, that is, we assume a fixed universe of types. In practice, programming languages have *module systems*, which allow distinct modules to have distinct, partial views of the universe of types. Then, it becomes possible for each module to come with its own data type definitions. Interestingly, it is even possible, in principle, to split the definition of a single data type over several modules, yielding *extensible* algebraic data types. For instance, module $A$ might declare the existence of a parameterized variant type $D\,\vec{X}$, without giving its components. Later on, module $B$ might define a component $\ell : \mathtt{T}$, where $ftv\,(\mathtt{T}) \subseteq \bar{\mathtt{X}}$. Such a definition makes $\ell$ a unary constructor with type scheme $\forall\bar{\mathtt{X}}.\mathtt{T} \to D\,\vec{X}$, as before. It becomes impossible, however, to introduce a destructor $\mathtt{case}_D$, because the definition of an extensible variant type can never be assumed to be complete—other, unknown modules might extend it further. To compensate for its absence, one may supplement every constructor $\ell$ with a destructor $\ell^{-1}$, whose semantics is given by $\ell^{-1}\,(\ell\,\mathtt{v})\,\mathtt{v}_1\,\mathtt{v}_2 \xrightarrow{\delta} \mathtt{v}_1\,\mathtt{v}$ and $\ell^{-1}\,(\ell'\,\mathtt{v})\,\mathtt{v}_1\,\mathtt{v}_2 \xrightarrow{\delta} \mathtt{v}_2\,(\ell'\,\mathtt{v})$ when $\ell \neq \ell'$, and whose type scheme is $\forall\bar{\mathtt{X}}\mathtt{Z}.D\,\vec{X} \to (\mathtt{T} \to \mathtt{Z}) \to (D\,\vec{X} \to \mathtt{Z}) \to \mathtt{Z}$. When

pattern matching is available, $\ell^{-1}$ may in fact be defined in the language. ML-the-programming-language does not offer extensible algebraic data types as a language feature, but does have one built-in extensible variant type, namely the type `exn` of exceptions. Thus, it is possible to define new constructors for the type `exn` within any module. The price of this extra flexibility is that no exhaustive case analysis on values of type `exn` is possible.

One significant drawback of algebraic data type definitions resides in the fact that a label $\ell$ cannot be *shared* by two distinct variant or record type definitions. Indeed, every algebraic data type definition extends the calculus with new constants. Strictly speaking, our presentation does not allow a single constant `c` to be associated with two distinct definitions. Even if we did allow such a collision, the initial environment would contain two bindings for `c`, one of which would then become inaccessible. This phenomenon arises in actual implementations of ML-the-programming-language, where a new algebraic data type definition may hide some of the data constructors or record labels introduced by a previous definition. An elegant solution to this lack of expressiveness is discussed in Section 1.11.

### Pattern matching

Our presentation of products, sums and algebraic data types has remained within the setting of ML-the-calculus: that is, data structures have been built out of constructors, while the case analysis and record access operations have been viewed as destructors. Some syntactic sugar has been used to recover standard notations. The language is now expressive enough to allow defining and manipulating complex data structures, such as lists and trees. Yet, experience shows that programming in such a language is still somewhat cumbersome. Indeed, case analysis and record access are low-level operations: the former allows inspecting a tag and branching, while the latter allows dereferencing a pointer. In practice, one often needs to carry out more complex tasks, such as determining whether a data structure has a certain shape or whether two data structures have comparable shapes. Currently, the only way to carry out these tasks is to program an explicit sequence of low-level operations. It would be much preferable to extend the language so that it becomes directly possible to describe shapes, called *patterns*, and so that checking whether a patterns *matches* a value becomes an elementary operation. ML-the-programming-language offers this feature, called *pattern matching*. Although pattern matching may be added to ML-the-calculus by introducing a family of destructors, we rather choose to extend the calculus with a new `match` construct, which subsumes the existing `let` construct. This approach appears somewhat simpler and more powerful. We now carry out this

$$
\begin{array}{lll}
\text{p} & ::= & & \textit{Patterns:} \\
& & \_ & \textit{Wildcard} \\
& & \text{z} & \textit{Variable} \\
& & \text{c p}_1 \ \cdots \ \text{p}_k & \textit{Data} \\
& & \qquad \text{c} \in \mathcal{Q}^+ \wedge k = a(\text{c}) & \\
& & \text{p} \wedge \text{p} & \textit{Conjunction} \\
& & \text{p} \vee \text{p} & \textit{Disjunction}
\end{array}
$$

*Pattern matching*

$$[\_ \mapsto \text{v}] = \varnothing$$

$$\big[\text{c p}_1 \ \cdots \ \text{p}_k \mapsto \text{c v}_1 \ \cdots \ \text{v}_k\big]$$
$$= [\text{p}_1 \mapsto \text{v}_1] \otimes \ldots \otimes [\text{p}_k \mapsto \text{v}_k]$$
$$[\text{p}_1 \wedge \text{p}_2 \mapsto \text{v}] = [\text{p}_1 \mapsto \text{v}] \otimes [\text{p}_2 \mapsto \text{v}]$$
$$[\text{p}_1 \vee \text{p}_2 \mapsto \text{v}] = [\text{p}_1 \mapsto \text{v}] \oplus [\text{p}_2 \mapsto \text{v}]$$

**Figure 1-13: Patterns and pattern matching**

extension.

Let us first define the syntax of patterns (Figure 1-13) and describe (informally, for now) which values they match. To a pattern $\text{p}$, we associate a set of *defined program variables* $dpi(\text{p})$, whose definition appears in the text that follows. The pattern $\text{p}$ is well-formed if and only if $dpi(\text{p})$ is defined. To begin, the wildcard $\_$ is a pattern, which matches every value and binds no variables. We let $dpi(\_) = \varnothing$. Although the wildcard may be viewed as an anonymous variable, and we have done so thus far, it is now simpler to view it as a distinct pattern. A program variable $\text{z}$ is also a pattern, which matches every value and binds $\text{z}$ to the matched value. We let $dpi(\text{z}) = \{\text{z}\}$. Next, if $\text{c}$ is a constructor of arity $k$, then $\text{c p}_1 \ \ldots \ \text{p}_k$ is a pattern, which matches $\text{c v}_1 \ \ldots \ \text{v}_k$ when $\text{p}_i$ matches $\text{v}_i$ for every $i \in \{1, \ldots, k\}$. We let $dpi(\text{c p}_1 \ \ldots \ \text{p}_k) = dpi(\text{p}_1) \uplus \ldots \uplus dpi(\text{p}_k)$. That is, the pattern $\text{c p}_1 \ \ldots \ \text{p}_k$ is well-formed when $\text{p}_1, \ldots, \text{p}_k$ define *disjoint* sets of variables. This condition rules out *nonlinear* patterns such as $(\text{z}, \text{z})$. Defining the semantics of such a pattern would require a notion of equality at every type, which introduces various complications, so it is commonly considered ill-formed. The pattern $\text{p}_1 \wedge \text{p}_2$ matches all values that both $\text{p}_1$ and $\text{p}_2$ match. It is commonly used with $\text{p}_2$ a program variable: then, it allows examining the shape of a value and binding a name to it at the same time. Again, we define $dpi(\text{p}_1 \wedge \text{p}_2) = dpi(\text{p}_1) \uplus dpi(\text{p}_2)$. The pattern $\text{p}_1 \vee \text{p}_2$ matches all values that either $\text{p}_1$ or $\text{p}_2$ matches. We define $dpi(\text{p}_1 \vee \text{p}_2) = dpi(\text{p}_1) = dpi(\text{p}_2)$. That is, the pattern $\text{p}_1 \vee \text{p}_2$ is well-formed when $\text{p}_1$ and $\text{p}_2$ define the *same* variables. Thus, $(\text{inj}_1 \ \text{z}) \vee (\text{inj}_2 \ \text{z})$ is a well-formed pattern, which binds $\text{z}$ to the component of a binary sum, without regard for its tag. However, $(\text{inj}_1 \ \text{z}_1) \vee (\text{inj}_2 \ \text{z}_2)$ is ill-formed, because one cannot statically predict whether it defines $\text{z}_1$ or $\text{z}_2$.

Let us now formally define whether a pattern $\text{p}$ matches a value $\text{v}$ and how the variables in $dpi(\text{p})$ become bound to values in the process. This is done by introducing a *generalized substitution*, written $[\text{p} \mapsto \text{v}]$, which is either

$$
\begin{array}{ll}
\texttt{t} & ::= \quad \dots \\
& \quad \texttt{match t with } (\texttt{p}_i \,.\, \texttt{t}_i)_{i=1}^{k} \\
\mathcal{E} & ::= \quad \dots \\
& \quad \texttt{match } \mathcal{E} \texttt{ with } (\texttt{p}_i \,.\, \texttt{t}_i)_{i=1}^{k}
\end{array}
$$

*Expressions:*

*Evaluation Contexts:*

*Reduction rules*

$$
\texttt{match v with } (\texttt{p}_i \,.\, \texttt{t}_i)_{i=1}^{k} \longrightarrow \bigoplus_{i=1}^{k} [\texttt{p}_i \mapsto \texttt{v}] \texttt{t}_i
$$

(R-Match)

**Figure 1-14: Extended syntax and semantics of ML-the-calculus**

undefined or a substitution of values for the program variables in $dpi(\texttt{p})$. If the former, then $\texttt{p}$ does not match $\texttt{v}$. If the latter, then $\texttt{p}$ matches $\texttt{v}$ and, for every $\texttt{z} \in dpi(\texttt{p})$, the variable $\texttt{z}$ becomes bound to the value $[\texttt{p} \mapsto \texttt{v}]\texttt{z}$. Of course, when $\texttt{p}$ is a variable $\texttt{z}$, the generalized substitution $[\texttt{z} \mapsto \texttt{v}]$ is defined and coincides with the substitution $[\texttt{z} \mapsto \texttt{v}]$, which justifies our abuse of notation. To construct generalized substitutions, we use two simple combinators. First, when $dpi(\texttt{p}_1)$ and $dpi(\texttt{p}_2)$ are disjoint, $[\texttt{p}_1 \mapsto \texttt{v}_1] \otimes [\texttt{p}_2 \mapsto \texttt{v}_2]$ stands for the set-theoretic union of $[\texttt{p}_1 \mapsto \texttt{v}_1]$ and $[\texttt{p}_2 \mapsto \texttt{v}_2]$, if both are defined, and is undefined otherwise. We use this combinator to ensure that $\texttt{p}_1$ matches $\texttt{v}_1$ and $\texttt{p}_2$ matches $\texttt{v}_2$ and to combine the two corresponding sets of bindings. Second, when $o_1$ and $o_2$ are two possibly undefined mathematical objects that belong to the same space when defined, $o_1 \oplus o_2$ stands for $o_1$, if it is defined, and for $o_2$ otherwise—that is, $\oplus$ is an angelic choice operator with a left bias. In particular, when $dpi(\texttt{p}_1)$ and $dpi(\texttt{p}_2)$ coincide, $[\texttt{p}_1 \mapsto \texttt{v}_1] \oplus [\texttt{p}_2 \mapsto \texttt{v}_2]$ stands for $[\texttt{p}_1 \mapsto \texttt{v}_1]$, if it is defined, and for $[\texttt{p}_2 \mapsto \texttt{v}_2]$ otherwise. We use this combinator to ensure that $\texttt{p}_1$ matches $\texttt{v}_1$ or $\texttt{p}_2$ matches $\texttt{v}_2$ and to retain the corresponding set of bindings. The full definition of generalized substitutions, which relies on these combinators, appears in Figure 1-13. It reflects the informal presentation of the previous paragraph.

Once patterns and pattern matching are defined, it is straightforward to extend the syntax and operational semantics of ML-the-calculus. We enrich the syntax of expressions with a new construct, $\texttt{match t with } (\texttt{p}_i \,.\, \texttt{t}_i)_{i=1}^{k}$, where $k \geq 1$. It consists of a term $\texttt{t}$ and a nonempty, ordered list of clauses, each of which is composed of a pattern $\texttt{p}_i$ and a term $\texttt{t}_i$. The syntax of evaluation contexts is extended as well, so that the term $\texttt{t}$ that is being examined is first reduced to a value $\texttt{v}$. The operational semantics is extended with a new rule, R-Match, which states that $\texttt{match v with } (\texttt{p}_i \,.\, \texttt{t}_i)_{i=1}^{k}$ reduces to $[\texttt{p}_i \mapsto \texttt{v}]\texttt{t}_i$, where $i$ is the least element of $\{1, \dots, k\}$ such that $\texttt{p}_i$ matches $\texttt{v}_i$. Technically, $\bigoplus_{i=1}^{k} [\texttt{p}_i \mapsto \texttt{v}]\texttt{t}_i$ stands for $[\texttt{p}_1 \mapsto \texttt{v}]\texttt{t}_1 \oplus \dots \oplus [\texttt{p}_k \mapsto \texttt{v}]\texttt{t}_k$, so that the reduct is the first term that is defined in this sequence.

As far as semantics is concerned, the $\texttt{match}$ construct may be viewed as a

generalization of the `let` construct. Indeed, `let z = t₁ in t₂` may now be viewed as syntactic sugar for `match t₁ with z . t₂`, that is, a `match` construct with a single clause and a variable pattern. Then, R-LET becomes a special case of R-MATCH.

It is pleasant to introduce some more syntactic sugar. We write $\lambda(\mathbf{p}_i.\mathbf{t}_i)_{i=1}^k$ for `λz.match z with` $(\mathbf{p}_i \ . \ \mathbf{t}_i)_{i=1}^k$, where `z` is fresh for $(\mathbf{p}_i.\mathbf{t}_i)_{i=1}^k$. Thus, it becomes possible to define functions by cases—a common idiom in ML-the-programming-language.

1.9.15 EXAMPLE: Using pattern matching, a function that computes the length of a list (Example 1.9.11) may now be written as follows:

$$\texttt{letrec length} = \lambda(\texttt{Nil \_} . \ \hat{0} \mid \texttt{Cons (\_,z)} . \ \hat{1} \ \hat{+} \ \texttt{length z})$$

The second pattern matches a nonempty list and binds `z` to its tail at the same time, obviating the need for an explicit application of $\pi_2$. □

1.9.16 EXERCISE [★★, RECOMMENDED, ↠]: Under the above definition of `length`, consider an application of `length` to the list `Cons(0̂, Nil())`. After eliminating the syntactic sugar, determine by which reduction sequence this expression reduces to a value. □

Before we can proceed and extend the type system to deal with the new `match` construct, we must make two mild extensions to the syntax and meaning of constraints. First, if $\sigma$ is $\forall \bar{\mathtt{X}}[C].\mathtt{T}$, where $\bar{\mathtt{X}} \ \# \ ftv(\mathtt{T}')$, then $\mathtt{T}' \preceq \sigma$ stands for the constraint $\exists \bar{\mathtt{X}}.(C \wedge \mathtt{T}' \leq \mathtt{T})$. This relation is identical to the instance relation (Definition 1.3.3), except the direction of subtyping is reversed. We extend the syntax of constraints with instantiation constraints of the form $\mathtt{T} \preceq \mathtt{x}$ and define their meaning by adding a symmetric counterpart of CM-INSTANCE. We remark that, when subtyping is interpreted as equality, the relations $\sigma \preceq \mathtt{T}$ and $\mathtt{T} \preceq \sigma$ coincide, so this extension is unnecessary in that particular case. Second, we extend the syntax of environments so that several successive bindings may *share* a set of quantifiers and a constraint. That is, we allow writing $\forall \bar{\mathtt{X}}[C].(\mathtt{x}_1 : \mathtt{T}_1; \ldots; \mathtt{x}_k : \mathtt{T}_k)$ for $\mathtt{x}_1 : \forall \bar{\mathtt{X}}[C].\mathtt{T}_1; \ldots; \mathtt{x}_k : \forall \bar{\mathtt{X}}[C].\mathtt{T}_k$. From a theoretical standpoint, this is little more than syntactic sugar; however, in practice, it is useful to implement this new idiom literally, since it avoids unnecessary copying of the constraint $C$.

Let us now extend the type system. For the sake of brevity, we extend the constraint generation rules only. Of course, it would also be possible to define corresponding extensions of the rule-based type systems shown earlier, namely DM, HM($X$), and PCB($X$). We begin by defining a constraint $[\![\mathtt{T} : \mathtt{p}]\!]$ that represents a necessary and sufficient condition for values of type $\mathtt{T}$ to be acceptable inputs for the pattern $\mathtt{p}$. Its free type variables are a subset of

$$\begin{aligned}
[\![\mathtt{T} : \_]\!] &= \mathsf{true} \\
[\![\mathtt{T} : \mathtt{z}]\!] &= \mathtt{T} \preceq \mathtt{z} \\
[\![\mathtt{T} : \mathtt{c}\ \mathtt{p}_1\ \ldots\ \mathtt{p}_k]\!] &= \exists \bar{\mathtt{X}}.(\vec{\mathtt{X}} \to \mathtt{T} \preceq \mathtt{c} \wedge \textstyle\bigwedge_{i=1}^{k} [\![\mathtt{X}_i : \mathtt{p}_i]\!]) \\
[\![\mathtt{T} : \mathtt{p}_1 \wedge \mathtt{p}_2]\!] &= [\![\mathtt{T} : \mathtt{p}_1]\!] \wedge [\![\mathtt{T} : \mathtt{p}_2]\!] \\
[\![\mathtt{T} : \mathtt{p}_1 \vee \mathtt{p}_2]\!] &= [\![\mathtt{T} : \mathtt{p}_1]\!] \wedge [\![\mathtt{T} : \mathtt{p}_2]\!]
\end{aligned}$$

$$[\![\mathtt{match\ t\ with}\ (\mathtt{p}_i\ .\ \mathtt{t}_i)_{i=1}^{k} : \mathtt{T}]\!] = \bigwedge_{i=1}^{k} \mathsf{let}\ \forall \mathtt{X}\bar{\mathtt{X}}_i [[\![\mathtt{t} : \mathtt{X}]\!] \wedge \mathsf{let}\ \vec{\mathtt{z}}_i : \vec{\mathtt{X}}_i\ \mathsf{in}\ [\![\mathtt{X} : \mathtt{p}_i]\!]].(\vec{\mathtt{z}}_i : \vec{\mathtt{X}}_i)\ \mathsf{in}\ [\![\mathtt{t}_i : \mathtt{T}]\!]$$
$$\text{where } \vec{\mathtt{z}}_i = dpi(\mathtt{p}_i)$$

**Figure 1-15: Constraint generation for patterns and pattern matching**

$ftv(\mathtt{T})$, while its free program identifiers are either constructors or program variables bound by $\mathtt{p}$. It is defined in the upper part of Figure 1-15. The first rule states that a wildcard matches values of arbitrary type. The second and third rules govern program variables and constructor applications in *patterns*. They are identical to the rules that govern these constructs in *expressions* (page 59), except that the direction of subtyping is reversed. In the absence of subtyping, they would be entirely identical. We write $\vec{\mathtt{X}}$ for $\mathtt{X}_1 \ldots \mathtt{X}_k$ and $\vec{\mathtt{X}} \to \mathtt{T}$ for $\mathtt{X}_1 \to \ldots \to \mathtt{X}_k \to \mathtt{T}$. As usual, the type variables $\mathtt{X}_1, \ldots, \mathtt{X}_k$ must have kind $\star$ and must be distinct and fresh for the equation's left-hand side. The last two rules simply distribute the type $\mathtt{T}$ to both subpatterns. It is easy to check that $[\![\mathtt{T} : \mathtt{p}]\!]$ is *contravariant* in $\mathtt{T}$:

1.9.17   LEMMA: $\mathtt{T}' \leq \mathtt{T} \wedge [\![\mathtt{T} : \mathtt{p}]\!]$ entails $[\![\mathtt{T}' : \mathtt{p}]\!]$.                                □

This property reflects the fact that $\mathtt{T}$ represents the type of an *input* for the pattern $\mathtt{p}$. Compare it with Lemma 1.6.3.

1.9.18   EXAMPLE: Consider the pattern $\mathtt{Cons}\,(\_, \mathtt{z})$, which appears in Example 1.9.15. We have

$$\begin{aligned}
&\quad [\![\mathtt{T} : \mathtt{Cons}\,(\_, \mathtt{z})]\!] \\
&\equiv \exists \mathtt{Z}_1.([\![\mathtt{Z}_1 \to \mathtt{T} : \mathtt{Cons}]\!] \wedge [\![\mathtt{Z}_1 : (\_, \mathtt{z})]\!]) \\
&\equiv \exists \mathtt{Z}_1.(\mathtt{Z}_1 \to \mathtt{T} \preceq \mathtt{Cons} \wedge \exists \mathtt{Z}_2 \mathtt{Z}_3.([\![\mathtt{Z}_2 \to \mathtt{Z}_3 \to \mathtt{Z}_1 : (\cdot, \cdot)]\!] \wedge [\![\mathtt{Z}_2 : \_]\!] \wedge [\![\mathtt{Z}_3 : \mathtt{z}]\!])) \\
&\equiv \exists \mathtt{Z}_1 \mathtt{Z}_2 \mathtt{Z}_3.(\mathtt{Z}_1 \to \mathtt{T} \preceq \mathtt{Cons} \wedge \mathtt{Z}_2 \to \mathtt{Z}_3 \to \mathtt{Z}_1 \preceq (\cdot, \cdot) \wedge \mathtt{Z}_3 \preceq \mathtt{z})
\end{aligned}$$

where $\mathtt{Z}_1$, $\mathtt{Z}_2$, $\mathtt{Z}_3$ are fresh for $\mathtt{T}$. Let us now place this constraint within the scope of the initial environment, which assigns type schemes to the constructors $\mathtt{Cons}$ and $(\cdot, \cdot)$, and within the scope of a binding of $\mathtt{z}$ to some type $\mathtt{T}'$.

We find

$$\text{let } \Gamma_0 \text{ in let } \mathtt{z} : \mathtt{T}' \text{ in } [\![ \mathtt{T} : \mathtt{Cons}\ (\_, \mathtt{z}) ]\!]$$

$$\equiv \quad \exists \mathtt{Z}_1 \mathtt{Z}_2 \mathtt{Z}_3. (\exists \mathtt{X}. (\mathtt{Z}_1 \to \mathtt{T} \le \mathtt{X} \times \mathsf{list}\ \mathtt{X} \to \mathsf{list}\ \mathtt{X}) \wedge$$
$$\exists \mathtt{Y}_1 \mathtt{Y}_2. (\mathtt{Z}_2 \to \mathtt{Z}_3 \to \mathtt{Z}_1 \le \mathtt{Y}_1 \to \mathtt{Y}_2 \to \mathtt{Y}_1 \times \mathtt{Y}_2) \wedge \mathtt{Z}_3 \le \mathtt{T}')$$

$$\equiv \quad \exists \mathtt{X}. (\mathtt{T} \le \mathsf{list}\ \mathtt{X} \wedge \mathsf{list}\ \mathtt{X} \le \mathtt{T}')$$

where the final simplification relies mainly on C-ARROW, on the corresponding rule for products, and on C-EXTRANS, and is left as an exercise to the reader. Thus, the constraint states that the pattern matches values that have type $\mathsf{list}\ \mathtt{X}$ (equivalently, values whose type $\mathtt{T}$ is a subtype of $\mathsf{list}\ \mathtt{X}$), for some undetermined element type $\mathtt{X}$, and binds $\mathtt{z}$ to values of type $\mathsf{list}\ \mathtt{X}$ (equivalently, values whose type $\mathtt{T}'$ is a supertype of $\mathsf{list}\ \mathtt{X}$).             $\square$

The above example seems to indicate that the constraint generation rules for patterns make some sense. Still, the careful reader may be somewhat puzzled by the third rule, which, compared to its analogue for expressions, reverses the direction of subtyping, but does not reverse the direction of instantiation. Indeed, in order for this rule to make sense, and to be sound, we must formulate a requirement concerning the type schemes assigned to constructors.

1.9.19    DEFINITION: A constructor $\mathtt{c}$ is *invertible* if and only if, when $\vec{\mathtt{X}}$ and $\vec{\mathtt{X}}'$ have length $a(\mathtt{c})$, the constraint $\mathsf{let}\ \Gamma_0$ in $(\vec{\mathtt{X}}' \to \mathtt{T} \preceq \mathtt{c} \wedge \mathtt{c} \preceq \vec{\mathtt{X}} \to \mathtt{T})$ entails $\vec{\mathtt{X}} \le \vec{\mathtt{X}}'$. In the following, we assume patterns contain invertible constructors only.    $\square$

Intuitively, when $\mathtt{c}$ is invertible, it is possible to recover the type of every $\mathtt{v}_i$ from the type of $\mathtt{c}\ \mathtt{v}_1\ \ldots\ \mathtt{v}_k$, a crucial property for pattern matching to be possible. Please note that, if $\Gamma_0(\mathtt{c})$ is monomorphic, then $\mathtt{c}$ is invertible. The following lemma identifies another important class of invertible constructors.

1.9.20    LEMMA: The constructors of algebraic data types are invertible.       $\square$

*Proof:*   Let $\mathtt{c}$ be a constructor introduced by the definition of an algebraic data type $\mathtt{D}$. Let $k = a(\mathtt{c})$. Then, the type scheme $\Gamma_0(\mathtt{c})$ is of the form $\forall \bar{\mathtt{Y}}. \vec{\mathtt{T}} \to \mathtt{D}\ \vec{\mathtt{Y}}$, where $\vec{\mathtt{Y}}$ are the parameters of the definition and $\vec{\mathtt{T}}$, a vector of length $k$, consists of *some of* the definition's components. (More precisely, $\vec{\mathtt{T}}$ contains just one component in the case of variant types and contains all components in the case of record types.) Let $\vec{\mathtt{X}}$ and $\vec{\mathtt{X}}'$ have length $k$. Let $\forall \bar{\mathtt{Y}}_1. \vec{\mathtt{T}}_1 \to \mathtt{D}\ \vec{\mathtt{Y}}_1$ and $\forall \bar{\mathtt{Y}}_2. \vec{\mathtt{T}}_2 \to \mathtt{D}\ \vec{\mathtt{Y}}_2$ be two $\alpha$-equivalent forms of the type scheme $\Gamma_0(\mathtt{c})$, with $\bar{\mathtt{Y}}_1 \mathbin{\#} \bar{\mathtt{Y}}_2$ and $\bar{\mathtt{Y}}_1 \bar{\mathtt{Y}}_2 \mathbin{\#} \mathit{ftv}(\bar{\mathtt{X}}, \bar{\mathtt{X}}', \mathtt{T})$. The constraint $\mathsf{let}\ \Gamma_0$ in $(\vec{\mathtt{X}}' \to \mathtt{T} \preceq \mathtt{c} \wedge \mathtt{c} \preceq \vec{\mathtt{X}} \to \mathtt{T})$ is, by definition, equivalent to $\vec{\mathtt{X}}' \to \mathtt{T} \preceq \Gamma_0(\mathtt{c}) \wedge \Gamma_0(\mathtt{c}) \preceq \vec{\mathtt{X}} \to \mathtt{T}$, that is, $\exists \bar{\mathtt{Y}}_1. (\vec{\mathtt{X}}' \to \mathtt{T} \le \vec{\mathtt{T}}_1 \to \mathtt{D}\ \vec{\mathtt{Y}}_1) \wedge \exists \bar{\mathtt{Y}}_2. (\vec{\mathtt{T}}_2 \to \mathtt{D}\ \vec{\mathtt{Y}}_2 \le \vec{\mathtt{X}} \to \mathtt{T})$. By C-EXAND and C-ARROW, this may be written $\exists \bar{\mathtt{Y}}_1 \bar{\mathtt{Y}}_2. (\mathtt{D}\ \vec{\mathtt{Y}}_2 \le \mathtt{T} \le \mathtt{D}\ \vec{\mathtt{Y}}_1 \wedge \vec{\mathtt{X}} \le \vec{\mathtt{T}}_2 \wedge \vec{\mathtt{T}}_1 \le \vec{\mathtt{X}}')$. Now,

by Definition 1.9.8, $\mathtt{D}\,\vec{\overline{Y}}_2 \leq \mathtt{D}\,\vec{\overline{Y}}_1$ entails $\vec{T}_2 \leq \vec{T}_1$, so the previous constraint entails $\exists \vec{\overline{Y}}_1 \vec{\overline{Y}}_2.(\vec{X} \leq \vec{X}')$, that is, $\vec{X} \leq \vec{X}'$. $\qquad\qquad\qquad\square$

An important class of *noninvertible* constructors are those associated with existential type definitions (page 118), where not all quantifiers of the type scheme $\Gamma_0(\mathtt{c})$ are parameters of the type constructor $\mathtt{D}$. For instance, under the definition $\mathtt{D} \approx \ell : \exists \mathtt{X}.\mathtt{X}$, the type scheme associated with $\ell$ is $\forall \mathtt{X}.\mathtt{X} \to \mathtt{D}$. Then, it is easy to check that $\ell$ is not invertible. This reflects the fact that it is not possible to recover the type of $\mathtt{v}$ from the type of $\ell\,\mathtt{v}$—which must be $\mathtt{D}$ in any case—and explains why existential types require special treatment.

We are now ready to associate a constraint generation rule with the `match` construct. It is given in the lower part of Figure 1-15. In the rule's right-hand side, we write $\vec{z}_i$ for the program variables bound by the pattern $\mathtt{p}_i$, and we write $\vec{X}_i$ for a vector of type variables of the same length. The type variables $\mathtt{X}\vec{\overline{X}}_i$ must have kind $\star$, must be pairwise distinct and must not appear free in the rule's left-hand side. Let us now explain the rule. Its right-hand side is a conjunction, where each conjunct deals with one clause of the `match` construct, requiring $\mathtt{t}_i$ to have type $\mathtt{T}$ under certain assumptions about the program variables $\vec{z}_i$ bound by the pattern $\mathtt{p}_i$. There remains to explain how these assumptions are built. First, as in the case of a `let` construct, we summon a fresh type variable $\mathtt{X}$ and produce $[\![\mathtt{t} : \mathtt{X}]\!]$, the least specific constraint that guarantees $\mathtt{t}$ has type $\mathtt{X}$. Then, reflecting the operational semantics, which feeds (the value produced by) $\mathtt{t}$ into the pattern $\mathtt{p}_i$, we feed the type $\mathtt{X}$ into $\mathtt{p}_i$ and produce $\mathsf{let}\ \vec{z}_i : \vec{X}_i\ \mathsf{in}\ [\![\mathtt{X} : \mathtt{p}_i]\!]$, a constraint that guarantees that $\vec{X}_i$ is a correct vector of type assumptions for the program variables $\vec{z}_i$ (see Example 1.9.18). This explains why we may place $[\![\mathtt{T} : \mathtt{t}_i]\!]$ within the scope of $(\vec{z}_i : \vec{X}_i)$. There remains to point out that, as in the case of the `let` construct, *every* assignment of ground types to $\mathtt{X}\vec{\overline{X}}_i$ that satisfies the constraint $[\![\mathtt{t} : \mathtt{X}]\!] \wedge \mathsf{let}\ \vec{z}_i : \vec{X}_i\ \mathsf{in}\ [\![\mathtt{X} : \mathtt{p}_i]\!]$ is acceptable, so it is valid to universally quantify these type variables. This allows the program variables $\vec{z}_i$ to receive polymorphic type schemes when $\mathtt{t}$ itself has polymorphic type.

1.9.21   EXERCISE [★, RECOMMENDED]: We have previously suggested viewing `let z = t`$_1$` in t`$_2$ as syntactic sugar for `match t`$_1$` with z . t`$_2$, and shown that the operational semantics validates this view. Check that it is also valid from a typing perspective. $\qquad\qquad\square$

The `match` constraint generation rule, if implemented literally, takes $k$ copies of the constraint $[\![\mathtt{t} : \mathtt{X}]\!]$. When $k$ is greater than 1, this compromises the linear time and space complexity of constraint generation. To remedy this problem, one may modify the rule as follows: replace every copy of $[\![\mathtt{t} : \mathtt{X}]\!]$ with $\mathtt{z} \preceq \mathtt{X}$ and place the constraint within the context $\mathsf{let}\ \mathtt{z} : \forall \mathtt{X}[\![\mathtt{t} : \mathtt{X}]\!].\mathtt{X}\ \mathsf{in}\ []$, where $\mathtt{z}$ is

a fresh program variable. It is not difficult to check that the logical meaning of
the constraint is not affected and that a linear behavior is recovered. In prac-
tice, solving the new constraint requires taking instances of the type scheme
$\forall \mathtt{X}[\llbracket \mathtt{t} : \mathtt{X} \rrbracket].\mathtt{X}$, which essentially requires copying $\llbracket \mathtt{t} : \mathtt{X} \rrbracket$ again—however, an
efficient solver may now *simplify* this subconstraint before duplicating it.

The following lemma is a key to establishing subject reduction for R-
Match. It relies on the requirement that constructors be invertible.

1.9.22   Lemma: Assume $[\mathtt{p} \mapsto \mathtt{v}]$ is defined and maps $\vec{z}$ to $\vec{w}$, where $\bar{z} = dpi(\mathtt{p})$.
Let $\vec{z} : \vec{\mathtt{T}}$ be an arbitrary monomorphic environment of domain $\bar{z}$. Then,
$\mathsf{let}\ \Gamma_0\ \mathsf{in}\ (\llbracket \mathtt{v} : \mathtt{T} \rrbracket \wedge \mathsf{let}\ \vec{z} : \vec{\mathtt{T}}\ \mathsf{in}\ \llbracket \mathtt{T} : \mathtt{p} \rrbracket)$ entails $\mathsf{let}\ \Gamma_0\ \mathsf{in}\ \llbracket \vec{w} : \vec{\mathtt{T}} \rrbracket$.   □

We now prove that our extension of ML-the-calculus with pattern match-
ing enjoys subject reduction. We only state that R-Match preserves types,
and leave the new subcase of R-Context, where the evaluation context in-
volves a `match` construct, to the reader. For this subcase to succeed, the value
restriction (Definition 1.7.7) must be extended to require that either all con-
stants have pure semantics or all `match` constructs are in fact of the form
`match v with` $(\mathtt{p}_i\ .\ \mathtt{t}_i)_{i=1}^k$.

1.9.23   Theorem [Subject reduction]: (R-Match) $\subseteq$ ($\sqsubseteq$).   □

1.9.24   Exercise [★★★, ↛]: For the sake of simplicity, we have omitted the pro-
duction `ref p` from the syntax of patterns. The pattern `ref p` matches every
memory location whose content (with respect to the current store) is matched
by `p`. Determine how the previous definitions and proofs must be extended in
order to accommodate this new production.   □

The progress property does *not* hold in general: for instance,
`match Nil with (Cons z . z)` is well-typed (with type $\forall \mathtt{X}.\mathtt{X}$) but is stuck.
In actual implementations of ML-the-programming-language, such errors are
dynamically detected. This may be considered a weakness of ML-the-type-
system. Fortunately, however, it is often possible to statically prove that a
particular `match` construct is *exhaustive* and cannot go wrong. Indeed, if
`match v with` $(\mathtt{p}_i\ .\ \mathtt{t}_i)_{i=1}^k$ is well-typed, then for every $i \in \{1, \ldots, k\}$, the
constraint $\mathsf{let}\ \Gamma_0\ \mathsf{in}\ (\llbracket \mathtt{v} : \mathtt{X} \rrbracket \wedge \exists \bar{\mathtt{x}}.\mathsf{let}\ \vec{z}_i : \vec{\mathtt{X}}\ \mathsf{in}\ \llbracket \mathtt{X} : \mathtt{p}_i \rrbracket)$, where $\bar{z}_i$ are the program
variables bound by $\mathtt{p}_i$, must be satisfiable; that is, $\mathtt{v}$ must have some type
that is an acceptable input for $\mathtt{p}_i$. This fact yields information about $\mathtt{v}$, from
which it may be possible to derive that $\mathtt{v}$ must match one of the patterns $\mathtt{p}_i$.

1.9.25   Example: Let $k = 2$, $\mathtt{p}_1 = \mathtt{Nil}\ \_$, and $\mathtt{p}_2 = \mathtt{Cons}\ (\mathtt{z}_1, \mathtt{z}_2)$. Then, the con-
straints $\mathsf{let}\ \Gamma_0\ \mathsf{in}\ \exists \bar{\mathtt{x}}.\mathsf{let}\ \vec{z}_i : \vec{\mathtt{X}}\ \mathsf{in}\ \llbracket \mathtt{X} : \mathtt{p}_i \rrbracket$, for $i \in \{1, 2\}$, are both equivalent
(after simplification, when $i = 2$) to $\exists \mathtt{Z}.\mathtt{X} \leq \mathsf{list}\ \mathtt{Z}$. Because the type construc-
tor $\mathsf{list}$ is isolated, every closed value $\mathtt{v}$ whose type $\mathtt{X}$ satisfies this constraint

must be an application of `Nil` or `Cons`. If the latter, because `Cons` has type $\forall X. X \times \mathsf{list}\, X \to \mathsf{list}\, X$, and because the type constructor $\times$ is isolated, the argument to `Cons` must be a pair. We conclude that `v` must match either $p_1$ or $p_2$, which guarantees that this `match` construct is exhaustive and its evaluation cannot go wrong. $\qquad\square$

It is beyond the scope of this chapter to give more details about the check for exhaustiveness. The reader is referred to (Sekar, Ramesh, and Ramakrishnan, 1995; Le Fessant and Maranget, 2001).

### Type annotations

So far, we have been interested in a very pure, and extreme, form of type inference. Indeed, in ML-the-calculus, expressions contain no explicit type information whatsoever: it is entirely inferred. In practice, however, it is often useful to insert *type annotations* within expressions, because they provide a form of machine-checked documentation. Type annotations are also helpful when attempting to trace the cause of a type error: by supplying the type-checker with (supposedly) correct type information, one runs a better chance of finding a type inconsistency near an actual programming mistake.

When type annotations are allowed to contain type variables, one must be quite careful about *where* (at which program point) and *how* (existentially or universally) these variables are bound. Indeed, the meaning of type annotations cannot be made precise without settling these issues. In what follows, we first explain how to introduce type annotations whose type variables are bound *locally* and existentially. We show that extending ML-the-calculus with such limited type annotations is again a simple matter of introducing new constants. Then, we turn to a more general case, where type variables may be explicitly existentially introduced *at any program point*. We defer the discussion of *universally* bound type variables to Section 1.10.

Let a *local existential type annotation* $\exists \bar{X}.T$ be a pair of a set of type variables $\bar{X}$ and a type $T$, where $T$ has kind $\star$, $\bar{X}$ is considered bound within $T$, and $\bar{X}$ contains $ftv(T)$. For every such annotation, we introduce a new unary destructor $(\cdot : \exists \bar{X}.T)$. Such a definition is valid only because a type annotation must be *closed*, that is, does not have any free type variables. We write $(t : \exists \bar{X}.T)$ for the application $((\cdot : \exists \bar{X}.T))\, t$. Since a type annotation does not affect the meaning of a program, the new destructor has identity semantics:

$$(v : \exists \bar{X}.T) \xrightarrow{\delta} v \qquad\qquad \text{(R-ANNOTATION)}$$

Its type scheme, however, is not that of the identity, namely $\forall X. X \to X$: instead, it is less general, so that annotating an expression *restricts* its type. Indeed,

we extend the initial environment $\Gamma_0$ with the binding

$$(\cdot : \exists \bar{\mathtt{X}}.\mathtt{T}) : \forall \bar{\mathtt{X}}.\mathtt{T} \to \mathtt{T}$$

1.9.26   EXERCISE [★]: Check that $\forall \bar{\mathtt{X}}.\mathtt{T} \to \mathtt{T}$ is an *instance* of $\forall \mathtt{X}.\mathtt{X} \to \mathtt{X}$ in Damas and Milner's sense, that is, the former is obtained from the latter via the rule DM-INST' given in Exercise 1.2.23. Does this allow arguing that the type scheme assigned to $(\cdot : \exists \bar{\mathtt{X}}.\mathtt{T})$ is sound? Check that the above definitions meet the requirements of Definition 1.7.6.   □

Although inserting a type annotation does not change the semantics of the program, it does affect constraint generation, hence type inference. We let the reader check that, assuming $\bar{\mathtt{X}} \mathbin{\#} ftv(\mathtt{t}, \mathtt{T}')$, the following derived constraint generation rule holds:

$$\mathsf{let}\ \Gamma_0\ \mathsf{in}\ [\![(\mathtt{t} : \exists \bar{\mathtt{X}}.\mathtt{T}) : \mathtt{T}']\!] \equiv \mathsf{let}\ \Gamma_0\ \mathsf{in}\ \exists \bar{\mathtt{X}}.([\![\mathtt{t} : \mathtt{T}]\!] \wedge \mathtt{T} \leq \mathtt{T}')$$

So far, expressions cannot have free type variables, so the hypothesis $\bar{\mathtt{X}} \mathbin{\#} ftv(\mathtt{t})$ may seem superfluous. However, we shall soon allow expressions to contain type annotations with free type variables, so we prefer to make this condition explicit now. According to this rule, the effect of the type annotation is to force the expression $\mathtt{t}$ to have type $\mathtt{T}$, for *some* choice of the type variables $\bar{\mathtt{X}}$. As usual in type systems with subtyping, the expression's final type $\mathtt{T}'$ may then be an arbitrary supertype of this particular instance of $\mathtt{T}$. When subtyping is interpreted as equality, $\mathtt{T}'$ and $\mathtt{T}$ are equated by the constraint, so this constraint generation rule may be read: *a valid type for* $(\mathtt{t} : \exists \bar{\mathtt{X}}.\mathtt{T})$ *must be of the form* $\mathtt{T}$, *for some choice of the type variables* $\bar{\mathtt{X}}$.

1.9.27   EXAMPLE: In DM extended with integers, the expression $(\lambda \mathtt{z}.\mathtt{z} : \mathtt{int} \to \mathtt{int})$ has most general type $\mathtt{int} \to \mathtt{int}$, even though the underlying identity function has most general type $\forall \mathtt{X}.\mathtt{X} \to \mathtt{X}$, so the annotation restricts its type. The expression $(\lambda \mathtt{z}.\mathtt{z} \mathbin{\hat{+}} \hat{1} : \exists \mathtt{X}.\mathtt{X} \to \mathtt{X})$ has type $\mathtt{int} \to \mathtt{int}$, which is also the most general type of the underlying function, so the annotation acts merely as documentation in this case. Note that the type variable $\mathtt{X}$ is instantiated to $\mathtt{int}$ by the constraint solver. The expression $(\lambda \mathtt{z}.(\mathtt{z}\ \hat{1}) : \exists \mathtt{X}.\mathtt{X} \to \mathtt{int})$ has type $(\mathtt{int} \to \mathtt{int}) \to \mathtt{int}$ because the underlying function has type $(\mathtt{int} \to \mathtt{Y}) \to \mathtt{Y}$, which successfully unifies with $\mathtt{X} \to \mathtt{int}$ by instantiating $\mathtt{X}$ to $\mathtt{int} \to \mathtt{int}$ and $\mathtt{Y}$ to $\mathtt{int}$. Last, the expression $(\lambda \mathtt{z}.(\mathtt{z}\ \hat{1}) : \exists \mathtt{X}.\mathtt{int} \to \mathtt{X})$ is ill-typed— even though the underlying expression is well-typed—because the equation $(\mathtt{int} \to \mathtt{Y}) \to \mathtt{Y} = \mathtt{int} \to \mathtt{X}$ is unsatisfiable.   □

1.9.28    EXAMPLE: In DM extended with pairs, the expression $\lambda z_1.\lambda z_2.((z_1 :$ $\exists X.X), (z_2 : \exists X.X))$ has most general type $\forall XY.X \to Y \to X \times Y$. In other words, the two occurrences of X do not represent the same type. Indeed, one could just as well have written $\lambda z_1.\lambda z_2.((z_1 : \exists X.X), (z_2 : \exists Y.Y))$. If one wishes $z_1$ and $z_2$ to receive the same type, one must lift the type annotations and merge them above the pair constructor, as follows: $\lambda z_1.\lambda z_2.((z_1, z_2) : \exists X.X \times X)$. In the process, the type constructor $\times$ has appeared in the annotation, causing its size to increase.                                                                    □

The above example reveals a limitation of this style of type annotations: by requiring every type annotation to be closed, we lose the ability for two separate annotations to *share* a type variable. Yet, such a feature is sometimes desirable. If the two annotations where sharing is desired are distant in the code, it may be awkward to lift and merge them into a single annotation; so, more expressive power is sometimes truly needed.

Thus, we are lead to consider more general type annotations, of the form $(t : T)$, where T has kind $\star$, and where the type variables that appear within T are considered *free*, so that distinct type annotations may refer to shared type variables. For this idea to make sense, however, it is still necessary to specify where these type variables are bound. We do so using expressions of the form $\exists \bar{X}.t$. Such an expression binds the type variables $\bar{X}$ within the expression t, so that all free occurrences of X (where $X \in \bar{X}$) in type annotations inside t stand for the same type. Thus, we break the simple type annotation construct $(\cdot : \exists \bar{X}.T)$ into two more elementary constituents, namely *existential type variable introduction* $\exists \bar{X}.\cdot$ and *type constraint* $(\cdot : T)$. Note that both are new forms of expressions; neither can be encoded by adding new constants to the calculus, because it is not possible to assign *closed* type schemes to them.

Technically, allowing expressions to contain type variables requires some care. Several constraint generation rules employ auxiliary type variables, which become bound in the generated constraint. These type variables may be chosen in an arbitrary way, provided they do not appear free in the rule's left-hand side—a side-condition intended to avoid inadvertent capture. So far, this side-condition could be read: the auxiliary type variables used to form the constraint $[\![t : T]\!]$ must not appear free within T. Now, since type annotations may contain free type variables, the side-condition becomes: *the auxiliary type variables used to form $[\![t : T]\!]$ must not appear free within t or T.*

With this extended side-condition in mind, our original constraint generation rules remain unchanged. We add two new rules to describe how the new expression forms affect constraint generation:

$$
\begin{aligned}
[\![\exists \bar{X}.t : T]\!] &= \exists \bar{X}.[\![t : T]\!] \qquad \text{provided } \bar{X} \mathbin{\#} \mathit{ftv}(T)\\
[\![(t : T) : T']\!] &= [\![t : T]\!] \wedge T \leq T'
\end{aligned}
$$

The effect of these rules is simple. The construct $\exists \bar{\mathtt{X}}.\mathtt{t}$ is an indication to the constraint generator that the type variables $\bar{\mathtt{X}}$, which may occur free within type annotations inside $\mathtt{t}$, should be existentially bound at this point. The side-condition $\bar{\mathtt{X}} \,\#\, ftv(\mathtt{T})$ ensures that quantifying over $\bar{\mathtt{X}}$ in the generated constraint does not capture type variables in the expected type $\mathtt{T}$. It can always be satisfied by $\alpha$-conversion of the expression $\exists \bar{\mathtt{X}}.\mathtt{t}$. The construct $(\mathtt{t} : \mathtt{T})$ is an indication to the constraint generator that the expression $\mathtt{t}$ should have type $\mathtt{T}$, and it is treated as such by generating the subconstraint $[\![\mathtt{t} : \mathtt{T}]\!]$. The expression's type may be an arbitrary supertype of $\mathtt{T}$, hence the auxiliary constraint $\mathtt{T} \leq \mathtt{T}'$.

1.9.29   EXAMPLE: In DM extended with pairs, the expression $\lambda \mathtt{z}_1.\lambda \mathtt{z}_2.\exists \mathtt{X}.((\mathtt{z}_1 : \mathtt{X}), (\mathtt{z}_2 : \mathtt{X}))$ has most general type $\forall \mathtt{X}.\mathtt{X} \rightarrow \mathtt{X} \rightarrow \mathtt{X} \times \mathtt{X}$. Indeed, the constraint generated for this expression contains the pattern $\exists \mathtt{X}.([\![\mathtt{z}_1 : \mathtt{X}]\!] \wedge [\![\mathtt{z}_2 : \mathtt{X}]\!] \wedge \ldots)$, which causes $\mathtt{z}_1$ and $\mathtt{z}_2$ to receive the same type. Note that this style is more flexible than that employed in Example 1.9.28, where we were forced to use a single, monolithic type annotation to express this sharing constraint.   □

1.9.30   REMARK: In practice, a type variable is usually represented as a memory cell in the typechecker's heap. So, one cannot say that the source code contains type variables; rather, it contains *names* that are meant to stand for type variables. Let us write $X$ for such a name, and $T$ for a type made of type constructors and names, rather than of type constructors and type variables. Then, our new expression forms are really $\exists \bar{X}.\mathtt{t}$ and $(\mathtt{t} : T)$. When the constraint generator enters the scope of an introduction form $\exists \bar{X}.\mathtt{t}$, it allocates a vector of fresh type variables $\bar{\mathtt{X}}$, and augments an internal environment with the bindings $\bar{X} \mapsto \bar{\mathtt{X}}$. Because the type variables are fresh, the side-condition of the first constraint generation rule above is automatically satisfied. When the constraint generator finds a type annotation $(\mathtt{t} : T)$, it looks up the internal environment to translate the type annotation $T$ into an internal type $\mathtt{T}$—which fails if $T$ contains a name that is not in scope—and applies the second constraint generation rule above.   □

1.9.31   EXERCISE [★★, ⇸]: Let $\bar{\mathtt{X}} \supseteq ftv(\mathtt{T})$ and $\bar{\mathtt{X}} \,\#\, ftv(\mathtt{t})$. Check that the constraints $[\![(\mathtt{t} : \exists \bar{\mathtt{X}}.\mathtt{T}) : \mathtt{T}']\!]$ and $[\![\exists \bar{\mathtt{X}}.(\mathtt{t} : \mathtt{T}) : \mathtt{T}']\!]$ are equivalent. In other words, the *local* type annotations introduced earlier may be expressed in terms of the more complex constructs described above.   □

1.9.32   EXERCISE [★★, ⇸]: One way of giving identity semantics to our new type annotation constructs is to *erase* them altogether prior to execution. Give an inductive definition of $\lfloor \mathtt{t} \rfloor$, the expression obtained by removing all type annotation constructs from the expression $\mathtt{t}$. Check that $[\![\mathtt{t} : \mathtt{T}]\!]$ entails $[\![\lfloor \mathtt{t} \rfloor : \mathtt{T}]\!]$ and explain why this is sufficient to ensure type soundness.   □

It is interesting to study how explicit introduction of existentially quantified type variables interacts with let-polymorphism. The source of their interaction lies in the difference between the constraints let $z : \forall \bar{\mathtt{X}}[\exists \mathtt{X}.C_1].\mathtt{T}$ in $C_2$ and $\exists \mathtt{X}.$let $z : \forall \bar{\mathtt{X}}[C_1].\mathtt{T}$ in $C_2$, which was explained in Example 1.3.28. In the former constraint, every free occurrence of z inside $C_2$ causes a copy of $\exists \mathtt{X}.C_1$ to be taken, thus creating its own fresh copy of X. In the latter constraint, on the other hand, every free occurrence of z inside $C_2$ produces a copy of $C_1$. All such copies *share* references to X, because its quantifier was not duplicated. In the former case, one may say that the type scheme assigned to z is *polymorphic* with respect to X, while in the latter case it is *monomorphic*. As a result, the placement of type variable introduction expressions with respect to let bindings in the source code is meaningful: introducing a type variable *outside* of a let construct *prevents* it from being generalized.

1.9.33 EXAMPLE: In DM extended with integers and Booleans, the program let f = $\exists \mathtt{X}.\lambda z.(z : \mathtt{X})$ in $(f\,0, f\,\text{true})$ is well-typed. Indeed, the type scheme assigned to f is $\forall \mathtt{X}.\mathtt{X} \to \mathtt{X}$. However, the program $\exists \mathtt{X}.$let f = $\lambda z.(z : \mathtt{X})$ in $(f\,0, f\,\text{true})$ is ill-typed. Indeed, the type scheme assigned to f is $\mathtt{X} \to \mathtt{X}$; then, no value of X satisfies the constraints associated with the applications $f\,0$ and $f\,\text{true}$. The latter behavior is observed in Objective Caml, where type variables are *implicitly* introduced at the outermost level of expressions:

```
# let f z = (z:'a) in (f 0, f true);;
This expression has type bool but is here used with type int
```

More details about the treatment of type annotations in Standard ML, Objective Caml, and Haskell are given on page 113. □

1.9.34 EXERCISE [★, ↛]: Determine which constraints are generated for the two programs in Example 1.9.33. Check that the former is indeed well-typed, while the latter is ill-typed. □

### Recursive types

We have shown that specializing HM($X$) with an equality-only syntactic model yields HM($=$), a constraint-based formulation of Damas and Milner's type system. Similarly, it is possible to specialize HM($X$) with an equality-only free *regular* tree model, yielding a constraint-based type system that may be viewed as an extension of Damas and Milner's type discipline with recursive types. This flavor of recursive types is sometimes known as *equirecursive*, since cyclic *equations*, such as $\mathtt{X} = \mathtt{X} \to \mathtt{X}$, are then satisfiable. Our theorems about type inference and type soundness, which are independent of the model, remain valid. The constraint solver described in Section 1.8 may be used in

the setting of an equality-only free regular tree model: the only difference with the syntactic case is that the occurs check is no longer performed.

Please note that, although *ground types* are regular, *types* remain finite objects: their syntax is unchanged. The $\mu$ notation commonly employed to describe recursive types may be emulated using type equations: for instance, the notation $\mu \mathtt{X}.\mathtt{X} \to \mathtt{X}$ corresponds, in our constraint-based approach, to the type scheme $\forall \mathtt{X}[\mathtt{X} = \mathtt{X} \to \mathtt{X}].\mathtt{X}$.

Although recursive types come for free, as explained above, they have not been adopted in mainstream programming languages based on ML-the-type-system. The reason is pragmatic: experience shows that many nonsensical expressions are well-typed in the presence of recursive types, whereas they are not in their absence. Thus, the gain in expressiveness is offset by the fact that many programming mistakes are detected later than otherwise possible. Consider, for instance, the following `OCaml` session:

```
ocaml -rectypes
# let rec map f = function
    | [] → []
    | x :: l → (map f x) :: (map f l);;
val map : 'a → ('b list as 'b) → ('c list as 'c) = <fun>
```

This nonsensical variant of `map` is essentially useless, yet well-typed. Its principal type scheme, in our notation, is $\forall \mathtt{XYZ}[\mathtt{Y} = \mathsf{list}\,\mathtt{Y} \wedge \mathtt{Z} = \mathsf{list}\,\mathtt{Z}].\mathtt{X} \to \mathtt{Y} \to \mathtt{Z}$. In the absence of recursive types, it is ill-typed, since the constraint $\mathtt{Y} = \mathsf{list}\,\mathtt{Y} \wedge \mathtt{Z} = \mathsf{list}\,\mathtt{Z}$ is then false.

The need for equirecursive types is usually suppressed by the presence of algebraic data types, which offer *isorecursive* types, in the language. Yet, they are still necessary in some situations, such as in Objective Caml's object-oriented extension (Rémy and Vouillon, 1998), where recursive object types are commonly inferred. In order to allow recursive object types while still rejecting the above variant of `map`, Objective Caml's constraint solver implements a selective occurs check, which forbids cycles unless they involve the type constructor $\langle \cdot \rangle$ associated with objects. The corresponding model is a tree model where every infinite path down a tree must encounter the type constructor $\langle \cdot \rangle$ infinitely often.

## 1.10   Universal quantification in constraints

The constraint logic studied so far allows a set of variables $\bar{\mathtt{X}}$ to be existentially quantified within a formula $C$. The resulting formula $\exists \bar{\mathtt{X}}.C$ receives its standard meaning: it requires $C$ to hold for *some* $\bar{\mathtt{X}}$. However, we currently have no way of requiring a formula $C$ to hold for *all* $\bar{\mathtt{X}}$. Is it possible to extend our

logic with universal quantification? If so, what are the new possibilities offered by this extension, in terms of type inference? The present section proposes some answers to these questions.

It is worth noting that, although the standard notation for type schemes involves the symbol ∀, type scheme introduction and instantiation constraints do not allow an encoding of universal quantification. Indeed, a universal quantifier in a type scheme is very much like an existential quantifier in a constraint: this is suggested, for instance, by Definition 1.3.3 and by C-LETEX.

### Constraints

We extend the syntax of constraints as follows:

$$C ::= \ldots \mid \forall \bar{\mathrm{X}}.C$$

Universally quantified variables are often referred to as *rigid*, while existentially quantified variables are known as *flexible*. The logical interpretation of constraints (Figure 1-5) is extended as follows:

$$\frac{\forall \vec{t} \quad \phi[\vec{\mathrm{X}} \mapsto \vec{t}\,] \vdash \mathsf{def}\ \Gamma\ \mathsf{in}\ C \qquad \bar{\mathrm{X}} \mathrel{\#} \mathit{ftv}(\Gamma)}{\phi \vdash \mathsf{def}\ \Gamma\ \mathsf{in}\ \forall \bar{\mathrm{X}}.C} \qquad \text{(CM-FORALL)}$$

We let the reader check that none of the results established in Section 1.3 are affected by this addition. Furthermore, the extended constraint language enjoys the following properties.

1.10.1   LEMMA: $\forall \bar{\mathrm{X}}.C \Vdash C$. Conversely, $\bar{\mathrm{X}} \mathrel{\#} \mathit{ftv}(C)$ implies $C \Vdash \forall \bar{\mathrm{X}}.C$.        □

1.10.2   LEMMA: $\bar{\mathrm{X}} \mathrel{\#} \mathit{ftv}(C_2)$ implies $\forall \bar{\mathrm{X}}.(C_1 \wedge C_2) \equiv (\forall \bar{\mathrm{X}}.C_1) \wedge C_2$.        □

1.10.3   LEMMA: $\forall \bar{\mathrm{X}}.\forall \bar{\mathrm{Y}}.C \equiv \forall \bar{\mathrm{X}}\bar{\mathrm{Y}}.C$.        □

1.10.4   LEMMA: Let $\bar{\mathrm{X}} \mathrel{\#} \bar{\mathrm{Y}}$. Then, $\exists \bar{\mathrm{X}}.\forall \bar{\mathrm{Y}}.C$ entails $\forall \bar{\mathrm{Y}}.\exists \bar{\mathrm{X}}.C$. Conversely, if $\exists \bar{\mathrm{Y}}.C$ determines $\bar{\mathrm{X}}$, then $\forall \bar{\mathrm{Y}}.\exists \bar{\mathrm{X}}.C$ entails $\exists \bar{\mathrm{X}}.\forall \bar{\mathrm{Y}}.C$.        □

### Constraint solving

We briefly explain how to extend the constraint solver described in Section 1.8 with support for universal quantification. (Thus, we again assume an equality-only free tree model.) Constraint solving in the presence of equations and of existential and universal quantifiers is known as *unification under a mixed prefix*. It is a particular case of the decision problem for the first-order theory of equality on trees; see *e.g.* (Comon and Lescanne, 1989). Extending our solver is straightforward: in fact, the treatment of universal quantifiers turns

$$S; U; \forall \bar{\mathtt{X}}.C \quad \rightarrow \quad S[\forall \bar{\mathtt{X}}.[]]; U; C \qquad \text{(S-SOLVE-ALL)}$$
$$\text{if } \bar{\mathtt{X}} \mathbin{\#} \mathit{ftv}(U)$$

$$S[\forall \bar{\mathtt{X}}.\exists \bar{\mathtt{Y}} \bar{\mathtt{Z}}.[]]; U; \mathsf{true} \quad \rightarrow \quad S[\exists \bar{\mathtt{Y}}.\forall \bar{\mathtt{X}}.\exists \bar{\mathtt{Z}}.[]]; U; \mathsf{true} \qquad \text{(S-ALLEX)}$$
$$\text{if } \bar{\mathtt{X}} \mathbin{\#} \bar{\mathtt{Y}} \wedge \exists \bar{\mathtt{X}} \bar{\mathtt{Z}}.U \text{ determines } \bar{\mathtt{Y}}$$

$$S[\forall \bar{\mathtt{X}} \mathtt{X}.\exists \bar{\mathtt{Y}}.[]]; U; \mathsf{true} \quad \rightarrow \quad \mathsf{false} \qquad \text{(S-ALL-FAIL-1)}$$
$$\text{if } \mathtt{X} \notin \bar{\mathtt{Y}} \wedge \mathtt{X} \prec^{\star}_{U} \mathtt{Z} \wedge \mathtt{Z} \notin \mathtt{X}\bar{\mathtt{Y}}$$

$$S[\forall \bar{\mathtt{X}} \mathtt{X}.\exists \bar{\mathtt{Y}}.[]]; \mathtt{X} = \mathtt{T} = \epsilon \wedge U; \mathsf{true} \quad \rightarrow \quad \mathsf{false} \qquad \text{(S-ALL-FAIL-2)}$$
$$\text{if } \mathtt{X} \notin \bar{\mathtt{Y}} \wedge \mathtt{T} \notin \mathcal{V}$$

$$S[\forall \bar{\mathtt{X}}.\exists \bar{\mathtt{Y}}.[]]; U_1 \wedge U_2; \mathsf{true} \quad \rightarrow \quad S; U_1; \mathsf{true} \qquad \text{(S-POP-ALL)}$$
$$\text{if } \bar{\mathtt{X}} \bar{\mathtt{Y}} \mathbin{\#} \mathit{ftv}(U_1) \wedge \exists \bar{\mathtt{Y}}.U_2 \equiv \mathsf{true}$$

**Figure 1-16: Solving universal constraints**

out to be surprisingly analogous to that of let constraints. To begin, we extend the syntax of stacks with so-called *universal frames*:

$$S ::= \dots \mid S[\forall \bar{\mathtt{X}}.[]]$$

Because existential quantifiers cannot, in general, be hoisted out of universal quantifiers, rules S-EX-1 to S-EX-4 now allow floating them up to the nearest enclosing let *or universal* frame, if any, or to the outermost level, otherwise. Thus, in our machine representation of stacks, where rules S-EX-1 to S-EX-4 are applied in an eager fashion, every universal frame carries a list of the type variables that are existentially bound immediately after it, and integer ranks count not only let frames, but also universal frames.

The solver's specification is extended with the rules in Figure 1-16. S-SOLVE-ALL, a forward rule, discovers a universal constraint and enters it, creating a new universal frame to record its existence. S-ALLEX exploits Lemma 1.10.4 to hoist existential quantifiers out of the universal frame. It is analogous to S-LETALL, and its implementation may rely upon the same procedure (Exercise 1.8.8). The next two rules detect failure conditions. S-ALL-FAIL-1 states that the constraint $\forall \mathtt{X}.\exists \bar{\mathtt{Y}}.U$ is false if the rigid variable X is directly or indirectly *dominated* by a *free* variable Z. Indeed, the value of X is then determined by that of Z—but a universally quantified variable ranges over *all* values, so this is a contradiction. In such a case, X is commonly said to *escape its scope*. S-ALL-FAIL-2 states that the same constraint is false if X is equated with a *nonvariable* term. Indeed, the value of X is then

partially determined, since its head constructor is known, which again contradicts its universal status. Last, S-Pop-All splits the current unification constraint into two components $U_1$ and $U_2$, where $U_1$ is made up entirely of *old* variables and $U_2$ constrains *young* variables only. This decomposition is analogous to that performed by S-Pop-Let. Then, it is not difficult to check that $\forall \bar{\mathtt{X}}.\exists \bar{\mathtt{Y}}.(U_1 \wedge U_2)$ is equivalent to $U_1$. So, the universal frame, as well as $U_2$, are discarded, and the solver proceeds by examining whatever remains on top of the stack $S$.

It is possible to further extend the treatment of universal frames with two rules analogous to S-Compress and S-UnName. In practice, this improves the solver's efficiency, and makes it easier to share code between the treatment of let frames and that of universal frames.

It is interesting to remark that, as far as the underlying unification algorithm is concerned, there is no difference between existentially and universally quantified type variables. The algorithm solves whatever equations are presented to it, without inquiring about the status of their variables. Equations that lead to failure, because a rigid variable escapes its scope or is equated with a nonvariable term, are detected only when the universal frame is exited. A perhaps more common approach is to *mark* rigid variables as such, allowing the unification algorithm to signal failure as soon as one of the two error conditions is encountered. In this approach, a rigid variable may successfully unify only with itself or with flexible variables fresher than itself. It is often called a *Skolem constructor* in the literature (Läufer and Odersky, 1994; Shields and Peyton Jones, 2002). An interesting variant of this approach appears in Dowek, Hardin, Kirchner and Pfenning's treatment of (higher-order) unification (1995; 1998), where flexible variables are represented as ordinary variables, while rigid variables are encoded using De Bruijn indices.

The properties of our constraint solver are preserved by this extension: it is possible to prove that Lemmas 1.8.9, 1.8.10, and 1.8.11 remain valid.

### Type annotations, continued

In Section 1.9, we introduced the expression form $(\mathtt{t} : \exists \bar{\mathtt{X}}.\mathtt{T})$, allowing an expression $\mathtt{t}$ to be annotated with a type $\mathtt{T}$ whose free variables $\bar{\mathtt{X}}$ are *locally* and *existentially* bound. It is now natural to introduce the symmetric expression form $(\mathtt{t} : \forall \bar{\mathtt{X}}.\mathtt{T})$, where $\mathtt{T}$ has kind $\star$, $\bar{\mathtt{X}}$ is bound within $\mathtt{T}$, and $\bar{\mathtt{X}}$ contains $\mathit{ftv}(\mathtt{T})$, as before. Its constraint generation rule is as follows:

$$[\![ (\mathtt{t} : \forall \bar{\mathtt{X}}.\mathtt{T}) : \mathtt{T}' ]\!] \quad = \quad \forall \bar{\mathtt{X}}.[\![ \mathtt{t} : \mathtt{T} ]\!] \wedge \exists \bar{\mathtt{X}}.(\mathtt{T} \leq \mathtt{T}') \quad \text{provided } \bar{\mathtt{X}} \mathrel{\#} \mathit{ftv}(\mathtt{t}, \mathtt{T}')$$

The first conjunct requires $\mathtt{t}$ to have type $\mathtt{T}$ for *all* values of $\bar{\mathtt{X}}$. Here, the type variables $\bar{\mathtt{X}}$ are *universally* bound, as expected. The second conjunct requires

T′ to be *some* instance of the universal annotation ∀X̄.T. Since T′ is only a monotype, it seems difficult to think of another sensible way of constraining T′. For this reason, the type variables X̄ are still *existentially* bound in the second conjunct. This makes the interpretation of the universal quantifier in type annotations a bit more complex than that of the existential quantifier. For instance, when subtyping is interpreted as equality, the constraint generation rule may be read: *a valid type for* (t : ∀X̄.T) *is of the form* T, *for* some *choice of the type variables* X̄, *provided* t *has type* T *for* all *choices of* X̄.

We remark that (t : ∀X̄.T) must be a new expression form: it cannot be encoded by adding new constants to the calculus—whereas (t : ∃X̄.T) could—because none of the existing constraint generation rules produce universally quantified constraints. Like all type annotations, it has identity semantics.

What is the use of universal type annotations, compared with existential type annotations? When a type variable is existentially bound, the typechecker is free to assign it whatever value makes the program well-typed. As a result, the expressions ($\lambda$z.z $\hat{+}$ $\hat{1}$ : ∃X.X → X) and ($\lambda$z.z : ∃X.X → X) are both well-typed: X is assigned int in the former case, and remains undetermined in the latter. However, it is sometimes useful to be able to insist that an expression should be polymorphic. This effect is naturally achieved by using a universally bound type variable. Indeed, ($\lambda$z.z$\hat{+}$$\hat{1}$ : ∀X.X → X) is ill-typed, because ∀X.(X = int) is false, while ($\lambda$z.z : ∀X.X → X) is well-typed.

1.10.5 EXERCISE [★]: Write down the constraints ∃Z.⟦($\lambda$z.z $\hat{+}$ $\hat{1}$ : ∀X.X → X) : Z⟧ and ∃Z.⟦($\lambda$z.z : ∀X.X → X) : Z⟧, which tell whether these expressions are well-typed. Check that the former is false, while the latter is satisfiable. □

A universal type annotation, as defined above, is nothing but a (closed) Damas-Milner type scheme. Thus, the new construct (t : ∀X̄.T) gives us the ability to ensure that the expression t admits the type scheme ∀X̄.T. This feature is exploited at the module level in ML-the-programming-language, where it is necessary to check that the inferred type for a module component t is more general than the type scheme S that appears in the module's signature. In our view, this process simply consists in ensuring that (t : S) is well-typed.

In Section 1.9, we have pointed out that local (that is, closed) type annotations offer limited expressiveness, because they cannot share type variables. To lift this limitation, we have introduced the expression forms ∃X̄.t and (t : T). The former binds the type variables X̄ within t, making them available for use in type annotations, and instructs the constraint generator to existentially quantify them at this point. The latter requires t to have T. It is natural to proceed in the same manner in the case of universal type annotations. We now introduce the expression form ∀X̄.t, which also binds X̄ within t, but comes

with a different constraint generation rule:

$$[\![\forall\bar{\mathtt{X}}.\mathtt{t} : \mathtt{T}]\!] \quad = \quad \forall\bar{\mathtt{X}}.\exists\mathtt{Z}.[\![\mathtt{t} : \mathtt{Z}]\!] \wedge \exists\bar{\mathtt{X}}.[\![\mathtt{t} : \mathtt{T}]\!] \quad \text{provided } \bar{\mathtt{X}} \mathrel{\#} \mathit{ftv}(\mathtt{T}) \wedge \mathtt{Z} \notin \mathit{ftv}(\mathtt{t})$$

This rule is a bit more complex than that associated with the expression form $\exists\bar{\mathtt{X}}.\mathtt{t}$. Again, this is due to the fact that we do not wish to overconstrain $\mathtt{T}$. The first exercise below shows that a more naïve version of the rule does not yield the desired behavior. The second exercise shows that this version does. The third exercise clarifies an efficiency concern.

1.10.6   EXERCISE [★]: Assume that $[\![\forall\bar{\mathtt{X}}.\mathtt{t} : \mathtt{T}]\!]$ is defined as $\forall\bar{\mathtt{X}}.[\![\mathtt{t} : \mathtt{T}]\!]$, provided $\bar{\mathtt{X}} \mathrel{\#} \mathit{ftv}(\mathtt{T})$. Write down the constraint $[\![\forall\mathtt{X}.(\lambda\mathtt{z}.\mathtt{z} : \mathtt{X} \to \mathtt{X}) : \mathtt{Z}]\!]$. Can you describe its solutions? Does it have the intended meaning?    □

1.10.7   EXERCISE [★★]: Let $\bar{\mathtt{X}} \supseteq \mathit{ftv}(\mathtt{T})$ and $\bar{\mathtt{X}} \mathrel{\#} \mathit{ftv}(\mathtt{t})$. Check that the constraints $[\![(\mathtt{t} : \forall\bar{\mathtt{X}}.\mathtt{T}) : \mathtt{T}']\!]$ and $[\![\forall\bar{\mathtt{X}}.(\mathtt{t} : \mathtt{T}) : \mathtt{T}']\!]$ are equivalent. In other words, *local* universal type annotations may also be expressed in terms of the more complex constructs described above.    □

1.10.8   EXERCISE [★★★★, ↛]: The constraint generation rule that appears above compromises the linear time and space complexity of constraint generation, because it duplicates the term $\mathtt{t}$. It is possible to avoid this problem, but this requires a slight generalization of the constraint language. Let us write $\mathsf{let}\ \mathtt{x} : \forall\underline{\bar{\mathtt{X}}}\bar{\mathtt{Y}}[C_1].\mathtt{T}\ \mathsf{in}\ C_2$ for $\forall\bar{\mathtt{X}}.\exists\bar{\mathtt{Y}}.C_1 \wedge \mathsf{def}\ \mathtt{x} : \forall\bar{\mathtt{X}}\bar{\mathtt{Y}}[C_1].\mathtt{T}\ \mathsf{in}\ C_2$. In this extended $\mathsf{let}$ form, the underlined variables $\bar{\mathtt{X}}$ are interpreted as *rigid*, instead of *flexible*, while checking that $C_1$ is satisfiable. However, the type scheme associated with $\mathtt{x}$ is not affected. *Check that the above constraint generation rule may now be written as follows:*

$$[\![\forall\bar{\mathtt{X}}.\mathtt{t} : \mathtt{T}]\!] \quad = \quad \mathsf{let}\ \mathtt{x} : \forall\underline{\bar{\mathtt{X}}}\mathtt{Z}[[\![\mathtt{t} : \mathtt{Z}]\!]].\mathtt{Z}\ \mathsf{in}\ \mathtt{x} \preceq \mathtt{T} \quad \text{provided } \mathtt{Z} \notin \mathit{ftv}(\mathtt{t})$$

Roughly speaking, the new rule forms a most general type scheme for $\mathtt{t}$, ensures that the type variables $\bar{\mathtt{X}}$ are unconstrained in it, and checks that $\mathtt{T}$ is an instance of it. Furthermore, it does not duplicate $\mathtt{t}$. To complete the exercise, *extend the specification of the constraint solver (Figures 1-12 and 1-16), as well as its implementation,* to deal with this extension of the constraint language.    □

    To conclude, let us once again stress that, if $\mathtt{T}$ has free type variables, the effect of the type annotation $(\mathtt{t} : \mathtt{T})$ depends on *how* and *where* they are bound. The effect of *how* stems from the fact that binding a type variable universally, rather than existentially, leads to a stricter constraint. Indeed, we let the reader check that $[\![\forall\bar{\mathtt{X}}.\mathtt{t} : \mathtt{T}]\!]$ entails $[\![\exists\bar{\mathtt{X}}.\mathtt{t} : \mathtt{T}]\!]$, while the converse

does not hold in general. The effect of *where* has been illustrated, in the case of existentially bound type variables, in Section 1.9. It is due, in that case, to the fact that let and ∃ do not commute. In the case of universally bound type variables, it may be imputed to the fact that ∀ and ∃ do not commute. For instance, $\lambda z.\forall X.(z : X)$ is ill-typed, because *inside the $\lambda$-abstraction*, the program variable z cannot be said to have every type. However, $\forall X.\lambda z.(z : X)$ is well-typed, because the identity function does have type $X \rightarrow X$ for every X.

1.10.9   EXERCISE [★]: Write down the constraints $\exists Z.[\![\lambda z.\forall X.(z : X) : Z]\!]$ and $\exists Z.[\![\forall X.\lambda z.(z : X) : Z]\!]$, which tell whether these expressions are well-typed. Is the former satisfiable? Is the latter?   □

   In Standard ML and Objective Caml, the type variables that appear in type annotations are *implicitly* bound. That is, there is no syntax in the language for the constructs $\exists \bar{X}.t$ and $\forall \bar{X}.t$. When a type annotation (t : T) contains a free type variable X, a fixed convention tells how and where X is bound. In Standard ML, X is *universally* bound at the nearest val binding that encloses all related occurrences of X (Milner, Tofte, and Harper, 1990). The 1997 revision of Standard ML (Milner, Tofte, Harper, and MacQueen, 1997b) slightly improves on this situation by allowing type variables to be *explicitly* introduced at val bindings. However, they still must be universally bound. In Objective Caml, X is *existentially* bound at the nearest enclosing toplevel let binding; this behavior seems to be presently undocumented. We argue that (i) allowing type variables to be implicitly introduced is confusing; and (ii) for expressiveness, both universal and existential quantifiers should be made available to programmers. Surprisingly, these language design and type inference issues seem to have received little attention in the literature, although they have most likely been "folklore" for a long time. Peyton Jones and Shields (2003) study these issues in the context of Haskell, and concur with (i). Concerning (ii), they seem to think that the language designer must choose between existential and universal type variable introduction forms—which they refer to as "type-sharing" and "type-lambda"—whereas we point out that they may and should coexist.

### Polymorphic recursion

Example 1.2.10 explains how the letrec construct found in ML-the-programming-language may be viewed as an application of the constant fix, wrapped inside a normal let construct. Exercise 1.9.6 shows that this gives rise to a somewhat restrictive constraint generation rule: generalization occurs only *after* the application of fix is typechecked. In other words, in letrec f = $\lambda z.t_1$ in $t_2$, all occurrences of f within $t_1$ must have the same

(monomorphic) *type*. This restriction is sometimes a nuisance, and seems un-warranted: if the function that is being defined is polymorphic, it should be possible to use it at different types even inside its own definition. Indeed, My-croft (1984) extended Damas and Milner's type system with a more liberal treatment of recursion, commonly known as *polymorphic recursion*. The idea is to only request occurrences of f within $t_1$ to have the same *type scheme*. Hence, they may have different *types*, all of which are instances of a common type scheme. It was later shown that well-typedness in Mycroft's extended type system is undecidable (Henglein, 1993; Kfoury, Tiuryn, and Urzyczyn, 1993). To work around this stumbling block, one solution is to use a semi-algorithm, falling back to monomorphic recursion if it does not succeed or fail in reasonable time. Although such a solution might be appealing in the setting of an automated program analysis, it is less so in the setting of a programmer-visible type system, because it may become difficult to under-stand why a program is ill-typed. Thus, we describe a simpler solution, which consists in requiring the programmer to explicitly supply a type scheme for f. This is an instance of a *mandatory* type annotation.

To begin, we must change the status of fix, because if fix remains a constant, then f must remain $\lambda$-bound and cannot receive a polymorphic type scheme. We turn fix into a language construct, which binds a program variable f, and annotates it with a DM type scheme. The syntax of values and expressions is thus extended as follows:

$$v ::= \ldots \mid \text{fix } f : S.\lambda z.t \qquad t ::= \ldots \mid \text{fix } f : S.\lambda z.t$$

Please note that f is bound within $\lambda z.t$. The operational semantics is extended as follows.

$$(\text{fix } f : S.\lambda z.t)\, v \longrightarrow (\text{let } f = \text{fix } f : S.\lambda z.t \text{ in } \lambda z.t)\, v \qquad (\text{R-Fix'})$$

The type annotation S plays no essential role in the reduction; it is merely preserved. It is now possible to define $\text{letrec } f : S = \lambda z.t_1 \text{ in } t_2$ as syntactic sugar for $\text{let } f = \text{fix } f : S.\lambda z.t_1 \text{ in } t_2$.

We now give a constraint generation rule for fix:

$$[\![\text{fix } f : S.\lambda z.t : T]\!] \quad = \quad \text{let } f : S \text{ in } [\![\lambda z.t : S]\!] \wedge S \preceq T$$

The left-hand conjunct requires the function $\lambda z.t$ to have type scheme S, under the assumption that f has type S. Thus, it is now possible for different occur-rences of f within t to receive different types. If S is $\forall \bar{X}.T$, where $\bar{X} \mathbin{\#} ftv(t)$, then we write $[\![t : S]\!]$ for $\forall \bar{X}.[\![t : T]\!]$. Indeed, checking the validity of a poly-morphic type annotation—be it mandatory, as is the case here, or optional, as was previously the case—requires a universally quantified constraint. The right-hand conjunct merely constrains T to be an instance of S.

Given the definition of `letrec f : S = λz.t₁ in t₂` as syntactic sugar, the above rule leads to the following derived constraint generation rule for `letrec`:

$$[\![\texttt{letrec f} : \texttt{S} = \lambda\texttt{z.t}_1 \texttt{ in t}_2 : \texttt{T}]\!] \quad = \quad \texttt{let f} : \texttt{S in } ([\![\lambda\texttt{z.t}_1 : \texttt{S}]\!] \wedge [\![\texttt{t}_2 : \texttt{T}]\!])$$

This rule is arguably quite natural. The program variable `f` is assigned the type scheme `S` throughout its scope, that is, both inside and outside of the function's definition. The function `λz.t₁` must itself have type scheme `S`. Last, `t₂` must have type `T`, as in every `let` construct.

1.10.10   EXERCISE [★★]: Prove that the derived constraint generation rule above is indeed valid.                                                     □

It is straightforward to prove that the extended language still enjoys subject reduction. The proof relies on the following lemma: if `t` has type scheme `S`, then every instance of `S` is also a valid type for `t`.

1.10.11   LEMMA: $[\![\texttt{t} : \texttt{S}]\!] \wedge \texttt{S} \preceq \texttt{T} \Vdash [\![\texttt{t} : \texttt{T}]\!]$.                        □

1.10.12   THEOREM [SUBJECT REDUCTION]: (R-FIX') ⊆ (⊑).                               □

The programming language Haskell (Hudak, Peyton Jones, Wadler, Boutel, Fairbairn, Fasel, Guzman, Hammond, Hughes, Johnsson, Kieburtz, Nikhil, Partain, and Peterson, 1992) offers polymorphic recursion. Interesting details about its typing rules may be found in (Jones, 1999).

It is worth pointing out that some restricted instances of type inference in the presence of polymorphic recursion are decidable. This is typically the case in certain program analyses, where a type derivation for the program is already available, and the goal is only to infer extra atomic annotations, such as binding time or strictness properties. Several papers that exploit this idea are (Dussart, Henglein, and Mossin, 1995a; Jensen, 1998; Rehof and Fähndrich, 2001).

## Universal types

ML-the-type-system enforces a strict stratification between types and type schemes, or, in other words, allows only prenex universal quantifiers inside types. We have pointed out earlier that there is good reason to do so: type inference for ML-the-type-system is decidable, while type inference for System F, which has no such restriction, is undecidable. Yet, this restriction comes at a cost in expressiveness: it prevents higher-order functions from accepting polymorphic function arguments, and forbids storing polymorphic functions inside data structures. Fortunately, it is in fact possible to circumvent the problem by requiring the programmer to supply additional type information.

The approach that we are about to describe is reminiscent of the way algebraic data type definitions allow circumventing the problems associated with equirecursive types (Section 1.9). Because we do not wish to extend the syntax of types with universal types of the form $\forall \bar{\mathtt{Y}}.\mathtt{T}$, we instead allow *universal type definitions*, of the form

$$\mathtt{D}\,\vec{\mathtt{X}} \approx \forall \bar{\mathtt{Y}}.\mathtt{T}$$

where $\mathtt{D}$ still ranges over data types. If $\mathtt{D}$ has signature $\vec{\kappa} \Rightarrow \star$, then the type variables $\vec{\mathtt{X}}$ must have kind $\vec{\kappa}$. The type $\mathtt{T}$ must have kind $\star$. The type variables $\bar{\mathtt{X}}$ and $\bar{\mathtt{Y}}$ are considered bound within $\mathtt{T}$, and the definition must be closed, that is, $\mathit{ftv}(\mathtt{T}) \subseteq \bar{\mathtt{X}}\bar{\mathtt{Y}}$ must hold. Last, the variance of the type constructor $\mathtt{D}$ must match its definition—a requirement stated as follows:

1.10.13    DEFINITION: Let $\mathtt{D}\,\vec{\mathtt{X}} \approx \forall \bar{\mathtt{Y}}.\mathtt{T}$ and $\mathtt{D}\,\vec{\mathtt{X}}' \approx \forall \bar{\mathtt{Y}}'.\mathtt{T}'$ be two $\alpha$-equivalent instances of a single universal type definition, such that $\bar{\mathtt{Y}} \mathbin{\#} \mathit{ftv}(\mathtt{T}')$ and $\bar{\mathtt{Y}}' \mathbin{\#} \mathit{ftv}(\mathtt{T})$. Then, $\mathtt{D}\,\vec{\mathtt{X}} \leq \mathtt{D}\,\vec{\mathtt{X}}' \Vdash \forall \bar{\mathtt{Y}}'.\exists \bar{\mathtt{Y}}.\mathtt{T} \leq \mathtt{T}'$ must hold.                    □

This requirement is analogous to that found in Definition 1.9.8. The idea is, if $\mathtt{D}\,\vec{\mathtt{X}}$ and $\mathtt{D}\,\vec{\mathtt{X}}'$ are comparable, then their unfoldings $\forall \bar{\mathtt{Y}}.\mathtt{T}$ and $\forall \bar{\mathtt{Y}}'.\mathtt{T}'$ should be comparable as well. The comparison between them is expressed by the constraint $\forall \bar{\mathtt{Y}}'.\exists \bar{\mathtt{Y}}.\mathtt{T} \leq \mathtt{T}'$, which may be read: *every instance of $\forall \bar{\mathtt{Y}}'.\mathtt{T}'$ is (a supertype of) an instance of $\forall \bar{\mathtt{Y}}.\mathtt{T}$*. Again, when subtyping is interpreted as equality, the requirement of Definition 1.10.13 is always satisfied; it becomes nontrivial only in the presence of true subtyping.

The effect of the universal type definition $\mathtt{D}\,\vec{\mathtt{X}} \approx \forall \bar{\mathtt{Y}}.\mathtt{T}$ is to enrich the programming language with a new construct:

$$\mathtt{v} ::= \ldots \mid \mathtt{pack}_{\mathtt{D}}\ \mathtt{v} \qquad \mathtt{t} ::= \ldots \mid \mathtt{pack}_{\mathtt{D}}\ \mathtt{t} \qquad \mathcal{E} ::= \ldots \mid \mathtt{pack}_{\mathtt{D}}\ \mathcal{E}$$

and with a new unary destructor $\mathtt{open}_{\mathtt{D}}$. Their operational semantics is as follows:

$$\mathtt{open}_{\mathtt{D}}\ (\mathtt{pack}_{\mathtt{D}}\ \mathtt{v}) \xrightarrow{\delta} \mathtt{v} \qquad\qquad \text{(R-OPEN-ALL)}$$

Intuitively, $\mathtt{pack}_{\mathtt{D}}$ and $\mathtt{open}_{\mathtt{D}}$ are the two coercions that witness the isomorphism between $\mathtt{D}\,\vec{\mathtt{X}}$ and $\forall \bar{\mathtt{Y}}.\mathtt{T}$. The value $\mathtt{pack}_{\mathtt{D}}\ \mathtt{v}$ behaves exactly like $\mathtt{v}$, except it is marked, as a hint to the typechecker. As a result, the mark must be removed using $\mathtt{open}_{\mathtt{D}}$ before the value can be used.

What are the typing rules for $\mathtt{pack}_{\mathtt{D}}$ and $\mathtt{open}_{\mathtt{D}}$? In System F, they would receive types $\forall \bar{\mathtt{X}}.(\forall \bar{\mathtt{Y}}.\mathtt{T}) \to \mathtt{D}\,\vec{\mathtt{X}}$ and $\forall \bar{\mathtt{X}}.\mathtt{D}\,\vec{\mathtt{X}} \to \forall \bar{\mathtt{Y}}.\mathtt{T}$, respectively. However, neither of these is a valid type scheme: both exhibit a universal quantifier under an arrow.

In the case of $\mathtt{pack}_{\mathtt{D}}$, which has been made a language construct rather than a constant, we work around the problem by embedding this universal

quantifier in the constraint generation rule:

$$\llbracket \mathtt{pack_D}\ \mathtt{t} : \mathtt{T'} \rrbracket \quad = \quad \exists \bar{X}.(\llbracket \mathtt{t} : \forall \bar{Y}.\mathtt{T} \rrbracket \wedge \mathtt{D}\,\vec{X} \le \mathtt{T'})$$

The rule implicitly requires that $\bar{X}$ be fresh for the left-hand side and that $\mathtt{D}\,\vec{X} \approx \forall \bar{Y}.\mathtt{T}$ be (an $\alpha$-variant of) the definition of $\mathtt{D}$. The left-hand conjunct requires $\mathtt{t}$ to have type scheme $\forall \bar{Y}.\mathtt{T}$. The notation $\llbracket \mathtt{t} : \mathtt{S} \rrbracket$ was defined on page 114. The right-hand conjunct states that a valid type for $\mathtt{pack_D}\ \mathtt{t}$ is (a supertype of) $\mathtt{D}\,\vec{X}$.

We deal with $\mathtt{open_D}$ as follows. Provided $\bar{X} \mathrel{\#} \bar{Y}$, we extend the initial environment $\Gamma_0$ with the binding $\mathtt{open_D} : \forall \bar{X}\bar{Y}.\mathtt{D}\,\vec{X} \to \mathtt{T}$. We have simply hoisted the universal quantifier outside of the arrow—a valid isomorphism in System F.

The proof of the subject reduction theorem must be extended with the following new case:

1.10.14   THEOREM [SUBJECT REDUCTION]:  (R-OPEN-ALL) $\subseteq$ ($\sqsubseteq$).                      □

*Proof:*   We have

$$
\begin{aligned}
&\quad\ \ \mathsf{let}\ \Gamma_0\ \mathsf{in}\ \llbracket \mathtt{open_D}\ (\mathtt{pack_D}\ \mathtt{v}) : \mathtt{T_0} \rrbracket \\
&\equiv\ \ \mathsf{let}\ \Gamma_0\ \mathsf{in}\ \exists \mathtt{Z}.(\mathtt{open_D} \preceq \mathtt{Z} \to \mathtt{T_0} \wedge \llbracket \mathtt{pack_D}\ \mathtt{v} : \mathtt{Z} \rrbracket) &\textbf{(1)} \\
&\equiv\ \ \mathsf{let}\ \Gamma_0\ \mathsf{in}\ \exists \mathtt{Z}.(\exists \bar{X}'\bar{Y}'.(\mathtt{D}\,\vec{X}' \to \mathtt{T'} \le \mathtt{Z} \to \mathtt{T_0}) \wedge \exists \bar{X}.(\llbracket \mathtt{v} : \forall \bar{Y}.\mathtt{T} \rrbracket \wedge \mathtt{D}\,\vec{X} \le \mathtt{Z})) &\textbf{(2)} \\
&\equiv\ \ \mathsf{let}\ \Gamma_0\ \mathsf{in}\ \exists \bar{X}\bar{X}'\bar{Y}'.(\llbracket \mathtt{v} : \forall \bar{Y}.\mathtt{T} \rrbracket \wedge \mathtt{D}\,\vec{X} \le \mathtt{D}\,\vec{X}' \wedge \mathtt{T'} \le \mathtt{T_0}) &\textbf{(3)} \\
&\Vdash\ \ \mathsf{let}\ \Gamma_0\ \mathsf{in}\ \exists \bar{X}\bar{Y}\bar{X}'\bar{Y}'.(\llbracket \mathtt{v} : \forall \bar{Y}.\mathtt{T} \rrbracket \wedge \mathtt{T} \le \mathtt{T'} \wedge \mathtt{T'} \le \mathtt{T_0}) &\textbf{(4)} \\
&\Vdash\ \ \mathsf{let}\ \Gamma_0\ \mathsf{in}\ \exists \bar{X}\bar{Y}\bar{X}'\bar{Y}'.\llbracket \mathtt{v} : \mathtt{T_0} \rrbracket &\textbf{(5)} \\
&\equiv\ \ \mathsf{let}\ \Gamma_0\ \mathsf{in}\ \llbracket \mathtt{v} : \mathtt{T_0} \rrbracket &\textbf{(6)}
\end{aligned}
$$

where (1) is by definition of constraint generation for applications and for constants; $\mathtt{Z}$ is fresh; (2) is by definition of constraint generation for $\mathtt{pack_D}$ and $\mathtt{open_D}$, where $\mathtt{D}\,\vec{X} \approx \forall \bar{Y}.\mathtt{T}$ and $\mathtt{D}\,\vec{X}' \approx \forall \bar{Y}'.\mathtt{T'}$ are two $\alpha$-equivalent instances of the definition of $\mathtt{D}$; $\bar{X}$, $\bar{Y}$, $\bar{X}'$, and $\bar{Y}'$ are fresh and satisfy $\bar{Y} \mathrel{\#} ftv(\mathtt{T'})$ and $\bar{Y}' \mathrel{\#} ftv(\mathtt{T})$; (3) is by C-EXAND, C-ARROW, and C-EXTRANS, which allows eliminating $\mathtt{Z}$; (4) is by Definition 1.10.13, Lemma 1.10.1, and C-EXAND; (5) is by Lemmas 1.10.11 and 1.6.3; (6) is by C-EX*.                      □

The proof of (R-CONTEXT) $\subseteq$ ($\sqsubseteq$) must also be extended with a new subcase, corresponding the new production $\mathcal{E} ::= \ldots \mid \mathtt{pack_D}\ \mathcal{E}$. If the language is pure, this is straightforward. In the presence of side effects, however, this subcase fails, because universal and existential quantifiers in constraints do not commute. The problem is then avoided by restricting $\mathtt{pack_D}$ to values, as in Definition 1.7.7.

This approach to extending ML-the-type-system with universal (or existential—see below) types has been studied in (Läufer and Odersky, 1994;

Rémy, 1994; Odersky and Läufer, 1996; Shields and Peyton Jones, 2002). Läufer and Odersky have suggested combining universal or existential type declarations with algebraic data type definitions. This allows suppressing the cumbersome $\mathtt{pack}_\mathtt{D}$ and $\mathtt{open}_\mathtt{D}$ constructs; instead, one simply uses the standard syntax for constructing and deconstructing variants and records.

### Existential types

Existential types (TAPL Chapter 24) are close cousins of universal types, and may be introduced into ML-the-type-system in the same manner. Actually, existential types have been introduced in ML-the-type-system before universal types. We give a brief description of this extension, insisting mainly on the differences with the case of universal types.

We now allow *existential type definitions*, of the form $\mathtt{D}\,\vec{\mathtt{X}} \approx \exists\bar{\mathtt{Y}}.\mathtt{T}$. The conditions required of a well-formed definition are unchanged, except the variance requirement, which is dual:

1.10.15   DEFINITION: Let $\mathtt{D}\,\vec{\mathtt{X}} \approx \exists\bar{\mathtt{Y}}.\mathtt{T}$ and $\mathtt{D}\,\vec{\mathtt{X}}' \approx \exists\bar{\mathtt{Y}}'.\mathtt{T}'$ be two $\alpha$-equivalent instances of a single existential type definition, such that $\bar{\mathtt{Y}}\ \#\ ftv(\mathtt{T}')$ and $\bar{\mathtt{Y}}'\ \#\ ftv(\mathtt{T})$. Then, $\mathtt{D}\,\vec{\mathtt{X}} \leq \mathtt{D}\,\vec{\mathtt{X}}' \Vdash \forall\bar{\mathtt{Y}}.\exists\bar{\mathtt{Y}}'.\mathtt{T} \leq \mathtt{T}'$ must hold.                          □

The effect of this existential type definition is to enrich the programming language with a new unary constructor $\mathtt{pack}_\mathtt{D}$ and with a new construct: $\mathtt{t} ::= \ldots \mid \mathtt{open}_\mathtt{D}\ \mathtt{t}\ \mathtt{t}$ and $\mathcal{E} ::= \ldots \mid \mathtt{open}_\mathtt{D}\ \mathcal{E}\ \mathtt{t} \mid \mathtt{open}_\mathtt{D}\ \mathtt{v}\ \mathcal{E}$. Their operational semantics is as follows:

$$\mathtt{open}_\mathtt{D}\ (\mathtt{pack}_\mathtt{D}\mathtt{v}_1)\ \mathtt{v}_2 \longrightarrow \mathtt{v}_2\ \mathtt{v}_1 \qquad\qquad \text{(R-OPEN-EX)}$$

In the literature, the second argument of $\mathtt{open}_\mathtt{D}$ is often required to be a $\lambda$-abstraction $\lambda\mathtt{z}.\mathtt{t}$, so the construct becomes $\mathtt{open}_\mathtt{D}\ \mathtt{t}\ (\lambda\mathtt{z}.\mathtt{t})$, often written $\mathtt{open}_\mathtt{D}\ \mathtt{t}\ \mathtt{as}\ \mathtt{z}\ \mathtt{in}\ \mathtt{t}$.

Provided $\bar{\mathtt{X}}\ \#\ \bar{\mathtt{Y}}$, we extend the initial environment $\Gamma_0$ with the binding $\mathtt{pack}_\mathtt{D} : \forall\bar{\mathtt{X}}\bar{\mathtt{Y}}.\mathtt{T} \to \mathtt{D}\,\vec{\mathtt{X}}$. The constraint generation rule for $\mathtt{open}_\mathtt{D}$ is as follows:

$$[\![\mathtt{open}_\mathtt{D}\ \mathtt{t}_1\ \mathtt{t}_2 : \mathtt{T}']\!] \quad = \quad \exists\bar{\mathtt{X}}.([\![\mathtt{t}_1 : \mathtt{D}\,\vec{\mathtt{X}}]\!] \wedge [\![\mathtt{t}_2 : \forall\bar{\mathtt{Y}}.\mathtt{T} \to \mathtt{T}']\!])$$

The rule implicitly requires that $\bar{\mathtt{X}}$ be fresh for the left-hand side, that $\bar{\mathtt{Y}}$ be fresh for $\mathtt{T}'$, and that $\mathtt{D}\,\vec{\mathtt{X}} \approx \forall\bar{\mathtt{Y}}.\mathtt{T}$ be (an $\alpha$-variant of) the definition of $\mathtt{D}$. The left-hand conjunct simply requires $\mathtt{t}_1$ to have type $\mathtt{D}\,\vec{\mathtt{X}}$. The right-hand conjunct states that the function $\mathtt{t}_2$ must be prepared to accept an argument of type $\mathtt{T}$, for *any* $\bar{\mathtt{Y}}$, and produce a result of the expected type $\mathtt{T}'$. In other words, $\mathtt{t}_2$ must be a polymorphic function.

The type scheme of existential $\mathtt{pack}_\mathtt{D}$ resembles that of universal $\mathtt{open}_\mathtt{D}$, while the constraint generation rule for existential $\mathtt{open}_\mathtt{D}$ is a close cousin

of that for universal $\mathtt{pack_D}$. Thus, the duality between universal and existential types is rather strong. The main difference lies in the fact that the existential $\mathtt{open_D}$ construct is *binary*, rather than unary, so as to limit the scope of the newly introduced type variables $\bar{\mathtt{Y}}$. The duality may be better understood by studying the encoding of existential types in terms of universal types (Reynolds, 1983b).

As expected, R-OPEN-EX preserves types.

1.10.16  THEOREM [SUBJECT REDUCTION]: $(\text{R-OPEN-EX}) \subseteq (\sqsubseteq)$. ☐

1.10.17  EXERCISE [★★, ↛]: Prove Theorem 1.10.16. The proof is analogous, although not identical, to that of Theorem 1.10.14. ☐

In the presence of side effects, the new production $\mathcal{E} ::= \ldots \mid \mathtt{open_D} \; \mathtt{v} \; \mathcal{E}$ is problematic. The standard workaround is to restrict the second argument to $\mathtt{open_D}$ to be a value.

## 1.11 Rows

In Section 1.9, we have shown how to extend ML-the-programming-language with algebraic data types, that is, variant and record type definitions, which we now refer to as *simple*. This mechanism has a severe limitation: two distinct definitions must define incompatible types. As a result, one cannot hope to write code that uniformly operates over variants or records of different shapes, because the type of such code is not even expressible.

For instance, it is impossible to express the type of the *polymorphic record access* operation, which retrieves the value stored at a particular field $\ell$ inside a record, *regardless* of which other fields are present. Indeed, if the label $\ell$ appears with type $\mathtt{T}$ in the definition of the simple record type $\mathtt{D} \; \vec{\mathtt{X}}$, then the associated record access operation has type $\forall \bar{\mathtt{X}}.\mathtt{D} \; \vec{\mathtt{X}} \to \mathtt{T}$. If $\ell$ appears with type $\mathtt{T}'$ in the definition of another simple record type, say $\mathtt{D}' \; \vec{\mathtt{X}}'$, then the associated record access operation has type $\forall \bar{\mathtt{X}}'.\mathtt{D}' \; \vec{\mathtt{X}}' \to \mathtt{T}'$; and so on. The most precise type scheme that subsumes all of these incomparable type schemes is $\forall \mathtt{XY}.\mathtt{X} \to \mathtt{Y}$. It is, however, not a sound type scheme for the record access operation. Another powerful operation whose type is currently not expressible is *polymorphic record extension*, which copies a record and stores a value at field $\ell$ in the copy, possibly creating the field if it did not previously exist, again *regardless* of which other fields are present. (If $\ell$ was known to previously exist, the operation is known as *polymorphic record update*.)

In order to assign types to polymorphic record operations, we must do away with record type definitions: we must replace *named* record types, such as $\mathtt{D} \; \vec{\mathtt{X}}$, with *structural* record types that provide a direct description of the record's

domain and contents. (Following the analogy between a record and a partial function from labels to values, we use the word *domain* to refer to the set of fields that are defined in a record.) For instance, a product type is structural: the type $T_1 \times T_2$ is the (undeclared) type of pairs whose first component has type $T_1$ and whose second component has type $T_2$. Thus, we wish to design record types that behave very much like product types. In doing so, we face two orthogonal difficulties. First, as opposed to pairs, records may have different domains. Because the type system must statically ensure that no undefined field is accessed, information about a record's domain must be made part of its type. Second, because we suppress record type definitions, labels must now be predefined. However, for efficiency and modularity reasons, it is impossible to explicitly list *every* label in existence in every record type.

In what follows, we explain how to address the first difficulty in the simple setting of a finite set of labels. Then, we introduce *rows*, which allow dealing with an infinite set of labels, and address the second difficulty. We define the syntax and logical interpretation of rows, study the new constraint equivalence laws that arise in their presence, and extend the first-order unification algorithm with support for rows. Then, we review several applications of rows, including polymorphic operations on records, variants, and objects, and discuss alternatives to rows.

### Records with finite carrier

Let us temporarily assume that $\mathcal{L}$ is finite. In fact, for the sake of definiteness, let us assume that $\mathcal{L}$ is $\{\ell_a, \ell_b, \ell_c\}$.

To begin, let us consider only *full* records, whose domain is exactly $\mathcal{L}$—in other words, tuples indexed by $\mathcal{L}$. To describe them, it is natural to introduce a type constructor record of signature $\star \otimes \star \otimes \star \Rightarrow \star$. The type record $T_a$ $T_b$ $T_c$ represents all records where the field $\ell_a$ (resp. $\ell_b$, $\ell_c$) contains a value of type $T_a$ (resp. $T_b$, $T_c$). Please note that record is nothing but a product type constructor of arity $3$. The basic operations on records, namely *creation* of a record out of a default value, which is stored into every field, *update* of a particular field (say, $\ell_b$), and *access* to a particular field (say, $\ell_b$), may be assigned the following type schemes:

$$
\begin{aligned}
\{\cdot\} : &\quad \forall X.X \rightarrow \text{record } X \ X \ X \\
\{\cdot \text{ with } \ell_b = \cdot\} : &\quad \forall X_a X_b X_b' X_c.\text{record } X_a \ X_b \ X_c \rightarrow X_b' \rightarrow \text{record } X_a \ X_b' \ X_c \\
\cdot.\{\ell_b\} : &\quad \forall X_a X_b X_c.\text{record } X_a \ X_b \ X_c \rightarrow X_b
\end{aligned}
$$

Here, polymorphism allows updating or accessing a field without knowledge of the types of the other fields. This flexibility is made possible by the property that all record types are formed using a single record type constructor.

This is fine, but in general, the domain of a record is not necessarily $\mathcal{L}$: it may be a subset of $\mathcal{L}$. How may we deal with this fact, while maintaining the above key property? A naïve approach consists in encoding arbitrary records in terms of full records, using the standard algebraic data type option, whose definition is option X $\approx$ pre X + abs. We use pre for *present* and abs for *absent*: indeed, a field that is defined with value v is encoded as a field with value pre v, while an undefined field is encoded as a field with value abs. Thus, an arbitrary record whose fields, *if present*, have types $T_a$, $T_b$, and $T_c$, respectively, may be encoded as a full record of type record (option $T_a$) (option $T_b$) (option $T_c$). This naive approach suffers from a serious drawback: record types still contain no domain information. As a result, field access must involve a dynamic check, so as to determine whether the desired field is present: in our encoding, this corresponds to the use of case$_\mathsf{option}$.

To avoid this overhead and increase programming safety, we must move this check from runtime to compile time. In other words, we must make the type system aware of the difference between pre and abs. To do so, we replace the definition of option by two separate algebraic data type definitions, namely pre X $\approx$ pre X and abs $\approx$ abs. In other words, we introduce a unary type constructor pre, whose only associated data constructor is pre, and a nullary type constructor abs, whose only associated data constructor is abs. Record types now contain domain information: for instance, a record of type record abs (pre $T_b$) (pre $T_c$) must have domain $\{\ell_b, \ell_c\}$. Thus, the type of a field tells whether it is defined. Since the type pre has no data constructors other than pre, the accessor pre$^{-1}$, whose type is $\forall$X.pre X $\rightarrow$ X, and which allows retrieving the value stored in a field, cannot fail. Thus, the dynamic check has been eliminated.

To complete the definition of our encoding, we now define operations on arbitrary records in terms of operations on full records. To distinguish between the two, we write the former with angle braces, instead of curly braces. The *empty record* $\langle\rangle$, where all fields are undefined, may be defined as {abs}. *Extension* at a particular field (say, $\ell_b$) $\langle\cdot$ with $\ell_b = \cdot\rangle$ is defined as $\lambda$r.$\lambda$z.{r with $\ell_b = $ pre z}. *Access* at a particular field (say, $\ell_b$) $\cdot.\langle\ell_b\rangle$ is defined as $\lambda$z.pre$^{-1}$z.{$\ell_b$}. It is straightforward to check that these operations have the following principal type schemes:

$$
\begin{aligned}
\langle\rangle : \quad & \text{record abs abs abs} \\
\langle\cdot \text{ with } \ell_b = \cdot\rangle : \quad & \forall X_a\, X_b\, X_b'\, X_c.\text{record } X_a\ X_b\ X_c \rightarrow X_b' \rightarrow \text{record } X_a\ (\text{pre } X_b')\ X_c \\
\cdot.\langle\ell_b\rangle : \quad & \forall X_a\, X_b\, X_c.\text{record } X_a\ (\text{pre } X_b)\ X_c \rightarrow X_b
\end{aligned}
$$

It is important to notice that the type schemes associated with extension and access at $\ell_b$ are polymorphic in $X_a$ and $X_c$, which now means that *these operations are insensitive not only to the type, but also to the presence or*

*absence of the fields $\ell_a$ and $\ell_c$.* Furthermore, extension is polymorphic in $\mathtt{X}_b$, which means that it is insensitive to the presence or absence of the field $\ell_b$ in its argument. The subterm $\mathsf{pre}\ \mathtt{X}'_b$ in its result type reflects the fact that $\ell_b$ is defined in the extended record. Conversely, the subterm $\mathsf{pre}\ \mathtt{X}_b$ in the type of the access operation reflects the requirement that $\ell_b$ be defined in its argument.

Our encoding of arbitrary records in terms of full records was carried out for pedagogical purposes. In practice, no such encoding is necessary: the *data* constructors $\mathsf{pre}$ and $\mathsf{abs}$ have no machine representation, and the compiler is free to lay out records in memory in an efficient manner. The encoding is interesting, however, because it provides a natural way of introducing the *type* constructors $\mathsf{pre}$ and $\mathsf{abs}$, which play an important role in our treatment of polymorphic record operations.

We remark that, in our encoding, the arguments of the type constructor $\mathsf{record}$ are expected to be either type variables or formed with $\mathsf{pre}$ or $\mathsf{abs}$, while, conversely, the type constructors $\mathsf{pre}$ and $\mathsf{abs}$ are not intended to appear anywhere else. It is possible to enforce this invariant using kinds. In addition to $\star$, let us introduce the kind $\diamond$ of *field types*. Then, let us adopt the following signatures: $\mathsf{pre} \colon \star \Rightarrow \diamond$, $\mathsf{abs} \colon \diamond$, and $\mathsf{record} \colon \diamond \otimes \diamond \otimes \diamond \Rightarrow \star$.

1.11.1   EXERCISE [$\bigstar$, RECOMMENDED, $\nrightarrow$]: Check that the three type schemes given above are well-kinded. What is the kind of each type variable?   □

1.11.2   EXERCISE [$\bigstar\bigstar$, RECOMMENDED, $\nrightarrow$]: Our $\mathsf{record}$ types contain information about every field, regardless of whether it is defined: we encode definedness information within the type of each field, using the type constructors $\mathsf{pre}$ and $\mathsf{abs}$. A perhaps more natural approach would be to introduce a family of record type constructors, indexed by the subsets of $\mathcal{L}$, so that the types of records with different domains are formed with different constructors. For instance, the empty record would have type $\{\}$; a record that defines the field $\ell_a$ only would have a type of the form $\{\ell_a : \mathtt{T}_a\}$; a record that defines the fields $\ell_b$ and $\ell_c$ only would have a type of the form $\{\ell_b : \mathtt{T}_b; \ell_c : \mathtt{T}_c\}$; and so on. Assuming that the type discipline is Damas and Milner's (that is, assuming an equality-only syntactic model), would it be possible to assign satisfactory type schemes to polymorphic record access and extension? Would it help to equip record types with a nontrivial subtyping relation?   □

### Records with infinite carrier

Finite records are insufficient both from practical and theoretical points of view. In practice, the set of labels could become very large, making the type of every record as large as the set of labels itself, even if only a few labels are

actually defined. In principle, the set of labels could even be infinite. Actually, in modular programs the whole set of labels may not be known in advance, which amounts in some way to working with an infinite set of labels. Thus, records must be drawn from an infinite set of labels—whether their domains are finite or infinite. Still, we can restrict our attention to records that are almost constant, that is, records where only a finite number of fields differ. With this restriction, full records (defined everywhere) can always be built by giving explicit definitions for a finite number of fields and a default value for all other fields, as in the finite case. For instance, the record $\{\{\{\texttt{false}\}$ with $\ell = 1\}$ with $\ell' = \texttt{true}\}$ is the record equal to $\texttt{true}$ on field $\ell'$, to 1 on field $\ell$, and to $\texttt{false}$ on any other field.

   Types of records are functions from labels to types, called *rows*. However, for sake of generality, we use a unary type constructor, say $\Pi$, as an indirection between rows and record types. Moreover, we further restrict our attention to the case where rows are also almost constant. (The fact that the property holds for record values does not imply that it also holds for record types, for the default value of some record could have a polymorphic type, and one could wish to see each field with a different instance of this polymorphic type. So this is a true restriction, but a reasonable one.) Thus, rows can also be represented by giving explicit types for a finite number of fields and a default type for all other fields. We write $\partial\texttt{T}$ the row whose type is $\texttt{T}$ on every field, and $(\ell : \texttt{T} \,;\, \texttt{T}')$ the row whose type is $\texttt{T}$ on field $\ell$ and $\texttt{T}'$ on other fields. Formally, $\partial$ is a unary type constructor and $\ell$ is a family of binary type constructors, written with syntactic sugar $(\ell : \cdot \,;\, \cdot)$. For example, $\Pi(\ell : \texttt{bool} \,;\, (\ell' : \texttt{int} \,;\, \partial\texttt{bool}))$ is a record type that describes records whose field $\ell$ carries a value of type $\texttt{bool}$, field $\ell'$ carries a value of type $\texttt{int}$, and all other fields carry values of type $\texttt{bool}$. In fact, this is a sound type for the record defined above. In fact, the type $\Pi(\ell' : \texttt{int} \,;\, (\ell : \texttt{bool} \,;\, \partial\texttt{bool}))$ should also be a sound type for this record, since the order in which fields are specified should not matter. We actually treat both types as equivalent. Furthermore, the row $(\ell : \texttt{bool} \,;\, \partial\texttt{bool})$, which stands for $\texttt{bool}$ on field $\ell$ and $\partial\texttt{bool}$ everywhere else, must also be equivalent to $\partial\texttt{bool}$, which stands for $\texttt{bool}$ everywhere.

   A record type may also contain type variables. For instance, the record $\lambda\texttt{z}.\{\texttt{z}\}$ that maps any value $\texttt{v}$ to a record with the default value $\texttt{v}$ has type $\texttt{X} \to \Pi(\partial\texttt{X})$. Projections of this record on any field will return a value of the same type $\texttt{X}$. By comparison, the function that reads some field $\ell$ of its (record) argument has type $\Pi(\ell : \texttt{X} \,;\, \texttt{Y}) \to \texttt{X}$: this says that the argument must be a record where field $\ell$ has type $\texttt{X}$ and other fields may have *any* type. Variable $\texttt{Y}$ is called a *row variable*, since it can be instantiated to any row. For instance, $\texttt{Y}$ can be instantiated to $(\ell' : \texttt{int} \,;\, \texttt{Y}')$ and as a result this function can be applied to the record above. Conversely, the row $\partial\texttt{X}$, which is equal to $(\ell' : \texttt{X} \,;\, \texttt{X})$, can

only be instantiated to rows of the form $\partial\mathtt{T}$, which are equal to $(\ell' : \mathtt{T} \; ; \; \mathtt{T})$, that is, to constant rows.

### Syntax of row types

Let $\mathcal{L}$ be a denumerable collection of labels. We write $\ell.L$ for $\{\ell\} \uplus L$, which implies $\ell \notin L$. We first introduce kinds, so as to distinguish rows such as $(\ell : \mathtt{int} \; ; \; \partial\mathtt{bool})$ from *basic types*, such as $\mathtt{int}$ or $\mathtt{int} \to \mathtt{int}$.

1.11.3 DEFINITION [ROW KINDS]: Let *row kinds* be composed of a particular kind *Type* and the collection of kinds $Row(L)$ where $L$ ranges over finite subsets of $\mathcal{L}$. We use letter $s$ to range over row kinds. □

Intuitively, a row of kind $Row(L)$ is a function of domain $\mathcal{L} \setminus L$ to types. That is, $L$ specifies the set of labels that the row *must not* define. For instance, the (basic) type $\Pi(\ell : \mathtt{int} \; ; \; \mathtt{X})$ has kind *Type*, the row $(\ell : \mathtt{int} \; ; \; \mathtt{X})$ has kind $Row(\emptyset)$ provided $\mathtt{X}$ has kind $Row(\{\ell\})$.

To remain abstract the definition of rows is parameterized by a signature $\mathcal{S}_0$ for building basic types and a signature $\mathcal{S}_1$ for building rows. From those two signatures, we then define a new signature $\mathcal{S}$ that completely specifies the set of types. However, the signature $\mathcal{S}$ must superimposed row kinds on top of the (basic) kinds of the two input signatures $\mathcal{S}_0$ and $\mathcal{S}_1$. We use product signatures to enlighten this process. More precisely, we build a product signature from two signatures $K \Rightarrow \kappa$ and $K' \Rightarrow \kappa'$ with the following notations: we write $\kappa.\kappa'$ for the pair $(\kappa, \kappa')$; $K.\kappa$ for the mapping $(d \mapsto K(d).\kappa)^{d \in dom(K)}$; $(K \Rightarrow \kappa).\kappa'$ for the kind signature $K.\kappa \Rightarrow \kappa.\kappa'$; and symmetrically, we write $\kappa.K'$ and $\kappa.(K' \Rightarrow \kappa')$. The signature $\mathcal{S}$ reuses the same input type constructors as $\mathcal{S}_0$ and $\mathcal{S}_1$, but at different row kinds. We use superscripts to provide copies of type constructors at different kinds, and thus avoid overloading of kinds.

1.11.4 DEFINITION [ROW EXTENSION OF A SIGNATURE]: Let $\mathcal{S}_0$ and $\mathcal{S}_1$ be signatures where all symbols of $\mathcal{S}_1$ are unary. The row extension of $\mathcal{S}_0$ with $\mathcal{S}_1$ is the signature $\mathcal{S}$ defined as follows where $\kappa$ ranges over basic kinds (those used in $\mathcal{S}_0$ and $\mathcal{S}_1$) and $s$ ranges over row kinds:

| $F \in dom(\mathcal{S})$ | Signature | Conditions |
|---|---|---|
| $G^s$ | $(K \Rightarrow \kappa).s$ | $(G : K \Rightarrow \kappa) \in \mathcal{S}_0$ |
| $H$ | $K.Row(\emptyset) \Rightarrow \kappa.Type$ | $(H : K \Rightarrow \kappa) \in \mathcal{S}_1$ |
| $\partial^{\kappa,L}$ | $\kappa.(Type \Rightarrow Row(L))$ | |
| $\ell^{\kappa,L}$ | $\kappa.(Type \otimes Row(\ell.L) \Rightarrow Row(L))$ | $\ell \notin L$ |

□

We usually write $\ell^{\kappa,L} : \mathtt{t}_1 \; ; \; \mathtt{t}_2$ instead of $\ell^{\kappa,L} \; \mathtt{t}_1 \; \mathtt{t}_2$ and let this symbol be right-associative. We often drop the superscripts of type constructors since,

for any type expression T, superscripts can be unambiguously recovered from
the kind of T.

1.11.5   Example:  Let us assume there is a single basic kind $\star$ and that $\mathcal{S}_1$ contain a
unique type constructor $\Pi$ (hence of kind $\star \Rightarrow \star$). An example of row type is
X$_0 \rightarrow \Pi(\ell_1 : G$ ; (Y $\rightarrow \partial$X$_0$)). With all superscripts annotations, this type is

$$\text{X}_0 \rightarrow^{\star, \, Type} \Pi(\ell_1^{\,\star, Row(\emptyset)} : G^{\,Type} \text{ ; } (\text{Y} \rightarrow^{\star, Row(\{\ell_1\})} \partial^{\star, Row(\{\ell_1\})} \text{X}_0)).$$

Intuitively, this is the type of a function that takes a value of type X$_0$ and
returns a record where field $\ell_1$ has type $G$ and all other fields are functions
that given a value of an arbitrary type would returns a value of (the same) type
X$_0$. An instance of this type is X$_0 \rightarrow \Pi(\ell_1 : G$ ; (($\ell_2$ : Y$_2$ ; Y$'$) $\rightarrow (\ell_2$ : X$_0$ ; $\partial$X$_0$))),
obtained by instantiating row variable Y and by expanding the constant row
$\partial$X$_0$. As shown below, this type is actually equivalent to X$_0 \rightarrow \Pi(\ell_1 : G$ ;
$\ell_2$ : Y$_2 \rightarrow$ X$_0$ ; (Y$' \rightarrow \partial$X$_0$)), by distributivity of type constructor $\rightarrow$ other type
constructor $\ell_2$. Please, note again the difference between Y, which is a row
variable that can expand to different type variables on different labels, and
$\partial$X, which is a constant row that expands to the same type variable X on all
labels.                                                                                            □

1.11.6   Example [ Ill-kinded expression ]:  Under the assumptions of the previ-
ous example, the expression X $\rightarrow \Pi($X$)$ is not a row type, since variable X
cannot simultaneously be of row kind *Type* and $Row(\emptyset)$ as required by its two
occurrences, from left to right respectively. The expression $(\ell$ : X ; $\ell$ : X$'$ ; X$'')$ is
also ill-kinded, for the inner expression $(\ell$ : X$'$ ; X$'')$ of row kind $Row(L)$ with
$\ell \notin L$ cannot also have row kind $Row(\{\ell\})$, as required by its occurrence in
the whole expression. Indeed, row kinds prohibit multiple definitions of the
same label, as well as using rows in place of basic types and conversely. Notice
that $\Pi(\Pi($X$))$ is also ill-formed, since type constructors of $\mathcal{S}_1$ are not lifted to
row kinds and thus cannot appear in rows, except under the type constructor
$\partial$, hence as basic types.                                                                   □

1.11.7   Exercise [★★★,↛]:  Design an algorithm that infers superscripts of type
constructors of a type expression from its kind. Can the kind of the expression
be inferred as well? Can you give an algorithm to check that type expressions
are well-kinded when both the superscripts of type constructors and the kinds
of the whole type expressions are omitted?                                                         □

### Meaning of rows

As mentioned above, a row of kind $Row(L)$ stands for a function from $\mathcal{L} \setminus L$
to types. Actually, it is simpler to represent this function explicitly as an

infinitely branching tree in the model. For this purpose, we use a collection of constructors $L$ of (infinite but denumerable) arity $\mathcal{L} \setminus L$.

1.11.8    DEFINITION [ROW MODEL]: Let $\mathcal{S}$ be the row extension of a signature $\mathcal{S}_0$ with a signature $\mathcal{S}_1$. Let $\mathcal{S}_\mathcal{M}$ be the following signature, where $\kappa$ ranges over basic kinds and $L$ ranges over finite subsets of $\mathcal{L}$:

| $F \in dom(\mathcal{S}_\mathcal{M})$ | Signature | Conditions |
|---|---|---|
| $G$ | $(K \Rightarrow \kappa).\, Type$ | $(G : K \Rightarrow \kappa) \in \mathcal{S}_0$ |
| $H$ | $K.Row(\emptyset) \Rightarrow \kappa.\, Type$ | $(H : K \Rightarrow \kappa) \in \mathcal{S}_1$ |
| $L^\kappa$ | $\kappa.(\, Type^{\mathcal{L} \setminus L} \Rightarrow Row(L))$ | |

Let $\mathcal{M}_\kappa$ consist of the regular trees $t$ built over the signature $\mathcal{S}_\mathcal{M}$ such that $t(\epsilon)$ has image kind $\kappa$. We interpret a type constructor $F$ of signature $K \Rightarrow \kappa.s$ as a function that maps $T \in \mathcal{M}_K$ to $t \in \mathcal{M}_{\kappa.s}$ defined by cases on $F$ and on the basic kind $\kappa$.

| $F \in \mathcal{S}$ | $t(\epsilon)$ | For $d \in dom(K)$ and $\ell \in \mathcal{L} \setminus L, \ell \neq \ell_0$. |
|---|---|---|
| $G^{\,Type}$ | $G$ | $t/d = T(d)$ |
| $H$ | $H$ | $t/1 = T(1)$ |
| $G^{\,Row(L)}$ | $L^\kappa$ | $t(\ell) = G \wedge t/(\ell \cdot d) = T(d)/\ell$ |
| $\partial^{\kappa,L}$ | $L^\kappa$ | $t/\ell = T(1)$ |
| $\ell_0^{\kappa,L}$ | $L^\kappa$ | $t/\ell_0 = T(1) \wedge t/\ell = T(2)/\ell$ |

In the presence of subtyping, we let type constructors $G$ and $H$ behave in $\mathcal{S}_\mathcal{M}$ as in $\mathcal{S}_0$ and $\mathcal{S}_1$ and type constructors $L^\kappa$ be isolated and covariant in every position. Models that define ground types and interpret type constructors in this manner are referred to as *row* models.                                    □

### Reasoning with row types

In this section, we assume a subtyping model. All reasoning principles also apply to the equality-only model, which is a subcase of the subtyping model.

The meaning of rows has been carefully defined so as to be independent of some syntactical choices. In particular, the order in which the types of significant fields have been declared leaves the meaning of rows unchanged. This is formally stated by the following Lemma.

1.11.9    LEMMA:  The equations of Figure 1-17 are equivalent to `true`.                     □

*Proof:*   Each equation can be considered independently. It suffices to see that any ground assignment $\phi$ sends both sides of the equation to the same element

$$(\ell_1 : \mathtt{T}_1 \ ; \ \ell_2 : \mathtt{T}_2 \ ; \ \mathtt{T}_3) = (\ell_2 : \mathtt{T}_2 \ ; \ \ell_1 : \mathtt{T}_1 \ ; \ \mathtt{T}_3) \hspace{2em} \text{(C-Row-LL)}$$

$$\partial \mathtt{T} = (\ell : \mathtt{T} \ ; \ \partial \mathtt{T}) \hspace{2em} \text{(C-Row-DL)}$$

$$\partial(G \ \mathtt{T}_1 \ \ldots \ \mathtt{T}_n) = G \ \partial \mathtt{T}_1 \ \ldots \ \partial \mathtt{T}_n \hspace{2em} \text{(C-Row-DG)}$$

$$G \ (\ell : \mathtt{T}_1 \ ; \ \mathtt{T}'_1) \ \ldots \ (\ell : \mathtt{T}_n \ ; \ \mathtt{T}'_n) = (\ell : G \ \mathtt{T}_1 \ \ldots \ \mathtt{T}_n \ ; \ G \ \mathtt{T}'_1 \ \ldots \ \mathtt{T}'_n) \hspace{2em} \text{(C-Row-GL)}$$

**Figure 1-17: Equational reasoning for row types.**

$$(\ell_1 : \mathtt{T}_1 \ ; \ \mathtt{T}'_1) = (\ell_2 : \mathtt{T}_2 \ ; \ \mathtt{T}'_2) \quad \equiv \quad \exists \mathtt{X}. (\mathtt{T}'_1 = (\ell_2 : \mathtt{T}_2 \ ; \ \mathtt{X}) \wedge \mathtt{T}'_2 = (\ell_1 : \mathtt{T}_1 \ ; \ \mathtt{X})) \hspace{1em} \text{(C-Mute-LL)}$$
$$\text{if } \mathtt{X} \ \# \ ftv(\mathtt{T}_1, \mathtt{T}'_1, \mathtt{T}_2, \mathtt{T}'_2) \wedge \ell_1 \neq \ell_2$$

$$(\ell : \mathtt{T} \ ; \ \mathtt{T}') = G \ \mathtt{T}^I_i \quad \equiv \quad \exists (\mathtt{X}_i, \mathtt{X}'_i)^I. (\mathtt{T} = G \ \mathtt{X}^I_i \wedge \mathtt{T}' = G \ {\mathtt{X}'}^I_i \wedge \mathtt{T}_i = (\ell : \mathtt{X}_i \ ; \ \mathtt{X}'_i))$$
$$\text{if } (\mathtt{X}_i, \mathtt{X}_i)^I \ \# \ ftv(\mathtt{T}, \mathtt{T}', \mathtt{T}^I_i) \hspace{1em} \text{(C-Mute-LG)}$$

$$\partial \mathtt{T} = G \ \mathtt{T}^I_i \quad \equiv \quad \exists \mathtt{X}^I_i. (\mathtt{T} = G \ \mathtt{X}^I_i \wedge (\mathtt{T}_i = \partial \mathtt{X}_i)^I) \hspace{2em} \text{(C-Mute-DG)}$$
$$\text{if } \mathtt{X}^I_i \ \# \ ftv(\mathtt{T}, \mathtt{T}^I_i)$$

$$\partial \mathtt{T} = (\ell : \mathtt{T}' \ ; \ \mathtt{T}'') \quad \equiv \quad \mathtt{T} = \mathtt{T}' \wedge \partial \mathtt{T} = \mathtt{T}'' \hspace{2em} \text{(C-Mute-DL)}$$

**Figure 1-18: Constraint equivalence laws for rows.**

in the model, which follows directly from the meaning of row types. Notice that this fact only depends on the semantics of types and not on the semantics of the subtyping predicate. □

It follows from those equations that type constructors $\ell$, $\partial$, and $G$ are never isolated, each equation exhibiting a pair of compatible top symbols. Variances and incompatible pairs of type constructors depend on the signature $\mathcal{S}_0 \uplus \mathcal{S}_1$. However, it is not difficult to see that type constructors $\partial$ and $\ell$ are logically covariant in all directions and that the logical variance of types constructors $G$ of $dom(\mathcal{S}_0 \uplus \mathcal{S}_1)$ correspond to their syntactic variance, which, in most cases, will allow the decomposition of equations with the same top symbols. Moreover, an equation between two terms whose top symbols form one of the four compatible pairs derived from the equations of Figure 1-17 holds only if immediate subexpressions can be "conciliated" in some way. There is a transformation quite similar to decomposition, called *mutation*, that mimics the equations for rows (Figure 1-17) and described by the rules of Figure 1-18. For sake of readability and conciseness, we write $\mathtt{T}^I_i$ instead of $\mathtt{T}^{i \in I}_i$.

1.11.10   Lemma [Mutation]:  All equivalence laws in Figure 1-18 hold. □

*Proof:*

∘ *Case* C-Mute-LL: Let $\mathtt{X} \mathrel{\#} ftv(\mathtt{T}_1, \mathtt{T}_1', \mathtt{T}_2, \mathtt{T}_2')$ **(1)** and $\ell_1 \neq \ell_2$. Let $Row(L)$ be the row kind of this equation. Let $\phi$ be a ground assignment that validates the constraint $(\ell_1 : \mathtt{T}_1 \;;\; \mathtt{T}_1') = (\ell_2 : \mathtt{T}_2 \;;\; \mathtt{T}_2')$. That is, $\phi$ sends all terms of the multi-equation to the same ground type $t$ of row kind $Row(L)$. Moreover, the row-term semantics implies that $t$ satisfies $t(\epsilon) = L$, $t/\ell_1 = \phi(\mathtt{T}_1) = \phi(\mathtt{T}_2')/\ell_1$, $t/\ell_2 = \phi(\mathtt{T}_1')/\ell_2 = \phi(\mathtt{T}_2)$, and $t/\ell = \phi(\mathtt{T}_2')/\ell = \phi(\mathtt{T}_1')/\ell$ for all $\ell \in \mathcal{L} \setminus \ell_1.\ell_2.L$ **(2)**. Let $t'$ be the tree defined by $t'(\epsilon) = \ell_1.\ell_2.L$ and $t'/\ell = t/\ell$ for all $\ell \in \mathcal{L} \setminus \ell_1.\ell_2.L$. By construction and (2), $\phi[\mathtt{X} \mapsto t']$ satisfies both equations $\mathtt{T}_1' = (\ell_2 : \mathtt{T}_2 \;;\; \mathtt{X})$ and $\mathtt{T}_2' = (\ell_1 : \mathtt{T}_1 \;;\; \mathtt{X})$. Thus by CM-Exists and (1), $\phi$ satisfies $\exists \mathtt{X}.\mathtt{T}_1' = (\ell_2 : \mathtt{T}_2 \;;\; \mathtt{X}) \wedge \mathtt{T}_2' = (\ell_1 : \mathtt{T}_1 \;;\; \mathtt{X})$. Conversely, we have the entailment:

$$\exists \mathtt{X}.(\mathtt{T}_1' = (\ell_2 : \mathtt{T}_2 \;;\; \mathtt{X}) \wedge \mathtt{T}_2' = (\ell_1 : \mathtt{T}_1 \;;\; \mathtt{X}))$$
$$\equiv \exists \mathtt{X}.((\ell_1 : \mathtt{T}_1 \;;\; \mathtt{T}_1') = (\ell_1 : \mathtt{T}_1 \;;\; \ell_2 : \mathtt{T}_2 \;;\; \mathtt{X}) \wedge$$
$$\qquad (\ell_2 : \mathtt{T}_2 \;;\; \mathtt{T}_2') = (\ell_2 : \mathtt{T}_2 \;;\; \ell_1 : \mathtt{T}_1 \;;\; \mathtt{X})) \qquad \mathbf{(3)}$$
$$\Vdash \exists \mathtt{X}.(\ell_1 : \mathtt{T}_1 \;;\; \mathtt{T}_1') = (\ell_2 : \mathtt{T}_2 \;;\; \mathtt{T}_2') \qquad \mathbf{(4)}$$
$$\equiv (\ell_1 : \mathtt{T}_1 \;;\; \mathtt{T}_1') = (\ell_2 : \mathtt{T}_2 \;;\; \mathtt{T}_2') \qquad \mathbf{(5)}$$

(3) follows by covariance of $\ell_1$ and $\ell_2$; (4) by C-Row-LL and transitivity of equivalence; (5) by C-Ex* and (1).

∘ *Cases* C-Mute-LG, C-Mute-DG, and C-Mute-DL: The reasoning is similar to the case C-Mute-LL. □

### Solving row constraints in an equality model

In this section, we extend the constraint solver for the equality-only free tree model (Figure 1-11), so as to handle rows. We thus assume an equality-only model.

Mutation is a common technique to solve equations in a large class of non-free algebras that are described by *syntactic theories* (Kirchner and Klay, 1990). The equations of Figure 1-17 happen to be a syntactic presentation of an equational theory, from which a unification algorithm could be automatically derived (Rémy, 1993). We recover the same transformation rules directly, without using results on syntactic theories.

The following lemma shows that all pairs of distinct type constructors for which there is no mutation rule are in fact incompatible.

1.11.11  Lemma: All symbols $H \in \mathcal{S}_1$ are isolated. Furthermore for every pair of distinct type constructors $G_1, G_2 \in dom(\mathcal{S}_0 \uplus \mathcal{S}_1)$, and every row kind $s$, we have $G_1^s \bowtie G_2^s$. □

*Proof:* Terms of the form $H\vec{\mathtt{T}}$ are interpreted by ground types with $H$ at occurrence $\epsilon$, and conversely the only interpretations of types with $H$ at occurrence $\epsilon$ are terms of the form $H\vec{\mathtt{T}}$. Hence, no ground assignment can ever

$$(\ell_1 : \mathtt{X}_1 \; ; \; \mathtt{X}_1') = (\ell_2 : \mathtt{T}_2 \; ; \; \mathtt{T}_2') = \epsilon \quad \rightarrow \quad \exists \mathtt{Y}. (\mathtt{X}_1' = (\ell_2 : \mathtt{T}_2 \; ; \; \mathtt{Y}) \wedge \mathtt{T}_2' = (\ell_1 : \mathtt{X}_1 \; ; \; \mathtt{Y}))$$
$$\wedge \; (\ell_1 : \mathtt{X}_1 \; ; \; \mathtt{X}_1') = \epsilon \qquad\qquad \text{(S-\textsc{Mute}-LL)}$$
$$\text{if } \mathtt{Y} \; \# \; ftv(\mathtt{X}_1, \mathtt{X}_1', \mathtt{T}_2, \mathtt{T}_2') \wedge \ell_1 \neq \ell_2$$

$$(\ell : \mathtt{X} \; ; \; \mathtt{X}') = G \; \mathtt{T}_i^{i \in I} = \epsilon \quad \rightarrow \quad \exists (\mathtt{Y}_i, \mathtt{Y}_i')^{i \in I}. (\mathtt{X} = G \; \mathtt{Y}_i^{i \in I} \wedge \mathtt{X}' = G \; {\mathtt{Y}'}_i^{i \in I} \wedge \mathtt{T}_i = (\ell : \mathtt{Y}_i \; ; \; \mathtt{Y}_i'))$$
$$\wedge \; (\ell : \mathtt{X} \; ; \; \mathtt{X}') = \epsilon \qquad\qquad \text{(S-\textsc{Mute}-LG)}$$
$$\text{if } (\mathtt{Y}_i, \mathtt{Y}_i')^{i \in I} \; \# \; ftv(\mathtt{X}, \mathtt{X}', \mathtt{T}_i^{i \in I})$$

$$\partial \mathtt{X} = G \; \mathtt{T}_i^{i \in I} = \epsilon \quad \rightarrow \quad \exists \mathtt{Y}_i^{i \in I}. (\mathtt{X} = G \; \mathtt{Y}_i^{i \in I} \wedge (\mathtt{T}_i = \partial \mathtt{Y}_i)^{i \in I})$$
$$\wedge \; \partial \mathtt{X} = \epsilon \qquad\qquad \text{(S-\textsc{Mute}-DG)}$$
$$\text{if } \mathtt{Y}_i^{i \in I} \; \# \; ftv(\mathtt{X}, \mathtt{T}_i^{i \in I})$$

$$\partial \mathtt{X} = (\ell : \mathtt{T} \; ; \; \mathtt{T}') = \epsilon \quad \rightarrow \quad \mathtt{X} = \mathtt{T} \wedge \partial \mathtt{X} = \mathtt{T}' \wedge \partial \mathtt{X} = \epsilon \qquad\qquad \text{(S-\textsc{Mute}-DL)}$$

$$G \; \vec{\mathtt{T}} = G' \; \vec{\mathtt{T}}' = \epsilon \quad \rightarrow \quad \mathsf{false} \qquad\qquad\qquad\qquad\qquad\qquad \text{(S-\textsc{Clash}-I)}$$
$$\text{if } F \bowtie F'$$

**Figure 1-19: Unification addendum for row types**

send $H\vec{\mathtt{T}}$ and $F\vec{\mathtt{T}}'$ to the same ground term when $F \neq H$ and, as a result, $H$ is isolated.

Let $G_1$ and $G_2$ be two type constructors of $\mathcal{S}_0$. For $s = \mathit{Type}$, the interpretations of terms of the form $G_1^s\vec{\mathtt{T}}$ and $G_2^s\vec{\mathtt{T}}'$ are ground types with symbols $G_1$ and $G_2$ at occurrence $\epsilon$, respectively. Hence they cannot be made equal under any ground assignment. For $s = \mathit{Row}(L)$, the interpretations of terms of the form $G_1^s\vec{\mathtt{T}}$ and $G_2^s\vec{\mathtt{T}}'$ are ground types with constructor $L$ at occurrence $\epsilon$ and constructors $G_1$ and $G_2$ at occurrence 1, respectively. Hence they cannot be made equal under any ground assignment. $\qquad\qquad\qquad\square$

Any other combination of type constructors forms a compatible pair, as illustrated by equations of Figure 1-17 and can be transformed by mutation rules of Figure 1-18. The constraint solver for row-terms is the relation $\rightarrow^\dagger$ defined by the rewriting rules of Figure 1-11, except Rule S-\textsc{Clash}, which is replaced by the set of rules of Figure 1-19.

1.11.12    L\textsc{emma}:  The rewriting system $\rightarrow^\dagger$ is strongly normalizing. $\qquad\qquad\square$

Please, note that the termination of $\rightarrow^\dagger$ relies on types being well-kinded. In particular, $\rightarrow^\dagger$ would not terminate on the ill-kinded system of equations $\mathtt{X} = \ell : \mathtt{T} \; ; \; \mathtt{X}' \wedge \mathtt{X}' = \ell' : \mathtt{T} \; ; \; \mathtt{X}$.

1.11.13    L\textsc{emma}:  If $U \rightarrow^\dagger U'$, then $U \equiv U'$. $\qquad\qquad\qquad\qquad\qquad\square$

*Proof:* It suffices to check the property independently for each rule defining $\to^\dagger$. The proof for rules of Figure 1-11 but S-CLASH remain valid for row terms. For S-DECOMPOSE, it follows by the invariance of all type constructors, which is preserved for row terms. For rule S-CLASS-I it follows by Lemma 1.11.11 and for mutation rules, it follows by Lemma 1.11.10. □

Although reductions $\to$ are not sound for row types, hence $\to$ cannot be used for computation over row types, it has some interest. In particular, the following property shows that normal forms for row types are identical to normal forms for regular types.

1.11.14 LEMMA: A system $U$ in normal form for $\to^\dagger$ is also in normal form for $\to$. □

*Proof:* The only rule of $\to$ that is not in $\to^\dagger$ is S-CLASH. Thus, it suffices to observe that if Rule S-CLASH would be applicable, then either Rule S-CLASS-I or a mutation rule would be applicable as well. □

As a corollary, Lemma 1.8.6 extends to row types.

### Operations on records

We now illustrate the use of rows for typechecking operations on records. The simplest application of rows are full records with infinite carrier. Records types are expressed with rows instead of a simple product type. The basic operations are the same as in the finite case, that is, *creation*, *polymorphic update*, and *polymorphic access*, but labels are now taken from an infinite set. However, creation and polymorphic update, which where destructors are now taken as constructors and used to represent records as association lists. The access of a record $\mathtt{v}$ at a field $\ell$ is obtained by linearly searching $\mathtt{v}$ for a definition of field $\ell$ and returning this definition, or returning the default value if no definition has been found for $\ell$.

1.11.15 EXAMPLE [FULL RECORDS]: We assume a unique basic kind $\star$ and a unique covariant isolated type constructor $\Pi$ in $\mathcal{S}_1$. Let $\{\cdot\}$ be a unary constructor, $(\{\cdot \text{ with } \cdot = \ell\})^{\ell \in \mathcal{L}}$ a collection of binary constructors, and $(\ell.\{\cdot\})^{\ell \in \mathcal{L}}$ a collection of unary destructors with the following reduction rules:

$$
\begin{aligned}
\{\mathtt{v}\}.\{\ell\} &\xrightarrow{\delta} \mathtt{v} &&& \text{(RD-DEFAULT)} \\
\{\mathtt{w} \text{ with } \ell = \mathtt{v}\}.\{\ell\} &\xrightarrow{\delta} \mathtt{v} &&& \text{(RD-FOUND)} \\
\{\mathtt{w} \text{ with } \ell' = \mathtt{v}\}.\{\ell\} &\xrightarrow{\delta} \mathtt{w}.\{\ell\} && \text{if } \ell \neq \ell' & \text{(RD-FOLLOW)}
\end{aligned}
$$

Let the initial environment $\Gamma_0$ contain the following bindings

$$
\begin{aligned}
\{\cdot\} : &\quad \forall \mathtt{X}.\ \mathtt{X} \to \Pi(\partial \mathtt{X}) \\
\{\cdot \text{ with } \ell = \cdot\} : &\quad \forall \mathtt{X}\mathtt{X}'\mathtt{Y}.\ \Pi(\ell : \mathtt{X} \ ; \ \mathtt{Y}) \to \mathtt{X}' \to \Pi(\ell : \mathtt{X}' \ ; \ \mathtt{Y}) \\
\cdot.\{\ell\} : &\quad \forall \mathtt{X}\mathtt{Y}.\ \Pi(\ell : \mathtt{X} \ ; \ \mathtt{Y}) \to \mathtt{X}
\end{aligned}
$$

□

1.11.16    Exercise [Full records, ★★★, ↛]: Check that these definitions meet
the requirements of Definition 1.7.6.                                                 □

1.11.17    Exercise [Field Exchange, ★★]: Add an operation to permute two fields
of a record: give the reduction rules and the typing assumptions and check
that the requirements of Definition 1.7.6 are preserved.                    □

1.11.18    Exercise [Normal forms for records, ★★★]: Record values may con-
tain repetitions. For instance, $\{\{w \text{ with } \ell = v\} \text{ with } \ell = v'\}$ is a value that is
in fact observationally equivalent to $\{w \text{ with } \ell = v'\}$. So are the two record
values $\{\{w \text{ with } \ell = v\} \text{ with } \ell' = v'\}$ and $\{\{w \text{ with } \ell' = v'\} \text{ with } \ell = v\}$ when
$\ell' \neq \ell$. Modify the semantics, so that two record values of the same type have
similar structure and records do not contain inaccessible values.         □

1.11.19    Exercise [Map Apply, ★★]: Add a binary operation map such that the
expressions $(\text{map } v \, w).\{\ell\}$ and $v.\{\ell\} \, w.\{\ell\}$ reduce to the same value for every
label $\ell$.                                                                              □

1.11.20    Exercise [★, ↛]: So far, full records are almost constants. This condition
is not necessary for values, but only for types. As an example, introduce a
primitive record, that is a nullary record constructor, that maps every label
to a distinct integer. Give its typing assumption and review the semantics of
records.                                                                                □

As opposed to full records, standard records are partial and their domains
are finite (but with infinite carrier) and statically determined from their types.
Standard records can be built by extending the empty record on a finite num-
ber of fields. We refer to such records as records with *polymorphic extension*.
Record with polymorphic extension can be obtained by means of encoding
into full records, much as in the finite case.

1.11.21    Example [Encoding of polymorphic extension]: Reusing the two type
definitions abs and pre that have introduced for the finite case, we let abs
encodes an undefined field and prev encode a field with value v. We use the
following syntactic sugar with their meaning and principal types:

$$\langle\rangle \stackrel{\text{def}}{=} \{\text{abs}\}$$
$$: \Pi(\partial\text{abs})$$
$$\langle\cdot \text{ with } \ell = \cdot\rangle \stackrel{\text{def}}{=} \lambda v.\lambda w.\{w \text{ with } \ell = \text{pre } v\}$$
$$: \forall \text{XX}'\text{Y. } \Pi(\ell : \text{X} \mathbin{;} \text{Y}) \to \text{X}' \to \Pi(\ell : \text{pre } \text{X}' \mathbin{;} \text{Y})$$
$$\cdot.\langle\ell\rangle \stackrel{\text{def}}{=} \lambda v.\text{pre}^{-1} \, (v.\{\ell\})$$
$$: \forall \text{XY. } \Pi(\ell : \text{pre } \text{X} \mathbin{;} \text{Y}) \to \text{X}$$

□

**1.11.22**    EXERCISE [RECOMMENDED, ★]: Extension may actually override an existing field. Can you define a version polymorphic extension that prevents this situation from happening? Add an operation that hide some particular field of a record.      □

Extensible records can also be implemented directly, without encoding into full records. In fact, this requires only a tiny variation on full records.

**1.11.23**    EXAMPLE [RECORDS WITH POLYMORPHIC EXTENSION]: Let $\star$ and $\diamond$ be two basic kinds. Let the basic signature $\mathcal{S}_0$ contain (in addition to $\rightarrow$) the covariant isolated type constructors $\mathsf{pre}$ of kind $\star \Rightarrow \diamond$ and $\mathsf{abs}$ of kind $\diamond$. In the presence of subtyping, we may assume $\mathsf{pre} \leqslant \mathsf{abs}$. Let $\mathcal{S}_1$ contain the unique covariant isolated type constructor $\Pi$ of kind $\diamond \Rightarrow \star$. Let $\langle\rangle$ be a unary constructor, $(\{\cdot \text{ with } \cdot = \ell\})^{\ell \in \mathcal{L}}$ be a binary constructor, and $(\ell.\{\cdot\})^{\ell \in \mathcal{L}}$ be a unary destructor with the following reduction rules:

$$\langle\mathtt{w} \text{ with } \ell = \mathtt{v}\rangle.\langle\ell\rangle \quad \xrightarrow{\delta} \quad \mathtt{v} \quad\quad\quad\quad\quad\quad \text{(ER-FOUND)}$$
$$\langle\mathtt{w} \text{ with } \ell' = \mathtt{v}\rangle.\langle\ell\rangle \quad \xrightarrow{\delta} \quad \mathtt{w}.\{\ell\} \quad\quad \text{if } \ell \neq \ell' \quad\quad \text{(ER-FOLLOW)}$$

Let $\Gamma_0$ contain the following typing assumptions:

$$\langle\rangle : \quad \Pi(\partial\mathsf{abs})$$
$$\langle\cdot \text{ with } \ell = \cdot\rangle : \quad \forall \mathtt{XX'Y}.\, \Pi(\ell : \mathtt{X} \mathbin{;} \mathtt{Y}) \rightarrow \mathtt{X'} \rightarrow \Pi(\ell : \mathsf{pre}\ \mathtt{X'} \mathbin{;} \mathtt{Y})$$
$$\cdot.\langle\ell\rangle : \quad \forall \mathtt{XY}.\, \Pi(\ell : \mathsf{pre}\ \mathtt{X} \mathbin{;} \mathtt{Y}) \rightarrow \mathtt{X}$$

□

Notice that the typing assumptions obtained from the direct approach are identical to those obtained via the encoding into full records in Example 1.11.21.

**1.11.24**    EXERCISE [★★★★, ↛]: Prove the equivalence between the direct semantics and the semantics via the encoding into records with a default.    □

**1.11.25**    EXERCISE [RECOMMENDED, ★★, ↛]: Prove that type soundness for extensible records hold in both the subtyping model and equality-only models. □

**1.11.26**    EXERCISE [RECOMMENDED, ★, ↛]: Check that in the subtyping a record with more fields can be used in place of records with fewer fields. Check that this is not the case in the equality-only model.    □

1.11.27    Example [Refinement of record types]:  In  an  equality-only  model,
records with more fields cannot be used in place of records with fewer fields.
However, this may be partially recovered by a small refinement of the struc-
ture of types. The presence of fields can actually be split form their types, thus
enabling some polymorphism over the presence of fields while type of fields
themselves remains fixed. Let ∘ be a new basic kind. Let type constructors
abs and pre be both of kind ∘ and let · be a new type constructor of kind
$\circ \otimes \star \Rightarrow \diamond$. Let $\Gamma_0$ contain the following typing assumptions:

$$\langle\rangle : \quad \forall \mathtt{X}.\Pi(\partial(\mathsf{abs} \cdot \mathtt{X}))$$
$$\langle \cdot \text{ with } \ell = \cdot\rangle : \quad \forall \mathtt{ZXX'Y}.\ \Pi(\ell : \mathtt{X} \ ; \ \mathtt{Y}) \to \mathtt{X}' \to \Pi(\ell : \mathtt{Z} \cdot \mathtt{X}' \ ; \ \mathtt{Y})$$
$$\therefore\langle\ell\rangle : \quad \forall \mathtt{XY}.\ \Pi(\ell : \mathsf{pre} \cdot \mathtt{X} \ ; \ \mathtt{Y}) \to \mathtt{X}$$

The semantics of records remain unchanged. The new signature strictly gener-
alizes the previous one (strictly more programs can be typed) while preserving
type soundness. Here is a program that can now be typed and that could not
be typed before:

$$(\text{if } a \text{ then } \langle\langle\langle\rangle \text{ with } \ell' = \mathtt{true}\rangle \text{ with } \ell = 1\rangle \text{ else } \langle\langle\rangle \text{ with } \ell = 2\rangle).\langle\ell\rangle$$

Notice however, that when a present field is forgotten, the type of the field
is not. Therefore two records defining the same field but with incompatible
types can still not be mixed, which is possible in the subtyping model.      □

1.11.28    Example [Refined subtyping]:  The previous refinement for an equality-
only model is not much interesting in the case of a subtyping model.
     The subtyping assumption pre ⩽ abs makes abs play the role of ⊤ for fields.
That is, abs encodes the absence of information and not the information of
absence. In other words, a value whose field $\ell$ has type abs may either be
undefined or defined on field $\ell$; in the latter case, the fact that field $\ell$ is
actually defined has just been forgotten. Thus, types only provides a lower
approximation of the actual domain of records. This is a lost of accuracy by
comparison with the equality-only model, where a record domain is known
from its type. As a result, some optimizations in the representation of records
that are only possible when the exact domain of a record is statically known
are lost.
     Fortunately, there is a way to recover such accuracy. A conservative solution
could of course to drop the inequality pre ⩽ abs. Notice that this would still
be more expressive than using an equality model since, for instance $\Pi(\ell :$
pre $(\mathtt{T}_1 \to \mathtt{T}_2) \ ; \ \mathtt{T}) \leq \Pi(\ell : \mathsf{pre} \ \top \ ; \ \mathtt{T})$ would still hold, as long as $\to\ \leq\ \top$ does
hold. This solution is known as depth-only subtyping for records, while the
previous one provided both depth and width record subtyping. Conversely, one
could also keep width subtyping and disallow depth subtyping, by preserving

the relation $\mathsf{pre} \leqslant \mathsf{abs}$ while requiring $\mathsf{pre}$ to be invariant; in this case, presence of fields can be forgotten as a whole, but the types of fields cannot be weakened as long as fields remain visible.

Another more interesting solution consists in introducing another type constructor $\mathsf{either}$ of signature $\diamond$ and assuming that $\mathsf{pre} \leqslant \mathsf{either}$ and $\mathsf{abs} \leqslant \mathsf{either}$ (but $\mathsf{pre} \not\leqslant \mathsf{abs}$). Here, $\mathsf{either}$ plays the role of $\top$ for fields and means either *present* (and forgotten) or *absent*. while $\mathsf{abs}$ really means *absent*. The accuracy of typechecking can be formally stated as the fact that a record value of type $\Pi(\ell : \mathsf{abs} \,;\, \mathtt{T})$ cannot define field $\ell$. $\qquad\qquad\qquad\square$

1.11.29  EXAMPLE [MIXED SUBTYPING]: It is tempting to mix all variations of Example 1.11.28 together. As a first attempt, we may assume that the basic signature $\mathcal{S}_0$ contains covariant type constructors $\mathsf{pre}$ and $\mathsf{maybe}$ and invariant type constructors $\mathsf{pre}_=$ and $\mathsf{maybe}_=$, all of kind $\star \Rightarrow \diamond$ and two type constructors $\mathsf{abs}$ and $\mathsf{either}$ of kind $\diamond$, and that the subtype ordering $\leqslant$ is defined by the following diagram:



Intuitively, we wish that $\mathsf{pre}_=$ and $\mathsf{maybe}_=$ be *logically* invariant, $\mathsf{pre}$ and $\mathsf{maybe}$ be *logically* covariant, and the equivalences $\mathsf{pre}_= \, \mathtt{T} \leq \mathsf{maybe}_= \, \mathtt{T}' \equiv \mathtt{T} = \mathtt{T}'$ and

$$\mathsf{pre}_= \, \mathtt{T} \leq \mathsf{pre} \, \mathtt{T}' \equiv \mathsf{pre} \, \mathtt{T} \leq \mathsf{maybe} \, \mathtt{T}' \equiv \mathsf{maybe}_= \, \mathtt{T} \leq \mathsf{maybe} \, \mathtt{T}' \equiv \mathtt{T} \leq \mathtt{T}' \qquad (\mathbf{1})$$

simultaneously hold. However, (1) requires, for instance, type constructors $\mathsf{pre}_=$ and $\mathsf{pre}$ to have the same direction, which is not currently possible since they do not have the same variance. Interestingly, this restriction may be relaxed by assigning variances of directions on a per type constructor basis and define structural subtyping accordingly (See Exercise 1.11.30). Then, replacing all occurrences of $\mathsf{pre}$ by $\mathsf{pre}_=$ in $\Gamma_0$ preserves type soundness and allows for both accurate record types and flexible subtyping in the same setting. $\qquad\square$

1.11.30  EXERCISE [RELAXED VARIANCES, ★★★, $\nrightarrow$]: Let $\emptyset$ be allowed as a new variance, let extend the composition of variances defined in Example 1.3.9 with $\nu\emptyset = \emptyset$, and let $\leqslant^\emptyset$ stands for the full relation on type constructors. Let each type constructor $F$ of signature $K \Rightarrow \kappa$ now come with a mapping $\vartheta(F)$ from $dom(K)$ to variances. Let $\vartheta(t, t', \pi)$ be the variance of two

ground types $t$ and $t'$ at a path $\pi$ recursively defined by $\vartheta(t, t', d \cdot \pi) = \left( \vartheta(t(\epsilon))(d) \cap \vartheta(t'(\epsilon))(d) \right) \vartheta(t/d, t'/d, \pi)$ and $\vartheta(t, t', \epsilon) = +$. Then define the interpretation of subtyping as follows: if $t, t' \in \mathcal{M}_\kappa$, let $t \leq t'$ hold if and only if for all path $\pi \in dom(t) \cap dom(t')$, $t(\pi) \leqslant^{\vartheta(t, t', \pi)} t'(\pi)$ holds.

Check that the relation $\leq$ remains a partial ordering. Check that a type constructor whose direction $d$ has been syntactically declared covariant (respectively contravariant, invariant) is still logically covariant (respectively contravariant, invariant) in $d$.                                      □

### Record concatenation

Record concatenation takes two records and combines them into a new record whose fields are taken from whatever argument defines them. Of course, there is an ambiguity when the two records do not have disjoint domains and a choice should be made to disambiguate such cases. *Symmetric* concatenation let concatenation be undefined in this case (Harper and Pierce, 1991), while *asymmetric* concatenation let one-side (usually the right side) always take priority. Despite a rather simple semantics, record concatenation remains hard to type (with either a strict or a priority semantics).  Solutions to type inference for record concatenation may be found, for instance, in (Wand, 1989; Rémy, 1992; Pottier, 2000).

### Polymorphic variants

Variants can be defined via algebraic data-type definitions. However, as fields for records, variant tags are taken from a relatively small, finite collection of labels and two variant definitions will have incompatible types. Thus, to remain compatible, two variants must chose their tag among a larger collection that is a superset of all the possible tags of either variant. In general, this reduces the accuracy of types and forces useless dynamic checks for tags that could otherwise be known not to occur. Extensible variants (page 93) allow to work with an arbitrary large collection of tags, but do not improve accuracy. Polymorphic variants refers to a more precise typechecking mechanism for variants, where types more accurately describes the tags that may actually occur. They allow to build values of sum types out of a large, potentially infinite predefined set of tags and call polymorphic functions to explore them. As for record, this problem could be tackled by first considering polymorphic operations over variants built from a finite set of tags and total variants with an infinite set of tag independently and then by combining both approaches together. We propose a direct solution by a simple analogy with records.

Indeed, type constructor pre can be used to distinguish a (finite) set of tags

that the variant may actually carry, from other tags that are certain not to occur and typed with abs. For example, a variant $\ell.\text{v}$, built from a value v with a constructor tag $\ell$ of arity one. may be assigned the principal type scheme $\forall \text{X}.\Sigma(\ell : \text{pre T}; \text{X})$ where T is the type of v. The unary type constructor $\Sigma$ is used to coerce rows to variant types—thus, variant types and record types may share the same inner row structure and be simply distinguished by their top symbol. An instance of this polymorphic type is $\forall \text{X}.\Sigma(\ell : \text{pre T}; \text{abs})$, which tells that the variant must have been built with tag $\ell$ and no other tag, thus retaining exact information about the shape of the value. Another instance of the variant polymorphic type is $\Sigma(\ell : \text{pre T}; \ell' : \text{pre T}'; \text{abs})$. Indeed, it is sound to assume that the value might also have been built with some other tag $\ell'$, even if we know that this is not actually the case. Interestingly, both values $\ell.\text{v}$ and $\ell'.\text{v}$ have this type and can be mixed at this type.

We use filters to explore variants. A filter $[\ell : \text{v} \mid \text{v}']$ is a function that expects a variant argument, thus of the form $\ell'.\text{w}$. It then proceeds with either v w, if $\ell' = \ell$, or v' w otherwise. The type of this filter is $\Sigma(\ell : \text{pre T}; \text{T}') \to \text{T}''$ where T is the type of values accepted by v, $\Sigma(\ell : \text{T}'''; \text{T}')$ is the type of values accepted by v', and $\text{T}''$ is the type of values returned by both v and v'. Any type $\text{T}'''$ would do, including, in particular, abs. Indeed, when w is passed to v', it is known not to have tag $\ell$, so the behavior of v' on $\ell$ does not matter. The null filter $[]$ can be used for v'. This filter should actually never be applied, which we ensure by assigning $[]$ the type $\forall \text{X}.\Sigma(\partial \text{abs}) \to \text{X}$, for no variant value has type $\Sigma(\partial \text{abs})$. For instance, the filter $[\ell : \text{v}_\ell \mid [\ell' : \text{v}_{\ell'} \mid []]]$, which may be abbreviated as $[\ell : \text{v}_\ell \mid \ell' : \text{v}_{\ell'}]$ can be applied to either $\ell.\text{v}$ or $\ell'.\text{v}'$. The following example formalizes polymorphic variants.

1.11.31 EXAMPLE [POLYMORPHIC VARIANTS]: Let $\star$ and $\diamond$ be two basic kinds. Let $\mathcal{S}$ contain in addition to the arrow type constructor the two type constructors pre of kind $\star \Rightarrow \diamond$ and abs of kind $\diamond$. In the presence of subtyping we may assume $\text{abs} \leqslant \text{pre}$. Let $\mathcal{S}_1$ contain the unique covariant isolated type constructor $\Sigma$ of kind $\diamond \Rightarrow \star$. Let $\Gamma_0$ be composed of unary constructors $(\ell.\cdot)^{\ell \in \mathcal{L}}$ and primitives $[]$ of arity 0 and $([\ell : \cdot \mid \cdot]\cdot)^{\ell \in \mathcal{L}}$ of arity 3, given with the following reduction rules:

$$[\ell : \text{v} \mid \text{v}'] \; \ell.\text{w} \quad \xrightarrow{\delta} \quad \text{v w} \qquad\qquad\qquad\qquad\qquad\qquad (\text{EV-FOUND})$$
$$[\ell : \text{v} \mid \text{v}'] \; \ell'.\text{w} \quad \xrightarrow{\delta} \quad \text{v}' \text{ w} \qquad\qquad \text{if } \ell \neq \ell' \qquad\qquad (\text{EV-FOLLOW})$$

and contain the following typing assumptions:

$$
\begin{aligned}
\ell.\cdot : & \quad \forall \text{XY}. \; \text{X} \to \Sigma(\ell : \text{pre X}; \text{Y}) \\
[] : & \quad \forall \text{X}.\Sigma(\partial \text{abs}) \to \text{X} \\
[\ell : \cdot \mid \cdot] : & \quad \forall \text{XX}'\text{YY}'. \; (\text{X} \to \text{Y}) \to (\Sigma(\ell : \text{X}'; \text{Y}') \to \text{Y}) \to \Sigma(\ell : \text{pre X}; \text{Y}') \to \text{Y}
\end{aligned}
$$

$\square$

1.11.32   EXERCISE [SOUNDNESS FOR EXTENSIBLE VARIANTS, ★★★,↛]:  Prove type
soundness for extensible variants in both equality-only and subtyping models.
□

### Other applications of rows

Polymorphic records and variants are the most well-known applications of
rows. Besides the many variations on their presentations—we have only il-
lustrated some of them—there are several other interesting applications of
rows.

Since objects can be viewed as record-of-functions, at least from a typing
point of view, rows can also be used to type structural objects (Wand, 1994;
Rémy, 1994; Rémy and Vouillon, 1998) and provide, in particular, polymor-
phic method invocation. This is the key to typechecking objects in Objec-
tive Caml (Rémy and Vouillon, 1998). First-class messages (Nishimura, 1998;
Müller and Nishimura, 1998; Pottier, 2000) combine records and variants in
an interesting way: while filters over variant types enforce all branches to have
the same return type, first-class messages treat filters as records of functions
(also called objects) rather than functions from a variant type to a shared
return type. A message is an element of a variant type. The application of an
object to a message, that is of a record of functions to a variant type, selects
from the record the branch labeled with the same tag as the message and
applies it to the content of the message, much as pattern matching. However,
these applications are typechecked more accurately by first restricting the do-
main of the record to the set of tags that the message may possibly carry, and
thus other branches and in particular their return type are left unconstrained.

Row types may also represent set of properties within types or type refine-
ments and be used in type systems for program analysis. Two examples worth
mentioning are their application to soft-typing (Cartwright and Fagan, 1991;
Wright and Cartwright, 1994) and typechecking of uncaught exceptions (Leroy
and Pessaux, 2000).

The key to rows is to decompose the set of row labels into a class of fi-
nite partitions that is closed by some operations. Here, those partitions are
composed of singleton labels and co-finite sets of labels; the operations are
merging (or conversely splitting) a singleton label and a co-finite set of la-
bels. Other decompositions are possible, for instance, one could imagine to
consider labels in a two-dimensional space. More generally, labels might also
be given internal structure, for instance, one might consider automatons as
labels. Notice also that record types are stratified, since rows, that is, expres-
sions of kind $Row(L)$, may not themselves contain records —constructors of
$\mathcal{S}_1$ are only given the image row kind *Type*. This restriction can be partially

relaxed leading to rows of increasing degrees (Rémy, 1992b) . . . and complexity! Yet more intriguing are typed-indexed rows where labels are themselves types (Shields and Meijer, 2001).

### Alternatives to rows

The original idea of using rows to describe types of extensible records is due to Wand (Wand, 1987, 1988). A key simplification to row types is to make them total functions from labels to types and encode definiteness explicitly in the structure of fields, for instance with pre and abs type constructors, as presented here. This decomposition reduces the resolution of unification constraints to a simple equational reasoning (Rémy, 1993, 1992a). Other approaches that do not treat rows as total functions seem more *ad hoc* and have often hard-wired restrictions (Jategaonkar and Mitchell, 1988; Ohori and Buneman, 1989; Berthomieu, 1993; Ohori, 1999). Among these partial solutions, (Ohori, 1999) is quite interesting for its overall simplicity in the case where polymorphic access alone is required. Rows and fields may also be represented within ad-hoc type constraints rather than terms and equality (or subtyping) constraints. For example, qualified types use the predicates (T has $\ell$ : T$'$) and (T lacks $\ell$) to mean that field $\ell$ of row T is defined with type T$'$ or undefined, respectively (Jones, 1994b; Odersky, Sulzmann, and Wehr, 1999b). These constraints are in fact equivalent in our equality-model to $\exists$X.T = ($\ell$:pre  T$'$ ; X) and $\exists$X.T = ($\ell$:abs ; X), respectively. Record typechecking has also been widely studied in the presence of subtyping. Usually, record subtyping is given meaning directly and not via rows. While these solutions are quite expressive, thanks to subtyping, they still suffer from their nonstructural treatment of record types and cannot type row extension. Thus, even in subtyping models the use of rows increases expressiveness, and is usually a simplification as well. The subtyping model can then also take advantage of the possibility of enriching type constructors pre and abs with more structure and relate them via subtyping (Pottier, 2000). Notice, that even though rows have been introduced for type inference, they seem to be beneficial to explicitly typed languages as well since even other advanced solutions (Cardelli and Mitchell, 1991; Cardelli, 1992) are limited.

Rules of Figure 1-19 are *one way* of solving row type constraints. In a model with subtyping constraints, a more direct closure-based resolution may be more appropriate (Pottier, 2003).

# B  *Solutions to Selected Exercises*

1.2.6    SOLUTION: The definition does not behave as expected, because `if` is a destructor, whose arguments—according to the call-by-value semantics of ML-the-calculus—are evaluated *before* R-TRUE or R-FALSE is allowed to fire. As a result, the semantics of the expression `if t`$_0$` then t`$_1$` else t`$_2$ is to evaluate *both* `t`$_1$ and `t`$_2$ before choosing one of them. Since these expressions may have side effects (for instance, they may fail to terminate, or update a reference), this semantics is undesirable. The desired evaluation order can be obtained by placing `t`$_1$ and `t`$_2$ within closures, which delays their evaluation, then invoking the closure returned by the conditional, forcing its body to be evaluated. In other words, the expression `if t`$_0$` then t`$_1$` else t`$_2$ should now be viewed as syntactic sugar for `if t`$_0$ $(\lambda z.t_1)$ $(\lambda z.t_2)$ $\hat{0}$. The choice of the constant $\hat{0}$ is arbitrary, since it is discarded; any value would do.

1.2.21   SOLUTION: Within Damas and Milner's type system, we have:

$$\cfrac{\cfrac{\overline{z_1 : X \vdash z_1 : X}\ \text{DM-VAR} \qquad \overline{z_1 : X; z_2 : X \vdash z_2 : X}\ \text{DM-VAR}}{z_1 : X \vdash \texttt{let } z_2 = z_1 \texttt{ in } z_2 : X}\ \text{DM-LET}}{\varnothing \vdash \lambda z_1.\texttt{let } z_2 = z_1 \texttt{ in } z_2 : X \to X}\ \text{DM-ABS}$$

Please note that, because X occurs free within the environment $z_1 : X$, it is impossible to apply DM-GEN to the judgement $z_1 : X \vdash z_1 : X$ in a nontrivial way. For this reason, $z_2$ cannot receive the type scheme $\forall X.X$, and the whole expression cannot receive type $X \to Y$, where X and Y are distinct.

1.2.22   SOLUTION: It is straightforward to prove that the identity function has type `int → int`:

$$\cfrac{\overline{\Gamma_0; z : \texttt{int} \vdash z : \texttt{int}}\ \text{DM-VAR}}{\Gamma_0 \vdash \lambda z.z : \texttt{int} \to \texttt{int}}\ \text{DM-ABS}$$

In fact, nothing in this type derivation depends on the choice of `int` as the type of `z`. Thus, we may just as well use a type variable X instead. Furthermore, after forming the arrow type $X \to X$, we may employ DM-GEN to quantify universally over X, since it no longer appears in the environment.

$$\cfrac{\cfrac{\overline{\Gamma_0; z : X \vdash z : X}\ \text{DM-VAR}}{\Gamma_0 \vdash \lambda z.z : X \to X}\ \text{DM-ABS} \qquad X \notin \mathit{ftv}(\Gamma_0)}{\Gamma_0 \vdash \lambda z.z : \forall X.X \to X}\ \text{DM-GEN}$$

It is worth noting that, although the type derivation employs an arbitrary type variable X, the final typing judgement has no free type variables. It is

thus independent of the choice of $\mathtt{X}$. In the following, we refer to the above type derivation as $\Delta_0$.

Next, we prove that the successor function has type $\mathtt{int} \to \mathtt{int}$ under the initial environment $\Gamma_0$. We write $\Gamma_1$ for $\Gamma_0; \mathtt{z} : \mathtt{int}$, and make uses of DM-VAR implicit.

$$
\cfrac{
\cfrac{
\cfrac{
\cfrac{\Gamma_1 \vdash \hat{+} : \mathtt{int} \to \mathtt{int} \to \mathtt{int} \qquad \Gamma_1 \vdash \mathtt{z} : \mathtt{int}}{\Gamma_1 \vdash \hat{+} \, \mathtt{z} : \mathtt{int} \to \mathtt{int}} \text{ DM-APP} \qquad \Gamma_1 \vdash \hat{1} : \mathtt{int}
}{\Gamma_1 \vdash \mathtt{z} \mathbin{\hat{+}} \hat{1} : \mathtt{int}} \text{ DM-APP}
}{\Gamma_0 \vdash \lambda\mathtt{z}.\mathtt{z} \mathbin{\hat{+}} \hat{1} : \mathtt{int} \to \mathtt{int}} \text{ DM-ABS}
}
$$

In the following, we refer to the above type derivation as $\Delta_1$. We may now build a derivation for the third typing judgement. We write $\Gamma_2$ for $\Gamma_0; \mathtt{f} : \mathtt{int} \to \mathtt{int}$.

$$
\cfrac{
\Delta_1 \qquad \cfrac{\Gamma_2 \vdash \mathtt{f} : \mathtt{int} \to \mathtt{int} \qquad \Gamma_2 \vdash \hat{2} : \mathtt{int}}{\Gamma_2 \vdash \mathtt{f} \, \hat{2} : \mathtt{int}} \text{ DM-APP}
}{\Gamma_0 \vdash \mathtt{let}\ \mathtt{f} = \lambda\mathtt{z}.\mathtt{z} \mathbin{\hat{+}} \hat{1}\ \mathtt{in}\ \mathtt{f} \, \hat{2} : \mathtt{int}} \text{ DM-LET}
$$

To derive the fourth typing judgement, we re-use $\Delta_0$, which proves that the identity function has polymorphic type. We write $\Gamma_3$ for $\Gamma_0; \mathtt{f} : \forall\mathtt{X}.\mathtt{X} \to \mathtt{X}$. By DM-VAR and DM-INST, we have $\Gamma_3 \vdash \mathtt{f} : (\mathtt{int} \to \mathtt{int}) \to (\mathtt{int} \to \mathtt{int})$ and $\Gamma_3 \vdash \mathtt{f} : \mathtt{int} \to \mathtt{int}$. Thus, we may build the following derivation:

$$
\cfrac{
\Delta_0 \qquad \cfrac{
\cfrac{
\cfrac{\Gamma_3 \vdash \mathtt{f} : (\mathtt{int} \to \mathtt{int}) \to (\mathtt{int} \to \mathtt{int}) \qquad \Gamma_3 \vdash \mathtt{f} : \mathtt{int} \to \mathtt{int}}{\Gamma_3 \vdash \mathtt{f}\,\mathtt{f} : \mathtt{int} \to \mathtt{int}} \text{ DM-APP} \qquad \Gamma_3 \vdash \hat{2} : \mathtt{int}
}{\Gamma_3 \vdash \mathtt{f}\,\mathtt{f}\,\hat{2} : \mathtt{int}} \text{ DM-APP}
}{\Gamma_0 \vdash \mathtt{let}\ \mathtt{f} = \lambda\mathtt{z}.\mathtt{z}\ \mathtt{in}\ \mathtt{f}\,\mathtt{f}\,\hat{2} : \mathtt{int}} \text{ DM-LET}
$$

The first and third judgements are valid in the simply-typed $\lambda$-calculus, because they use neither DM-GEN nor DM-INST, and use DM-LET only to introduce the *monomorphic* binding $\mathtt{f} : \mathtt{int} \to \mathtt{int}$ into the environment. The second judgement, of course, is not: because it involves a nontrivial type scheme, it is not even a well-formed judgement in the simply-typed $\lambda$-calculus. The fourth judgement is well-formed, but *not* derivable, in the simply-typed $\lambda$-calculus. This is because $\mathtt{f}$ is used at two incompatible types, namely $(\mathtt{int} \to \mathtt{int}) \to (\mathtt{int} \to \mathtt{int})$ and $\mathtt{int} \to \mathtt{int}$, inside the expression $\mathtt{f}\,\mathtt{f}\,\hat{2}$. Both of these types are *instances* of $\forall\mathtt{X}.\mathtt{X} \to \mathtt{X}$, the type scheme assigned to $\mathtt{f}$ in the environment $\Gamma_3$.

By inspection of DM-VAR, DM-GEN, and DM-INST, it is straightforward to see that, if $\Gamma_0 \vdash \hat{1} : T$ is derivable, then $T$ must be $\texttt{int}$. Since $\texttt{int}$ is not an arrow type, the application $\hat{1}\ \hat{2}$ cannot be well-typed under $\Gamma_0$. In fact, because this expression is stuck, it cannot be well-typed in a sound type system.

The expression $\lambda\texttt{f}.(\texttt{f f})$ is ill-typed in the simply-typed $\lambda$-calculus, because no type $T$ may coincide with a type of the form $T \to T'$. Indeed, $T$ would then be a subterm of itself. For the same reason, this expression is ill-typed in DM as well. Indeed, it is not difficult to check that the presence of DM-GEN and DM-INST makes no difference: DM-GEN cannot generalize $T$ as long as the binding $\texttt{f} : T$ appears in the environment, and DM-INST can only instantiate $T$ to $T$ itself. Thus, the self-application $\texttt{f f}$ is well-typed in DM *only* if $\texttt{f}$ is $\texttt{let}$-bound, as opposed to $\lambda$-bound. The argument crucially relies on the fact that $\texttt{f}$ must be assigned a *monotype*. Indeed, the expression $\lambda\texttt{f}.(\texttt{f f})$ is well-typed in an implicitly-typed variant of System F: one of its types is $(\forall\texttt{X}.\texttt{X} \to \texttt{X}) \to (\forall\texttt{X}.\texttt{X} \to \texttt{X})$. It also relies on the fact that types are *finite*: indeed, this expression is well-typed in an extension of the simply-typed $\lambda$-calculus with recursive types, where the equation $T = T \to T'$ has a solution.

1.2.23    SOLUTION: It is clear that the effect of DM-GEN may be obtained by a series of successive applications of DM-GEN'. Conversely, consider an instance of DM-GEN', whose premises are $\Gamma \vdash \texttt{t} : S$ **(1)** and $\texttt{X} \notin ftv(\Gamma)$ **(2)**. Let us write $S = \forall\bar{\texttt{X}}.T$, where $\bar{\texttt{X}} \mathbin{\#} ftv(\Gamma)$ **(3)**. Applying DM-INST to (1) and to the identity substitution yields $\Gamma \vdash \texttt{t} : T$ **(4)**. Applying DM-GEN to (4), (2) and (3) yields $\Gamma \vdash \texttt{t} : \forall\texttt{X}\bar{\texttt{X}}.T$, that is, $\Gamma \vdash \texttt{t} : \forall\texttt{X}.S$. Thus, the effect of DM-GEN' may be obtained by DM-INST and DM-GEN.

It is clear that DM-INST is a particular case of DM-INST' where $\bar{\texttt{Y}}$ is empty. Conversely, consider an instance of DM-INST', whose premises are $\Gamma \vdash \texttt{t} : \forall\bar{\texttt{X}}.T$ **(1)** and $\bar{\texttt{Y}} \mathbin{\#} ftv(\forall\bar{\texttt{X}}.T)$ **(2)**. Let $\rho$ be a renaming that exchanges $\bar{\texttt{Y}}$ with $\bar{\texttt{Z}}$, where $\bar{\texttt{Z}} \mathbin{\#} ftv(\forall\bar{\texttt{Y}}.[\vec{\texttt{X}} \mapsto \vec{\texttt{T}}]T)$ **(3)** and $\bar{\texttt{Z}} \mathbin{\#} ftv(\Gamma)$ **(4)**. Applying DM-INST to (1) yields $\Gamma \vdash \texttt{t} : [\vec{\texttt{X}} \mapsto \rho\vec{\texttt{T}}]T$ **(5)**. Applying DM-GEN to (5) and (4) yields $\Gamma \vdash \texttt{t} : \forall\bar{\texttt{Z}}.[\vec{\texttt{X}} \mapsto \rho\vec{\texttt{T}}]T$, that is, $\Gamma \vdash \texttt{t} : \forall\rho\bar{\texttt{Y}}.[\vec{\texttt{X}} \mapsto \rho\vec{\texttt{T}}]T$ **(6)**. Now, by (2) and (3), we have $[\vec{\texttt{X}} \mapsto \rho\vec{\texttt{T}}]T = \rho([\vec{\texttt{X}} \mapsto \vec{\texttt{T}}]T)$, so (6) may be written $\Gamma \vdash \texttt{t} : \forall\rho\bar{\texttt{Y}}.\rho([\vec{\texttt{X}} \mapsto \vec{\texttt{T}}]T)$, that is, $\Gamma \vdash \texttt{t} : \rho(\forall\bar{\texttt{Y}}.[\vec{\texttt{X}} \mapsto \vec{\texttt{T}}]T)$ **(7)**. By (3), this is exactly $\Gamma \vdash \texttt{t} : \forall\bar{\texttt{Y}}.[\vec{\texttt{X}} \mapsto \vec{\texttt{T}}]T$. Thus, the effect of DM-INST' may be obtained by DM-INST and DM-GEN.

1.4.4    SOLUTION: Let us recall that a program $\texttt{t}$ is well-typed if and only if a judgement of the form $C, \Gamma \vdash \texttt{t} : \sigma$, where $C$ is satisfiable, holds. Let us show that it is in fact possible, without loss of generality, to require $\sigma$ to be a monotype.

Assume $C, \Gamma \vdash \texttt{t} : \sigma$ **(1)** is derivable within HM($X$). Let us write $\sigma = \forall\bar{\texttt{X}}[D].T$, where $\bar{\texttt{X}} \mathbin{\#} ftv(C)$ **(2)**. Applying Lemma 1.4.1 to (1) yields $C \Vdash \exists\bar{\texttt{X}}.D$ **(3)**. By HM-INST, (1) implies $C \wedge D, \Gamma \vdash \texttt{t} : T$ **(4)**. By (3), we have $C \equiv C \wedge \exists\bar{\texttt{X}}.D \equiv \exists\bar{\texttt{X}}.(C \wedge D)$. Because $C$ is satisfiable, this implies that $C \wedge D$

is satisfiable as well. Thus, the judgement (4), which involves the monotype T, witnesses that t is well-typed.

We have shown that a program t is well-typed if and only if a judgement of the form $C, \Gamma \vdash$ t : T, where $C$ is satisfiable, holds. Thus, by Theorems **??** and **??**, well-typedness is the same for both rule sets.

1.4.5   SOLUTION: By Theorem **??**, every rule in Figure 1-8 is admissible in HM($X$). Of course, so is HM-GEN. So, every judgement that is derivable via the rules of Figure 1-8 and HM-GEN is a valid HM($X$) judgement.

Conversely, assume $C, \Gamma \vdash$ t : $\sigma$ **(1)** holds in HM($X$). We must show that it is derivable via the rules of Figure 1-8 and HM-GEN. Let us write $\sigma = \forall \bar{\text{X}}[D].\text{T}$, where $\bar{\text{X}} \mathrel{\#} \mathit{ftv}(C, \Gamma)$ **(2)**. By HM-INST and (1), the judgement $C \wedge D, \Gamma \vdash$ t : T **(3)** holds in HM($X$). This judgement involves a monotype, so, by Theorem **??**, it is derivable via the rules of Figure 1-8. Furthermore, from (3) and (2), HM-GEN allows deriving $C \wedge \exists \sigma, \Gamma \vdash$ t : $\sigma$ **(4)**. Applying Lemma 1.4.1 to (1) yields $C \Vdash \exists \sigma$, so the judgement (4) may be written $C, \Gamma \vdash$ t : $\sigma$. We have shown that (1) is derivable via the rules of Figure 1-8 and HM-GEN. In fact, it is possible to apply HM-GEN only once, at the end of the derivation.

1.5.1   SOLUTION: Within the type system PCB($X$), we have

$$
\cfrac{
  \cfrac{\rule{0pt}{0pt}}{\text{z}_1 \preceq \text{Z} \vdash \text{z}_1 : \text{Z}} \text{VAR} \qquad
  \cfrac{\rule{0pt}{0pt}}{\text{z}_2 \preceq \text{Y} \vdash \text{z}_2 : \text{Y}} \text{VAR}
}{
  \cfrac{\text{let } \text{z}_2 : \forall \text{Z}[\text{z}_1 \preceq \text{Z}].\text{Z in } \text{z}_2 \preceq \text{Y} \vdash \text{let } \text{z}_2 = \text{z}_1 \text{ in } \text{z}_2 : \text{Y}}{\text{let } \text{z}_1 : \text{X}; \text{z}_2 : \forall \text{Z}[\text{z}_1 \preceq \text{Z}].\text{Z in } \text{z}_2 \preceq \text{Y} \vdash \lambda \text{z}_1.\text{let } \text{z}_2 = \text{z}_1 \text{ in } \text{z}_2 : \text{X} \to \text{Y}} \text{ABS}
} \text{LET}
$$

The type variable Z, which occurs free in the left-hand instance of VAR, is generalized. However, $\text{z}_2$ does *not* receive the type scheme $\forall \text{Z}.\text{Z}$, which, as suggested earlier, is unsound; instead, it receives the *constrained* type scheme $\forall \text{Z}[\text{z}_1 \preceq \text{Z}].\text{Z}$. The latter is more restrictive than the former: indeed, the former claims that $\text{z}_2$ has every type, while the latter only claims that every valid type for $\text{z}_1$ is also a valid type for $\text{z}_2$. Let us now examine the constraint let $\text{z}_1 : \text{X}; \text{z}_2 : \forall \text{Z}[\text{z}_1 \preceq \text{Z}].\text{Z in } \text{z}_2 \preceq \text{Y}$, which appears at the root of the derivation. By C-INID and C-IN*, it is equivalent to let $\text{z}_1 : \text{X in } \exists \text{Z}.(\text{z}_1 \preceq \text{Z} \wedge \text{Z} \le \text{Y})$ and to $\exists \text{Z}.(\text{X} \le \text{Z} \wedge \text{Z} \le \text{Y})$, which by C-EXTRANS is equivalent to $\text{X} \le \text{Y}$. Thus, the judgement at the root of the above derivation may be written $\text{X} \le \text{Y} \vdash \lambda \text{z}_1.\text{let } \text{z}_2 = \text{z}_1 \text{ in } \text{z}_2 : \text{X} \to \text{Y}$. In other words, the expression let $\text{z}_2 = \text{z}_1 \text{ in } \text{z}_2$ has type $\text{X} \to \text{Y}$ only under the assumption that X is a subtype of Y, which is sound. Even though LET allows unrestricted generalization of type variables, it remains sound, because the type scheme that it produces typically has *free program identifiers*, such as $\forall \text{Z}[\text{z}_1 \preceq \text{Z}].\text{Z}$ above.

1.7.10   SOLUTION: Let $\mathcal{E} = \mathtt{let}\ \mathtt{z} = \mathcal{E}_1\ \mathtt{in}\ \mathtt{t}_1$ and $\mathcal{E}_1[\mathtt{t}]/\mu \sqsubseteq \mathcal{E}_1[\mathtt{t}']/\mu'$ **(1)**. Then,

$$\mathsf{let}\ \Gamma_0;\mathsf{ref}\ M\ \mathsf{in}\ [\![\mathcal{E}[\mathtt{t}]/\mu : \mathtt{T}/M]\!]$$

$\equiv\quad \mathsf{let}\ \Gamma_0;\mathsf{ref}\ M\ \mathsf{in}\ ((\mathsf{let}\ \mathtt{z} : \forall\mathtt{X}[[\![\mathcal{E}_1[\mathtt{t}] : \mathtt{X}]\!]].\mathtt{X}\ \mathsf{in}\ [\![\mathtt{t}_1 : \mathtt{T}]\!]) \wedge [\![\mu : M]\!])$     **(2)**

$\equiv\quad \mathsf{let}\ \Gamma_0;\mathsf{ref}\ M;\mathtt{z} : \forall\mathtt{X}[[\![\mathcal{E}_1[\mathtt{t}]/\mu : \mathtt{X}/M]\!]].\mathtt{X}\ \mathsf{in}\ [\![\mathtt{t}_1 : \mathtt{T}]\!]$     **(3)**

$\equiv\quad \mathsf{let}\ \Gamma_0;\mathsf{ref}\ M;\mathtt{z} : \forall\mathtt{X}[\mathsf{let}\ \Gamma_0;\mathsf{ref}\ M\ \mathsf{in}\ [\![\mathcal{E}_1[\mathtt{t}]/\mu : \mathtt{X}/M]\!]].\mathtt{X}\ \mathsf{in}\ [\![\mathtt{t}_1 : \mathtt{T}]\!]$     **(4)**

$\Vdash\quad \mathsf{let}\ \Gamma_0;\mathsf{ref}\ M;\mathtt{z} : \forall\mathtt{X}\bar{\mathtt{Y}}[\mathsf{let}\ \Gamma_0;\mathsf{ref}\ M'\ \mathsf{in}\ [\![\mathcal{E}_1[\mathtt{t}']/\mu' : \mathtt{X}/M']\!]].\mathtt{X}\ \mathsf{in}\ [\![\mathtt{t}_1 : \mathtt{T}]\!]$     **(5)**

where (2) is by definition of constraint generation, where $\mathtt{X} \not\in ftv(\mathtt{T}, M)$ **(6)**; (3) is by (6), C-LETAND, and by definition of constraint generation; (4) is by (6) and C-LETDUP; (5) follows from (1) and C-LETEX, for some $\bar{\mathtt{Y}}$ and $M'$ such that $\bar{\mathtt{Y}} \mathbin{\#} ftv(\mathtt{X}, M)$ **(7)** and $ftv(M') \subseteq \bar{\mathtt{Y}} \cup ftv(M)$ **(8)** and $dom(M') = dom(\mu')$ and $M'$ extends $M$. Note that (6), (7) and (8) imply $\mathtt{X} \not\in ftv(M')$ **(9)**.

At this point, the type variables $\bar{\mathtt{Y}}$, which determine the types of the newly allocated store cells, are universally quantified in the type scheme assigned to $\mathtt{z}$, which is undesirable. We are stuck, because we cannot in general apply C-LETALL to hoist $\exists\bar{\mathtt{Y}}$ out of the let constraint. Let us now assume that, by some external means, we are guaranteed $\bar{\mathtt{Y}} = \varnothing$ **(10)**. Then, we may proceed as follows:

$\equiv\quad \mathsf{let}\ \Gamma_0;\mathsf{ref}\ M';\mathtt{z} : \forall\mathtt{X}[\mathsf{let}\ \Gamma_0;\mathsf{ref}\ M'\ \mathsf{in}\ [\![\mathcal{E}_1[\mathtt{t}']/\mu' : \mathtt{X}/M']\!]].\mathtt{X}\ \mathsf{in}\ [\![\mathtt{t}_1 : \mathtt{T}]\!]$     **(11)**

$\equiv\quad \mathsf{let}\ \Gamma_0;\mathsf{ref}\ M'\ \mathsf{in}\ [\![\mathcal{E}[\mathtt{t}']/\mu' : \mathtt{T}/M']\!]$     **(12)**

where (11) follows from the fact the the memory locations that appear free in $[\![\mathtt{t}_1 : \mathtt{T}]\!]$ are members of $dom(\mu)$, thus are not members of $dom(M') \setminus dom(M)$; (12) is obtained by performing the steps that lead to (4) in reverse.

The requirement that $\bar{\mathtt{Y}}$ be empty, that is, $ftv(M) = ftv(M')$, is classic (Tofte, 1988). How is it enforced? Assume that the left-hand side of every let construct is required to be a non-expansive expression. By assumptions (ii) and (iii), this invariant is preserved by reduction. So, $\mathcal{E}_1[\mathtt{t}]$ must be non-expansive, which, by assumption (i), guarantees that the reduction step does not allocate new memory cells. Then, $\mu'$ is $\mu$, so $M'$ is $M$.

1.9.1   SOLUTION: We must first ensure that R-ADD respects $\sqsubseteq$ (Definition 1.7.5). Since the rule is pure, it is sufficient to establish that $\mathsf{let}\ \Gamma_0\ \mathsf{in}\ [\![\hat{k}_1 \mathbin{\hat{+}} \hat{k}_2 : \mathtt{T}]\!]$ entails $\mathsf{let}\ \Gamma_0\ \mathsf{in}\ [\![\widehat{k_1 + k_2} : \mathtt{T}]\!]$. In fact, we have

$$\mathsf{let}\ \Gamma_0\ \mathsf{in}\ [\![\hat{k}_1 \mathbin{\hat{+}} \hat{k}_2 : \mathtt{T}]\!]$$

$\equiv\quad \mathsf{let}\ \Gamma_0\ \mathsf{in}\ \exists\mathtt{XY}.(\hat{+} \preceq \mathtt{X} \to \mathtt{Y} \to \mathtt{T} \wedge \hat{k}_1 \preceq \mathtt{X} \wedge \hat{k}_2 \preceq \mathtt{Y})$     **(1)**

$\equiv\quad \exists\mathtt{XY}.(\mathtt{int} \to \mathtt{int} \to \mathtt{int} \leq \mathtt{X} \to \mathtt{Y} \to \mathtt{T} \wedge \mathtt{int} \leq \mathtt{X} \wedge \mathtt{int} \leq \mathtt{Y})$     **(2)**

$\equiv\quad \exists\mathtt{XY}.(\mathtt{X} = \mathtt{int} \wedge \mathtt{Y} = \mathtt{int} \wedge \mathtt{int} \leq \mathtt{T})$     **(3)**

$\equiv\quad \mathtt{int} \leq \mathtt{T}$     **(4)**

$\equiv\quad \mathsf{let}\ \Gamma_0\ \mathsf{in}\ [\![\widehat{k_1 + k_2} : \mathtt{T}]\!]$     **(5)**

where (1) is by definition of constraint generation; (2) is by definition of $\Gamma_0$, by C-InId and C-In*; (3) is by C-Arrow and by antisymmetry of subtyping; (4) is by C-ExAnd and C-Name; (5) is again by definition of $\Gamma_0$, by C-InId and C-In*, and by definition of constraint generation.

Second, we must check that if the configuration c $v_1$ ... $v_k/\mu$ (where $k \geq 0$) is well-typed, then either it is reducible, or c $v_1$ ... $v_k$ is a value.

We begin by checking that every value that is well-typed with type int is of the form $\hat{k}$. Indeed, suppose that let $\Gamma_0$; ref $M$ in $[\![v : \text{int}]\!]$ is satisfiable. Then, v cannot be a program variable, for a well-typed value must be closed. v cannot be a memory location $m$, for otherwise ref $M(m) \leq$ int would be satisfiable—but the type constructors ref and int are incompatible. v cannot be $\hat{+}$ or $\hat{+}$ $v'$, for otherwise int $\to$ int $\to$ int $\leq$ int or int $\to$ int $\leq$ int would be satisfiable—but the type constructors $\to$ and int are incompatible. Similarly, v cannot be a $\lambda$-abstraction. Thus, v must be of the form $\hat{k}$, for it is the only case left.

Next, we note that, according to the constraint generation rules, if the configuration c $v_1$ ... $v_k/\mu$ is well-typed, then a constraint of the form let $\Gamma_0$; ref $M$ in (c $\preceq$ $X_1$ $\to$ ... $\to$ $X_k$ $\to$ T $\wedge$ $[\![v_1 : X_1]\!]$ $\wedge$ ... $\wedge$ $[\![v_k : X_k]\!]$) is satisfiable. We now reason by cases on c.

$\circ$ *Case* c is $\hat{k}$. Then, $\Gamma_0$(c) is int. Because the type constructors int and $\to$ are incompatible with each other, this implies $k = 0$. Since $\hat{k}$ is a constructor, the expression is a value.

$\circ$ *Case* c is $\hat{+}$. We may assume $k \geq 2$, because otherwise the expression is a value. Then, $\Gamma_0$(c) is int $\to$ int $\to$ int, so, by C-Arrow, the above constraint entails let $\Gamma_0$; ref $M$ in ($X_1 \leq$ int $\wedge X_2 \leq$ int $\wedge [\![v_1 : X_1]\!] \wedge [\![v_2 : X_2]\!]$), which, by Lemma 1.6.3, entails let $\Gamma_0$; ref $M$ in ($[\![v_1 : \text{int}]\!] \wedge [\![v_2 : \text{int}]\!]$). Thus, $v_1$ and $v_2$ are well-typed with type int. By the remark above, they must be integer literals $\hat{k}_1$ and $\hat{k}_2$. As a result, the configuration is reducible by R-Add.

1.9.5    Solution: We must first ensure that R-Ref, R-Deref and R-Assign respect $\sqsubseteq$ (Definition 1.7.5).

$\circ$ *Case* R-Ref. The reduction is ref $v/\varnothing \longrightarrow m/(m \mapsto v)$, where $m \notin fpi(v)$ **(1)**. Let T be an arbitrary type. According to Definition 1.7.5, the goal is to show that there exist a set of type variables $\bar{Y}$ and a store type $M'$ such that $\bar{Y} \# ftv(\text{T})$ and $ftv(M') \subseteq \bar{Y}$ and $dom(M') = \{m\}$ and let $\Gamma_0$ in $[\![\text{ref } v : \text{T}]\!]$

entails $\exists \bar{Y}.\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M'\ \mathsf{in}\ [\![m/(m \mapsto \mathtt{v}) : \mathtt{T}/M']\!]$. Now, we have

$$
\begin{aligned}
&\mathsf{let}\ \Gamma_0\ \mathsf{in}\ [\![\mathtt{ref}\ \mathtt{v} : \mathtt{T}]\!] \\
\equiv\ &\mathsf{let}\ \Gamma_0\ \mathsf{in}\ \exists \mathtt{XY}.(\mathtt{Y} \to \mathsf{ref}\,\mathtt{Y} \leq \mathtt{X} \to \mathtt{T} \wedge [\![\mathtt{v} : \mathtt{X}]\!]) && \mathbf{(2)} \\
\equiv\ &\exists \mathtt{Y}.\mathsf{let}\ \Gamma_0\ \mathsf{in}\ (\mathsf{ref}\,\mathtt{Y} \leq \mathtt{T} \wedge [\![\mathtt{v} : \mathtt{Y}]\!]) && \mathbf{(3)} \\
\equiv\ &\exists \mathtt{Y}.\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M'\ \mathsf{in}\ (m \preceq \mathtt{T} \wedge [\![\mathtt{v} : M'(m)]\!]) && \mathbf{(4)} \\
\equiv\ &\exists \mathtt{Y}.\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M'\ \mathsf{in}\ [\![m/(m \mapsto \mathtt{v}) : \mathtt{T}/M']\!] && \mathbf{(5)}
\end{aligned}
$$

where (2) is by definition of constraint generation and by definition of $\Gamma_0(\mathtt{ref})$; (3) is by C-ARROW, Lemma 1.6.4, and C-INEX; (4) assumes $M'$ is defined as $m \mapsto \mathtt{Y}$, and follows from (1), C-INID and C-IN*; and (5) is by definition of constraint generation.

    *Subcase* R-DEREF. The reduction is $!m/(m \mapsto \mathtt{v}) \longrightarrow \mathtt{v}/(m \mapsto \mathtt{v})$. Let $\mathtt{T}$ be an arbitrary type and let $M$ be a store type of domain $\{m\}$. We have

$$
\begin{aligned}
&\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ [\![!m/(m \mapsto \mathtt{v}) : \mathtt{T}/M]\!] \\
\equiv\ &\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ \exists \mathtt{XY}.(\mathsf{ref}\,\mathtt{Y} \to \mathtt{Y} \leq \mathtt{X} \to \mathtt{T} \wedge m \preceq \mathtt{X} \wedge [\![\mathtt{v} : M(m)]\!]) && \mathbf{(1)} \\
\equiv\ &\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ \exists \mathtt{XY}.(\mathsf{ref}\,M(m) \leq \mathtt{X} \leq \mathsf{ref}\,\mathtt{Y} \wedge \mathtt{Y} \leq \mathtt{T} \wedge [\![\mathtt{v} : M(m)]\!]) && \mathbf{(2)} \\
\equiv\ &\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ \exists \mathtt{Y}.(M(m) = \mathtt{Y} \wedge \mathtt{Y} \leq \mathtt{T} \wedge [\![\mathtt{v} : M(m)]\!]) && \mathbf{(3)} \\
\equiv\ &\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ (M(m) \leq \mathtt{T} \wedge [\![\mathtt{v} : M(m)]\!]) && \mathbf{(4)} \\
\Vdash\ &\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ ([\![\mathtt{v} : \mathtt{T}]\!] \wedge [\![\mathtt{v} : M(m)]\!]) && \mathbf{(5)} \\
\equiv\ &\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ [\![\mathtt{v}/(m \mapsto \mathtt{v}) : \mathtt{T}/M]\!] && \mathbf{(6)}
\end{aligned}
$$

where (1) is by definition of constraint generation and by definition of $\Gamma_0(!)$; (2) is by C-ARROW and C-INID; (3) follows from C-EXTRANS and from the fact that $\mathsf{ref}$ is an invariant type constructor; (4) is by C-NAMEEQ; (5) is by Lemma 1.6.3 and C-DUP; and (6) is again by definition of constraint generation.

    ◦ *Case* R-ASSIGN. The reduction is $m := \mathtt{v}/(m \mapsto \mathtt{v}_0) \longrightarrow \mathtt{v}/(m \mapsto \mathtt{v})$. Let $\mathtt{T}$ be an arbitrary type and let $M$ be a store type of domain $\{m\}$. We have

$$
\begin{aligned}
&\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ [\![m := \mathtt{v}/(m \mapsto \mathtt{v}_0) : \mathtt{T}/M]\!] \\
\Vdash\ &\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ [\![m := \mathtt{v} : \mathtt{T}]\!] && \mathbf{(1)} \\
\equiv\ &\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ \exists \mathtt{XYZ}.(\mathsf{ref}\,\mathtt{Z} \to \mathtt{Z} \to \mathtt{Z} \leq \mathtt{X} \to \mathtt{Y} \to \mathtt{T} \wedge m \preceq \mathtt{X} \wedge [\![\mathtt{v} : \mathtt{Y}]\!]) && \mathbf{(2)} \\
\equiv\ &\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ \exists \mathtt{XYZ}.(\mathsf{ref}\,M(m) \leq \mathtt{X} \leq \mathsf{ref}\,\mathtt{Z} \wedge \mathtt{Z} \leq \mathtt{T} \wedge [\![\mathtt{v} : \mathtt{Y}]\!] \wedge \mathtt{Y} \leq \mathtt{Z}) && \mathbf{(3)} \\
\equiv\ &\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ \exists \mathtt{Z}.(M(m) = \mathtt{Z} \wedge \mathtt{Z} \leq \mathtt{T} \wedge [\![\mathtt{v} : \mathtt{Z}]\!]) && \mathbf{(4)} \\
\equiv\ &\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ (M(m) \leq \mathtt{T} \wedge [\![\mathtt{v} : M(m)]\!]) && \mathbf{(5)} \\
\Vdash\ &\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ [\![\mathtt{v}/(m \mapsto \mathtt{v}) : \mathtt{T}/M]\!] && \mathbf{(6)}
\end{aligned}
$$

where (1) and (2) are by definition of constraint generation; (3) is by C-ARROW and C-INID; (4) is by C-EXTRANS, Lemma 1.6.4, and from the fact that $\mathsf{ref}$ is an invariant type constructor; (5) is by C-NAMEEQ; and (6) is obtained as in the previous case.

Second, we must check that if the configuration c $v_1$ ... $v_k/\mu$ (where $k \geq 0$) is well-typed, then either it is reducible, or c $v_1$ ... $v_k$ is a value. We only give a sketch of this proof; see the solution to Exercise 1.9.1 for details of a similar proof.

We begin by checking that every value that is well-typed with a type of the form ref T is a memory location. This assertion relies on the fact that the type constructor ref is isolated.

Next, we note that, according to the constraint generation rules, if the configuration c $v_1$ ... $v_k/\mu$ is well-typed, then a constraint of the form let $\Gamma_0$; ref $M$ in (c $\preceq$ $X_1$ $\to$ ... $\to$ $X_k$ $\to$ T $\wedge$ $[\![v_1 : X_1]\!]$ $\wedge$ ... $\wedge$ $[\![v_k : X_k]\!]$) is satisfiable. We now reason by cases on c.

∘ *Case* c is ref. If $k = 0$, then the expression is a value; otherwise, it is reducible by R-REF.

∘ *Case* c is !. We may assume $k \geq 1$, because otherwise the expression is a value. Then, by definition of $\Gamma_0(!)$, the above constraint entails let $\Gamma_0$; ref $M$ in $\exists Y.$(ref Y $\to$ Y $\leq$ $X_1$ $\to$ ... $\to$ $X_k$ $\to$ T $\wedge$ $[\![v_1 : X_1]\!]$), which, by C-ARROW, Lemma 1.6.3, and C-INEX, entails $\exists Y.$let $\Gamma_0$; ref $M$ in $[\![v_1 : \text{ref } Y]\!]$. Thus, $v_1$ is well-typed with a type of the form ref Y. By the remark above, $v_1$ must be a memory location $m$. Furthermore, because every well-typed configuration is closed, $m$ must be a member of $dom(\mu)$. As a result, the configuration ref $v_1$ ... $v_k/\mu$ is reducible by R-DEREF.

∘ *Case* c is :=. We may assume $k \geq 2$, because otherwise the expression is a value. As above, we check that $v_1$ must be a memory location and a member of $dom(\mu)$. Thus, the configuration is reducible by R-ASSIGN.

1.9.6   SOLUTION: We must first ensure that R-FIX respects $\sqsubseteq$ (Definition 1.7.5). Since the rule is pure, it is sufficient to establish that let $\Gamma_0$ in $[\![\text{fix } v_1 \ v_2 : T]\!]$ entails let $\Gamma_0$ in $[\![v_1 \ (\text{fix } v_1) \ v_2 : T]\!]$. Let $C$ stand for the constraint fix $\preceq$ $((X \to Y) \to (X \to Y)) \to X \to Y \wedge Y \leq T \wedge [\![v_1 : (X \to Y) \to (X \to Y)]\!] \wedge [\![v_2 : X]\!]$. We have

$$\begin{aligned}
&\text{let } \Gamma_0 \text{ in } [\![\text{fix } v_1 \ v_2 : T]\!] \\
\equiv\ &\text{let } \Gamma_0 \text{ in } \exists X_1 X_2.(\text{fix} \preceq X_1 \to X_2 \to T \wedge [\![v_1 : X_1]\!] \wedge [\![v_2 : X_2]\!]) &(\mathbf{1})\\
\equiv\ &\text{let } \Gamma_0 \text{ in } \exists X_1 X_2 XY.(((X \to Y) \to (X \to Y)) \to X \to Y \leq X_1 \to X_2 \to T \\
&\wedge [\![v_1 : X_1]\!] \wedge [\![v_2 : X_2]\!]) &(\mathbf{2})\\
\equiv\ &\text{let } \Gamma_0 \text{ in } \exists XY.(Y \leq T \wedge [\![v_1 : (X \to Y) \to (X \to Y)]\!] \wedge [\![v_2 : X]\!]) &(\mathbf{3})\\
\equiv\ &\text{let } \Gamma_0 \text{ in } \exists XY.C &(\mathbf{4})
\end{aligned}$$

where (1) is by definition of constraint generation; (2) is by definition of $\Gamma_0(\text{fix})$; (3) is by C-ARROW and Lemma 1.6.4; (4) is by definition of $\Gamma_0(\text{fix})$. By Theorem 1.6.2 and WEAKEN, the judgements $C \vdash v_1 : (X \to Y) \to$

$(X \to Y)$ and $C \vdash v_2 : X$ hold. By VAR, WEAKEN, APP, and SUB, it follows that $C \vdash v_1 \, (\text{fix} \, v_1) \, v_2 : T$ holds. By Theorem 1.6.6, this implies $C \Vdash [\![v_1 \, (\text{fix} \, v_1) \, v_2 : T]\!]$. By congruence of entailment and by C-EX*, (4) entails let $\Gamma_0$ in $[\![v_1 \, (\text{fix} \, v_1) \, v_2 : T]\!]$.

Second, we must check that if the configuration $\text{fix} \, v_1 \, \ldots \, v_k / \mu$ (where $k \geq 0$) is well-typed, then either it is reducible, or $\text{fix} \, v_1 \, \ldots \, v_k$ is a value. This is immediate, for it is a value when $k < 2$, and it is reducible by R-FIX when $k \geq 2$.

We now recall that the construct $\text{letrec} \, f = \lambda z.t_1 \, \text{in} \, t_2$ provided by ML-the-programming-language may be viewed as syntactic sugar for $\text{let} \, f = \text{fix} \, (\lambda f.\lambda z.t_1) \, \text{in} \, t_2$, and set forth to discover the constraint generation rule that arises out of such a definition. We have

$$
\begin{aligned}
&\phantom{\equiv} \ \text{let} \, \Gamma_0 \, \text{in} \, [\![\text{fix} \, (\lambda f.\lambda z.t_1) : T]\!] \\
&\equiv \ \text{let} \, \Gamma_0 \, \text{in} \, \exists Z.(\text{fix} \preceq Z \to T \wedge [\![\lambda f.\lambda z.t_1 : Z]\!]) && (1) \\
&\equiv \ \text{let} \, \Gamma_0 \, \text{in} \, \exists XY.(X \to Y \leq T \wedge [\![\lambda f.\lambda z.t_1 : (X \to Y) \to (X \to Y)]\!]) && (2) \\
&\equiv \ \text{let} \, \Gamma_0 \, \text{in} \, \exists XY.(X \to Y \leq T \wedge \text{let} \, f : X \to Y; z : X \, \text{in} \, [\![t_1 : Y]\!]) && (3)
\end{aligned}
$$

where (1) is by definition of constraint generation; (2) is by definition of $\Gamma_0(\text{fix})$, by C-ARROW, and by Lemma 1.6.4; and (3) follows from Lemma 1.6.5. This allows us to write

$$
\begin{aligned}
&\phantom{\equiv} \ \text{let} \, \Gamma_0 \, \text{in} \, [\![\text{let} \, f = \text{fix} \, (\lambda f.\lambda z.t_1) \, \text{in} \, t_2 : T]\!] \\
&\equiv \ \text{let} \, \Gamma_0; f : \forall Z[[\![\text{fix} \, (\lambda f.\lambda z.t_1) : Z]\!]].Z \, \text{in} \, [\![t_2 : T]\!] && (4) \\
&\equiv \ \text{let} \, \Gamma_0; f : \forall Z[\exists XY.(X \to Y \leq Z \wedge \text{let} \, f : X \to Y; z : X \, \text{in} \, [\![t_1 : Y]\!])].Z \, \text{in} \, [\![t_2 : T]\!] && (5) \\
&\equiv \ \text{let} \, \Gamma_0; f : \forall XY[\text{let} \, f : X \to Y; z : X \, \text{in} \, [\![t_1 : Y]\!]].X \to Y \, \text{in} \, [\![t_2 : T]\!] && (6)
\end{aligned}
$$

where (4) is by definition of constraint generation; (5) follows from C-LETDUP and from the previous series of equivalences; (6) is by C-LETEX, C-EXTRANS and Lemma 1.3.22.

1.9.21   SOLUTION: We have

$$
\begin{aligned}
&\phantom{\equiv} \ [\![\text{match} \, t_1 \, \text{with} \, z \, . \, t_2 : T]\!] \\
&\equiv \ \text{let} \, \forall XX'[[\![t_1 : X]\!] \wedge \text{let} \, z : X' \, \text{in} \, [\![X : z]\!]].(z : X') \, \text{in} \, [\![t_2 : T]\!] && (1) \\
&\equiv \ \text{let} \, z : \forall X'[\exists X.([\![t_1 : X]\!] \wedge X \leq X')].X' \, \text{in} \, [\![t_2 : T]\!] && (2) \\
&\equiv \ \text{let} \, z : \forall X'[[\![t_1 : X']\!]].X' \, \text{in} \, [\![t_2 : T]\!] && (3) \\
&\equiv \ [\![\text{let} \, z = t_1 \, \text{in} \, t_2 : T]\!] && (4)
\end{aligned}
$$

where (1) is by definition of constraint generation for match; (2) is by definition of constraint generation for patterns, by C-INID, C-IN*, and C-LETEX; (3) is by Lemma 1.6.4; (4) is by definition of constraint generation for let.

1.9.26   SOLUTION: The type scheme $\forall \bar{X}.T \to T$ may be written $\forall \bar{X}.[X \mapsto T](X \to X)$. Furthermore, $\bar{X} \, \# \, \forall X.X \to X$ holds. Thus, $\forall \bar{X}.T \to T$ is an instance of $\forall X.X \to X$

in the sense of DM-INST'. Since DM-INST' is an admissible rule for the type system DM, and since it is clear that the identity function $\lambda z.z$ has type $\forall X.X \to X$, it must also have type $\forall \bar{X}.T \to T$. (A more direct proof of this fact would not be difficult.) So, the destructor $(\cdot : \exists \bar{X}.T)$ has not only identity semantics, but also an identity type. This shows that our definitions are sound.

Let us now check requirement (i) of Definition 1.7.6. Since R-ANNOTATION is pure, it suffices to show that let $\Gamma_0$ in $[\![(v : \exists \bar{X}.T) : T']\!]$ entails let $\Gamma_0$ in $[\![v : T']\!]$. Now, we have

$$
\begin{aligned}
&\quad \text{let } \Gamma_0 \text{ in } [\![(v : \exists \bar{X}.T) : T']\!] \\
&\equiv \text{ let } \Gamma_0 \text{ in } \exists X \bar{X}. (T \to T \leq X \to T' \wedge [\![v : X]\!]) \quad (1) \\
&\equiv \text{ let } \Gamma_0 \text{ in } \exists X \bar{X}. (X \leq T \leq T' \wedge [\![v : X]\!]) \qquad\quad (2) \\
&\Vdash \text{ let } \Gamma_0 \text{ in } [\![v : T']\!] \qquad\qquad\qquad\qquad\qquad (3)
\end{aligned}
$$

where (1) is by definition of constraint generation and by definition of $\Gamma_0((\cdot : \exists \bar{X}.T))$; (2) is by C-ARROW; and (3) follows from Lemma 1.6.3 and C-EX*.

1.10.5   SOLUTION: We have

$$
\begin{aligned}
&\quad \text{let } \Gamma_0 \text{ in } \exists Z. [\![(\lambda z.z \mathbin{\hat{+}} \hat{1} : \forall X.X \to X) : Z]\!] \\
&\equiv \text{ let } \Gamma_0 \text{ in } \exists Z. (\forall X. [\![\lambda z.z \mathbin{\hat{+}} \hat{1} : X \to X]\!] \wedge \exists X. (X \to X \leq Z)) \quad (\mathbf{1}) \\
&\equiv \text{ let } \Gamma_0 \text{ in } \forall X. \text{let } z : X \text{ in } [\![z \mathbin{\hat{+}} \hat{1} : X]\!] \qquad\qquad\qquad\qquad (\mathbf{2}) \\
&\equiv \forall X. (\texttt{int} \to \texttt{int} \to \texttt{int} \leq X \to \texttt{int} \to X) \qquad\qquad\qquad (\mathbf{3}) \\
&\equiv \forall X. (X = \texttt{int}) \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad (\mathbf{4}) \\
&\equiv \text{false} \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad (\mathbf{5})
\end{aligned}
$$

where (1) is by definition of constraint generation for universal type annotations; (2) is obtained by restricting the scope of $\exists Z$ to the second conjunct, then dropping the latter altogether, since it is equivalent to $\texttt{true}$, and by Lemma 1.6.5; (3) is obtained by definition of constraint generation, by definition of $\Gamma_0(\hat{+})$ and of $\Gamma_0(\hat{1})$, and by a few simple equivalence laws; (4) follows from C-ARROW and antisymmetry of subtyping; (5) follows from the fact that $\texttt{int}$ and (say) $\texttt{int} \to \texttt{int}$ have distinct interpretations, since the type constructors $\texttt{int}$ and $\to$ are incompatible. On the other hand, we have

$$
\begin{aligned}
&\quad \text{let } \Gamma_0 \text{ in } \exists Z. [\![(\lambda z.z : \forall X.X \to X) : Z]\!] \\
&\equiv \text{ let } \Gamma_0 \text{ in } \forall X. \text{let } z : X \text{ in } [\![z : X]\!] \qquad (\mathbf{1}) \\
&\equiv \forall X. (X \leq X) \qquad\qquad\qquad\qquad\qquad\quad (\mathbf{2}) \\
&\equiv \text{true} \qquad\qquad\qquad\qquad\qquad\qquad\quad (\mathbf{3})
\end{aligned}
$$

where (1) is obtained as above; (2) by definition of constraint generation, C-INID and C-IN*; (3) is by reflexivity of subtyping.

1.10.6   SOLUTION: Under the naïve constraint generation rule for universal type variable introduction, the constraint $[\![\forall X.(\lambda z.z : X \rightarrow X) : Z]\!]$ is equivalent to $\forall X.([\![\lambda z.z : X \rightarrow X]\!] \wedge X \rightarrow X \leq Z)$. Since the first conjunct is a tautology, this is in turn equivalent to $\forall X.(X \rightarrow X \leq Z)$. In a nondegenerate free term model where subtyping is interpreted as equality, this constraint is unsatisfiable. In a non-structural subtyping model equipped with a least type $\bot$ and a greatest type $\top$, it is equivalent to $\bot \rightarrow \top \leq Z$. This is a pretty restrictive constraint: since no value has type $\bot$, a function whose type is (a supertype of) $\bot \rightarrow \top$ cannot ever be invoked at runtime. This situation is clearly unsatisfactory. Checking that $\forall X.[\![\lambda z.z : X \rightarrow X]\!]$ holds was indeed part of our intent, but constraining $Z$ to be a supertype of $X \rightarrow X$ for *every* $X$ was not.

1.10.7   SOLUTION: Let $\bar{X} \supseteq ftv(T)$ **(1)** and $\bar{X} \# ftv(t)$ **(2)**. We may assume, *w.l.o.g.*, $\bar{X} \# ftv(T')$ **(3)**. By (1), (2), (3), and by definition of constraint generation for local universal type annotations, $[\![(t : \forall \bar{X}.T) : T']\!]$ is well-defined and is $\forall \bar{X}.[\![t : T]\!] \wedge \exists \bar{X}.(T \leq T')$ **(4)**. By (3) and by definition of constraint generation for introduction of universal type variables and for general type annotations, $[\![\forall \bar{X}.(t : T) : T']\!]$ is $\forall \bar{X}.\exists Z.([\![t : T]\!] \wedge T \leq Z) \wedge \exists \bar{X}.([\![t : T]\!] \wedge T \leq T')$, where $Z$ is fresh, which we may immediately simplify to $\forall \bar{X}.[\![t : T]\!] \wedge \exists \bar{X}.([\![t : T]\!] \wedge T \leq T')$ **(5)**. Using C-EXAND and Lemma 1.10.1, it is straightforward to check that (4) and (5) are equivalent.

1.10.9   SOLUTION: We have

$$
\begin{aligned}
&\exists Z.[\![\lambda z.\forall X.(z : X) : Z]\!] \\
\Vdash \quad &\exists Z_1 Z_2.\mathsf{let}\ z : Z_1\ \mathsf{in}\ [\![\forall X.(z : X) : Z_2]\!] \quad &\textbf{(1)} \\
\Vdash \quad &\exists Z_1.\forall X.(Z_1 \leq X) \quad &\textbf{(2)}
\end{aligned}
$$

where (1) is by definition of constraint generation for $\lambda$-abstractions, dropping the constraint that relates $Z$, $Z_1$, and $Z_2$; (2) is by definition of constraint generation for universal type variable introduction, this time dropping information about $Z_2$. Now, in a nondegenerate equality model, the constraint (2) is equivalent to $\mathsf{false}$. In fact, for (2) to be satisfiable, the interpretation of subtyping must admit a least element $\bot$. We now see that $[\![\lambda z.\forall X.(z : X) : Z]\!]$ is a very restrictive constraint. Indeed, it requires $z$ to have every type at once. Because $z$ is $\lambda$-bound—hence monomorphic—it must in fact have type $\bot$. On the other hand, we have

$$
\begin{aligned}
&\exists Z.[\![\forall X.\lambda z.(z : X) : Z]\!] \\
\equiv \quad &\forall X.\exists Z.[\![\lambda z.(z : X) : Z]\!] \quad &\textbf{(1)} \\
\equiv \quad &\forall X.\exists Z Z_1 Z_2.(Z_1 \leq X \wedge X \leq Z_2 \wedge Z_1 \rightarrow Z_2 \leq Z) \quad &\textbf{(2)} \\
\equiv \quad &\mathsf{true} \quad &\textbf{(3)}
\end{aligned}
$$

where (1) is by definition of constraint generation for universal type variable introduction, dropping the second conjunct, which is entailed by the first; (2) is by Lemma 1.6.5, by definition of constraint generation for general type annotations, and by a few simple equivalence laws; (3) follows from C-NAMEEQ and the witness substitution $[Z_1 \mapsto X, Z_2 \mapsto X, Z \mapsto (X \to X)]$.

1.10.10   SOLUTION: We have

$$
\begin{array}{cll}
& [\![\texttt{letrec } f : S = \lambda z.t_1 \texttt{ in } t_2 : T]\!] & \\
\equiv & \texttt{let } f : \forall X[[\![\texttt{fix } f : S.\lambda z.t_1 : X]\!]].X \texttt{ in } [\![t_2 : T]\!] & \textbf{(1)} \\
\equiv & \texttt{let } f : \forall X[\texttt{let } f : S \texttt{ in } [\![\lambda z.t_1 : S]\!] \wedge S \preceq X].X \texttt{ in } [\![t_2 : T]\!] & \textbf{(2)} \\
\equiv & \texttt{let } f : S \texttt{ in } [\![\lambda z.t_1 : S]\!] \wedge \texttt{let } f : \forall X[S \preceq X].X \texttt{ in } [\![t_2 : T]\!] & \textbf{(3)} \\
\equiv & \texttt{let } f : S \texttt{ in } ([\![\lambda z.t_1 : S]\!] \wedge [\![t_2 : T]\!]) & \textbf{(4)}
\end{array}
$$

where (1) is by definition of the `letrec` syntactic sugar and by the definition of constraint generation for `let` constructs; we have $X \notin \mathit{ftv}(S, t_1)$; (2) is by definition of constraint generation for `fix`; (3) is by C-LETAND; (4) follows from the equivalence between the type schemes $\forall X[S \preceq X].X$ and $S$—which itself is a direct consequence of C-EXTRANS—and from C-INAND.

1.11.16   SOLUTION: We reason simultaneously in both the subtyping model or the equal-only model, that is, we only rely on properties that are valid in both models.

We must first ensure that rules RD-DEFAULT, RD-FOUND, and RD-FOLLOW respect (Definition 1.7.5).

○ *Case* RD-DEFAULT. The reduction is $\{v\}.\{\ell\} \xrightarrow{\delta} v$, which is pure. Therefore, it is sufficient to establish that $\texttt{let } \Gamma_0 \texttt{ in } [\![\{v\}.\{\ell\} : T]\!]$ entails $\texttt{let } \Gamma_0 \texttt{ in } [\![v : T]\!]$. In fact, we have:

$$
\begin{array}{cll}
& \texttt{let } \Gamma_0 \texttt{ in } [\![\{v\}.\{\ell\} : T]\!] & \\
\equiv & \texttt{let } \Gamma_0 \texttt{ in } \exists XY.(\cdot.\{\ell\} \preceq X \to T \wedge \{\cdot\} \preceq Y \to X \wedge [\![v : Y]\!]) & \textbf{(1)} \\
\equiv & \texttt{let } \Gamma_0 \texttt{ in } \exists XY.(\exists X_1 X_2.(\Pi(\ell : X_1 ; X_2) \to X_1 \leq X \to T) & \textbf{(2)} \\
& \qquad\qquad\qquad\quad \wedge \exists Y_1.(Y_1 \to \Pi(\partial Y_1) \leq Y \to X) \wedge [\![v : Y]\!]) & \\
\Vdash & \texttt{let } \Gamma_0 \texttt{ in } \exists X_2 Y.(\partial Y \leq (\ell : X_1 ; X_2) \wedge X_1 \leq T \wedge [\![v : Y]\!]) & \textbf{(3)} \\
\Vdash & \texttt{let } \Gamma_0 \texttt{ in } \exists Y.(Y \leq X_1 \wedge X_1 \leq T \wedge [\![v : Y]\!]) & \textbf{(4)} \\
\Vdash & \texttt{let } \Gamma_0 \texttt{ in } [\![v : T]\!] & \textbf{(5)}
\end{array}
$$

where (1) is by definition of constraint generation; (2) is by definition of $\Gamma_0$, C-INID; (3) by variances of $\Pi$, $\ell$, and $\to$, C-AND, C-EX*, C-EXAND; (4) by C-ROW-DL and covariance of $\ell$; (5) by Lemma 1.6.3.

○ *Case* RD-FOUND: The reduction is $\{w \texttt{ with } \ell = v\}.\{\ell\} \xrightarrow{\delta} v$. It suffices to establish $\texttt{let } \Gamma_0 \texttt{ in } [\![\{w \texttt{ with } \ell = v\}.\{\ell\} : T]\!]$ entails $\texttt{let } \Gamma_0 \texttt{ in } [\![v : T]\!]$. In fact, we

have:

$$\mathsf{let}\ \Gamma_0\ \mathsf{in}\ [\![\{\mathtt{w}\ \mathsf{with}\ \ell = \mathtt{v}\}.\{\ell\} : \mathtt{T}]\!]$$

$\equiv$   $\mathsf{let}\ \Gamma_0\ \mathsf{in}\ \exists \mathtt{XYY}'.(\cdot.\{\ell\} \preceq \mathtt{X} \to \mathtt{T} \wedge \{\cdot\ \mathsf{with}\ \ell = \cdot\} \preceq \mathtt{Y} \to \mathtt{Y}' \to \mathtt{X}\ \wedge$

$\hspace{6cm} \wedge\ [\![\mathtt{w} : \mathtt{Y}]\!] \wedge [\![\mathtt{v} : \mathtt{Y}']\!])$   **(1)**

$\equiv$   $\mathsf{let}\ \Gamma_0\ \mathsf{in}\ \exists \mathtt{XYY}'.(\exists \mathtt{X}_1 \mathtt{X}_2.(\Pi(\ell : \mathtt{X}_1\ ;\ \mathtt{X}_2) \to \mathtt{X}_1 \leq \mathtt{X} \to \mathtt{T})$

$\hspace{1.3cm} \wedge\ \exists \mathtt{Y}_1 \mathtt{Y}_2 \mathtt{Y}_3.(\Pi(\ell : \mathtt{Y}_1\ ;\ \mathtt{Y}_3) \to \mathtt{Y}_2 \to \Pi(\ell : \mathtt{Y}_2\ ;\ \mathtt{Y}_3) \leq \mathtt{Y} \to \mathtt{Y}' \to \mathtt{X})$

$\hspace{6cm} \wedge\ [\![\mathtt{w} : \mathtt{Y}]\!] \wedge [\![\mathtt{v} : \mathtt{Y}']\!])$   **(2)**

$\Vdash$   $\mathsf{let}\ \Gamma_0\ \mathsf{in}\ \exists \mathtt{Y}'\mathtt{X}_1 \mathtt{Y}_2.(\mathtt{Y}' \leq \mathtt{Y}_2 \wedge \mathtt{Y}_2 \leq \mathtt{X}_1 \wedge \mathtt{X}_1 \leq \mathtt{T} \wedge [\![\mathtt{v} : \mathtt{Y}']\!])$   **(3)**

$\Vdash$   $\mathsf{let}\ \Gamma_0\ \mathsf{in}\ [\![\mathtt{v} : \mathtt{T}]\!]$   **(4)**

where (1) is by definition of constraint generation; (2) is by definition of $\Gamma_0$, C-INID; (3) by variances of $\Pi$, $\ell$, and $\to$, C-AND, C-EX\*, C-EXAND; (4) by Lemma 1.6.3.

   ∘ *Case* RD-FOLLOW The proof is similar to the previous case.

We must now check that if the configuration $F\ \mathtt{v}_1\ \ldots\ \mathtt{v}_k/\mu$ is is well-typed, then either it is reducible, or it is a value.

   We begin by checking that every value that is well-typed with type $\Pi\ \mathtt{T}$ is a record value, that is, either of the form $\{\mathtt{v}'\}$ or $\{\mathtt{v}''\ \mathsf{with}\ \ell' = \mathtt{v}'\}$. Indeed, suppose that $\mathsf{let}\ \Gamma_0\ \mathsf{in}\ [\![\mathtt{v} : \Pi\ \mathtt{T}]\!]$ is satisfiable. Then, $v$ cannot be a program variable, for a well-typed value must be closed; $v$ cannot be a memory location $m$, for otherwise $\mathsf{ref}\ M(m) \leq \Pi\ \mathtt{T}$ would be satisfiable—but the top type constructors $\mathsf{ref}$ and $\Pi$ are incompatible (since $\Pi$ is isolated); $\mathtt{v}$ cannot be a partial application of a constructor or a primitive, nor a $\lambda$-abstraction, since otherwise $\mathtt{T}' \to \mathtt{T}'' \leq \Pi\ \mathtt{T}$ would be satisfiable but the top type constructors $\to$ and $\Pi$ are incompatible (since they are both isolated); thus $\mathtt{v}$ must either be of the form $\{\mathtt{v}\}$ or $\{\mathtt{w}\ \mathsf{with}\ \ell = \mathtt{v}\}$, for these are the only left cases.

   Next, we note that, according to the constraint generation rules, if the configuration $\mathtt{c}\ \mathtt{v}_1\ \ldots\ \mathtt{v}_k/\mu$ is well-typed, then a constraint of the form $\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ (\mathtt{c} \preceq \mathtt{X}_1 \to \ldots \to \mathtt{X}_k \to \mathtt{T} \wedge [\![\mathtt{v}_1 : \mathtt{X}_1]\!] \wedge \ldots \wedge [\![\mathtt{v}_k : \mathtt{X}_k]\!])$ is satisfiable. We now reason by cases on $\mathtt{c}$.

   ∘ *Case* $\mathtt{c}$ is $\{\cdot\}$. We may asume $k \geq 2$, since otherwise, the expression is a value. Then $\Gamma_0(\mathtt{c})$ is $\forall \mathtt{XY}.\mathtt{X} \to \Pi(\partial \mathtt{X})$, so by C-INID and C-ARROW the above constraint entails $\exists \mathtt{X}.(\Pi(\partial \mathtt{X}) \leq \mathtt{X}_2 \to \ldots \to \mathtt{T})$, which by C-Class-I entails $\mathtt{false}$ since $\to$ and $\Pi$ are imcompatible. Thus, this case cannot occur.

   ∘ *Case* $\mathtt{c}$ is $\{\cdot\ \mathsf{with}\ \ell = \cdot\}$. Similar to the previous case.

   ∘ *Case* $\mathtt{c}$ is $\cdot.\{\ell\}$. We may asume $k \geq 1$, since otherwise, the expression is a value. Then $\Gamma_0(\mathtt{c})$ is $\forall \mathtt{XY}.\Pi(\ell : \mathtt{X}\ ;\ \mathtt{Y}) \to \mathtt{X}$, so by C-INID and C-ARROW the above constraint entails $\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ (\exists \mathtt{XY}.(\mathtt{X}_1 \leq \Pi(\ell : \mathtt{X}\ ;\ \mathtt{Y})) \wedge [\![\mathtt{v}_1 : \mathtt{X}_1]\!])$, which by lemma 1.6.3 entails $\mathsf{let}\ \Gamma_0; \mathsf{ref}\ M\ \mathsf{in}\ \exists \mathtt{XY}.[\![\mathtt{v}_1 : \Pi(\ell : \mathtt{X}\ ;\ \mathtt{Y})]\!]$. Thus $\mathtt{v}_1$ is a record value, that is, either of the form $\{\mathtt{v}'\}$ and the configuration is reducible

to $v'$ or of the form $\{v''$ with $\ell' = v'\}$ and the configuration is reducible to either $v'$ or $v''.\{\ell\}$.

1.11.17   SOLUTION:  We add a collection of destructors $\cdot[\ell_1 \leftrightarrow \ell_2]$ of arity 1 for all pairs of distinct labels, with the following semantics:

$$\{v\}[\ell_1 \leftrightarrow \ell_2] \xrightarrow{\delta} v$$
$$\{w \text{ with } \ell = v\}[\ell_1 \leftrightarrow \ell_2] \xrightarrow{\delta} \{w[\ell_1 \leftrightarrow \ell_2] \text{ with } \ell = v\} \qquad\qquad \text{if } \ell \notin \{\ell_1, \ell_2\}$$
$$\{w \text{ with } \ell = v\}[\ell_1 \leftrightarrow \ell_2] \xrightarrow{\delta} \{w[\ell_1 \leftrightarrow \ell_2] \text{ with } \bar{\ell} = v\} \qquad \text{if } \{\ell, \bar{\ell}\} = \{\ell_1, \ell_2\}$$

The initial environment $\Gamma_0$ must be extended with the following typing asumption:

$$\cdot[\ell_1 \leftrightarrow \ell_2] : \quad \forall X_1 X_2 Y.\ \Pi(\ell_1 : X_1\ ;\ \ell_2 : X_2\ ;\ Y) \to \Pi(\ell_1 : X_2\ ;\ \ell_2 : X_1\ ;\ Y)$$

We must then check subjection reduction for the new primitive. Since we only added a constructor, it sufficies to check progress for the new primitive, that is, that well-typed expressions of the form  $[\ell_1 \leftrightarrow \ell_2] v_1\ \ldots\ v_n$ are either value or can be further reduced. Both parts are easy and similar to the corresponding parts in Exercice 1.11.16.

1.11.18   SOLUTION: There are several solutions. One of them is to asume a fixed total ordering on row-labels, and to retain as constructors only $\ell^{\kappa, L}$ such that $\ell < L$, that is $\ell < \ell'$ for all $\ell' \in L$; other constants $\ell^{\kappa, L}$ such that $\ell \not< L$ are moved from constructors to the status of destructors with the following collection of reduction rules:

$$\{\{w \text{ with } \ell' = v'\} \text{ with } \ell = v\} \xrightarrow{\delta} \{\{w \text{ with } \ell = v\} \text{ with } \ell' = v'\}$$
$$\text{(RD-TRANSPOSE)}$$

for all labels $\ell$ and $\ell'$ such that $\ell' < \ell$ and

$$\{\{w \text{ with } \ell = v'\} \text{ with } \ell = v\} \xrightarrow{\delta} \{w \text{ with } \ell = v\} \qquad \text{(RD-DISCARD)}$$

for all labels $\ell$. It is now obvious that values are in normal forms, in the sense that explicit fields are never repeated and are always listed in order. Typing rules need not be changed, so requirement (i) of Definition 1.7.6 still holds. Requirement (ii) need to be check, in particular, for the new primitives $\ell^L$, which we leave to the reader (the proof for $\cdot.\{\ell\}$ should hold unchanged).

1.11.19   SOLUTION:  Let map have type $\Pi(X \to Y) \to \Pi(X) \to \Pi(Y)$, and the following reduction rules in the semantics with normal forms:

$$\text{map } \{v' \text{ with } \ell = v\}\ w \xrightarrow{\delta} \{\text{map } v'\ w \text{ with } \ell = v\ (w.\{\ell\})\}$$
$$\text{map } v\ \{w' \text{ with } \ell = w\} \xrightarrow{\delta} \{\text{map } v\ w' \text{ with } \ell = (v.\{\ell\})\ w\}$$
$$\text{map } \{v\}\ \{w\} \xrightarrow{\delta} \{v\ w\}$$

1.11.22  SOLUTION: To ensure that the field is not present in the argument of extension, it sufficies to restrict its the typing asumptions as follows:

$$\langle \cdot \text{ with } \ell = \cdot \rangle : \forall \text{XX}'\text{Y}. \; \Pi(\ell : \text{abs} \; ; \; \text{Y}) \to \text{X}' \to \Pi(\ell : \text{pre } \; \text{X}' \; ; \; \text{Y}).$$

To remove an existing field, we can use the following syntactic sugar:

$$\cdot \setminus \ell \stackrel{\text{def}}{=} \lambda \text{v}.\{\text{v with } \ell = \text{abs}\}$$
$$: \forall \text{XY}. \; \Pi(\ell : \text{X} \; ; \; \text{Y}) \to \Pi(\ell : \text{abs} \; ; \; \text{Y})$$

The following weaker typing asumption could also be used to ensure that the field is always present before removal:

$$\forall \text{XY}. \; \Pi(\ell : \text{pre } \; \text{X} \; ; \; \text{Y}) \to \Pi(\ell : \text{abs} \; ; \; \text{Y})$$

1.11.25  SOLUTION: The proof is similar to 1.11.16 but slightly more complex because we must also check that labels are defined when accessed, and with subtyping.

We reason simultaneously in both the subtyping model or the equal-only model, that is, we only rely on properties that are valid in both models.

We must first ensure that rules RE-FOUND and RE-FOLLOW respect (Definition 1.7.5).

○ *Case* RE-FOUND: See Exercice **??**. In line **??**, field $\ell$ is pre $\text{X}_1$ instead of $\text{X}_1$ and pre $\text{Y}_2$ instead of $\text{Y}_2$ and step **??** also uses covariance of pre.

○ *Case* RE-FOLLOW The proof is similar.

We must then check that if the configuration $F \; \text{v}_1 \; \ldots \; \text{v}_k/\mu$ is is well-typed, then either it is reducible, or it is a value.

We begin by checking that every value that is well-typed with type $\Pi \, \text{T}$ is a record value, that is, either of the form $\langle \rangle$ or $\langle \text{v}'' \text{ with } \ell' = \text{v}' \rangle$. See Exercice 1.11.16.

Next, we note that, according to the constraint generation rules, if the configuration $\text{c } \text{v}_1 \; \ldots \; \text{v}_k/\mu$ is well-typed, then a constraint of the form let $\Gamma_0; \text{ref } M$ in $(\text{c} \preceq \text{X}_1 \to \ldots \to \text{X}_k \to \text{T} \land [\![\text{v}_1 : \text{X}_1]\!] \land \ldots \land [\![\text{v}_k : \text{X}_k]\!])$ is satisfiable. We now reason by cases on $\text{c}$.

○ *Case* $\text{c}$ is $\langle \rangle$ or $\langle \cdot \text{ with } \ell = \cdot \rangle$. See Exercice 1.11.16.

○ *Case* $\text{c}$ is $\cdot.\langle \ell \rangle$. We may asume $k \geq 1$, since otherwise, the expression is a value. Then $\Gamma_0(\text{c})$ is $\forall \text{XY}.\Pi(\ell : \text{pre } \; \text{X} \; ; \; \text{Y}) \to \text{X}$, so by C-INID and C-ARROW the above constraint entails let $\Gamma_0; \text{ref } M$ in $(\exists \text{XY}.(\text{X}_1 \leq \Pi(\ell : \text{pre } \; \text{X} \; ; \; \text{Y})) \land [\![\text{v}_1 : \text{X}_1]\!])$, which by lemma 1.6.3 entails let $\Gamma_0; \text{ref } M$ in $\exists \text{XY}.[\![\text{v}_1 : \Pi(\ell : \text{pre } \; \text{X} \; ; \; \text{Y})]\!]$. Thus $\text{v}_1$ is a record value, that is, either of the form $\langle \rangle$ or $\langle \text{v}'' \text{ with } \ell = \text{v}' \rangle$. In fact, the former case cannot occur, since let $\Gamma_0; \text{ref } M$ in $\exists \text{XY}.[\![\langle \rangle : \Pi(\ell : \text{pre } \; \text{X} \; ; \; \text{Y})]\!]$ entails $\exists \text{XY}\Pi(\partial\text{abs}) \leq \Pi(\ell : \text{pre } \text{X} \; ; \; \text{Y})$ by C-INID and C-IN\*, which in turns

entails $\exists \mathtt{X}.\mathsf{abs} \leq \mathsf{pre}\ \mathtt{X}$ by C-Row-DL and covariance of $\Pi$ and $\ell$. However, this constraint is equivalent to false, because $\phi(\mathsf{abs}) \leq \phi(\mathsf{pre}\ \mathtt{X})$ does not hold in any ground assignment $\phi$. Thus $\mathtt{v}_1$ is $\langle \mathtt{v}'' \ \mathsf{with}\ \ell' = \mathtt{v}' \rangle$ and the configuration is reducible to $\mathtt{v}'$ if $\ell'$ is $\ell$ or $\mathtt{v}''$ otherwise.