

# 1

## *Substructural Type Systems*

*David Walker*

Advanced type systems make it possible to restrict access to data structures and to limit the use of newly-defined operations. Oftentimes, this sort of access control is achieved through the definition of new abstract types under control of a particular module. For example, consider the following simplified file system interface.

```
type file

val open   : string → file option
val read   : file → string * file
val append : file * string → file
val write  : file * string → file
val close  : file → unit
```

By declaring that the type `file` is abstract, the implementer of the module can maintain strict control over the representation of files. A client has no way to accidentally (or maliciously) alter any of the file's representation invariants. Consequently, the implementer may assume that the invariants that he or she establishes upon opening a file hold before any `read`, `append`, `write` or `close`.

While abstract types are a powerful means of controlling the structure of data, they are not sufficient to limit the *ordering* and *number of uses* of functions in an interface. Try as we might, there is no (static) way to prevent a file from being read after it has been closed. Likewise, we cannot stop a client from closing a file twice or forgetting to close a file.

This chapter introduces *substructural* type systems, which augment standard type abstraction mechanisms with the ability to control the number and order of uses of a data structure or operation. Substructural type systems are particularly useful for constraining interfaces that provide access to system

resources such as files, locks and memory. Each of these resources undergoes a series of changes of state throughout its lifetime. Files, as we have seen, may be open or closed; locks may be held or not; and memory may be allocated or deallocated. Substructural type systems provide sound static mechanisms for keeping track of just these sorts of state changes and preventing operations on objects in an invalid state.

The bulk of this chapter will focus on applications of substructural type systems to the control of memory resources. Memory is a pervasive resource that must be managed carefully in any programming system so it makes an excellent target of study. However, the general principles that we establish can be applied to other sorts of resources as well.

## 1.1 Structural Properties

Most of the type systems in this book allow *unrestricted* use of variables in the type checking context. For instance, each variable may be used once, twice, three times, or not at all. A precise analysis of the properties of such variables will suggest a whole new collection of type systems.

To begin our exploration, we will analyze the simply-typed lambda calculus, which is reviewed in Figure 1-1. In this discussion, we are going to be particularly careful when it comes to the form of the type-checking context  $\Gamma$ . We will consider such contexts to be simple lists of variable-type pairs. The "," operator appends a pair to the end of the list. We also write  $(\Gamma_1, \Gamma_2)$  for the list that results from appending  $\Gamma_2$  onto the end of  $\Gamma_1$ . As usual, we allow a given variable to appear at most once in a context and to maintain this invariant, we implicitly alpha-convert bound variables before entering them into the context.

We are now in position to consider three basic *structural* properties satisfied by our simply-typed lambda calculus. The first property, *exchange*, indicates that the order in which we write down variables in the context is irrelevant. A corollary of exchange is that if we can type check a term with the context  $\Gamma$ , then we can type check that term with any permutation of the variables in  $\Gamma$ . The second property, *weakening*, indicates that adding extra, unneeded assumptions to the context, does not prevent a term from type checking. Finally, the third property, *contraction*, states that if we can type check a term using two identical assumptions ( $x_2 : T_1$  and  $x_3 : T_1$ ) then we can check the same term using a single assumption.

- 1.1.1 LEMMA [EXCHANGE]: If  $\Gamma_1, x_1 : T_1, x_2 : T_2, \Gamma_2 \vdash t : T$  then  
 $\Gamma_1, x_2 : T_2, x_1 : T_1, \Gamma_2 \vdash t : T$  □
- 1.1.2 LEMMA [WEAKENING]: If  $\Gamma_1, \Gamma_2 \vdash t : T$  then  $\Gamma_1, x_1 : T_1, \Gamma_2 \vdash t : T$  □

Syntax		Typing	
$b ::=$			$\Gamma \vdash t : T$
true	<i>booleans:</i>	$\frac{}{\Gamma_1, x:T, \Gamma_2 \vdash x : T}$	(T-VAR)
false	true	$\frac{}{\Gamma \vdash b : \text{Bool}}$	(T-BOOL)
$t ::=$	<i>terms:</i>	$\frac{\Gamma \vdash t_1 : \text{Bool} \quad \Gamma \vdash t_2 : T \quad \Gamma \vdash t_3 : T}{\Gamma \vdash \text{if } t_1 \text{ then } t_2 \text{ else } t_3 : T}$	(T-IF)
x	false	$\frac{\Gamma, x:T_1 \vdash t_2 : T_2}{\Gamma \vdash \lambda x:T_1. t_2 : T_1 \rightarrow T_2}$	(T-ABS)
b	variable	$\frac{\Gamma \vdash t_1 : T_{11} \rightarrow T_{12} \quad \Gamma \vdash t_2 : T_{11}}{\Gamma \vdash t_1 t_2 : T_{12}}$	(T-APP)
if t then t else t	boolean		
$\lambda x:T. t$	conditional		
t t	abstraction		
$T ::=$	<i>types:</i>		
Bool	booleans		
$T \rightarrow T$	type of functions		
$\Gamma ::=$	<i>contexts:</i>		
$\emptyset$	empty context		
$\Gamma, x:T$	term variable binding		

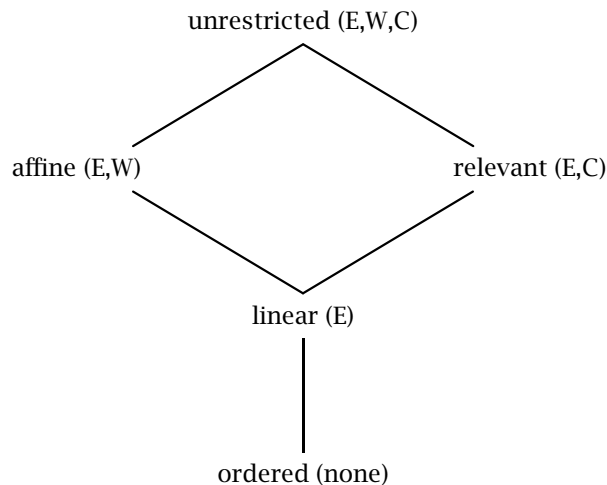
Figure 1-1: Simply-typed lambda calculus with booleans

- 1.1.3 LEMMA [CONTRACTION]: If  $\Gamma_1, x_2:T_1, x_3:T_1, \Gamma_2 \vdash t : T_2$  then  $\Gamma_1, x_1:T_1, \Gamma_2 \vdash [x_2 \mapsto x_1][x_3 \mapsto x_1]t : T_2$  □
- 1.1.4 EXERCISE [RECOMMENDED, ★]: Prove that exchange, weakening and contraction lemmas hold for the simply-typed lambda calculus. □

A *substructural type system* is any type system that is designed so that one or more of the structural properties do not hold. Different substructural type systems arise when different properties are withheld.

- *Linear* type systems ensure that every variable is used exactly once by allowing exchange but not weakening or contraction.
- *Affine* type systems ensure that every variable is used at most once by allowing exchange and weakening, but not contraction.
- *Relevant* type systems ensure that every variable is used at least once by allowing exchange and contraction, but not weakening.
- *Ordered* type systems ensure that every variable is used exactly once and in the order in which it is introduced. Ordered type systems do not allow any of the structural properties.

The picture below can serve as a mnemonic for the relationship between these systems. The system at the bottom of the diagram (the ordered type system) admits no structural properties. As we proceed upwards in the diagram, we add structural properties: E stands for exchange; W stands for weakening; and C stands for contraction. It might be possible to define type systems containing other combinations of structural properties, such as contraction only or weakening only, but so far researchers have not found applications for such combinations. Consequently, we have excluded them from the diagram.



The diagram can be realized as a relation between the systems. We say system  $q_1$  is more restrictive than system  $q_2$  and write  $q_1 \sqsubseteq q_2$  when system  $q_1$  exhibits fewer structural rules than system  $q_2$ . Figure 1-2 specifies the relation, which we will find useful in the coming sections of this chapter.

## 1.2 A Linear Type System

In order to safely deallocate data, we need to know that the data we deallocate is never used in the future. Unfortunately, we cannot, in general, deduce whether data will be used after execution passes a certain program point: The problem is clearly undecidable. However, there are a number of sound, but useful approximate solutions. One such solution may be implemented using a *linear type system*. Linear type systems ensure that objects are used exactly once, so it is completely obvious that after the use of an object, it may be safely deallocated.

$q ::=$	<i>system:</i>	$\text{ord} \sqsubseteq \text{lin}$	(Q-ORDLIN)
ord	<i>ordered</i>	$\text{lin} \sqsubseteq \text{rel}$	(Q-LINREL)
lin	<i>linear</i>	$\text{lin} \sqsubseteq \text{aff}$	(Q-LINAFF)
rel	<i>relevant</i>	$\text{rel} \sqsubseteq \text{un}$	(Q-RELUN)
aff	<i>affine</i>	$\text{aff} \sqsubseteq \text{un}$	(Q-AFFUN)
un	<i>unrestricted</i>	$q \sqsubseteq q$	(Q-REFLEX)
		$\frac{q_1 \sqsubseteq q_2 \quad q_2 \sqsubseteq q_3}{q_1 \sqsubseteq q_3}$	(Q-TRANS)

Figure 1-2: A relation between substructural type systems

### Syntax

Figure 1-3 presents the syntax of our linear language, which is an extension of the simply-typed lambda calculus. The main addition to be aware of, at this point, are the type qualifiers  $q$  that annotate the introduction forms for all data structures. The linear qualifier ( $\text{lin}$ ) indicates that the data structure in question will be *used* (i.e., appear in the appropriate elimination form) exactly once in the program. Operationally, we deallocate these linear values immediately after they are used. The unrestricted qualifier ( $\text{un}$ ) indicates that the data structure behaves as in the standard simply-typed lambda calculus. In other words, unrestricted data can be used as many times as desired and its memory resources will be automatically recycled by some extra-linguistic mechanism (a conventional garbage collector).

Apart from the qualifiers, the only slightly unusual syntactic form is the elimination form for pairs. The term  $\text{split } t_1 \text{ as } x, y \text{ in } t_2$  projects the first and second components from the pair  $t_1$  and calls them  $x$  and  $y$  in  $t_2$ . This  $\text{split}$  operation allows us to extract two components while only counting a single use of a pair. Extracting two components using the more conventional projections  $\pi_1 t_1$  and  $\pi_2 t_1$  requires two uses of the pair  $t_1$ . (It is also possible, but a bit tricky, to provide the conventional projections.)

To avoid dealing with an unnecessarily heavy syntax, we adopt a couple abbreviations in our examples in this section. First, we omit all unrestricted qualifiers and only annotate programs with the linear ones. Second, we freely use  $n$ -ary tuples (triples, quadruples, unit, etc.) in addition to pairs and also allow multi-argument functions. The latter may be defined as single-argument functions that take linear pairs (triples, etc) as arguments and immediately split them upon entry to the function body. Third, we often use ML-style type

<i>Syntax</i>			
$q ::=$	<i>qualifiers:</i>	$\text{split } t \text{ as } x, y \text{ in } t$	<i>split</i>
$\text{lin}$	<i>linear</i>	$q \lambda x:T. t$	<i>abstraction</i>
$\text{un}$	<i>unrestricted</i>	$t \ t$	<i>application</i>
$b ::=$	<i>booleans:</i>	$P ::=$	<i>pretypes:</i>
$\text{true}$	<i>true</i>	$\text{Bool}$	<i>booleans</i>
$\text{false}$	<i>false</i>	$T * T$	<i>pairs</i>
$t ::=$	<i>terms:</i>	$T \rightarrow T$	<i>functions</i>
$x$	<i>variable</i>	$T ::=$	<i>types:</i>
$q \ b$	<i>boolean</i>	$q \ P$	<i>qualified pretype</i>
$\text{if } t \ \text{then } t \ \text{else } t$	<i>conditional</i>	$\Gamma ::=$	<i>contexts:</i>
$q \langle t, t \rangle$	<i>pair</i>	$\emptyset$	<i>empty context</i>
		$\Gamma, x:T$	<i>term variable binding</i>

Figure 1-3: Linear lambda calculus: Syntax

declarations, value declarations and let expressions where convenient; they all have the obvious meanings.

### Typing

To ensure that linear objects are used exactly once, our type system maintains two important invariants.

1. Linear variables are used exactly once along every control-flow path.
2. Unrestricted data structures may not contain linear data structures. More generally, data structures with less restrictive type may not contain data structures with more restrictive type.

To understand why these invariants are useful, consider what could happen if either invariant is broken. When considering the first invariant, assume we have constructed a function `free` that uses its argument and then deallocates it. Now, if we allow a linear variable (say `x`) to appear twice, a programmer might write `<free x, free x>`, or, slightly more deviously,

$$(\lambda z. \lambda y. \langle \text{free } z, \text{free } y \rangle) \ x \ x.$$

In either case, the program ends up attempting to use and then free `x` after it has already been deallocated, causing the program to crash.

Now consider the second invariant and suppose we allow a linear data structure (call it `x`) to appear inside an unrestricted pair (`un <x, 3>`). We can

<p><i>Context Split</i></p> $\emptyset = \emptyset \circ \emptyset$ $\frac{\Gamma = \Gamma_1 \circ \Gamma_2}{\Gamma, x:\text{un } P = (\Gamma_1, x:\text{un } P) \circ (\Gamma_2, x:\text{un } P)} \quad (\text{M-UN})$	<div style="border: 1px solid black; display: inline-block; padding: 2px 5px;"><math>\Gamma = \Gamma_1 \circ \Gamma_2</math></div> (M-EMPTY)	$\frac{\Gamma = \Gamma_1 \circ \Gamma_2}{\Gamma, x:\text{lin } P = (\Gamma_1, x:\text{lin } P) \circ \Gamma_2} \quad (\text{M-LIN1})$ $\frac{\Gamma = \Gamma_1 \circ \Gamma_2}{\Gamma, x:\text{lin } P = \Gamma_1 \circ (\Gamma_2, x:\text{lin } P)} \quad (\text{M-LIN2})$
---	---	---

**Figure 1-4: Linear lambda calculus: Context splitting**

get exactly the same effect as above by using the unrestricted data structure multiple times:

```
let z = un <x,3> in
split z as x1,_ in
split z as x2,_ in
<free x1,free x2>
```

Fortunately, our type system ensures that none of these situations can occur.

We maintain the first invariant through careful context management. When type checking terms with two or more subterms, we pass all of the unrestricted variables in the context to each subterm. However, we split the linear variables between the different subterms to ensure each variable is used exactly once. Figure 1-4 defines a relation,  $\Gamma = \Gamma_1 \circ \Gamma_2$ , which describes how to split a single context in a rule conclusion ( $\Gamma$ ) into two contexts ( $\Gamma_1$  and  $\Gamma_2$ ) that will be used to type different subterms in a rule premise.

To check the second invariant, we define the predicate  $q(T)$  (and its extension to contexts  $q(\Gamma)$ ) to express the types  $T$  that can appear in a  $q$ -qualified data structure. These containment rules state that linear data structures can hold objects with linear or unrestricted type, but unrestricted data structures can only hold objects with unrestricted type.

- $q(T)$  if and only if  $T = q' P$  and  $q \sqsubseteq q'$
- $q(\Gamma)$  if and only if  $(x:T) \in \Gamma$  implies  $q(T)$

Recall, we have already defined  $q \sqsubseteq q'$  such that it is reflexive, transitive and  $\text{lin} \sqsubseteq \text{un}$ .

Now that we have defined the rules for containment and context splitting, we are ready for the typing rules proper, which appear in Figure 1-5. Keep in mind that these rules are constructed anticipating a call-by-value operational semantics.

It is often the case when designing a type system that the rules for the base cases, variables and constants, are hardly worth mentioning. However,

<i>Typing</i>	$\boxed{\Gamma \vdash \tau : T}$		
$\frac{\text{un}(\Gamma_1, \Gamma_2)}{\Gamma_1, x:T, \Gamma_2 \vdash x : T}$	(T-VAR)		$\frac{\Gamma_1 \vdash \tau_1 : T_1 \quad \Gamma_2 \vdash \tau_2 : T_2}{\Gamma_1 \circ \Gamma_2 \vdash q \langle \tau_1, \tau_2 \rangle : q(T_1 * T_2)}$ (T-PAIR)
$\frac{\text{un}(\Gamma)}{\Gamma \vdash q b : q \text{Bool}}$	(T-BOOL)		$\frac{\Gamma_1 \vdash \tau_1 : q(T_1 * T_2)}{\Gamma_2, x:T_1, y:T_2 \vdash \tau_2 : T}$ (T-SPLIT)
$\frac{\Gamma_1 \vdash \tau_1 : q \text{Bool} \quad \Gamma_2 \vdash \tau_2 : T \quad \Gamma_2 \vdash \tau_3 : T}{\Gamma_1 \circ \Gamma_2 \vdash \text{if } \tau_1 \text{ then } \tau_2 \text{ else } \tau_3 : T}$	(T-IF)		$\frac{q(\Gamma) \quad \Gamma, x:T_1 \vdash \tau_2 : T_2}{\Gamma \vdash q \lambda x:T_1. \tau_2 : q T_1 \rightarrow T_2}$ (T-ABS)
			$\frac{\Gamma_1 \vdash \tau_1 : q T_{11} \rightarrow T_{12} \quad \Gamma_2 \vdash \tau_2 : T_{11}}{\Gamma_1 \circ \Gamma_2 \vdash \tau_1 \tau_2 : T_{12}}$ (T-APP)

**Figure 1-5: Linear lambda calculus: Typing**

in substructural type systems these cases have a special role in defining the nature of the type system, and subtle changes can make all the difference. In our linear system, the base cases must ensure that no linear variable is discarded without being used. To enforce this invariant in rule (T-VAR), we explicitly check that  $\Gamma_1$  and  $\Gamma_2$  contain no linear variables using the condition  $\text{un}(\Gamma_1, \Gamma_2)$ . We make a similar check in rule (T-BOOL). Notice also that rule (T-VAR) is written carefully to allow the variable  $x$  to appear anywhere in the context, rather than just at the beginning or at the end.

1.2.1 EXERCISE [★]: What is the effect of rewriting the variable rule as follows?

$$\frac{\text{un}(\Gamma)}{\Gamma, x:T \vdash x : T} \quad (\text{T-BROKENVAR})$$

The inductive cases of the typing relation take care to use context splitting to partition linear variables between various subterms. For instance, rule (T-IF) splits the incoming context into two parts, one of which is used to check subterm  $\tau_1$  and the other which is used to check both  $\tau_2$  and  $\tau_3$ . As a result, a particular linear variable will occur once in  $\tau_2$  and once in  $\tau_3$ . However, the linear object bound to the variable in question will be used (and hence deallocated) exactly once at run time since only one of  $\tau_2$  or  $\tau_3$  will be executed.

The rules for creation of pairs and functions make use of the containment rules. In each case, the data structure's qualifier  $q$  is used in the premise of the typing rule to limit the sorts of objects it may contain. For example, in the rule (T-ABS), if the qualifier  $q$  is  $\text{un}$  then the variables in  $\Gamma$ , which will inhabit the function closure, must satisfy  $\text{un}(\Gamma)$ . In other words, they must all have



unrestricted type. If we omitted this constraint, we could write the following badly behaved functions. (For clarity, we have retained the unrestricted qualifiers in this example rather than omitting them.)

```

type T = un (un bool → lin bool)

val discard =
  lin λx:lin bool.
    (lin λf:T.lin true) (un λy:un bool.x)

val duplicate =
  lin λx:lin bool.
    (lin λf:T.lin <f (un true),f (un true)>)) (un λy:un bool.x)

```

The first function discards a linear argument  $x$  without using it and the second duplicates a linear argument and returns two copies of it in a pair. Hence, in the first case, we fail to deallocate  $x$  and in the second case, a subsequent function may project both elements of the pair and use  $x$  twice, which would result in a memory error as  $x$  would be deallocated immediately after the first use. Fortunately, the containment constraint disallows the linear variable  $x$  from appearing in the unrestricted function ( $\lambda y:\text{bool}. x$ ).

Now that we have defined our type system, we should verify our intended structural properties: exchange for all variables, and weakening and contraction for unrestricted variables.

- 1.2.2 LEMMA [EXCHANGE]: If  $\Gamma_1, x_1:T_1, x_2:T_2, \Gamma_2 \vdash t : T$  then  $\Gamma_1, x_2:T_2, x_1:T_1, \Gamma_2 \vdash t : T$ . □
- 1.2.3 LEMMA [UNRESTRICTED WEAKENING]: If  $\Gamma \vdash t : T$  then  $\Gamma, x_1:\text{un } P_1 \vdash t : T$ . □
- 1.2.4 LEMMA [UNRESTRICTED CONTRACTION]:  
 If  $\Gamma, x_2:\text{un } P_1, x_3:\text{un } P_1 \vdash t : T_3$  then  $\Gamma, x_1:\text{un } P_1 \vdash [x_2 \mapsto x_1][x_3 \mapsto x_1]t : T_3$ . □

*Proof:* The proofs of all three lemmas follow by induction on the structure of the appropriate typing derivation. □

### Algorithmic Linear Type Checking

The inference rules provided in the previous subsection give a clear, concise specification of the linearly-typed programs. However, these rules are also highly non-deterministic and cannot be implemented directly. The primary difficulty is that to implement the non-deterministic splitting operation,

<i>Algorithmic Typing</i>	$\boxed{\Gamma_{in} \vdash \tau : T; \Gamma_{out}}$	
$\Gamma_1, x : \text{un } P, \Gamma_2 \vdash x : \text{un } P; \Gamma_1, x : \text{un } P, \Gamma_2$	(A-UVAR)	
$\Gamma_1, x : \text{lin } P, \Gamma_2 \vdash x : \text{lin } P; \Gamma_1, \Gamma_2$	(A-LVAR)	
$\Gamma \vdash q \text{ b} : q \text{ Bool}; \Gamma$	(A-BOOL)	
$\Gamma_1 \vdash t_1 : q \text{ Bool}; \Gamma_2$		
$\frac{\Gamma_2 \vdash t_2 : T; \Gamma_3 \quad \Gamma_2 \vdash t_3 : T; \Gamma_3}{\Gamma_1 \vdash \text{if } t_1 \text{ then } t_2 \text{ else } t_3 : T; \Gamma_3}$	(A-IF)	
$\frac{\Gamma_1 \vdash t_1 : T_1; \Gamma_2 \quad \Gamma_2 \vdash t_2 : T_2; \Gamma_3}{\Gamma_1 \vdash q \langle t_1, t_2 \rangle : q (T_1 * T_2); \Gamma_3}$	(A-PAIR)	
		$\frac{\Gamma_1 \vdash t_1 : q (T_1 * T_2); \Gamma_2}{\Gamma_1 \vdash \text{split } t_1 \text{ as } x, y \text{ in } t_2 : T; \Gamma_3 \div (x : T_1, y : T_2)}$
		(A-SPLIT)
		$\frac{\Gamma_1, x : T_1 \Rightarrow \Gamma_1 = \Gamma_2 \div (x : T_1) \quad \Gamma_1, x : T_1 \vdash t_2 : T_2; \Gamma_2}{\Gamma_1 \vdash q \lambda x : T_1. t_2 : q T_1 \rightarrow T_2; \Gamma_2 \div (x : T_1)}$
		(A-ABS)
		$\frac{\Gamma_1 \vdash t_1 : q T_{11} \rightarrow T_{12}; \Gamma_2 \quad \Gamma_2 \vdash t_2 : T_{11}; \Gamma_3}{\Gamma_1 \vdash t_1 t_2 : T_{12}; \Gamma_3}$
		(A-APP)

**Figure 1-6: Linear lambda calculus: Algorithmic type checking**

$\Gamma = \Gamma_1 \circ \Gamma_2$ , we must guess how to split an input context  $\Gamma$  into two parts. Fortunately, it is relatively straightforward to restructure the type checking rules to avoid having to make these guesses. This restructuring leads directly to a practical type checking algorithm.

The central idea is that rather than splitting the context into parts before checking a complex expression composed of several subexpressions, we can pass the entire context as an input to the first subexpression and have it return the unused portion as an output. This output may then be used to check the next subexpression, which may also return some unused portions of the context as an output, and so on. Figure 1-6 makes these ideas concrete. It defines a new algorithmic type checking judgment with the form  $\Gamma_{in} \vdash \tau : T; \Gamma_{out}$ , where  $\Gamma_{in}$  is the input context, some portion of which will be consumed during type checking of  $\tau$ , and  $\Gamma_{out}$  is the output context, which will be synthesized alongside the type  $T$ .

There are several key changes in our reformulated system. First, the base cases for variables and constants allow any context to pass through the judgment rather than restricting the number of linear variables that appear. In order to ensure that linear variables are used, we move these checks to the rules where variables are introduced. For instance, consider the rule (A-SPLIT). The second premise has the form

$$\Gamma_2, x : T_1, y : T_2 \vdash t_2 : T; \Gamma_3$$

If  $T_1$  and  $T_2$  are linear, then they should be used in  $t_2$  and should not appear in  $\Gamma_3$ . Conversely,  $T_1$  and  $T_2$  are unrestricted, then they will always appear

in  $\Gamma_3$ , but we should delete them from the final outgoing context of the rule so that the ordinary scoping rules for the variables are enforced. To handle both the check that linear variables do not appear and the removal of unrestricted variables, we use a special “context difference” operator ( $\div$ ). Using this operator, the final outgoing context of the rule (A-SPLIT) is defined to be  $\Gamma_3 \div (\mathbf{x}:\mathsf{T}_1, \mathbf{y}:\mathsf{T}_2)$ . Formally, context difference is defined as follows.

$$\begin{aligned} \Gamma \div \emptyset &= \Gamma \\ \frac{\Gamma_1 \div \Gamma_2 = \Gamma_3 \quad (\mathbf{x}:\mathsf{lin} \mathsf{P}) \notin \Gamma_3}{\Gamma_1 \div (\Gamma_2, \mathbf{x}:\mathsf{lin} \mathsf{P}) = \Gamma_3} \\ \frac{\Gamma_1 \div \Gamma_2 = \Gamma_3 \quad \Gamma_3 = \Gamma_4, \mathbf{x}:\mathsf{un} \mathsf{P}, \Gamma_5}{\Gamma_1 \div (\Gamma_2, \mathbf{x}:\mathsf{un} \mathsf{P}) = \Gamma_4, \Gamma_5} \end{aligned}$$

Notice that this operator is undefined when we attempt to take the difference of two contexts,  $\Gamma_1$  and  $\Gamma_2$ , that contain bindings for the same linear variable ( $\mathbf{x}:\mathsf{lin} \mathsf{P}$ ). If the undefined quotient  $\Gamma_1 \div \Gamma_2$  were to appear anywhere in a typing rule, the rule itself would not be considered defined and could not be part of a valid typing derivation.

The rule for abstraction (A-ABS) also introduces a variable and hence it also uses context difference to manipulate the output context for the rule. Abstractions must also satisfy the appropriate containment conditions. In other words, rule (A-ABS) must check that unrestricted functions do not contain linear variables. We perform this last check by verifying that when the function qualifier is unrestricted, the input and output contexts from checking the function body are the same. This equivalence check is sufficient because if a linear variable was used in the body of an unrestricted function (and hence captured in the function closure), that linear variable would not show up in the outgoing context.

It is completely straightforward to check that every rule in our algorithmic system is syntax directed and that all our auxiliary functions including context membership tests and context difference are easily computable. Hence, we need only show that our algorithmic system is equivalent to the simpler and more elegant declarative system specified in the previous section. The proof of equivalence can be broken down into the two standard components: *soundness* and *completeness* of the algorithmic system with respect to the declarative system. However, before we can get to the main results, we will need to show that our algorithmic system satisfies some basic structural properties of its own. In the following lemmas, we use the notation  $\mathcal{L}(\Gamma)$  and  $\mathcal{U}(\Gamma)$  to refer to the list of linear and unrestricted assumptions in  $\Gamma$  respectively.

- 1.2.5 LEMMA [ALGORITHMIC MONOTONICITY]: If  $\Gamma \vdash \mathbf{t} : \mathbb{T}; \Gamma'$  then  $\mathcal{U}(\Gamma') = \mathcal{U}(\Gamma)$  and  $\mathcal{L}(\Gamma') \subseteq \mathcal{L}(\Gamma)$ .  $\square$
- 1.2.6 LEMMA [ALGORITHMIC EXCHANGE]: If  $\Gamma_1, \mathbf{x}_1 : \mathbb{T}_1, \mathbf{x}_2 : \mathbb{T}_2, \Gamma_2 \vdash \mathbf{t} : \mathbb{T}; \Gamma_3$  then  $\Gamma_1, \mathbf{x}_2 : \mathbb{T}_2, \mathbf{x}_1 : \mathbb{T}_1, \Gamma_2 \vdash \mathbf{t} : \mathbb{T}; \Gamma'_3$  and  $\Gamma_3$  is the same as  $\Gamma'_3$  up to transposition of the bindings for  $\mathbf{x}_1$  and  $\mathbf{x}_2$ .  $\square$
- 1.2.7 LEMMA [ALGORITHMIC WEAKENING]: If  $\Gamma \vdash \mathbf{t} : \mathbb{T}; \Gamma'$  then  $\Gamma, \mathbf{x} : \mathbb{T}' \vdash \mathbf{t} : \mathbb{T}; \Gamma', \mathbf{x} : \mathbb{T}'$ .  $\square$
- 1.2.8 LEMMA [ALGORITHMIC LINEAR STRENGTHENING]: If  $\Gamma, \mathbf{x} : \text{lin } \mathbb{P} \vdash \mathbf{t} : \mathbb{T}; \Gamma'$  then  $\Gamma \vdash \mathbf{t} : \mathbb{T}; \Gamma'$ .  $\square$

Each of these lemmas may be proven directly by induction on the initial typing derivation. The algorithmic system also satisfies a contraction lemma, but since it will not be necessary in the proofs of soundness and completeness, we have not stated it here.

- 1.2.9 THEOREM [ALGORITHMIC SOUNDNESS]: If  $\Gamma_1 \vdash \mathbf{t} : \mathbb{T}; \Gamma_2$  and  $\mathcal{L}(\Gamma_2) = \emptyset$  then  $\Gamma_1 \vdash \mathbf{t} : \mathbb{T}$ .  $\square$

*Proof:* As usual, the proof is by induction on the typing derivation. The structural lemmas we have just proven are required to push through the result, but it is mostly straightforward.  $\square$

- 1.2.10 THEOREM [ALGORITHMIC COMPLETENESS]: If  $\Gamma_1 \vdash \mathbf{t} : \mathbb{T}$  then  $\Gamma_1 \vdash \mathbf{t} : \mathbb{T}; \Gamma_2$  and  $\mathcal{L}(\Gamma_2) = \emptyset$ .  $\square$

*Proof:* The proof is by induction on the typing derivation.  $\square$

## Operational Semantics

To make the memory management properties of our language clear, we will evaluate terms in an abstract machine with an explicit store. As indicated in Figure 1-7, stores are a sequence of variable-value pairs. We will implicitly assume that any variable appears at most once on the left-hand side of a pair so the sequence may be treated as a finite partial map.

A value is a pair of a qualifier together with some data (a *prevalue*  $w$ ). For the sake of symmetry, we will also assume that all values are stored, even base types such as booleans. As a result, both components of any pair will be pointers (variables).

We define the operation of our abstract machine using a context-based, small-step semantics. Figure 1-7 defines the computational contexts  $E$ , which

$w ::=$	<i>prevalues:</i>	$E ::=$	<i>evaluation contexts:</i>
$b$	<i>boolean</i>	$[\ ]$	<i>context hole</i>
$\langle x, y \rangle$	<i>pair</i>	$\text{if } E \text{ then } t \text{ else } t$	<i>if context</i>
$\lambda x:T. t$	<i>abstraction</i>	$q \langle E, t \rangle$	<i>fst context</i>
$v ::=$	<i>values:</i>	$q \langle x, E \rangle$	<i>snd context</i>
$q w$	<i>qualified prevalue</i>	$\text{split } E \text{ as } x, y \text{ in } t$	<i>split context</i>
$S ::=$	<i>stores:</i>	$E t$	<i>fun context</i>
$\emptyset$	<i>empty context</i>	$x E$	<i>arg context</i>
$S, x \mapsto v$	<i>store binding</i>		

Figure 1-7: Linear lambda calculus: Run-time data

are terms with a single hole. Contexts define the order of evaluation of terms—they specify the places in a term where a computation can occur. In our case, evaluation is left-to-right since, for example, there is a context with the form  $E t$  indicating that we can reduce the term in the function position before reducing the term in the argument position. However, there is no context with the form  $t E$ . Instead, there is only the more limited context  $x E$ , indicating that we must reduce the term in the function position to a pointer  $x$  before proceeding to evaluate the term in the argument position. We use the notation  $E[t]$  to denote the term composed of the context  $E$  with its hole plugged by the computation  $t$ .

The operational semantics, defined in Figure 1-8, is factored into two relations. The first relation,  $(S; t) \rightarrow (S'; t')$ , picks out a subcomputation to evaluate. The second relation,  $(S; t) \rightarrow_{\beta} (S'; t')$ , does all the real work. In order to avoid creation of two sets of operational rules, one for linear data, which is deallocated when used, and one for unrestricted data, which is never deallocated, we define an auxiliary function,  $S \stackrel{q}{\sim} x$ , to manage the differences.

$$\begin{aligned} (S_1, x \mapsto v, S_2) \stackrel{1}{\sim} x &= S_1, S_2 \\ S \stackrel{un}{\sim} x &= S \end{aligned}$$

Aside from these details, the operational semantics is standard.

### Preservation and Progress

In order to prove the standard safety properties for our language, we need to be able to show that programs are well-formed after each step in evaluation. Hence, we will define typing rules for our abstract machine. Since these typing rules are only necessary for the proof of soundness, and have no place in an

<p><i>Top-level Evaluation</i>      <math>\boxed{(S; t) \rightarrow (S'; t')}</math></p> $\frac{(S; t) \rightarrow_{\beta} (S; t')}{(S; E[t]) \rightarrow (S; E[t'])} \quad (\text{E-CTXT})$ <p><i>Evaluation</i>      <math>\boxed{(S; t) \rightarrow_{\beta} (S'; t')}</math></p> $(S; q \ b) \rightarrow_{\beta} (S, x \mapsto q \ b; x) \quad (\text{E-BOOL})$ $\frac{S(x) = q \ \text{true}}{(S; \text{if } x \text{ then } t_1 \text{ else } t_2) \rightarrow_{\beta} (S \stackrel{q}{\sim} x; t_1)} \quad (\text{E-IF1})$	$\frac{S(x) = q \ \text{false}}{(S; \text{if } x \text{ then } t_1 \text{ else } t_2) \rightarrow_{\beta} (S \stackrel{q}{\sim} x; t_2)} \quad (\text{E-IF2})$ $(S; q \ \langle y, z \rangle) \rightarrow_{\beta} (S, x \mapsto q \ \langle y, z \rangle; x) \quad (\text{E-PAIR})$ $\frac{S(x) = q \ \langle y_1, z_1 \rangle}{(S; \text{split } x \text{ as } y, z \text{ in } t) \rightarrow_{\beta} (S \stackrel{q}{\sim} x; [y \mapsto y_1][z \mapsto z_1]t)} \quad (\text{E-SPLIT})$ $(S; q \ \lambda y:T. t) \rightarrow_{\beta} (S, x \mapsto q \ \lambda y:T. t; x) \quad (\text{E-FUN})$ $\frac{S(x_1) = q \ \lambda y:T. t}{(S; x_1 \ x_2) \rightarrow_{\beta} (S \stackrel{q}{\sim} x_1; [y \mapsto x_2]t)} \quad (\text{E-APP})$
--	---

Figure 1-8: Linear lambda calculus: Operational semantics

implementation, we will extend the declarative typing rules rather than the algorithmic typing rules.

Figure 1-9 presents the machine typing rules in terms of two judgments, one for stores and the other for programs. The store typing rules generate a context that describes the available bindings in the store. The program typing rule uses the generated bindings to check the expression that will be executed.

With this new machinery in hand, we are able to prove the standard progress and preservation theorems.

- 1.2.11 THEOREM [PRESERVATION]: If  $\vdash (S; t)$  and  $(S; t) \rightarrow (S'; t')$  then  $\vdash (S'; t')$ .  $\square$
- 1.2.12 THEOREM [PROGRESS]: If  $\vdash (S; t)$  then  $(S; t) \rightarrow (S'; t')$  or  $t$  is a value.  $\square$
- 1.2.13 EXERCISE [RECOMMENDED, ★]: You will need a substitution lemma to complete the proof of preservation. Is the following the right one?  
 Conjecture: Let  $\Gamma_3 = \Gamma_1 \circ \Gamma_2$ . If  $\Gamma_1, x:T \vdash t_1 : T_1$  and  $\Gamma_2 \vdash t : T$  then  $\Gamma_3 \vdash [x \mapsto t]t_1 : T_1$ .  $\square$
- 1.2.14 EXERCISE [★★★, →]: Prove progress and preservation using *TAPL*, Chapters 9 and 13, as an approximate guide.  $\square$

<p><i>Store Typing</i></p> $\frac{}{\vdash \emptyset : \emptyset} \quad \boxed{\vdash S : \Gamma} \quad \text{(T-EMPTYYS)}$ $\frac{\vdash S : \Gamma_1 \circ \Gamma_2 \quad \Gamma_1 \vdash \text{lin } w : T}{\vdash S, x \mapsto \text{lin } w : \Gamma_2, x:T} \quad \text{(T-NEXTLINS)}$	<p><i>Program Typing</i></p> $\frac{\vdash S : \Gamma_1 \circ \Gamma_2 \quad \Gamma_1 \vdash \text{un } w : T}{\vdash S, x \mapsto \text{un } w : \Gamma_2, x:T} \quad \text{(T-NEXTUNS)}$ $\boxed{\vdash (S; t)}$ $\frac{\vdash S : \Gamma \quad \Gamma \vdash t : T}{\vdash (S; t)} \quad \text{(T-PROG)}$
--	--

Figure 1-9: Linear lambda calculus: Program typing

### 1.3 Extensions and Variations

Most features found in modern programming languages can be defined to interoperate successfully with linear type systems, although some are trickier than others. In this section, we will consider a variety of practical extensions to our simple linear lambda calculus.

#### Sums and Recursive Types

Complex data structures, such as the recursive data types found in ML-like languages, pose little problem for linear languages. To demonstrate the central ideas involved, we extend the syntax for the linear lambda calculus with the standard introduction and elimination forms for sums and recursive types. The details are presented in Figure 1-10.

Values with sum type are introduced by injections  $q \text{ inl}_P t$  or  $q \text{ inr}_P t$ , where  $P$  is  $T_1 + T_2$ , the resulting pretype of the term. In the first instance, the underlying term  $t$  must have type  $T_1$ , and in the second instance, the underlying term  $t$  must have type  $T_2$ . The qualifier  $q$  indicates the linearity of the argument in exactly the same way as for pairs. The case expression will execute its first branch if its primary argument is a left injection and its second branch if its primary argument is a right injection. We assume that  $+$  binds more tightly than  $\rightarrow$  but less tightly than  $*$ .

Recursive types are introduced with a  $\text{roll}_P t$  expression, where  $P$  is the recursive pretype the expression will assume. Unlike all the other introduction forms, roll expressions are not annotated with a qualifier. Instead, they take on the qualifier of the underlying expression  $t$ . The reason for this distinction is that we will treat this introduction form as a typing coercion that has no real operational effect. Unlike functions, pairs or sums, recursive data types have no data of their own and therefore do not need a separate qualifier to control their allocation behavior. To simplify the notational overhead

$t ::=$	<i>terms:</i>	<i>Typing</i>	$\boxed{\Gamma \vdash t : T}$
...	<i>as before</i>		
$q \text{ inl}_P t$	<i>left inj.</i>	$\frac{\Gamma \vdash t : T_1 \quad q(T_1) \quad q(T_2)}{\Gamma \vdash q \text{ inl}_{T_1+T_2} t : q(T_1+T_2)}$	(T-INL)
$q \text{ inr}_P t$	<i>right inj.</i>		
$\text{case } t \text{ (inl } x \Rightarrow t \mid \text{inr } y \Rightarrow t)$	<i>case</i>	$\frac{\Gamma \vdash t : T_2 \quad q(T_1) \quad q(T_2)}{\Gamma \vdash q \text{ inr}_{T_1+T_2} t : q(T_1+T_2)}$	(T-INR)
$\text{roll}_P t$	<i>roll into rec type</i>		
$\text{unroll } t$	<i>unroll from rec type</i>		
$\text{fun } f(x:T_1):T_2.t$	<i>recursive fun</i>		
$P ::=$	<i>pretypes:</i>		
...	<i>as before</i>		
$a$	<i>pretype variables</i>		
$T_1+T_2$	<i>sum types</i>	$\frac{\Gamma_1 \circ \Gamma_2 \vdash \text{case } t \text{ (inl } x \Rightarrow t_1 \mid \text{inr } y \Rightarrow t_2) : T}{\Gamma \vdash t : [a \rightarrow P]q P_1 \quad P = \text{rec } a.q P_1}$	(T-CASE)
$\text{rec } a.T$	<i>recursive types</i>	$\frac{\Gamma \vdash t : [a \rightarrow P]q P_1 \quad P = \text{rec } a.q P_1}{\Gamma \vdash \text{roll}_P t : q P}$	(T-ROLL)
		$\frac{\Gamma \vdash t : P \quad P = \text{rec } a.q P_1}{\Gamma \vdash \text{unroll } t : [a \rightarrow P]q P_1}$	(T-UNROLL)
		$\frac{\text{un } (\Gamma) \quad \Gamma, f:\text{un } T_1 \rightarrow T_2, x:T_1 \vdash t : T_2}{\Gamma \vdash \text{fun } f(x:T_1):T_2.t : \text{un } T_1 \rightarrow T_2}$	(T-TFUN)

Figure 1-10: Linear lambda calculus: Sums and recursive types

of sums and recursive types, we will normally omit the typing annotations on their introduction forms in our examples.

In order to write computations that process recursive types, we add recursive function declarations to our language as well. Since the free variables in a recursive function closure will be used on each recursive invocation of the function, we cannot allow the closure to contain linear variables. Hence, all recursive functions are unrestricted data structures.

A simple but useful data structure is the linear list of Ts:

```
type T llist = rec a.lin (unit + lin (T * lin a))
```

Here, the entire spine (aside from the terminating value of unit type) is linear while the underlying T objects may be linear or unrestricted. To create a fully unrestricted list, we simply omit the linear qualifiers on the sum and pairs that make up the spine of the list:

```
type T list = rec a.unit + T * a
```



After defining the linear lists, the memory conscious programmer can write many familiar list-processing functions in a minimal amount of space. For example, here is how we map an unrestricted function across a linear list. Remember, multi-argument functions are abbreviations for functions that accept linear pairs as arguments.

```

fun nil(_:unit) : T2 llist =
  roll (lin inl ())

fun cons(hd:T2, tl:T2 llist) : T2 llist =
  roll (lin inr (lin <hd,tl>))

fun map(f:T1→T2, xs:T1 llist) : T2 llist =
  case unroll xs (
    inl _ ⇒ nil()
  | inr xs ⇒
    split xs as hd,tl in
    cons(f hd,map lin <f,tl>))

```

In this implementation of `map`, we can observe that on each iteration of the loop, it is possible to reuse the space deallocated by `split` or `case` operations for the allocation operations that follow in the body of the function (inside the calls to `nil` and `cons`).

Hence, at first glance, it appears that `map` will execute with only a constant space overhead. Unfortunately, however, there are some hidden costs as `map` executes. A typical implementation will store local variables and temporaries on the stack before making a recursive call. In this case, the result of `f hd` will be stored on the stack while `map` iterates down the list. Consequently, rather than having a constant space overhead, our `map` implementation will have an  $O(n)$  overhead, where  $n$  is the length of the list. This is not too bad, but we can do better.

In order to do better, we need to avoid implicit stack allocation of data each time we iterate through the body of a recursive function. Fortunately, many functional programming languages guarantee that if the last operation in a function is itself a function call then the language implementation will deallocate the current stack frame before calling the new function. We name such function calls *tail calls* and we say that any language implementation that guarantees that the current stack frame will be deallocated before a tail call is *tail-call optimizing*.

Assuming that our language is tail-call optimizing, we can now rewrite `map` so that it executes with only a constant space overhead. The main trick involved is that we will explicitly keep track of both the part of the input list we have yet to process and the output list that we have already processed. The

output list will wind up in reverse order, so we will reverse it at the end. Both of the loops in the code, `mapRev` and `reverse` are *tail-recursive* functions. That is, they end in a tail call and have a space-efficient implementation.

```

fun map(f:T1→T2, input:T1 llist) : T2 llist =
  reverse(mapRev(f,input,nil()),nil())

and mapRev(f:T1→T2,
           input:T1 llist,
           output:T2 llist) : T2 llist =
  case unroll input (
    inl _ ⇒ output
  | inr xs ⇒
    split xs as hd,tl in
    mapRev (f,tl,cons(f hd,output)))

and reverse(input:T2 llist, output:T2 llist)
  case unroll input (
    inl _ ⇒ output
  | inr xs ⇒
    split xs as hd,tl in
    reverse(tl,cons(hd,output)))

```

This *link reversal* algorithm is a well-known way of traversing a list in constant space. It is just one of a class of algorithms developed well before the invention of linear types. A similar algorithm was invented by Deutsch, Schorr, and Waite for traversing trees and graphs in constant space. Such constant space traversals are essential parts of mark-sweep garbage collectors—at garbage collection time there is no extra space for a stack so any traversal of the heap must be done in constant space.

- 1.3.1 EXERCISE [\*\*\*]: Define a recursive type that describes linear binary trees that hold data of type  $T$  in their internal nodes (nothing at the leaves). Write a constant-space function `treeMap` that produces an identically-shaped tree on output as it was given on input, modulo the action of the function  $f$  that is applied to each element of the tree. Feel free to use reasonable extensions to our linear lambda calculus including mutually recursive functions,  $n$ -ary tuples and  $n$ -ary sums. □

## Polymorphism

Parametric polymorphism is a crucial feature of almost any functional language, and our linear lambda calculus is no exception. The main function of polymorphism in our setting is to support two different sorts of code reuse.

1. Reuse of code to perform the same algorithm, but on data with different shapes.
2. Reuse of code to perform the same algorithm, but on data governed by different memory management strategies.

To support the first kind of polymorphism, we will allow quantification over pretypes. To support the second kind of polymorphism, we will allow quantification over qualifiers. A good example of both sorts of polymorphism arises in the definition of a polymorphic map function. In the code below, we use  $a$  and  $b$  to range over pretype variables as we did in the previous section, and  $p$  to range over qualifier variables.

```

type (p1,p2,a) list =
  rec a.p1 (unit + p1 (p2 a * (p1,p2,a) list))

map :
  ∀a,b.
  ∀pa,pb.
  lin ((pa a → pb b)*(lin,pa,a) list)→(lin,pb,b) list

```

The type definition in the first line defines lists in terms of three parameters. The first parameter,  $p_1$ , gives the usage pattern (linear or unrestricted) for the spine of the list, while the second parameter gives the usage pattern for the elements of the list. The third parameter is a pretype parameter, which gives the (pre)type of the elements of list. The map function is polymorphic in the argument ( $a$ ) and result ( $b$ ) element types of the list. It is also polymorphic (via parameters  $p_a$  and  $p_b$ ) in the way those elements are used. Overall, the function maps lists with linear spines to lists with linear spines.

Developing a system for polymorphic, linear type inference is a challenging research topic, beyond the scope of this book, so we will assume that, unlike in ML, polymorphic functions are introduced explicitly using the syntax  $\Lambda a. t$  or  $\Lambda p. t$ . Here,  $a$  and  $p$  are the type parameters to a function with body  $t$ . The body does not need to be a value, like in ML, since we will run the polymorphic function every time a pretype or qualifier is passed to the function as an argument. The syntax  $t' [P]$  or  $t' [q]$  applies the function  $t'$  to its pretype or qualifier argument. Figure 1-11 summarizes the syntactic extensions to the language.

Before we get to writing the map function, we will take a look at the polymorphic constructor functions for linear lists. These functions will take a pretype parameter and two qualifier parameters, just like the type definition for lists.

$q ::=$	<i>qualifiers:</i>	$q \Delta p.t$	<i>qualifier abstraction</i>
...	<i>as before</i>	$t [q]$	<i>qualifier application</i>
$p$	<i>polymorphic qualifier</i>	$P ::=$	<i>pretypes:</i>
$t ::=$	<i>terms:</i>	...	<i>as before</i>
...	<i>as before</i>	$\forall a.T$	<i>pretype polymorphism</i>
$q \Delta a.t$	<i>pretype abstraction</i>	$\forall p.T$	<i>qualifier polymorphism</i>
$t [P]$	<i>pretype application</i>		

**Figure 1-11: Linear lambda calculus: Polymorphism syntax**

```

val nil :  $\forall a, p_2. (\text{lin}, p_2, a) \text{ list} =$ 
   $\Lambda a, p_2. \text{roll } (\text{lin } \text{inl } ())$ 

val list :
   $\forall a, p_2. \text{lin } (p_2 \ a \ * \ (\text{lin}, p_2, a) \ \text{list}) \rightarrow (\text{lin}, p_2, a) \ \text{list} =$ 
   $\Lambda a, p_2.$ 
   $\lambda \text{cell} : \text{lin } (p_2 \ a \ * \ (\text{lin}, p_2, a) \ \text{list}).$ 
   $\text{roll } (\text{lin } \text{inr } (\text{lin } \text{cell}))$ 

```

Now our most polymorphic map function may be written as follows.

```

val map =
   $\Lambda a, b. \Lambda p_a, p_b.$ 
  fun aux(f:( $p_a \ a \ \rightarrow \ p_b \ b$ ),
    xs:( $\text{lin}, p_a, a$ ) list)) : ( $\text{lin}, p_b, b$ ) list =
  case unroll xs (
    inl _  $\Rightarrow$  nil [ $b, p_b$ ] ()
  | inr xs  $\Rightarrow$  split xs as hd,tl in
    cons [ $b, p_b$ ] ( $p_b \ <f \ \text{hd}, \text{map } (\text{lin } <f, \text{tl}>)>)$ )

```

In order to ensure that our type system remains sound in the presence of pretype polymorphism, we add the obvious typing rules, but change very little else. However, adding qualifier polymorphism, as we have done, is a little more involved. Before arriving at the typing rules themselves, we need to adapt some of our basic definitions to account for abstract qualifiers that may either be linear or unrestricted.

First, we need to ensure that we propagate contexts containing abstract qualifiers safely through the other typing rules in the system. Most importantly, we add additional cases to the context manipulation rules defined in the previous section. We need to ensure that linear hypotheses are not duplicated and therefore we cannot risk duplicating unknown qualifiers, which might turn out to be linear. Figure 1-12 specifies the details.

<p><i>Context Split</i></p> $\frac{\Gamma = \Gamma_1 \circ \Gamma_2}{\Gamma, x : p \ P = (\Gamma_1, x : p \ P) \circ \Gamma_2} \quad (\text{M-ABS1})$	<div style="border: 1px solid black; padding: 2px; display: inline-block;"><math>\Gamma = \Gamma_1 \circ \Gamma_2</math></div>	$\frac{\Gamma = \Gamma_1 \circ \Gamma_2}{\Gamma, x : p \ P = \Gamma_1 \circ (\Gamma_2, x : p \ P)} \quad (\text{M-ABS2})$
---	--	---

Figure 1-12: Linear context manipulation rules

<p><math>\Delta ::=</math></p> <ul style="list-style-type: none"> <li><math>\emptyset</math></li> <li><math>\Delta, a</math></li> <li><math>\Delta, p</math></li> </ul> <p><i>Typing</i></p> $\frac{q(\Gamma) \quad \Delta, a; \Gamma \vdash t : T}{\Delta; \Gamma \vdash q \ \Lambda a. t : q \ \forall a. T} \quad (\text{T-PABS})$	<p><i>type contexts:</i></p> <p><i>empty</i></p> <p><i>pretype var.</i></p> <p><i>qualifier var.</i></p> <div style="border: 1px solid black; padding: 2px; display: inline-block;"><math>\Delta; \Gamma \vdash t : T</math></div>	$\frac{\Delta; \Gamma \vdash t : q \ \forall a. T \quad FV(P) \subseteq \Delta}{\Delta; \Gamma \vdash t \ [P] : [a \mapsto P]T} \quad (\text{T-PAPP})$ $\frac{q(\Gamma) \quad \Delta, p; \Gamma \vdash t : T}{\Delta; \Gamma \vdash q \ \Lambda p. t : q \ \forall p. T} \quad (\text{T-QABS})$ $\frac{\Delta; \Gamma \vdash t : q_1 \ \forall p. T \quad FV(q) \subseteq \Delta}{\Delta; \Gamma \vdash t \ [q] : [p \mapsto q]T} \quad (\text{T-QAPP})$
---	--	--

Figure 1-13: Linear lambda calculus: Polymorphic typing

Second, we need to conservatively extend the relation on type qualifiers  $q_1 \sqsubseteq q_2$  so that it is sound in the presence of qualifier polymorphism. Since the linear qualifier is the least qualifier in the current system, the following rule should hold.

$$l \text{ in } \sqsubseteq p \quad (\text{Q-LINP})$$

Likewise, since  $un$  is the greatest qualifier in the system, we can be sure the following rule is sound.

$$p \sqsubseteq un \quad (\text{Q-PUN})$$

Aside from these rules, we will only be able to infer that an abstract qualifier  $p$  is related to itself via the general reflexivity rule. Consequently, linear data structures can contain abstract ones; abstract data structures can contain unrestricted data structures; and data structure with qualifier  $p$  can contain other data with qualifier  $p$ .

In order to define the typing rules for the polymorphic linear lambda calculus proper, we need to change the judgment form to keep track of the type variables that are allowed to appear free in a term. The new judgment uses the type context  $\Delta$  for this purpose. The typing rules for the introduction and elimination forms for each sort of polymorphism are fairly straightforward now and are presented in Figure 1-13.

The typing rules for the other constructs we have seen are almost unchanged. One relatively minor alteration is that the incoming type context  $\Delta$  will be propagated through the rules to account for the free type variables. Unlike term variables, type variables can always be used in an unrestricted fashion; it is difficult to understand what it would mean to restrict the use of a type variable to one place in a type or term. Consequently, all parts of  $\Delta$  are propagated from the conclusion of any rule to all premises. We also need the occasional side condition to check that whenever a programmer writes down a type, its free variables are contained in the current type context  $\Delta$ . For instance the rules for function abstraction and application will now be written as follows.

$$\frac{q(\Gamma) \quad FV(T_1) \subseteq \Delta \quad \Delta; \Gamma, x:T_1 \vdash t_2 : T_2}{\Delta; \Gamma \vdash q \lambda x:T_1. t_2 : q T_1 \rightarrow T_2} \quad (\text{T-ABS})$$

$$\frac{\Delta; \Gamma_1 \vdash t_1 : q T_1 \rightarrow T_2 \quad \Delta; \Gamma_2 \vdash t_2 : T_1}{\Delta; \Gamma_1 \circ \Gamma_2 \vdash t_1 t_2 : T_2} \quad (\text{T-APP})$$

The most important way to test our system for faults is to prove the type substitution lemma. In particular, the proof will demonstrate that we have made safe assumptions about how abstract type qualifiers may be used.

### 1.3.2 LEMMA [TYPE SUBSTITUTION]:

1. If  $\Delta, p; \Gamma \vdash t : T$  and  $FV(q) \subseteq \Delta$  then  $\Delta; [p \mapsto q]\Gamma \vdash [p \mapsto q]t : [p \mapsto q]T$
2. If  $\Delta, a; \Gamma \vdash t : T$  and  $FV(P) \subseteq \Delta$  then  $\Delta; [a \mapsto P]\Gamma \vdash [a \mapsto P]t : [a \mapsto P]T \quad \square$

### 1.3.3 EXERCISE [ $\star$ ]: Sketch the proof of the type substitution lemma. What structural rule(s) do you need to carry out the proof? $\square$

Operationally, we will choose to implement polymorphic instantiation using substitution. As a result, our operational semantics changes very little. We only need to specify the new computational contexts and to add the evaluation rules for polymorphic functions and application as in Figure 1-14.

## Arrays

Arrays pose a special problem for linearly typed languages. If we try to provide an operation fetches an element from an array in the usual way, perhaps using an array index expression  $a[i]$ , we would need to reflect the fact that the  $i^{th}$  element (and only the  $i^{th}$  element) of the array had been “used.” However, there is no simple way to reflect this change in the type of an array as the usual form of array types ( $\text{array}(T)$ ) provides no mechanism to distinguish between the properties of different elements of the array.

$E ::=$	<i>evaluation contexts:</i> $E [P]$ <i>pretype app context</i> $E [q]$ <i>qualifier app context</i>	$\frac{S(x) = q \Lambda a. t}{(S; x [P]) \rightarrow_{\beta} (S \stackrel{q}{\Delta} x; [a \mapsto P]t)} \quad (\text{E-PAPP})$ $(S; q \Lambda p. t) \rightarrow_{\beta} (S, x \mapsto q \Lambda p. t; x) \quad (\text{E-QFUN})$ $\frac{S(x) = q \Lambda p. t}{(S; x [q_1]) \rightarrow_{\beta} (S \stackrel{q}{\Delta} x; [p \mapsto q_1]t)} \quad (\text{E-QAPP})$
---------	---	--

**Figure 1-14: Linear lambda calculus: Polymorphic operational semantics**

We dodged this problem when we constructed our tuple operations by defining a pattern matching construct that simultaneously extracted all of the elements of a tuple. Unfortunately, we cannot follow the same path for arrays because in modern languages like Java and ML, the length of an array (and therefore the size of the pattern) is unknown at compile time.

Another non-solution to the problem is to add a special built-in iterator to process all the elements in an array at once. However, this last prevents programmers from using arrays as efficient, constant-time, random-access data structures; they might as well use lists instead.

One way out of this jam is to design the central array access operations so that, unlike the ordinary “get” and “set” operations, they *preserve* the number of pointers to the array and the number of pointers to each of its elements. We avoid our problem because there is no change to the array data structure that needs to be reflected in the type system. Using this idea, we will be able to allow programmers to define linear arrays that can hold a collection of arbitrarily many linear objects. Moreover, programmers will be able to access any of these linear objects, one at a time, using a convenient, constant-time, random-access mechanism.

So, what are the magic pointer-preserving array access operations? Actually, we need only one: a swap operation with the form  $\text{swap}(a[i], t)$ . The swap replaces the  $i^{\text{th}}$  element of the array  $a$  (call it  $t'$ ) with  $t$  and returns a (linear) pair containing the new array and  $t'$ . Notice the number of pointers to  $t$  and  $t'$  does not change during the operation. If there was one pointer to  $t$  (as an argument to  $\text{swap}$ ) before the call, then there is one pointer to  $t$  afterward (from within the array  $a$ ) and vice versa for  $t'$ . If, in addition, all of the elements of  $a$  had one pointer to them before the swap, then they will all have one pointer to them after the swap as well. Consequently, we will find it easy to type the swap operation, even when it works over linear arrays of linear objects.

In addition to `swap`, we provide functions to allocate an array given its list of elements (`array`), to determine array length (`length`) and to deallocate arrays (`free`). The last operation is somewhat unusual in that it takes two arguments `a` and `f`, where `a` is an array of type `lin array(T)` and `f` is a function with type `T → unit` that is run on each element of `T`. The function may be thought of as a finalizer for the elements; it may be used to deallocate any linear components of the array elements, thereby preserving the single pointer property.

Our definition of arrays is compatible with the polymorphic system from the previous subsection, but for simplicity, we formalize it in the context of the simply-typed lambda calculus (see Figure 1-15).

- 1.3.4 EXERCISE [RECOMMENDED, ★]: The typing rule for array allocation (`T-ARRAY`) contains the standard containment check to ensure that unrestricted arrays cannot contain linear objects. What kinds of errors can occur if this check is omitted? □
- 1.3.5 EXERCISE [★★, ↗]: With the presence of mutable data structures, it is possible to create cycles in the store. How should we modify the store typing rules to take this into account? □

The `swap` and `free` functions are relatively low-level operations. Fortunately, it is easy to build more convenient, higher-level abstractions out of them. For instance, the following code defines some simple functions for manipulating linear matrices of unrestricted integers.

```

type iArray = lin array(int)
type matrix = lin array(iArray)

fun dummy(x:unit):iArray = lin array()

fun freeElem(x:int):unit = ()
fun freeArray(a:iArray):unit = free(a,freeElem)
fun freeMatrix(m:matrix):unit = free(m,freeArray)

fun get(a:matrix,i:int,j:int):lin (matrix * int) =
  split swap(a[i],dummy()) as a,b in
  split swap(b[j],0) as b,k in
  split swap(b[j],k) as b,_ in
  split swap(a[i],b) as a,junk in
  freeArray(junk);
  lin <a,k>

```



<p>P ::=</p> <p>... array(T)</p> <p>t ::=</p> <p>... q array(t, ..., t) swap(t[t], t) length(t) free(t, t)</p> <p>w ::=</p> <p>... array[n, x, ..., x]</p> <p>E ::=</p> <p>... q array(v, ..., v, E, t, ..., t)</p> <p>swap(E(t), t) swap(v(E), t) swap(v(v), E) length(E) free(E, t) free(v, E)</p>	<p><i>pretypes:</i> as before</p> <p><i>array pretypes</i></p> <p><i>terms:</i> as before</p> <p><i>array creation</i> swap length deallocate</p> <p><i>prevalues:</i> as before</p> <p><i>array</i></p> <p><i>evaluation contexts:</i> as before</p> <p><i>array context</i> swap context swap context swap context length context free context free context</p>	<p><i>Typing</i> <span style="border: 1px solid black; padding: 2px;"><math>\Gamma \vdash t : T</math></span></p> $\frac{q(T) \quad \Gamma \vdash t_i : T \quad (\text{for } 1 \leq i \leq n)}{\Gamma \vdash q \text{ array}(t_1, \dots, t_n) : q \text{ array}(T)} \quad (\text{T-ARRAY})$ $\frac{\Gamma \vdash t_1 : q_1 \text{ array}(T_1) \quad \Gamma \vdash t_2 : q_2 \text{ int} \quad \Gamma \vdash t_3 : T_1}{\Gamma \vdash \text{ swap}(t_1[t_2], t_3) : \text{ lin } (q_1 \text{ array}(T_1) * T_1)} \quad (\text{T-SWAP})$ $\frac{\Gamma \vdash t : q \text{ array}(T)}{\Gamma \vdash \text{ length}(t) : \text{ lin } (q \text{ array}(T) * \text{int})} \quad (\text{T-LENGTH})$ $\frac{\Gamma \vdash t_1 : q \text{ array}(T) \quad \Gamma \vdash t_2 : T \rightarrow \text{unit}}{\Gamma \vdash \text{ free}(t_1, t_2) : \text{unit}} \quad (\text{T-FREE})$ <p><i>Evaluation</i> <span style="border: 1px solid black; padding: 2px;"><math>(S; t) \rightarrow_\beta (S'; t')</math></span></p> $\frac{(S; q \text{ array}(x_0, \dots, x_{n-1})) \rightarrow_\beta ((S, x \mapsto q \text{ array}[n, x_0, \dots, x_{n-1}]; x))}{S(x_i) = q_i j} \quad (\text{E-ARRAY})$ $\frac{S = S_1, x_a \mapsto q \text{ array}[n, \dots, x_j, \dots], S_2 \quad S' = S_1, x_a \mapsto q \text{ array}[n, \dots, x_e, \dots], S_2}{(S; \text{ swap}(x_a[x_i], x_e)) \rightarrow_\beta (S' \stackrel{q_i}{\sim} x_i; \text{ lin } \langle x_a, x_j \rangle)} \quad (\text{E-SWAP})$ $\frac{S(x) = q \text{ array}[n, x_0, \dots, x_{n-1}]}{(S; \text{ length}(x)) \rightarrow_\beta (S; \text{ lin } \langle x, \text{un } n \rangle)} \quad (\text{E-LENGTH})$ $\frac{S(x_a) = q \text{ array}[n, x_0, \dots, x_{n-1}]}{(S; \text{ free}(x_a, x_f)) \rightarrow_\beta (S \stackrel{q}{\sim} x_a; \text{ App}(x_f, x_0, \dots, x_{n-1}))} \quad (\text{E-FREE})$ <p>where</p> $\text{App}(x_f, \cdot) = ()$ $\text{App}(x_f, x_0, \dots) = x_f x_0; \text{App}(x_f, \dots)$
--	---	---

Figure 1-15: Linear lambda calculus: Arrays

```

fun set(a:matrix,i:int,j:int,e:int):matrix =
  split swap(a[i],dummy()) as a,b in
  split swap(b[j],e) as b,_ in
  split swap(a[i],b) as a,junk in
  freeArray(junk);
a

```

- 1.3.6 EXERCISE [**★★**, **→**]: Use the functions provided above to write matrix-matrix multiply. Your multiply function should return an integer and deallocate both arrays in the process. Use any standard integer operations necessary.  $\square$

In the examples above, we needed some sort of dummy value to swap into an array to replace the value we wanted to extract. For integers and arrays it was easy to come up with one. However, when dealing with polymorphic or abstract types, it may not be possible to conjure up a value of the right type. Consequently, rather than manipulating arrays with type `q array(a)` for some abstract type `a`, we may need to manipulate arrays of options with type `q array(a + unit)`. In this case, when we need to read out a value, we always have another value (`inr ()`) to swap in in its place. Normally such operations are called *destructive reads*; they are a common way to preserve the single pointer property when managing complex structured data.

## Reference Counting

Array swaps and destructive reads are dynamic techniques that can help overcome a lack of compile-time knowledge about the number of uses of a particular object. *Reference counting* is another dynamic technique that serves a similar purpose. Rather than restricting the number of pointers to an object to be exactly one, we can allow any number of pointers to the object and keep track of that number dynamically. Only when the last reference is used will the object be deallocated.

There are various ways to integrate reference counts into the current system. Here, we choose the simplest, which is to add a new qualifier `rc` for reference-counted data structures, and operations that allow the programmer to explicitly increment (`inc`) and decrement (`dec`) the counts (see Figure 1-16). More specifically, the increment operation takes a pointer argument, increments the reference count for the object pointed to, and returns two copies of the pointer in a (linear) pair. The decrement operation takes two arguments, a pointer and a function, and works as follows. In the case the object pointed to (call it `x`) has a reference count of 1 before the decrement, the function is executed with `x` as a linear argument. Since the function treats `x`

<p><i>Syntax</i></p> <p><math>q ::=</math></p> <p style="padding-left: 20px;">...</p> <p style="padding-left: 20px;"><math>rc</math></p> <p><math>t ::=</math></p> <p style="padding-left: 20px;">...</p> <p style="padding-left: 20px;"><math>inc(t)</math></p> <p style="padding-left: 20px;"><math>dec(t, t)</math></p>	<p><i>qualifiers:</i></p> <p style="padding-left: 20px;"><i>as before</i></p> <p style="padding-left: 20px;"><i>ref. count</i></p> <p><i>terms:</i></p> <p style="padding-left: 20px;"><i>as before</i></p> <p style="padding-left: 20px;"><i>increment count</i></p> <p style="padding-left: 20px;"><i>decrement count</i></p>	<p><i>Qualifier Relations</i></p> <p style="text-align: right;"><math>rc \sqsubseteq un</math> (Q-RCUN)</p> <p style="text-align: right;"><math>lin \sqsubseteq rc</math> (Q-LINRC)</p> <p><i>Typing</i></p> <div style="border: 1px solid black; display: inline-block; padding: 2px; margin-bottom: 10px;"><math>\Gamma \vdash t : T</math></div> <p style="text-align: right;"><math>\frac{\Gamma \vdash t : rc P}{\Gamma \vdash inc(t) : lin (rc P * rc P)}</math> (T-INC)</p> <p style="text-align: right;"><math>\frac{\Gamma \vdash t_1 : rc P \quad \Gamma \vdash t_2 : lin P \rightarrow unit}{\Gamma \vdash dec(t_1, t_2) : unit}</math> (T-DEC)</p>
--	---	--

**Figure 1-16: Linear lambda calculus: Reference counting syntax and typing**

linearly, it will deallocate  $x$  before it completes. In the other case, when  $x$  has a reference count greater than 1, the reference count is simply decremented and the function is not called; `unit` is returned as the result of the operation.

The main typing invariant in this system is that whenever a reference-counted variable appears in the static type-checking context, there is one dynamic reference count associated with it. Linear typing will ensure the number of references to an object is properly preserved.

The new `rc` qualifier should be treated in the same manner as the linear qualifier when it comes to context splitting. In other words, a reference-counted variable should be placed in exactly one of the left-hand context or the right-hand context (not both). In terms of containment, the `rc` qualifier sits between unrestricted and linear qualifiers: A reference-counted data structure may not be contained in unrestricted data structures and may not contain linear data structures. Figure 1-16 presents the appropriate qualifier relation and typing rules for our reference counting additions.

In order to define the execution behavior of reference-counted data structures, we will define a new sort of stored value with the form  $rc(n) w$ . The integer  $n$  is the reference count: it keeps track of the number of times the value is referenced elsewhere in the store or in the program.

The operational semantics for the new commands and reference-counted pairs and functions are summarized in Figure 1-17. Several new bits of notation show up here to handle the relatively complex computation that must go on to increment and decrement reference counts. First, in a slight abuse of notation, we allow  $q$  to range over static qualifiers `un`, `lin` and `rc` as well as dynamic qualifiers `un`, `lin` and `rc(n)`. Context will disambiguate the two

different sorts of uses. Second, we extend the notation  $S^q_x$  so that  $q$  may be  $rc(n)$  as well as  $lin$  and  $un$ . If  $n$  is 1 then  $S^{rc(n)}_x$  removes the binding  $x \mapsto rc(n)$   $w$  from  $S$ . Otherwise,  $S^{rc(n)}_x$  replaces the binding  $x \mapsto rc(n)$   $w$  with  $x \mapsto rc(n-1)$   $w$ . Finally, given a store  $S$  and a set of variables  $X$ , we define the function  $incr(S; X)$ , which produces a new store  $S'$  in which the reference count associated with any reference-counted variables  $x \in X$  is increased by 1.

To understand how the reference counting operational semantics works, we will focus on the rules for pairs. Allocation and use of linear and unrestricted pairs stays unchanged from before as in rules (E-PAIR') and (E-SPLIT'). Rule (E-PAIRRC) specifies that allocation of reference-counted pairs is similar to allocation of other data, except for the fact that the dynamic reference count must be initialized to 1. Use of reference-counted pairs is identical to use of other kinds of pairs when the reference count is 1: We remove the pair from the store via the function  $S^{rc(n)}_x$  as shown in rule and substitute the two components of the pair in the body of the term as shown in (E-SPLIT'). When the reference count is greater than 1, rule (E-SPLITRC) shows there are additional complications. More precisely, if one of the components of the pair, say  $y_1$ , is reference-counted then  $y_1$ 's reference count must be increased by 1 since an additional copy of  $y_1$  is substituted through the body of  $\tau$ . We use the  $incr$  function to handle the possible increase. In most respects, the operational rules for reference-counted functions follow the same principles as reference-counted pairs. Increment and decrement operations are also relatively straightforward.

In order to state and prove the progress and preservation lemmas for our reference-counting language, we must generalize the type system slightly. In particular, our typing contexts must be able specify the fact that a particular reference should appear exactly  $n$  times in the store or current computation. Reference-counted values in the store are described by these contexts and the context-splitting relation is generalized appropriately. Figure 1-18 summarizes the additional typing rules.

- 1.3.7 EXERCISE [\*\*\*, +]: State and prove progress and preservation lemmas for the simply-typed linear lambda calculus (functions and pairs) with reference counting. □

## 1.4 An Ordered Type System

Just as linear type systems provide a foundation for managing memory allocated on the heap, *ordered* type systems provide a foundation for managing memory allocated on the stack. The central idea is that by controlling the

$v ::= \dots$ <p style="text-align: right;"><i>values:</i> <i>as before</i></p> $rc(n) w$ <p style="text-align: right;"><i>ref-counted value</i></p> $E ::= \dots$ <p style="text-align: right;"><i>evaluation contexts:</i> <i>as before</i></p> $inc(E)$ <p style="text-align: right;"><i>inc context</i></p> $dec(E, t)$ <p style="text-align: right;"><i>dec context</i></p> $dec(x, E)$ <p style="text-align: right;"><i>dec context</i></p> <p><i>Evaluation</i> <span style="border: 1px solid black; padding: 2px;"><math>(S; t) \rightarrow_{\beta} (S'; t')</math></span></p>	$\frac{(q \in \{un, lin\})}{(S; q \lambda y:T. t) \rightarrow_{\beta} (S, x \mapsto q \lambda y:T. t; x)} \quad (E-FUN')$ $\frac{(S; rc \lambda y:T. t) \rightarrow_{\beta} (S, x \mapsto rc(1) \lambda y:T. t; x)}{(S; rc \lambda y:T. t) \rightarrow_{\beta} (S, x \mapsto rc(1) \lambda y:T. t; x)} \quad (E-FUNRC)$ $\frac{S(x_1) = q \lambda y:T. t \quad (q \in \{un, lin, rc(1)\})}{(S; x_1 x_2) \rightarrow_{\beta} (S \stackrel{q}{\sim} x_1; [y \mapsto x_2]t)} \quad (E-APP')$ $\frac{S(x_1) = rc(n) \lambda y:T. t \quad (n > 1 \text{ and } X = FV(\lambda y:T. t)) \quad incr(S; X) = S'}{(S; x_1 x_2) \rightarrow_{\beta} (S' \stackrel{rc(n)}{\sim} x_1; [y \mapsto x_2]t)} \quad (E-APPRC)$ $\frac{incr(S; \{x\}) = S'}{(S; inc(x)) \rightarrow_{\beta} (S'; lin \langle x, x \rangle)} \quad (E-INC)$ $\frac{S(x) = rc(n) w \quad (n > 1)}{(S; dec(x, x_f)) \rightarrow_{\beta} (S \stackrel{rc(n)}{\sim} x; un ())} \quad (E-DEC1)$ $\frac{S = S_1, x \mapsto rc(1) w, S_2 \quad S' = S_1, x \mapsto lin w, S_2}{(S; dec(x, x_f)) \rightarrow_{\beta} (S'; x_f x)} \quad (E-DEC2)$
$\frac{(q \in \{un, lin\})}{(S; q \langle y, z \rangle) \rightarrow_{\beta} (S, x \mapsto q \langle y, z \rangle; x)} \quad (E-PAIR')$ $\frac{(S; rc \langle y, z \rangle) \rightarrow_{\beta} (S, x \mapsto rc(1) \langle y, z \rangle; x)}{(S; rc \langle y, z \rangle) \rightarrow_{\beta} (S, x \mapsto rc(1) \langle y, z \rangle; x)} \quad (E-PAIRRC)$ $\frac{S(x) = q \langle y_1, z_1 \rangle \quad (q \in \{un, lin, rc(1)\})}{(S; split x as y, z in t) \rightarrow_{\beta} (S \stackrel{q}{\sim} x; [y \mapsto y_1][z \mapsto z_1]t)} \quad (E-SPLIT')$ $\frac{S(x) = rc(n) \langle y_1, z_1 \rangle \quad (n > 1) \quad incr(S; \{y_1, z_1\}) = S'}{(S; split x as y, z in t) \rightarrow_{\beta} ((S' \stackrel{rc(n)}{\sim} x); [y \mapsto y'_1][z \mapsto z'_1]t)} \quad (E-SPLITRC)$	

Figure 1-17: Linear lambda calculus: Reference counting operational semantics

exchange property, we are able to guarantee that certain values, those values allocated on the stack, are used in a first-in/last-out order.

To formalize this idea, we organize the store into two parts: a stack, which is a sequence of locations that can be accessed on one end (the “top”) and a heap, which is like the store described in previous sections of this chapter. Pairs, functions and other objects introduced with unrestricted or linear qualifiers are allocated on the heap as before. And as before, when a linear pair or function is used, it is deallocated. Also, we allow programmers to allocate simple data structures on the stack. Without the exchange property, an ordered object can only be used when it is at the top of the stack. When this happens, the ordered object is popped off the top of the stack.

<p><i>Syntax</i></p> $\Gamma ::=$ $\dots$ $\Gamma, x : rc(n)P$ <p><i>Store Typing</i></p> $\frac{\vdash S : \Gamma_1 \circ \Gamma_2 \quad \Gamma_1 \vdash rc\ w : rc\ P}{\vdash S, x \mapsto rc(n)\ w : \Gamma_2, x : rc(n)\ P} \text{ (T-NEXTRCS)}$	<p><i>Context Splitting</i></p> $\frac{\Gamma = \Gamma_1 \circ \Gamma_2 \quad n = i + j}{\Gamma, x : rc(n)P = (\Gamma_1, x : rc(i)P) \circ (\Gamma_2, x : rc(j)P)} \text{ (M-RC)}$ <p>(when <math>i</math> or <math>j</math> is 0, the corresponding binding is removed from the context)</p> <p><i>Variable Typing</i></p> $\frac{un(\Gamma_1, \Gamma_2)}{\Gamma_1, x : rc(1)P, \Gamma_2 \vdash x : rc\ P} \text{ (T-RCVAR)}$
--	--

**Figure 1-18: Linear lambda calculus: Reference counting run-time typing**

### Syntax

The overall structure and mechanics of the ordered type system are very similar to the linear type system developed in previous sections. Figure 1-19 presents the syntax. One key change from our linear type system is that we have introduced an explicit sequencing operation  $\text{let } x = t_1 \text{ in } t_2$  that first evaluates the term  $t_1$ , binds the result to  $x$ , and then continues with the evaluation of  $t_2$ . This sequencing construct gives programmers explicit control over the order of evaluation of terms, which is crucial now that we are introducing data that must be used in a particular order. Terms that normally can contain multiple nested subexpressions such as pair introduction and function application are syntactically restricted so that their primary subterms are variables and the order of evaluation is clear.

The other main addition is a new qualifier `ord` that marks data allocated on the stack. We only allow pairs and values with base type to be stack-allocated; functions are allocated on the unordered heap. Therefore, we declare types  $\text{ord } T_1 \rightarrow T_2$  and terms  $\text{ord } \lambda x : T. t$  to be syntactically ill-formed.

Ordered assumptions are tracked in the type checking context  $\Gamma$  like other assumptions. However, they are not subject to the exchange property. Moreover, the order that they appear in  $\Gamma$  mirrors the order that they appear on the stack, with the rightmost position representing the stack's top.

### Typing

The first step in the development of the type system is to determine how assumptions will be used. As before, unrestricted assumptions can be used

Syntax			
$q ::=$		<i>qualifiers:</i>	
ord		ordered	
lin		linear	
un		unrestricted	
$t ::=$		<i>terms:</i>	
x		variable	
q b		Boolean	
if t then t else t		conditional	
q <x,y>		pair	
split t as x,y in t		split	
q $\lambda x:T.t$		abstraction	
			$x y$ <i>application</i>
			$\text{let } x = t \text{ in } t$ <i>sequencing</i>
		$P ::=$	<i>pretypes:</i>
		Bool	booleans
		$T * T$	pairs
		$T \rightarrow T$	functions
		$T ::=$	<i>types:</i>
		q P	qualified pretype
		$\Gamma ::=$	<i>contexts:</i>
		$\emptyset$	empty context
		$\Gamma, x:T$	term variable binding

Figure 1-19: Ordered lambda calculus: Syntax

as often as the programmer likes but linear assumptions must be used exactly once along every control flow path. Ordered assumptions must be used exactly once along every control flow path, in the order in which they appear.

As before, the context splitting operator ( $\Gamma = \Gamma_1 \circ \Gamma_2$ ) helps propagate assumptions properly, separating the context  $\Gamma$  into  $\Gamma_1$  and  $\Gamma_2$ . Some sequence of ordered assumptions taken from the left-hand side of  $\Gamma$  are placed in  $\Gamma_1$  and the remaining ordered assumptions are placed in  $\Gamma_2$ . Otherwise, the splitting operator works the same as before. In the typing rules, the context  $\Gamma_2$  is used by the first subexpression to be evaluated (since the top of the stack is at the right) and  $\Gamma_1$  is used by the second subexpression to be evaluated. Formally, we define the "=" relation in terms of two subsidiary relations: " $=_1$ ," which places ordered assumptions in  $\Gamma_1$ , and " $=_2$ ," which places ordered assumptions in  $\Gamma_2$ . See Figure 1-20.

The second step in the development of the type system is to determine the containment rules for ordered data structures. Previously, we saw that if an unrestricted object can contain a linear object, a programmer can write functions that duplicate or discard linear objects, thereby violating the central invariants of the system. A similar situation arises if linear or unrestricted objects can contain stack objects; in either case, the stack object might be used out of order, after it has been popped off the stack. The typing rules use the qualifier relation  $q_1 \sqsubseteq q_2$ , which specifies that  $\text{ord} \sqsubseteq \text{lin} \sqsubseteq \text{un}$ , to ensure such problems do not arise.

The typing rules for the ordered lambda calculus appear in Figure 1-21. For the most part, the containment rules and context splitting rules encapsulate

<i>Context Split</i>	$\Gamma = \Gamma_1 \circ \Gamma_2$	
$\frac{\Gamma =_2 \Gamma_1 \circ \Gamma_2}{\Gamma = \Gamma_1 \circ \Gamma_2}$	(M-TOP)	$\frac{\Gamma =_1 \Gamma_1 \circ \Gamma_2}{\Gamma =_2 \Gamma_1 \circ \Gamma_2}$ (M-1TO2)
$\emptyset =_1 \emptyset \circ \emptyset$	(M-EMPTY)	$\frac{\Gamma =_{1,2} \Gamma_1 \circ \Gamma_2}{\Gamma, x:\text{lin } P =_{1,2} (\Gamma_1, x:\text{lin } P) \circ \Gamma_2}$ (M-LINA)
$\frac{\Gamma =_1 \Gamma_1 \circ \Gamma_2}{\Gamma, x:\text{ord } P =_1 (\Gamma_1, x:\text{ord } P) \circ \Gamma_2}$	(M-ORD1)	$\frac{\Gamma =_{1,2} \Gamma_1 \circ \Gamma_2}{\Gamma, x:\text{lin } P =_{1,2} \Gamma_1 \circ (\Gamma_2, x:\text{lin } P)}$ (M-LINB)
$\frac{\Gamma =_2 \Gamma_1 \circ \Gamma_2}{\Gamma, x:\text{ord } P =_2 \Gamma_1 \circ (\Gamma_2, x:\text{ord } P)}$	(M-ORD2)	$\frac{\Gamma =_{1,2} \Gamma_1 \circ \Gamma_2}{\Gamma, x:\text{un } P =_{1,2} (\Gamma_1, x:\text{un } P) \circ (\Gamma_2, x:\text{un } P)}$ (M-UN)

Figure 1-20: Ordered lambda calculus: Context splitting

the tricky elements of the type system. The rules for pairs illustrate how this is done. The rule for introducing pairs (T-OPAIR) splits the incoming context into two parts,  $\Gamma_1$  and  $\Gamma_2$ ; any ordered assumptions in  $\Gamma_2$  will represent data closer to the top of the stack than  $\Gamma_1$ . Therefore, if the pair  $(x)$  and its two components  $x_1$  and  $x_2$  are all allocated on the stack, then the pointer  $x$  will end up on top,  $x_2$  next and  $x_1$  on the bottom. The elimination rule for pairs (T-OSPLIT) is careful to maintain the proper ordering of the context. As above, the rule splits the context into  $\Gamma_1$  and  $\Gamma_2$ , where  $\Gamma_2$ , which represents data on top of the stack, is used in a computation  $t_1$  that generates a pair. The context  $\Gamma_1, x_1:T_1, x_2:T_2$  is used to check  $t_2$ . Notice that if both components of the pair,  $x_1$  and  $x_2$ , were allocated on the stack when the pair was introduced, they reappear back in the context in the appropriate order.

Consider the following function, taking a boolean and a pair allocated sequentially at the top of the stack. The boolean is at the very top of the stack and the integer pair is next (the top is to the right). If the boolean is true, it leaves the components of the pair (two unrestricted integers) in the same order as given; otherwise, it swaps them.

```

λx:ord (ord (int * int) * bool).
  split x as p,b in
  if b then
    p
  else
    split p as i1,i2 in
    ord <i2,i1>

```



<i>Typing</i>	$\boxed{\Gamma \vdash t : T}$	
$\frac{\text{un } (\Gamma_1, \Gamma_2)}{\Gamma_1, x:T, \Gamma_2 \vdash x : T}$	(T-OVAR)	$\frac{\Gamma_2 \vdash t_1 : q (T_1 * T_2) \quad \Gamma_1, x_1:T_1, x_2:T_2 \vdash t_2 : T}{\Gamma_1 \circ \Gamma_2 \vdash \text{split } t_1 \text{ as } x_1, x_2 \text{ in } t_2 : T}$
$\frac{\text{un } (\Gamma)}{\Gamma \vdash q \text{ b} : q \text{ Bool}}$	(T-OBOOL)	$\frac{q(\Gamma) \quad \Gamma, x:T_1 \vdash t_2 : T_2}{\Gamma \vdash q \lambda x:T_1. t_2 : q T_1 \rightarrow T_2}$
$\frac{\Gamma_2 \vdash t_1 : q \text{ Bool} \quad \Gamma_1 \vdash t_2 : T \quad \Gamma_1 \vdash t_3 : T}{\Gamma_1 \circ \Gamma_2 \vdash \text{if } t_1 \text{ then } t_2 \text{ else } t_3 : T}$	(T-OIF)	$\frac{\Gamma_1 \vdash x_1 : q T_{11} \rightarrow T_{12} \quad \Gamma_2 \vdash x_2 : T_{11}}{\Gamma_1 \circ \Gamma_2 \vdash x_1 x_2 : T_{12}}$
$\frac{\Gamma_1 \vdash x_1 : T_1 \quad \Gamma_2 \vdash x_2 : T_2 \quad q(T_1) \quad q(T_2)}{\Gamma_1 \circ \Gamma_2 \vdash q \langle x_1, x_2 \rangle : q (T_1 * T_2)}$	(T-OPAIR)	$\frac{\Gamma_2 \vdash t_1 : T_1 \quad \Gamma_1, x:T_1 \vdash t_2 : T_2}{\Gamma_1 \circ \Gamma_2 \vdash \text{let } x = t_1 \text{ in } t_2 : T_2}$
		(T-OSPLIT)
		(T-OABS)
		(T-OAPP)
		(T-OLET)

Figure 1-21: Ordered lambda calculus: Typing

### Operational Semantics

To define the operational semantics for our new ordered type system, we will divide our previous stores into two parts, a heap  $H$  and a stack  $K$ . Both are just a list of bindings as stores were before (see Figure 1-22). We also define a couple of auxiliary functions. The first says what it means to add a binding to the store. This is straightforward: unrestricted and linear bindings are added to the heap and ordered bindings are added to the top of the stack.

$$\begin{aligned} (H;K), x \mapsto \text{ord } w &= (H;K, x \mapsto \text{ord } w) \\ (H;K), x \mapsto \text{lin } w &= (H, x \mapsto \text{lin } w;K) \\ (H;K), x \mapsto \text{un } w &= (H, x \mapsto \text{un } w;K) \end{aligned}$$

The second function specifies how to remove a binding from the store. Notice that ordered deallocation will only remove the object at the top of the stack.

$$\begin{aligned} (H;K, x \mapsto v) \overset{\text{ord}}{\sim} x &= H;K \\ (H_1, x \mapsto v, H_2;K) \overset{\text{lin}}{\sim} x &= H_1, H_2;K \\ (H;K) \overset{\text{un}}{\sim} x &= H;K \end{aligned}$$

With these simple changes, the evaluation rules from previous sections can be reused essentially unchanged. However, we do need to add the evaluation context for sequencing ( $\text{let } x = E \text{ in } t$ ) and its evaluation rule:

$$(S; \text{let } x = x_1 \text{ in } t_2) \rightarrow_{\beta} (S; [x \mapsto x_1]t_2) \quad (\text{E-LET})$$

$S ::= H;K$ $H ::= \emptyset$ $H, x \mapsto \text{lin } w$ $H, x \mapsto \text{un } w$	$ $	$K ::= \emptyset$ $K, x \mapsto \text{ord } w$
<p style="text-align: right; margin-right: 20px;"><i>stores:</i></p> <p style="text-align: right; margin-right: 20px;"><i>complete store</i></p> <p style="text-align: right; margin-right: 20px;"><i>heap:</i></p> <p style="text-align: right; margin-right: 20px;"><i>empty heap</i></p> <p style="text-align: right; margin-right: 20px;"><i>linear heap binding</i></p> <p style="text-align: right; margin-right: 20px;"><i>unrestricted heap binding</i></p>		<p style="text-align: right; margin-right: 20px;"><i>stack:</i></p> <p style="text-align: right; margin-right: 20px;"><i>empty stack</i></p> <p style="text-align: right; margin-right: 20px;"><i>stack binding</i></p>

**Figure 1-22: Ordered lambda calculus: Operational semantics**

- 1.4.1 EXERCISE [RECOMMENDED, ★]: Write a program that demonstrates what can happen if the syntax of pair formation is changed to allow programmers to write nested subexpressions (i.e., we allow the term  $\text{ord } \langle t_1, t_2 \rangle$  rather than the term  $\text{ord } \langle x, y \rangle$ ). □
- 1.4.2 EXERCISE [RECOMMENDED, ★★]: Demonstrate the problem with allowing ordered functions (i.e., admitting the syntax  $\text{ord } \lambda x:T_1. \tau$  and  $\text{ord } T_1 \rightarrow T_2$ ) by writing a well-typed program that uses ordered functions and gets stuck. □
- 1.4.3 EXERCISE [★★★]: Modify the language so that programmers can use stack-allocated, ordered functions. There are many solutions to this problem, some more sophisticated than others. □

## 1.5 Further Applications

Memory management applications make good motivation for substructural type systems and provides a concrete framework for studying their properties. However, substructural types systems, and their power to control the number and order of uses of data and operations, have found many applications outside of this domain. In the following paragraphs, we informally discuss a few of them.

### Controlling Temporal Resources

We have studied several ways that substructural type systems can be used to control physical resources such as memory and files. What about controlling the temporal resources? Amazingly, substructural type systems can play a role here as well: Careful crafting of a language with an *affine* type system, where values are used at most once, can ensure that computations execute in polynomial time.

To begin, we will allow our polynomial time language to contain affine booleans, pairs and (non-recursive) functions. In addition, to make things interesting, we will add affine lists to our language, which have constructors `nil` and `cons` and a special iterator to recurse over objects with list type. Such iterators have the following form.

$$\text{iter } (\text{stop} \Rightarrow t_1 \mid x \text{ with } y \Rightarrow t_2)$$

If  $t_1$  has type  $T$  and  $t_2$  also has type  $T$  (under the assumption that  $x$  has type  $T_1$  and  $y$  has type  $T_1 \text{ list}$ ), our iterator defines a function from  $T_1$  lists to objects with type  $T$ . Operationally, the iterator does a case to see whether its input list is `nil` or `cons(hd, tl)` and executes the corresponding branch. We can define the operation of iterators using two simple rules.<sup>1</sup>

$$\begin{array}{l} \text{iter } (\text{stop} \Rightarrow t_1 \mid \text{hd with rest} \Rightarrow t_2) \text{ nil} \longrightarrow_{\beta} t_1 \quad (\text{E-ITERNIL}) \\ \hline \text{iter } (\text{stop} \Rightarrow t_1 \mid \text{hd with rest} \Rightarrow t_2) v_2 \xrightarrow{*}_{\beta} v'_2 \quad (\text{E-ITERCONS}) \\ \text{iter } (\text{stop} \Rightarrow t_1 \mid \text{hd with rest} \Rightarrow t_2) \text{ cons}(v_1, v_2) \longrightarrow_{\beta} \\ \quad [\text{hd} \mapsto v_1][\text{rest} \mapsto v'_2]t_2 \end{array}$$

In the second rule, the iterator is invoked inductively on  $v_2$ , giving the result  $v'_2$ , which is used in term  $t_2$ . The familiar `append` function below illustrates the use of iterators.

```
val append : T list → T list → T list =
  iter (
    stop ⇒ λ(l:T list).l
    | hd with rest ⇒ λ(l:T list).cons(hd, rest l))
```

When applied to a list  $l_1$ , `append` builds up a function that expects a second list  $l_2$  and concatenates  $l_2$  to the end of  $l_1$ . Clearly, `append` is a polynomial time function, a linear-time one in fact, but it is straightforward to write exponential time algorithms in the language as we have defined it so far. For instance:

```
val double : T list → T list =
  iter (stop ⇒ nil | hd with rest ⇒ cons(hd, cons(hd, rest)))

val exp : T list → T list =
  iter (stop ⇒ nil | hd with rest ⇒ double (cons(hd, rest)))
```

---

1. Since we are not interested in memory management here, we have simplified our operational semantics from previous parts of this chapter by deleting the explicit store and using substitution instead. The operational judgment has the form  $t \rightarrow_{\beta} t'$  and, in general, is defined similarly to the operational systems in *TAPL*.

The key problem here is that it is trivial to write iterators like `double` that increase the size of their arguments. After constructing one of these, we can use it as the inner loop of another, like `exp`, and cause an exponential blow-up in running time. But this is not the only problem. Higher-order functions make it even easier to construct exponential-time algorithms:

```
val compose =
  λ(fg:(T list→T list) * (T list→T list)).
  λ(x:T list).
    split fg as f,g in f (g x)

val junk : T

val exp2 : T list→T list→T list =
  iter (
    stop ⇒ λ(l:T list).cons(junk,l)
  | hd with rest ⇒ λ(l:T list).compose <rest,rest> l)
```

Fortunately, a substructural type system can be used to eliminate both problems by allowing us to define a class of *non-size-increasing* functions and by preventing the construction of troublesome higher-order functions, such as `exp2`.

The first step is to demand that all user-defined objects have affine type. They can be used zero or one times, but not more. This restriction immediately rules out programs such as `exp2`. System defined operators like `cons` can be used many times.

The next step is to put mechanisms in place to prevent iterators from increasing the size of their inputs. This can be achieved by altering the `cons` constructor so that it can only be applied when it has access to a special resource with type `R`.

```
operator cons : (R,T,T list) → T list
```

There is no constructor for resources with type `R` so they cannot be generated out of thin air; we can only apply `fcons` as many times as we have resources. We also adapt the syntax for iterators as follows.

```
iter (stop ⇒ t1 | hd with t1 and r ⇒ t2)
```

Inside the second clause of the iterator, we are only granted a single resource (`r`) with which to allocate data. Consequently, we can allocate at most one `cons` cell in `t2`. This provides us with the power to rebuild a list of the same size, but we cannot write a function such as `double` that doubles the length of the list or `exp` that causes an exponential increase in size. To ensure that

a single resource from an outer scope does not percolate inside the iterator and get reused on each iteration of the loop, we require that iterators be closed, mirroring the containment rules for recursive functions defined in earlier sections of this chapter.

Although restricted to polynomial time, our language permits us to write many useful functions in a convenient fashion. For instance, we can still write `append` much like we did before. The resource we acquire from destructuring the list during iteration can be used to rebuild the list later.

```
val append : T list → T list → T list =
  iter (
    stop ⇒ λ(l:T list).l
    | hd with rest and r ⇒ λ(l:T list). cons(r,hd,rest l))
```

We can also write `double` if our input list comes with appropriate credits, in the form of unused resources.

```
val double : (T*R) list → T list =
  iter (
    stop ⇒ nil
    | hd with rest and r1 ⇒
      split hd as x,r2 in cons(r1,hd,cons(r2,hd,rest)))
```

Fortunately, we will never be able to write `exp`, unless, of course, we are given an exponential number of credits in the size of the input list. In that case, our function `exp` would still only run in linear time with respect to our overall input (list and resources included).

The proof that all (first-order) functions we can define in this language run in polynomial time uses some substantial domain theory that lies outside the scope of this book. However, the avid reader should see Section 1.6 for references to the literature where these proofs can be found.

## Compiler Optimizations

Many compiler optimizations are enabled when we know that there will be *at most one use* or *at least one use* of a function, expression or data structure. If there is at most one use of an object then we say that object has *affine* type. If there is at least one use then we say the object has *relevant* (or *strict*) type. The following sorts of optimizations employ usage information directly; several of them have been implemented in the Glasgow Haskell Compiler.

- *Floating in bindings.* Consider the expression `let x = e in (λy...x...)`. Is it a good idea to float the binding inside the lambda and create the new

expression  $\lambda y. \text{let } x = e \text{ in } (\dots x \dots)$ ? The answer depends in part on how many times the resulting function is used. If it is used at most once, the optimization might be a good one: we may avoid computing  $e$  and will never compute it more than once.

- *Inlining expressions.* In the example above, if we have the further information that  $x$  itself is used at most once inside the body of the function, then we might want to substitute the expression  $e$  for  $x$ . This may give rise to further local optimizations at the site where  $e$  is used. Moreover, if it turns out that  $e$  is used zero times (as opposed to one time) we will have saved ourselves the trouble of computing it.
- *Thunk update avoidance.* In lazy functional languages such as Haskell, evaluation of function parameters is delayed until the parameter is actually used in the function body. In order to avoid recomputing the value of the parameter each time it is used, implementers make each parameter a *thunk*—a reference that may either hold the computation that needs to be run or the value itself. The first time the thunk is used, the computation will be run and will produce the necessary result. In general, this result is stored back in the thunk for all future uses of the parameter. However, if the compiler can determine that the data structure is used as most once, this thunk update can be avoided.
- *Eagerness.* If we can tell that a Haskell expression is used at least once, then we can evaluate it right away and avoid creating a thunk altogether.

The optimizations described above may be implemented in two phases. The first phase is a program analysis that may be implemented as affine and/or relevant type inference. After the analysis phase, the compiler uses the information to transform programs. Formulating compiler optimizations as type inference followed by type-directed translation has a number of advantages over other techniques. First, the language of types can be used to communicate optimization information across modular boundaries. This can facilitate the process of scaling intra-procedural optimizations to inter-procedural optimizations. Second, the type information derived in one optimization pass can be maintained and propagated to future optimization passes or into the back end of the compiler where it can be used to generate Typed Assembly Language or Proof-Carrying Code, as discussed in Chapters 4 and 5.

## 1.6 Notes

*Substructural logics* are very old, dating back to at least Orlov (1928), who axiomatized the implicational fragment of relevant logic. Somewhat later, Moh

(1950) and Church (1951) provided alternative axiomatizations of the relevant logic now known as R. In the same time period, Church was developing his theory of the lambda calculus at Princeton University, and his  $\lambda I$  calculus (1941), which disallowed abstraction over variables that did not appear free in the body of the term, was the first substructural lambda calculus. Lambek (1958) introduced the first “ordered logic,” and used it to reason about natural language sentence structure. More recently, Girard (1987) developed linear logic, which gives control over both contraction and weakening, and yet provides the full power of intuitionistic logic through the unrestricted modality “!”. O’Hearn and Pym (1999) show that the logic of bunched implications provides another way to recapture the power of intuitionistic logic while giving control over the structural rules.

For a comprehensive account of the history of substructural logics, please see Došen (1993), who is credited with coining the phrase “substructural logic,” or Restall (2005). Restall’s textbook on substructural logics (2000) provides good starting point to those looking to study the technical details of either the proof theory or model theory for these logics.

Reynolds pioneered the study of substructural type systems for programming languages with his development of syntactic control of interference (1978; 1989), which prevents two references from being bound to the same variable and thereby facilitates reasoning about Algol programs. Later, Girard’s development of linear logic inspired many researchers to develop functional languages with linear types. One of the main applications of these new type systems was to control effects and enable in-place update of arrays in pure functional languages.

Lafont (1988) was the one of the first to study programming languages with linear types, developing a linear abstract machine. He was soon followed by many other researchers, including Baker (1992) who informally showed how to compile Lisp into a linear assembly language in which all allocation, deallocation and pointer manipulation is completely explicit, yet safe. Another influential piece of work is due to Chirimar, Gunter, and Riecke (1996) who developed an interpretation of linear logic based on reference counting. The reference counting scheme described here is directly inspired by the work of Chirimar et al., but the technical setup is slightly different; we have explicit operations to increment and decrement reference counts whereas incrementing and decrementing counts in Chirimar’s system is done implicitly. Stephanie Weirich suggested the invariant for proving our reference counting system sound. Turner and Wadler (1999) summarize two computational interpretations that arise directly through the Curry-Howard isomorphism from Girard’s linear logic. They differ from the account given in this chapter as neither account has both shared, usable data structures and deallocation. Unfortunately, these two features together appear incompatible with a type

system derived directly from linear logic and its single unrestricted modality.

The development of practical linear type systems with two classes of type, one linear and one unrestricted, began with Wadler's work (1990) in the early nineties. The presentation given in this chapter is derived from Wadler's work and is also inspired by work from Wansbrough and Peyton Jones (1999) and Walker and Watkins (2001). Wansbrough and Peyton Jones included qualifier subtyping and bounded parametric polymorphism in their system in addition to many of the features described here. Walker and Watkins added reference counting features to a language with linear types and also memory regions. The idea of formulating the system with a generic context splitting operator was taken from Cervesato and Pfenning's presentation of Linear LF (2002).

The algorithmic type system described in section 1-5 solves what is commonly known in the linear logic programming and theorem proving literature, as the *resource management problem*. Many of the ideas for the current presentation came from work by Cervesato, Hodas, and Pfenning (2000), who solve the more general problem that arises when linear logic's additive connectives are considered. Hofmann takes a related approach when solving the type inference problem for a linearly-typed functional language (1997a).

The ordered type system developed here is derived from Polakow and Pfenning's ordered logic (1999), in the same way that the practical linear type systems mentioned above emerged from linear logic. It was also inspired by the ordered lambda calculus of Petersen, Harper, Crary, and Pfenning (2003), though there are some technical differences. Ahmed and Walker (2003) and Ahmed, Jia, and Walker (2003) use an ordered, modal logic to specify memory invariants and have integrated the logical specifications into a low-level typed language. Igarashi and Kobayashi (2002) have used ordered types to explore the more general problem of resource usage analysis. In addition, they have developed effective type inference algorithms for their type systems.

Recently, O'Hearn (2003) has proposed *bunched typing*, a new form of substructural typing, to control interference between mutable variables, generalizing Reynolds's earlier work on syntactic control of interference. These bunched types were derived from earlier work by O'Hearn and Pym (1999) on bunched logic. Together, Reynolds, Ishtiaq, and O'Hearn (Reynolds, 2000; Ishtiaq and O'Hearn, 2001) have used bunched logic to develop a system for verifying programs that explicitly allocate and deallocate data.

Analysis and reasoning about the time and space complexity of programs has always been an important part of computer science. However, the use of programming language technology, and type systems in particular, to automatically constrain the complexity of programs is somewhat more recent. For instance, Bellantoni and Cook (1992) and Leivant (1993) developed predicative systems that control the use and complexity of recursive functions.



It is possible to write all, and only, the polynomial-time functions in their system. However, it is not generally possible to compose functions and therefore many “obviously” polynomial-time algorithms cannot be coded naturally in their system. Girard (1998), Hofmann (2000; 1999), and Bellantoni, Niggl, and Schwichtenberg (2000) show how linear type systems can be used to alleviate some of these difficulties. The material presented in this chapter is derived from Hofmann’s work.

One of the most successful and extensive applications of substructural type systems in programming practice can be found in the Concurrent Clean programming language (Nöcker, Smetsers, van Eekelen, and Plasmeijer, 1991). Clean is a commercially developed, pure functional programming language. It uses *uniqueness types* (Barendsen and Smetsers, 1993), which are a variant of linear types, and strictness annotations (Nöcker and Smetsers, 1993) to help support concurrency, I/O and in-place update of arrays. The implementation is fast and is fully supported by a wide range of program development tools including an Integrated Development Environment for project management and GUI libraries, all developed in Clean itself.

Substructural type systems have also found gainful employment in the intermediate languages of the Glasgow Haskell Compiler. For instance, Turner, Wadler, and Mossin (1995) and Wansbrough and Peyton Jones (1999) showed how to use affine types and affine type inference to optimize programs as discussed earlier in this chapter. They also use extensive strictness analysis to avoid thunk creation.

Recently, researchers have begun to investigate ways to combine substructural type systems with dependent types and effect systems such as those described in Chapters 2 and 3. The combination of both dependent and substructural types provides a very powerful tool for enforcing safe memory management and more general resource-usage protocols. For instance, DeLine and Fähndrich developed Vault (2001; 2002), a programming language that uses static capabilities (Walker, Crary, and Morrisett, 2000) (a hybrid form of linear types and effects) to enforce a variety of invariants in Microsoft Windows device drivers including locking protocols, memory management protocols and others. Cyclone (Jim et al., 2002; Grossman et al., 2002), a completely type-safe substitute for C, also uses linear types and effects to grant programmers fine-grained control over memory allocation and deallocation. In each of these cases, the authors do not stick to the pure linear types described here. Instead, they add coercions to the language to allow linearly-typed objects to be temporarily aliased in certain contexts, following a long line of research on this topic (Wadler, 1990; Odersky, 1992; Kobayashi, 1999; Smith, Walker, and Morrisett, 2000; Aspinall and Hofmann, 2002; Foster, Terauchi, and Aiken, 2002; Aiken, Foster, Kodumal, and Terauchi, 2003).