# Lightweight Linear Types in System F°

Karl Mazurak Jianzhou Zhao

Steve Zdancewic

University of Pennsylvania {mazurak,jianzhou,stevez}@cis.upenn.edu

# Abstract

We present System  $F^{\circ}$ , an extension of System F that uses *kinds* to distinguish between linear and unrestricted types, simplifying the use of linearity for general-purpose programming. We demonstrate through examples how System  $F^{\circ}$  can elegantly express many useful protocols, and we prove that any protocol representable as a DFA can be encoded as an  $F^{\circ}$  type. We supply mechanized proofs of System  $F^{\circ}$ 's soundness and parametricity properties, along with a nonstandard operational semantics that formalizes common intuitions about linearity and aids in reasoning about protocols.

We compare System  $F^{\circ}$  to other linear systems, noting that the simplicity of our kind-based approach leads to a more explicit account of what linearity is meant to capture, allowing otherwise-conflicting interpretations of linearity (in particular, restrictions on *aliasing* versus restrictions on resource *usage*) to coexist peace-fully. We also discuss extensions to System  $F^{\circ}$  aimed at making the core language more practical, including the additive fragment of linear logic, algebraic datatypes, and recursion.

*Categories and Subject Descriptors* D.3.3 [*Programming Languages*]: Language Constructs and Features

General Terms Design, Languages, Theory

Keywords Linear logic, Polymorphism, Type systems

# 1. Introduction

*Linear logic* [15, 16] models resource usage by restricting the properties of contraction (the ability to duplicate a resource) and weakening (the ability to discard a resource). In the context of programming languages, ideas from linear logic were quickly adopted, at first to eliminate garbage collection [21] and shortly thereafter to handle mutable state [29].

Since their introduction, variants, refinements, and improvements on linear type systems have been proposed for many applications, including explicit memory management and control of aliasing [2, 13, 17, 28, 36], capabilities [9, 10], and tracking state changes for program analysis [11, 32]. Of particular interest is work on *typestates*, which ensure that a sequence of API calls is wellbehaved [13, 12], and on *session types*, which check that the endpoints of a channel agree on the next message to be sent or received [18, 25, 27]. Walker has a more comprehensive survey [33].

TLDI'10 January 23, 2010, Madrid, Spain.

Copyright © 2010 ACM 978-1-60558-891-9/10/01...\$5.00

Given these success stories, it is perhaps surprising that we have yet to see general linear types seriously considered for inclusion in a mainstream functional programming language.<sup>1</sup> But alas, linear types can easily lead to awkward programming models and potentially complicated language designs that are difficult both to implement and to program with. This paper seeks to address these issues by introducing System  $F^{\circ}$ —pronounced "F-pop"—a language that is intended to be a simple foundation for practical linear programming. Rather than aiming at one particular problem, System  $F^{\circ}$  lets programmers enforce their own protocol abstractions through the power of linearity and polymorphism, yet its typing discipline is lightweight enough to expose in a surface language.

System  $F^{\circ}$  is simply the Girard–Reynolds polymorphic  $\lambda$ -calculus [14, 23] extended with two base kinds:  $\star$ , classifying ordinary, unrestricted types, and  $\circ$ , classifying linear types. A sub-kinding relation  $\star \leq \circ$  makes explicit the observation that values of unrestricted types may safely be treated linearly; *i.e.*, since variables of unrestricted type may be used any number of times and linear variables must be used exactly once, it is always safe to use unrestricted variables as though they were linear. Any System  $F^{\circ}$  expression with kinds erased is a well-typed System F expression.

In introducing System  $F^{\circ}$ , this paper contributes the following:

- The design of System F°, and in particular its use of kinds and kind subsumption, which is lightweight and structured so as to integrate well with existing functional languages. (Section 2)
- Mechanized proofs of standard soundness and parametricity results for System F°, which ensure that the properties functional programmers rely on continue to hold. (Section 2.1)
- A second, linearity-aware semantics for System F°, which formalizes common intuitions about linearity and shows that these intuitions do indeed hold. (Section 4)
- Several examples—including all regular languages—along with extensions and proposals for compiler support, remarkable primarily for their simplicity, which showcase System F°'s potential usefulness. (Sections 3 and 5)

In the rest of this section we discuss the design of System  $F^{\circ}$  in the context of prior work, and we showcase System  $F^{\circ}$ 's ability to enforce programmer-defined protocols with a familiar example.

#### 1.1 Prior work: Linear type system design considerations

There are many variants on linear type systems in the literature [33], but the crucial design decision from our perspective is how linear and unrestricted variables are differentiated and how that mechanism interacts with polymorphism. Our use of kinds and kind subsumption is intended to capture the essence of linearity simply and generally while remaining faithful to the standard

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

<sup>&</sup>lt;sup>1</sup> Clean is often seen as the exception to this, but there are subtle differences between its uniqueness types, which concern aliasing, and standard linear types (see Section 3.1).

semantics and programming model of System F. This division between linear and unrestricted types is inspired by that proposed for *monomorphic* systems by Wadler [29] and Benton [5]. Here we contrast our approach with some alternatives from the literature.

Broadly speaking, there are two other approaches to distinguishing linear from unrestricted variables. The first approach, which closely follows linear logic, is to treat *all* types as linear and introduce the modal constructor ! to account for unrestricted terms, which must be closed with respect to linear variables. This interacts easily enough with polymorphism, since all types are treated uniformly—any type can be substituted for any type variable.

Unfortunately, assuming linearity by default requires very explicit handling of unrestricted terms; the standard approach uses a **let** ! construct to introduce unrestricted variables from suspended computations !e [6, 30]. Because unrestricted values are common in practice, this can prove quite cumbersome, and the burden extends even to base types—the constant 3 would be of type Int, not !Int. Further, full inference of !s has been shown to be impossible [22].

A second approach distinguishes linear from unrestricted types by means of *qualifiers* lin (for "linear") and un (for "unrestricted") applied to a collection of pre-types [1, 33]. These qualifiers (which do not nest) constrain usage (or aliasing), while the pre-types determine the introduction and elimination forms of values. This separation facilitates implicit copying and discarding for unrestricted types, yielding a less burdensome programming model.

Type qualifiers, however, have more complex interactions with polymorphism. To retain both soundness and expressivity, one is led to introduce quantification over qualifiers, pre-types, and types independently. This quickly leads to large and complex types; for instance, the type of plus has five qualifiers:

plus : 
$$(q_1 \text{ Int}) \xrightarrow{q_4} (q_2 \text{ Int}) \xrightarrow{q_5} (q_3 \text{ Int})$$

The relationships among such qualifiers are often nontrivial (*e.g.*  $q_5$  should be lin if  $q_1$  is), which can be captured (at the expense of additional complexity) by qualifier-level bounded quantification.<sup>2</sup> Qualifiers thus ease the use of unrestricted types but are too unwieldy for a polymorphic source language—indeed, others have argued against qualifiers even for intermediate languages [34].

Our use of kinds in lieu of ! or type qualifiers strikes a good balance on these issues. As with qualifiers, programming with unrestricted types is natural; as with !, polymorphism remains simple. Subkinding also plays well with base types: 3 has type lnt, which has kind  $\star$  (and, by subsumption,  $\circ$ ), while the type of plus is the simple lnt  $\stackrel{*}{\rightarrow}$  lnt. We thus have flexible polymorphic types without the need for bounded quantification or other complexities of subtyping in a higher-order setting.

Closest in design to System  $F^{\circ}$  is probably Ahmed, Fluet, and Morrisett's language for substructural state [1], though they use qualifiers as described above. Due to their focus on aliasing in mutable state, their language does not admit subtyping, which would be the analog in their setting of our subkinds. Our subkind relation  $* \leq \circ$  agrees with the interpretation of linearity as related to *usage*; it does not reflect linearity as *alias-freedom*, in which linear types are analogous to uniqueness types [17, 28]. Nevertheless, we show in Section 3.1 that matters of aliasing can indeed be addressed in System  $F^{\circ}$  given an appropriate representation of references.

Ahmed et al. also admit *affine types*, for which only weakening is permitted, and *relevant types*, which allow only contraction. Such concepts would fit well with our subkinding relation, but would increase the complexity of typing context management. Other systems have considered notions of usage much more finegrained than ours [19], but the types in such systems can quickly become overwhelming if exposed to the programmer. As we are incorporating the full power of System F, exposing at least some types is unavoidable; we also believe that System  $F^{\circ}$  convincingly demonstrates that just the simple distinction between linear and unrestricted types has much to offer.

# 1.2 System F° by example: Types for filehandles

Linearity lets us specify a filesystem interface that requires a filehandle to be closed exactly once and forbids its use thereafter. A first approximation<sup>3</sup> for an idealized linear filesystem might be

FileHandle : o

readLine : FileHandle  $\stackrel{\star}{\rightarrow}$  (String, FileHandle)

Here open, read, write, and close are intended to be primitive operations over the linear (note the kind ascription) type FileHandle, while readLine is one of many library functions defined to make file access more convenient. The  $\star$  decorating the arrow in a type like open : String  $\stackrel{\star}{\rightarrow}$  FileHandle indicates that open is an unrestricted function, which may be used more than once; each time open is invoked, it will return a new FileHandle value that must be used linearly. Unrestricted functions are also free to take arguments of linear type, as can be seen in the other operations. Similarly, a linear function of type  $\tau_1 \stackrel{\circ}{\rightarrow} \tau_2$  should be invoked exactly once, but this has no bearing on the kind of either  $\tau_1$  or  $\tau_2$ .

Because operations like read and write consume a FileHandle as input and return a linear FileHandle as output, and because such values cannot be duplicated, client programs are forced to sequence the calls to these functions, fixing their order of evaluation. Since FileHandle values cannot be discarded, the program—unless its overall type indicates that it contains a filehandle—must eventually use close to dispose of any FileHandles it has created. Linearity ensures that no aliased or duplicated FileHandles representing closed files remain to be improperly accessed. Clients of this interface are thus constrained, after opening a file, to access that file according to the regular protocol (read|write)\*close.

Unfortunately, today's operating systems do not understand linearity and instead provide unrestricted interfaces, which make no guarantees about correct filehandle usage:

UnsafeFH : \*

unsafeOpen	:	$String \xrightarrow{\star} UnsafeFH$
unsafeRead	:	$UnsafeFH\overset{\star}{\to}Char$
unsafeWrite	:	$Char \xrightarrow{\star} UnsafeFH \xrightarrow{\star} Unit$
unsafeClose	:	$UnsafeFH \xrightarrow{\star} Unit$

System  $F^{\circ}$  makes it easy to create a safe interface protecting the above from misuse. First, for  $\alpha$  of kind  $\circ$ —our abstract representation of an actual filehandle—we define a record of safe file operations:

$$\mathsf{File}(\alpha) = \{ \mathsf{read} : \alpha \stackrel{\star}{\to} (\mathsf{Char}, \alpha), \\ \mathsf{write} : \mathsf{Char} \stackrel{\star}{\to} \alpha \stackrel{\star}{\to} \alpha, \\ \mathsf{close} : \alpha \stackrel{\star}{\to} \mathsf{Unit} \}$$

<sup>&</sup>lt;sup>2</sup> Making Int a proper type rather than a pre-type simplifies the type of plus but prohibits certain polymorphic functions from accepting integers.

<sup>&</sup>lt;sup>3</sup> This example uses roughly the same interface given by DeLine and Fähndrich to motivate Vault [11] and discussed by Kiselyov and Shan [20] as an alternative to their filehandle regions.

$rac{\kappa}{ au}$	::= ::=	$ \begin{array}{c c} \star & \circ \\ \alpha & \tau \xrightarrow{\kappa} \tau & \forall \alpha : \kappa. \ \tau \end{array} $	kinds types
$e \\ v$	::= ::=	$\begin{array}{c c} x \mid \lambda^{\kappa} x{:}\tau. \; e \mid e \; e \mid \Lambda \alpha{:}\kappa. \; v \mid e \; [\tau] \\ \lambda^{\kappa} x{:}\tau. \; e \mid \Lambda \alpha{:}\kappa. \; v \end{array}$	expressions values
$\Gamma \Delta$	::= ::=	$ \begin{array}{c c} \cdot & \Gamma, \alpha {:} \kappa & \Gamma, x {:} \tau & \textit{unrestricted typ} \\ \cdot & \Delta, x {:} \tau & \textit{linear typ} \end{array} $	oing contexts oing contexts

riguic I. System r	Figure	1.	System	$\mathrm{F}^\circ$
--------------------	--------	----	--------	--------------------

Now we can define open to return both the hidden filehandle and its associated operations:

Note that, while the type UnsafeFH is unrestricted within the scope of open, the outside world sees its occurrences at the existentially bound linear type variable  $\alpha$ . If open treats filehandles correctly, then any use of an existential package created by open must treat them correctly as well.

We no longer read, write, and close as separate functions, but library functions like readLine, are still useful to have. We could simply replace the FileHandles in the first type proposed for readLine with open's return type,  $\exists \alpha : \circ$ . (File( $\alpha$ ),  $\alpha$ ), but a smarter choice is

readLine : 
$$\forall \alpha : \circ$$
. File $(\alpha) \xrightarrow{\star} \alpha \xrightarrow{\star} (\mathsf{String}, \alpha)$ 

This makes clear that  $File(\alpha)$  contains only the (unrestricted) file operations, not the filehandle itself, and by separating out **unpacks** and **packs** from calls to both primitive operations and library functions, our types can reflect the fact that the filehandle returned by readLine is the same one that it was given. Writing our functions this way allows for useful type coercions; for example, suppose we are also able to open files in read-only mode, resulting in restricted existential packages of the form

$$\mathsf{ROFile}(\alpha) = \{ \mathsf{read} : \alpha \stackrel{\star}{\to} (\mathsf{Char}, \alpha) \\ \mathsf{close} : \alpha \stackrel{\star}{\to} \mathsf{Unit} \}$$

Many functions, including readLine, may be defined over this weaker interface, and we can always construct a record of type  $\mathsf{ROFile}(\alpha)$  out of a record of type  $\mathsf{File}(\alpha)$  to allow a read-write filehandle to be treated as though it were read-only.

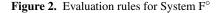
Of course, such a filesystem interface is of little use if it is too cumbersome for everyday programming. In Section 5.2 we discuss modest compiler support that ensures that this is not the case.

# **2.** System $F^{\circ}$ Defined

The syntax of System  $F^{\circ}$  is given in Figure 1; Figure 2 shows its call-by-value (and kind-agnostic) operational semantics, which is completely standard. Type variables are annotated with their kinds when bound (by  $\Lambda$  in expressions or by  $\forall$  in types); kinds also appear in functions ( $\lambda^{\kappa}x:\tau$ . e) and function types ( $\tau_1 \xrightarrow{\kappa} \tau_2$ ). F<sup> $\circ$ </sup> adopts the *value restriction* [35], permitting type abstraction only over values, for reasons that will become clear later.

Typing and kinding rules for System  $F^{\circ}$ , along with auxiliary judgments, are given in Figure 3. Our typing rules take a linear

$$\begin{split} & [\text{E-APPLAM}] \ (\lambda^{\kappa} x : \tau. \ e) \ v \longrightarrow \ \{x \mapsto v\} e \\ & [\text{E-TAPPLAM}] \ (\Lambda \alpha : \kappa. \ v) \ [\tau] \longrightarrow \ \{\alpha \mapsto \tau\} v \\ & [\text{E-APP1}] \ \frac{e_1 \longrightarrow e_1'}{e_1 \ e_2 \longrightarrow e_1' \ e_2} \qquad & [\text{E-APP2}] \ \frac{e \longrightarrow e'}{v \ e \longrightarrow v \ e'} \\ & [\text{E-TAPP}] \ \frac{e \longrightarrow e'}{e \ [\tau] \longrightarrow e' \ [\tau]} \end{split}$$



typing context  $\Delta$ , binding only term variables, in addition to the standard unrestricted context  $\Gamma$ , which binds both type and term variables; following Barber's DILL [3], this greatly simplifies the mechanization of our soundness proofs. At the kind level, note that the rule K-ARROW gives  $\tau_1 \stackrel{\kappa}{\to} \tau_2$  the kind  $\kappa$  regardless of the kinds of  $\tau_1$  and  $\tau_2$ . As described in Section 1.2, this decouples the notion of a linear function (i.e., of type  $\tau_1 \stackrel{\circ}{\to} \tau_2$ ), which must be used exactly once, from a function that takes a linear argument (i.e., where  $\tau_1$  has kind  $\circ$ ), which must use its argument exactly once. Rule K-ALL, by contrast, simply gives  $\forall \alpha: \kappa. \tau$  the same kind as  $\tau$ —this is a design choice made in the interests of keeping F° simple, as little seems to be gained by allowing these kinds to differ; this choice is also compatible with a type-erasure interpretation.

Type application must check that the kind of the supplied type argument is compatible with the kind of the variable for which it will be substituted. As mentioned in Section 1, we see linearity as constraining the permitted usage of a variable, which means that it should always be safe to replace a linear type variable with an unrestricted type; we thus allow subkinding via the rule K-SUB. Subsumption turns out to be key in many useful examples: for instance, it is what allowed the unrestricted type UnsafeFH to be protected by a linear interface.

Weakening—neglecting to use a variable—and contraction using a variable more than once—should only be possible with variables of unrestricted type. Rather than add explicit contraction and weakening rules, we have built these properties into the rules that require them. The separation of linear and unrestricted typing contexts makes this fairly straightforward; rules T-LVAR and T-UVAR permit weakening by allowing an arbitrary  $\Gamma$  at the leaves of typing derivations, while rule T-APP duplicates  $\Gamma$  but splits  $\Delta$ via the U relation. We can thus prove:<sup>4</sup>

**Lemma 1** (Weakening). If  $\Gamma_1, \Gamma_2; \Delta \vdash e : \tau'$  and  $\Gamma_1 \vdash \tau : \star$ , then  $\Gamma_1, x; \tau, \Gamma_2; \Delta \vdash e : \tau'$ .

**Lemma 2** (Contraction). If  $\Gamma_1, x:\tau, y:\tau, \Gamma_2; \Delta \vdash e : \tau'$  and  $\Gamma_1 \vdash \tau : \star$ , then  $\Gamma_1, x:\tau, \Gamma_2; \Delta \vdash \{y \mapsto x\}e : \tau'$ .

Linear variables must not inadvertently be captured by unrestricted function closures. To this end, rule T-LAM constrains its linear context  $\Delta$  according to the  $\lambda$ 's kind annotation  $\kappa$ :  $\Delta$  must be empty if  $\kappa$  is  $\star$ . We thus ensure that an unrestricted function cannot capture linear variables in its closure, even though it may well take an argument of linear type.

Only unrestricted function arguments should be placed in  $\Gamma$ , but subkinding allows any type to be considered as linear and hence any expression variable to be bound in  $\Delta$ —a fact which proves crucial in the proof of preservation. We write this potentially nondeterministic context extension as  $[\Gamma; \Delta], x:\tau \supseteq \Gamma'; \Delta'$ .

<sup>&</sup>lt;sup>4</sup> The lemmas in this section of the paper have all been proved in Coq; the source scripts of our proofs are available from the last author's web pages.

$$\begin{bmatrix} \mathbf{K} - \mathbf{SUB} \end{bmatrix} \frac{\Gamma \vdash \tau : \star}{\Gamma \vdash \tau : \circ} \qquad \begin{bmatrix} \mathbf{K} - \mathbf{ARR} \end{bmatrix} \frac{\Gamma \vdash \tau_1 : \kappa_1}{\Gamma \vdash \tau_1 \stackrel{\kappa}{\to} \tau_2 : \kappa} \qquad \begin{bmatrix} \mathbf{K} - \mathbf{TVAR} \end{bmatrix} \frac{\alpha : \kappa \in \Gamma}{\Gamma \vdash \alpha : \kappa} \qquad \begin{bmatrix} \mathbf{K} - \mathbf{ALL} \end{bmatrix} \frac{\Gamma, \alpha : \kappa \vdash \tau : \kappa'}{\Gamma \vdash \forall \alpha : \kappa, \tau : \kappa'} \stackrel{\alpha \notin \Gamma}{\Gamma \vdash \forall \alpha : \kappa, \tau : \kappa'}$$

$$\begin{bmatrix} \mathbf{U} - \mathbf{EMPTY} \end{bmatrix} \cdot \Downarrow \cdot = \cdot \qquad \begin{bmatrix} \mathbf{U} - \mathbf{LEFT} \end{bmatrix} \frac{\Delta_1 \Downarrow \Delta_2 = \Delta}{\Delta_1, x : \tau \Downarrow \Delta_2 = \Delta, x : \tau} \qquad \begin{bmatrix} \mathbf{U} - \mathbf{RIGHT} \end{bmatrix} \frac{\Delta_1 \Downarrow \Delta_2 = \Delta}{\Delta_1 \sqcup \Delta_2, x : \tau = \Delta, x : \tau}$$

$$\begin{bmatrix} \mathbf{B} - \mathbf{LIN} \end{bmatrix} \frac{\Gamma \vdash \tau : \circ}{[\Gamma; \Delta], x : \tau} \stackrel{\chi \notin \Gamma, \Delta}{\equiv \Gamma; (\Delta, x : \tau)} \qquad \begin{bmatrix} \mathbf{B} - \mathbf{UN} \end{bmatrix} \frac{\Gamma \vdash \tau : \star}{[\Gamma; \Delta], x : \tau \ni (\Gamma, x : \tau); \Delta} \qquad \begin{bmatrix} \mathbf{T} - \mathbf{LVAR} \end{bmatrix} \Gamma; x : \tau \vdash x : \tau \qquad \begin{bmatrix} \mathbf{T} - \mathbf{UVAR} \end{bmatrix} \frac{x : \tau \in \Gamma}{\Gamma; \cdot \vdash x : \tau}$$

$$\begin{bmatrix} \mathbf{T} - \mathbf{LAM} \end{bmatrix} \frac{[\Gamma; \Delta], x : \tau_1 \supseteq \Gamma'; \Delta'}{\Gamma; \Delta \vdash \lambda^{\kappa} x : \tau_1 . e : \tau_1 \stackrel{\kappa}{\to} \tau_2} \qquad \begin{bmatrix} \mathbf{T} - \mathbf{LAM} \end{bmatrix} \frac{\Gamma; \Delta \vdash e_1 : \tau_1 \stackrel{\kappa}{\to} \tau_2}{\Gamma; \Delta \vdash e_1 : e_2 : \tau_2} \qquad \Delta_1 \lor \Delta_2 = \Delta \qquad \begin{bmatrix} \mathbf{T} - \mathbf{TAPP} \end{bmatrix} \frac{\Gamma; \Delta \vdash e : \forall \alpha : \kappa, \tau'}{\Gamma; \Delta \vdash e : \tau_1 : (\alpha \mapsto \tau) \uparrow'}$$

Figure 3. Kinding and typing rules for System F<sup>o</sup>

let $x = e$ in $e'$	≜	$(\lambda^{\circ}x:\tau. e') e$ where e has type $\tau$
$Unit_{e_1; e_2}$	$\underline{\Delta}$	$\forall \alpha : \circ. \alpha \xrightarrow{\star} \alpha$ $\Lambda \alpha : \circ. \lambda^{\star} x : \alpha. x$ $\mathbf{let} \ _{-} = e_1 \ \mathbf{in} \ e_2$
$ \begin{aligned} &(\tau_1,\tau_2)\\ &(,) \end{aligned} \\ \mathbf{let}\ &(x,y)=e\ \mathbf{in}\ e' \end{aligned}$	≜	$ \begin{array}{l} \forall \alpha : \circ. \ (\tau_1 \stackrel{\circ}{\rightarrow} \tau_2 \stackrel{\circ}{\rightarrow} \alpha) \stackrel{\circ}{\rightarrow} \alpha \\ \Lambda \alpha : \circ. \ \Lambda \beta : \circ. \ \lambda^* x : \alpha. \ \lambda^\circ y : \beta. \\ \Lambda \gamma : \circ. \ \lambda^\circ f : \alpha \stackrel{\circ}{\rightarrow} \beta \stackrel{\circ}{\rightarrow} \gamma. \ f \ x \ y \\ e \ [\tau'] \ (\lambda^\circ x : \tau_1. \ \lambda^\circ y : \tau_2. \ e') \\ where \ e' \ has \ type \ \tau' \end{array} $
$\exists \alpha: \kappa. \ \tau'$ pack $\alpha: \kappa = \tau$ in $e: \tau'$ unpack $\alpha, x = e$ in $e'$	≜	$ \begin{array}{l} \forall \beta : \circ. \ (\forall \alpha : \kappa. \ \tau' \xrightarrow{\circ} \beta) \xrightarrow{\kappa'} \beta \\ \textbf{let} \ x = e \ \textbf{in} \\ \Lambda \beta : \circ. \ \lambda^{\kappa'} f : (\forall \alpha : \kappa. \ \tau' \xrightarrow{\circ} \beta). \ f \ [\tau] \ x \\ e \ [\tau'] \ (\Lambda \alpha : \kappa. \ \lambda^{\circ} x : \tau. \ e') \\ where \ e' \ has \ type \ \tau' \end{array} $

It is easy to see that, modulo the value restriction, System  $F^{\circ}$  is an extension of System F. With this in mind, Figure 4 gives several well-known System F encodings that we make use of<sup>5</sup>. Aside from kind annotations, these are all standard; the linear annotations on type variables are for generality—the pairs so encoded are linear, for instance—while those that are on arrows account for captured linear variables in the arguments.

## 2.1 Metatheory of System F°

**Type Soundness** We have verified in Coq that System  $F^{\circ}$  enjoys type safety—a crucial but unsurprising result, given its similarity to System F. As is standard, we define soundness in terms of two properties: progress and preservation. Progress, which states that a closed, well-typed non-value can always take an evaluation step, is no different than in ordinary System F:

**Lemma 3** (Progress). If  $\cdot; \cdot \vdash e : \tau$ , then either e is a value or there exists some e' such that  $e \longrightarrow e'$ .

*Proof.* Induction on typing derivations, completely standard.  $\Box$ 

Preservation, on the other hand, requires a bit of care. As usual, it depends on various substitution lemmas, and we must keep linearity in mind when formulating them.

## Lemma 4 (Substitution).

1. If 
$$\Gamma_1, \alpha:\kappa', \Gamma_2 \vdash \tau:\kappa$$
 and  $\Gamma_1 \vdash \tau':\kappa'$  then  
 $\Gamma_1, \{\alpha \mapsto \tau'\}\Gamma_2 \vdash \{\alpha \mapsto \tau'\}\tau:\kappa.$   
2. If  $\Gamma_1, \alpha:\kappa', \Gamma_2; \Delta \vdash e:\tau$  and  $\Gamma_1 \vdash \tau':\kappa'$  then  
 $\Gamma_1, \{\alpha \mapsto \tau'\}\Gamma_2; \{\alpha \mapsto \tau'\}\Delta \vdash \{\alpha \mapsto \tau'\}e: \{\alpha \mapsto \tau'\}\tau.$   
3. If  $\Gamma_1, x:\tau', \Gamma_2; \Delta \vdash e:\tau$  and  $\Gamma_1; \cdot \vdash e':\tau'$  then  
 $\Gamma_1, \Gamma_2; \Delta \vdash \{x \mapsto e'\}e:\tau.$   
4. If  $\Gamma; \Delta_1, x:\tau', \Delta_2 \vdash e:\tau$  and  $\Gamma; \Delta' \vdash e':\tau'$  then  
 $\Gamma; \Delta_1, \Delta', \Delta_2 \vdash \{x \mapsto e'\}e:\tau.$ 

*Proof.* Each case by induction on the first derivation. The result relies heavily on various strengthening, weakening, and permutation lemmas (with respect to typing, context well-formedness, and the  $\ensuremath{\mathbb{U}}$  relation) regarding the handling of typing contexts.

Note that Substitution (3) does *not* hold if e' is permitted to contain free linear variables. Call-by-value reduction allows us to consider only values, however, and—because we have adopted the value restriction—we can prove that unrestricted values contain no free linear variables:

**Lemma 5.** If  $\Gamma; \Delta \vdash v : \tau$  and  $\Gamma \vdash \tau : \star$  then  $\Delta = \cdot$ .

**Lemma 6** (Preservation). If  $\Gamma; \Delta \vdash e : \tau$  and  $e \longrightarrow e'$ , then  $\Gamma; \Delta \vdash e' : \tau$ .

From preservation and progress, soundness follows naturally:

**Theorem 7** (Type soundness). If  $\cdot; \cdot \vdash e:\tau$ , then it is never the case that  $e \longrightarrow^* e'$  where e' is not a value but cannot step further.

**Strong normalization and parametricity** System F also has other properties of interest: it is *strongly normalizing*—evaluation always eventually reaches a value—and it enjoys relational *parametricity* [24]. We are able to cheat somewhat in proving the former for  $F^{\circ}$  by observing that a well-typed System  $F^{\circ}$  expression becomes a well-typed System F expression upon erasure of kind annotations.

 $<sup>^{5}</sup>$  We also make use of records of unrestricted type; their encodings would generalize that of pairs, but with more  $\star$  annotations.

As the two systems have identical operational behavior, strong normalization follows immediately.

Parametricity for  $F^{\circ}$  cannot be proved by such an erasure, but, as one might hope from the similarity to System F, it is possible to directly adapt the standard logical relations proof with only minor syntactic differences due to our separation of unrestricted and linear contexts. We have thus proved (in Coq), for the standard relation between type substitutions  $\rho_1 \approx \rho_2$ :  $\Gamma$ , relation between term substitutions  $\rho \vdash \gamma_1 \approx \gamma_2$ :  $\Gamma$ ;  $\Delta$ , and computation closure  $C[\![\tau]\!]_{\rho}$ of relations induced by a type  $\tau$ , where  $\rho$  maps type variables to term relations and pairs of types and  $\Gamma$  and  $\Delta$  bind the variables in the domain of the various substitutions:

**Lemma 8** (Parametricity). If  $\Gamma; \Delta \vdash e : \tau, \rho_1 \approx \rho_2 : \Gamma$ , and  $\rho \vdash \gamma_1 \approx \gamma_2 : \Gamma; \Delta$ , then  $(\rho_1(\gamma_1(e)), (\rho_2(\gamma_2(e))) \in C[\![\tau]\!]_{\rho}$ .

Succinctly, this means that an expression e with type  $\tau$ , under appropriate closing substitutions, is related to itself by (the computation closure of) the relation induced logically by  $\tau$ .

Of course, this is only the simplest parametricity result we could provide; it does not take into account the extensions we propose in Section 5, nor does it take advantage of linearity in any way. The interactions between linearity and relational parametricity have been explored by Birkedal et al. [8] and Bierman et al. [7, 6] have explored program equivalences in the presence of !, both of which suggest avenues for future investigation in the context of System  $F^{\circ}$ . Our appeals to parametricity in Section 3.2, however, require nothing beyond what we have proved.

We view the ease with which we can adapt standard results from System F to System  $F^{\circ}$  as a significant benefit of this design. A direct correspondence between these metatheoretic properties and intuitions about what linearity provides is not immediately obvious, however. In Section 4 we will show, by means of elaborated operational semantics, that the restrictions required for our soundness proofs are exactly those needed to satisfy our intuitions.

#### 2.2 Comparison to traditional formulations

In contrast to our approach, linear type systems are more traditionally presented without kinds, assuming linearity by default and using the modality ! to allow unrestricted variables (see, for example [3, 6]). For the polymorphic  $\lambda$ -calculus, this gives us

$$\sigma ::= \alpha \mid \sigma \multimap \sigma \mid \forall \alpha. \sigma \mid !\sigma$$

$$t ::= a \mid x \mid \lambda a:\sigma. t \mid t t \mid \Lambda \alpha. t \mid t [\sigma]$$

$$\mid !t \mid \mathbf{let} !x = t \mathbf{in} t$$

$$\Phi ::= \cdot \mid \Phi, \alpha \mid \Phi, x:\sigma$$

$$\Psi ::= \cdot \mid \Psi, a:\sigma$$

For clarity, we distinguish between linear variables a, bound by  $\lambda$  terms, and non-linear variables x, bound by **let** !. The interesting typing rules concern the ! modality:

$$\frac{\Phi; \cdot \vdash t : \sigma}{\Phi; \cdot \vdash !t : !\sigma}$$

$$\frac{\Phi; \Psi_1 \vdash t_1 : !\sigma_1 \qquad \Phi, x:\sigma_1; \Psi_2 \vdash t_2 : \sigma_2 \qquad \Psi_1 \uplus \Psi_2 = \Psi}{\Phi; \Psi \vdash \mathbf{let} ! x = t_1 \mathbf{in} t_2 : \sigma_2}$$

In other words, the type  $!\sigma$  indicates a term of type  $\sigma$  which uses no linear variables—the same constraint we place on unrestricted functions—and such a term can be captured by the **let** ! operation and subsequently used in an unrestricted fashion. Note, however, that the type  $!\sigma$  itself is still linear, even though terms of that type allow for the introduction of unrestricted assumptions; while it is possible to formulate systems where terms of ! types can be duplicated or discarded directly, naive attempts to do so are unsound for precisely the same reasons that, as discussed in Section 2.1 and Section 4, we require call-by-value reduction and the value restriction in System  $F^{\circ}$ —and sound formulations end up being heavier than those that make this distinction [31, 4].

We can encode the above system in ours easily enough; we first define a translation on types  $[\sigma]$  as

$$\begin{array}{rcl} \llbracket \alpha \rrbracket &=& \alpha \\ \llbracket \sigma_1 \multimap \sigma_2 \rrbracket &=& \llbracket \sigma_1 \rrbracket \overset{\circ}{\to} \llbracket \sigma_2 \rrbracket \\ \llbracket \forall \alpha. \ \sigma \rrbracket &=& \forall \alpha{:}\circ. \llbracket \sigma \rrbracket \\ \llbracket ! \sigma \rrbracket &=& \operatorname{Unit} \overset{\star}{\to} \llbracket \sigma \rrbracket \end{array}$$

The corresponding translation on terms is straightforward. The only interesting cases involve unrestricted variables and the ! modality:

$$\begin{array}{rcl} \|x\| &=& x \text{ unit} \\ \|!t\| &=& \lambda^* \_: \text{Unit.} \|t\| \\ \|\text{let } !x = t_1 \text{ in } t_2\| &=& (\lambda^\circ x : \text{Unit} \stackrel{\star}{\to} \|\sigma_1\| . \|t_2\|) \|t_1\| \\ & \text{where } t_1 \text{ has type } !\sigma_1 \end{array}$$

Here we treat terms !t as suspended computations which may be evaluated more than once—or not at all—which is standard.

Translating in the other direction is much less straightforward and space constraints preclude us from including the translation here. The translation on types is kind-directed and, at the term level, the insertion of ! and **let** ! operations is not trivial: a function that takes an unrestricted type must now take a ! type and bind it so that the variable need not appear linearly, but such arguments must be repackaged under ! in order to be passed to any subsequent functions. Polymorphic types and expressions also need care—for example, there are four cases needed to translate type instantiation, one for each combination of linear/unrestricted for the polymorphic term and its type argument.

We see this asymmetry in the translations as evidence of System  $F^{\circ}$ 's expressive power and its ability to handle unrestricted terms (the bulk of any most programs) in a concise way, justifying our claim that it provides linearity in a lightweight fashion.

# 3. Examples

To demonstrate System  $F^{\circ}$ 's applicability, we now turn to two categories of examples. First, we demonstrate that, while System  $F^{\circ}$ does not build in notions of references and aliasing, protocols defining correct memory management can indeed be enforced by System  $F^{\circ}$  types. Second, we prove that *any* protocol expressible as a finite automaton has a corresponding  $F^{\circ}$  type; in addition to establishing that a rather large class of protocols can be encoded in our type system, this proof also highlights intuitions about linearity that still need formalization, which we will tackle in Section 4. (It is easy to see, however, that the regular languages are not an upper limit on System  $F^{\circ}$ 's expressivity; the classic non-regular parenthesis matching example has an obvious protocol type.)

#### 3.1 Reference cells

The filesystem interface in Section 1.2 can, under an appropriate renaming and abstraction over the contents type, also be seen as an interface for linear reference cells:<sup>6</sup>

$$\begin{aligned} \mathsf{Ref}[\tau](\alpha) &= \{ \mathsf{set} : \tau \xrightarrow{\star} \alpha \xrightarrow{\star} \alpha, \\ \mathsf{get} : \alpha \xrightarrow{\star} (\tau, \alpha), \\ \mathsf{free} : \alpha \xrightarrow{\star} \mathsf{Unit} \} \\ \mathsf{mkRef} &: \forall \beta : \star, \beta \xrightarrow{\star} \exists \alpha : \circ. \ (\alpha, \mathsf{Ref}[\beta](\alpha)) \end{aligned}$$

<sup>&</sup>lt;sup>6</sup> Here and elsewhere, we separate the "type parameters" like  $\tau$  from the linear "state parameters" like  $\alpha$  using the notation  $[\tau](\alpha)$ .

On its own this is not particularly interesting, as a linear reference cell simply encodes the practice of threading a value through a program. Indeed, as with Haskell's State monad, we could instantiate mkRef such that  $\alpha$  is  $\tau$ . Instead, however, let us consider variations on Ref that make more sense as safe interfaces wrapping true, potentially unsafe reference cells, much as File in Section 1.2 protected the primitive, unchecked filehandle calls.

While the above Ref could be such a safe interface, a more obvious one—which itself requires no linearity—is the type GCRef:

$$\begin{aligned} \mathsf{GCRef}[\tau](\alpha) \ = \ \{ \ \mathsf{set} : \tau \xrightarrow{\star} \alpha \xrightarrow{\star} \alpha, \\ \mathsf{get} : \alpha \xrightarrow{\star} (\tau, \alpha) \ \} \end{aligned}$$

mkGCRef : 
$$\forall \beta : \star, \beta \xrightarrow{\star} \exists \alpha : \star, (\alpha, \mathsf{GCRef}[\beta](\alpha))$$

The operations given by  $\mathsf{GCRef}[\tau](\alpha)$  are, of course, those of a garbage-collected reference cell of type  $\tau$ . By hiding the reference type behind  $\alpha$  we ensure that such garbage-collected references cannot be freed, and, in this case,  $\alpha$  can be unrestricted.

In System  $F^{\circ}$ , however, we can also define references that begin their lives as linear (and manually managed) but later are put under the garbage collector's control.  $\text{Ref}[\tau](\alpha)$  simply needs an additional function to serve as the appropriate coercion:

$$gc: \alpha \xrightarrow{\star} \exists \beta: \star. \ (\beta, \mathsf{GCRef}[\tau](\beta))$$

By consuming and not returning  $\alpha$ , gc prevents free from being called on the now garbage-collected (but unrestricted) reference. We thus have a coercion from alias-free to potentially aliased pointers, a fact that, had we conflated linearity with alias-freedom, would run counter to our subkinding relation of  $\star \leq \circ$ . (Whether gc needs to do any work at run time depends on the implementation of the memory management system.)

We can take other approaches to memory management as well. For instance, an intermediate point between a strictly linear and a garbage collected reference is a reference that must be explicitly aliased, and where aliases must be explicitly discarded. We can define this easily enough:

$$\begin{aligned} \mathsf{RCRef}[\tau](\alpha) \ = \ \{ \ \mathsf{set} : \tau \xrightarrow{\star} \alpha \xrightarrow{\star} \alpha, \\ \mathsf{get} : \alpha \xrightarrow{\star} (\tau, \alpha), \\ \mathsf{alias} : \alpha \xrightarrow{\star} (\alpha, \alpha), \\ \mathsf{drop} : \alpha \xrightarrow{\star} \mathsf{Unit} \ \} \end{aligned}$$

mkRCRef : 
$$\forall \beta : \star, \beta \stackrel{\star}{\to} \exists \alpha : \circ, (\alpha, \mathsf{RCRef}[\beta](\alpha))$$

A straightforward implementation of mkRCRef could return a pair of the desired reference cell and an additional cell to act as a counter, along with alias and drop operations that adjust this counter. Both the primary cell and this counter could safely be freed when the count reaches zero.

However, while our access capability is still linear with RCRef, we can never be certain that we possess the only reference to cell in question. To remedy this, as is often done in capability calculi [2, 9], we can give our cells both an exclusive capability  $\alpha$  and a shared capability  $\beta$ , with possession of the exclusive capability implying that no outstanding copies of the shared capability remain. If, for example, the cell contents should not be altered if any aliases exist, we can use

$$\begin{array}{ll} \mathsf{ShareRef}[\tau](\alpha,\beta) \ = \ \{ \ \mathsf{set} : \tau \stackrel{\star}{\to} \alpha \stackrel{\star}{\to} \alpha, \\ \mathsf{getE} : \alpha \stackrel{\star}{\to} (\tau, \alpha), & \mathsf{getS} : \beta \stackrel{\star}{\to} (\tau, \beta), \\ \mathsf{share} : \alpha \stackrel{\star}{\to} \beta, & \mathsf{claim} : \beta \stackrel{\star}{\to} \alpha \oplus \beta, \\ \mathsf{alias} : \beta \stackrel{\star}{\to} (\beta, \beta), & \mathsf{drop} : \beta \stackrel{\star}{\to} \mathsf{Unit}, \\ \mathsf{free} : \alpha \stackrel{\star}{\to} \mathsf{Unit} \ \} \end{array}$$

mkShareRef :  $\forall \gamma : \star. \gamma \xrightarrow{\star} \exists \alpha : \circ. \exists \beta : \circ. (\alpha, \mathsf{ShareRef}[\gamma](\alpha, \beta))$ 

Here  $\alpha \oplus \beta$  is a standard sum type; as mentioned in Section 2, these are not encodable in System F° as presented so far, but they are an easy extension and are discussed in Section 5.1. The claim function exchanges a shared  $\beta$  for an exclusive  $\alpha$  when the underlying counter indicates that only one alias exists; in all other cases it simply returns the supplied shared capability.

We can do still more if we extend System  $F^{\circ}$  with quantification over higher kinds, an extension which meshes well with Section 5.1's datatypes. For instance, we can define linear references that support *strong updates*—that is, updates that change the type contained in the reference cell—by abstracting the type of the access capability over the type of the contents, giving it kind  $\star \Rightarrow \circ$ :

$$\begin{aligned} \mathsf{SURef}[\tau](\alpha) \ = \ \{ \ \mathsf{set} : \forall \beta : \star. \ \forall \gamma : \star. \ \gamma \xrightarrow{\star} \alpha \ \beta \xrightarrow{\star} \alpha \ \gamma, \\ \mathsf{get} : \forall \beta : \star. \ \alpha \ \beta \xrightarrow{\star} (\beta, \alpha \ \beta), \\ \mathsf{free} : \forall \beta : \star. \ \alpha \ \beta \xrightarrow{\star} \mathsf{Unit} \ \end{aligned} \end{aligned}$$

mkSURef :  $\forall \beta : \star, \beta \xrightarrow{\star} \exists \alpha : \star \Rightarrow \circ. (\alpha \ \beta, \mathsf{SURef}[\beta](\alpha))$ 

Operation records of type SURef[ $\tau$ ]( $\alpha$ ) can easily be coerced to type Ref[ $\tau$ ]( $\alpha \tau$ ) by partial application of member functions. Augmenting SURef along the lines of ShareRef could further allow for sharable reference cells that support strong updates only when they are not shared, a fairly sophisticated feature.

### 3.2 Regular protocols

Apart from memory cells and file operations, what other protocols can System  $F^{\circ}$  enforce? Rather than present more individual examples, we will show how *any* protocol expressible as a regular language can be written in  $F^{\circ}$ .

We take the standard definition of a DFA as a tuple  $M = (Q, \Sigma, \delta, q_0, F)$ , where Q is a finite set of states,  $\Sigma$  is a finite set of alphabet symbols (in our context, protocol actions),  $\delta$  is a subset of  $\Sigma \times Q \times Q$  (that is, a ternary relation among actions, current states, and next states),  $q_0$  is a distinguished initial state, and F is a set of final states. The file access protocol from Section 1.2, for instance, can be thought of as a very simple automaton with states Open and Closed and an alphabet consisting of read, write, and close, with open as the creator of such an automaton.

For ease of notation, we extend our metavariable conventions to allow q and r as type variables. Taking  $Q = \{q_0, \ldots, q_n\}$ , we can now define the automaton type for the DFA M as

$$\tau_{M} = \exists q_{0}:\circ, \dots, q_{n}:\circ. (q_{0}, \{a\_takes\_q\_to\_q': q \stackrel{\star}{\to} q' \\ for \ every \ (a, q, q') \in \delta \\ \vdots \\ done\_at\_q: q \stackrel{\star}{\to} Unit \\ for \ every \ q \in F \})$$

Such type is trivially inhabited: one can supply an existential package of this type where all type variables are bound to Unit.

We say that the evaluation to a value of an expression e containing a subexpression of type  $\tau_M$  reflects a word  $w = a_0 \dots a_k$  if the sequence of record fields used is exactly

$$a_0\_takes\_q_0\_to\_r_0, \dots a_k\_takes\_r_{k-1}\_to\_r_k$$

for some states  $r_0$  through  $r_k$ . To prove that  $\tau_M$  is an accurate representation of M, we need to show, first, that each  $w = a_0 \dots a_k$  accepted by M corresponds to an expression of type  $\tau_M \stackrel{\star}{\to} \text{Unit}$  that reflects w and, second, that each f of type  $\tau_M \stackrel{\star}{\to} \text{Unit}$  reflects some w accepted by M. The former is fairly straightforward:

**Lemma 9.** For any  $w = a_0 \dots a_k$  accepted by M, there exists f such that  $\vdash f : \tau_M \stackrel{*}{\to} \text{Unit and for any } n : \tau_M, f n$  reflects w.

*Proof.* Recall that, if  $w = a_0 \dots a_k$  is accepted by M, we have a trace of M on w of the form  $q_0a_0r_0\dots r_{k-1}a_kr_k$ , where  $(a_0, q_0, r_0) \in \delta$ ,  $(a_i, r_{i-1}, r_i) \in \delta$ , and  $r_k \in F$ . Knowing this and the definition of  $\tau_M$ , we can construct

$$f = \lambda^* n: \tau_M. \text{ unpack } (q_0, \dots, q_n, p) = n \text{ in}$$
  
let  $(start, ops) = p$  in  
let  $s_0 = ops.a_0\_takes\_q_0\_to\_r_0 \text{ start in}$   
 $\vdots$   
let  $s_k = ops.a_k\_takes\_r_{k-1}\_to\_r_k \text{ } s_{k-1}$  in  
 $x.done\_at\_r_k \text{ } s_k$ 

This clearly reflects w, and, if w is indeed accepted by M, it will typecheck successfully.

The other direction depends on arguments from parametricity (which we have already proved) and the nature of linearity (which we will formalize and prove in Section 4).

**Lemma 10.** If  $\cdot \vdash f : \tau_M \stackrel{\star}{\to} \text{Unit}$ , then there exists some w such that M accepts w and f n reflects w for any n where  $\cdot \vdash n : \tau_M$ .

*Proof.* Because  $\tau_M$  is linear, f must eliminate n before it can return anything of type Unit. This will require first unpacking the existential and then pattern-matching against the pair, leaving the linear *start* and unrestricted *ops*.

If  $q_0 \in F$ , then f might immediately apply  $ops.done\_at\_q_0$ . This reflects the empty word  $\epsilon$ , and, indeed, if  $q_0 \in F$ , then M accepts  $\epsilon$ . Alternatively, regardless of whether  $q_0 \in F$ , start can become some other state type (or even  $q_0$  again) by repeated applications of the  $ops.a\_takes\_q\_to\_q'$  functions. If, after such applications, the result can be eliminated by some  $ops.done\_at\_q_j$ , then it must be that  $q_j \in F$ . Moreover, because of the construction of  $\tau_M$ , each  $ops.a\_takes\_q\_to\_q'$  application must represent a valid transition of  $\delta$ . Thus we are reflecting some w accepted by M.  $\Box$ 

In the above proof, familiar intuitions about parametricity justify the arguments that f behaves the same way for any argument of type  $\tau_M$  and that expressions whose types are variables can only be used in certain ways. The assurance what w is properly reflected, however, also depends on the connection between the static property of linear typing and the behavior of expressions at runtime. We would like type soundness to guarantee that our intuitions about the behavior of linear expressions are valid; in the next section we prove that this is indeed the case.

# 4. Linear Semantics

The preceding section makes clear that we have yet to formalize all of our intuitions about what linearity means for run-time behavior. This is not as straightforward as it might appear, because some of our intuitions about linearity turn out to be misleading. In essence, this is because linearity restricts *variables* of linear type, while we want to reason, at runtime, about the behavior of *expressions*.

For instance, linearity does not guarantee that subexpressions or even values—of linear type will be used exactly once. Nor would we want it to: recall that the encoding of **let** expressions in Figure 4 involves a linear function, and we certainly want to allow for **let** expressions of unrestricted type.

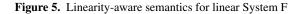
However, there are examples which, in a call-by-name or callby-need setting, do behave in ways we'd like to prevent. For example, if we have unpacked a reference cell from Section 3.1 as (r, ops), we can dispose of it and continue with the expression evia

$$(\lambda^{\star}$$
\_:Unit. e) (ops.free r)

But this call to free will only be evaluated in call-by-value—since its result is not used, other evaluation strategies would simply  $\begin{bmatrix} \text{L-APPLAM} \end{bmatrix} (\lambda^{\kappa} x : \tau. e) v \xrightarrow{(v:\tau)} \{x \mapsto (v:\tau)\}e$  $\begin{bmatrix} \text{L-TAPPLAM} \end{bmatrix} \Lambda \alpha : \kappa. v [\tau] \xrightarrow{\epsilon} \{\alpha \mapsto (\tau:\kappa)\}v$  $\begin{bmatrix} \text{L-TAG} \end{bmatrix} (v:\tau) \xrightarrow{(v:\tau)} v \qquad \begin{bmatrix} \text{L-TAPP} \end{bmatrix} \frac{e \frac{C}{D} \longrightarrow e'}{e [\tau] \frac{C}{D} \longrightarrow e' [\tau]}$ 

$$\begin{bmatrix} \text{L-APP1} \end{bmatrix} \xrightarrow{e_1 \ D} e_2' \xrightarrow{C} e_1' = e_1 \ e_2 \ D \xrightarrow{C} e_1' = e_2 \ e_2' = e_1' = e_2' \\ \begin{bmatrix} \text{L-APP2} \end{bmatrix} \xrightarrow{e_1 \ D} e_2' \xrightarrow{C} e_2' = e_2' \\ \hline v \ e_2 \ D \xrightarrow{C} v \ e_1' = e_2' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C} v \ e_1' = e_2' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C} v \ e_1' = e_2' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C} v \ e_1' = e_2' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C} v \ e_1' = e_2' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C} v \ e_1' = e_2' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C} v \ e_1' = e_2' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C} v \ e_1' = e_2' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C} v \ e_1' = e_2' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C} v \ e_1' = e_2' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C} v \ e_1' = e_2' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C} v \ e_1' = e_2' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C} v \ e_1' = e_2' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C} F \xrightarrow{C} v \ e_1' = e_2' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C} F \xrightarrow{C} v \ e_1' = e_2' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C} F \xrightarrow{C} F \xrightarrow{C} v \ e_1' \\ \hline r \ F \xrightarrow{C} F \xrightarrow{C}$$

$$[\text{T-TTAG1}] \frac{\Gamma; \Delta \vdash e : \tau}{\Gamma; \Delta \vdash e : (\tau:\kappa)} \qquad [\text{T-TTAG2}] \frac{\Gamma; \Delta \vdash e : (\tau:\kappa)}{\Gamma; \Delta \vdash e : \tau}$$



discard the subexpression, leaking memory, even though this is precisely the sort of error we hoped to rule out. Similar concerns for uniqueness types are known in the Clean community [28].

Call-by-value, with its assurance that function arguments are evaluated exactly once, seems well-suited to avoiding these problems, and indeed, in a call-by-value setting—with the addition of the value restriction—we are not able to construct such problematic examples. In other words, the restrictions already in place to ensure soundness in Section 2.1 are everything we need to support our intuitions about linearity. To prove that this is so, the remainder of this section provides an extended operational semantics with which we can reason about linearity at runtime.

#### 4.1 Annotated F°

What does it mean to "use" a subexpression? We might mean by this that an expression is evaluated to a value, that it is passed to a function, or that it is applied to an argument. In order to understand the run-time effects of linearity, however, we focus on the use of a value, which we divide into two phases:

- 1. A value is *conscripted* when it is substituted into a function body in place of a variable.
- 2. A previously conscripted value is *discharged* when it is used as a value to allow evaluation to continue—in a context where, were it replaced by a free variable, evaluation would be stuck.<sup>7</sup>

From this linearity's contribution is clear: in the course of evaluation every conscripted linear value not contained within the final result will be discharged exactly once. To formalize this, we define a linearity-aware operational semantics on expressions extended with tagged values (representing the arguments passed to functions) and tagged types (representing type arguments):

$$\begin{aligned} \tau & ::= & \dots & | & (\tau:\kappa) \\ e & ::= & \dots & | & (v:\tau) \\ C, D & ::= & \epsilon & | & (v:\tau) & | & C, C \end{aligned}$$

Our new semantics is given in Figure 5: we write  $e \stackrel{C}{_{D}} \longrightarrow e'$  if e steps to e' while conscripting the sequence of tagged values C

<sup>&</sup>lt;sup>7</sup> As variables are not values, a free variable on *either* side of an application results in a stuck expression by the semantics in Figure 2.

and discharging the sequence *D*. Values are conscripted in rule L-APPLAM, and, in a slight but critical departure from the usual operational semantics, it is the argument value tagged with the type expected by the function that is substituted into the function body. Similarly, while L-TAPPLAM does not conscript, it does tag its type argument with its expected kind, thus keeping a record of the kind at which said argument will be considered within its body.

Since a tagged value  $(v:\tau)$  is not itself a value, we have the rule L-TAG, which both discharges the value and removes its tag. Intuitively this means that the argument to a previous function application is either about to be used again—regardless of which side of an application it is on—or that it is being returned as the final result of the computation. Of course, discharging may be delayed for some time, as tagged values may linger beneath function closures. As with the presentation in Section 2, evaluation does not depend on typing—types and kinds are simply recorded—and, in particular, our new operational semantics does not treat expressions of non-linear type differently from those of linear type.

The kinding and typing rules in Figure 5 are straightforward; the ability to conclude  $e : \tau$  from  $e : (\tau:\kappa)$  regardless of  $\kappa$  makes kind annotations completely transparent, but the same is not true of type annotations. The soundness results in Section 2.1 carry over to the extended system without complication.

As an example, consider an application of the linear polymorphic identity function  $\Lambda \alpha :\circ. \lambda^{\circ} x :\alpha. x$  to the natural number 42 of type Nat. This expression steps as follows:

$$\begin{array}{ccc} & (\Lambda \alpha : \circ. \ \lambda^{\circ} x : \alpha . \ x) \ [\mathsf{Nat}] \ 42 \\ & \stackrel{\epsilon}{\underset{\epsilon}{\longrightarrow}} & (\lambda^{\circ} x : (\mathsf{Nat} : \circ) . \ x) \ 42 & \mathsf{L}\text{-TAPPLAM} \\ & \stackrel{(42 : (\mathsf{Nat} : \circ))}{\underset{\epsilon}{\longrightarrow}} & (42 : (\mathsf{Nat} : \circ)) & \mathsf{L}\text{-APPLAM} \\ & \stackrel{(42 : (\mathsf{Nat} : \circ))}{\underset{\epsilon}{\longrightarrow}} & 42 & \mathsf{L}\text{-TAG} \end{array}$$

Or, taking the obvious tag-concatenating closure  ${}_{D}^{C} \longrightarrow^{*}$ :

$$(\Lambda \alpha : \circ. \lambda^{\circ} x : \alpha. x)$$
 [Nat] 42  $(42:(\text{Nat}; \circ)) \longrightarrow *$  42

In other words, over the course of evaluation, the value 42 was passed to a function and used (in this case, returned as the final result) exactly once, and it was considered at type (Nat: $\circ$ ). Were it to appear more than once in the lower annotation, we would know that linearity had failed us—although 42 is of type Nat, here it is being considered as a linear type because  $\alpha$  was declared to be linear, so after being conscripted at this type it should be discharged exactly once. Were it to instead appear as (42:Nat), however, there should be no such restrictions on its use, as we can indeed write functions from Nat to Nat which use their arguments more than once or ignore them altogether.

This semantics is similar in scope to heap semantics for linear functional languages [26], and indeed we could have taken this approach by tracking, instead of the traces C and D, a single heap H. Such a semantics would perform substitution at application only if the argument type is unrestricted; for linear types, we would simply associate the variable with its argument in the heap, removing it when the variable is used. But, while we could still prove Theorem 12 with such a semantics, we need the ability to track unrestricted arguments—typically not placed in the heap—and to reason about the order of function calls in order to support the reasoning we used in Section 3, which we will return to in Section 4.3.

#### 4.2 Linearity at run time

Before we can formalize our intuitions about linearity, we need some auxiliary definitions. We write  $\mathfrak{T}(e)$  for the multiset of tagged values appearing within e, defined simply as:

$$\begin{array}{rcl} \mathfrak{T}(x) &=& \emptyset & \qquad \mathfrak{T}(\lambda^{\kappa}x:\sigma_{1}.\,e) &=& \mathfrak{T}(e) \\ \mathfrak{T}((v:\tau)) &=& \{(v:\tau)\} \uplus \mathfrak{T}(v) & \qquad \mathfrak{T}(\Lambda\alpha:\kappa.\,v) &=& \mathfrak{T}(v) \\ \mathfrak{T}(e_{1}\,e_{2}) &=& \mathfrak{T}(e_{1}) \uplus \mathfrak{T}(e_{2}) & \qquad \mathfrak{T}(e\left[\tau\right]) &=& \mathfrak{T}(e) \end{array}$$

We write  $\{C\}$  to denote the treatment of a sequence as a multiset, and, if S is a multiset of tagged values, we write  $S \setminus \kappa$  for the subset of S that omits any  $(v:\tau)$  where  $\cdot \vdash \tau : \kappa$ .

We define *proper* expressions as those in which no unrestricted value contains a value as a subexpression tagged with linear type—that is, a well-typed expression e is proper iff, for every value subexpression v in e, if v has some type  $\tau$  and  $\tau$  can be given kind  $\star$ , then  $\mathfrak{T}(v) \setminus \star = \emptyset$ . This property is preserved by evaluation:

**Lemma 11** (Proper expressions). If e is proper and  $e \stackrel{C}{_{D}} \rightarrow e'$ , then e' is also proper.

*Proof.* Values can only become tagged in a subexpression by function application. Because of the value restriction, in order for a linear value to be substituted into an unrestricted value, a linear variable bound at an outer scope would need to appear under an unrestricted  $\lambda$ . Rule T-LAM forbids this.

Thus, as long as we write our source expressions in the language described in Section 2, we will never evaluate to an improper result—this is essentially the runtime version of Lemma 5. If we add new constructs, this property will continue to hold as long as

- constructs under which evaluation can proceed (like conventional tuples) are given unrestricted types only if all of their components are also unrestricted, and
- 2. constructs that suspend computation either contain only values (like  $\Lambda$ ) or only typecheck at an unrestricted type given an empty linear context (like  $\lambda$ ).

(Of course, it also suffices to require that a new construct always be typed linearly, as we do for additive conjunction in Section 5.1.)

We are now equipped to prove our main result, that linearity guarantees a correspondence between values conscripted and discharged at linear type:

**Theorem 12** (Run time linearity). If e is proper and  $e \stackrel{C}{_{D}} \rightarrow e'$ , then  $\mathfrak{T}(e) \setminus \star \uplus \{C\} \setminus \star = \mathfrak{T}(e') \setminus \star \uplus \{D\} \setminus \star$ .

*Proof.* By straightforward induction over the annotated evaluation relation; L-APPLAM is the only interesting case. From T-UVAR and Lemma 11 we know that, unless the argument occurs exactly once in the function body, it will be of unrestricted type and contain no linear tags. We thus preserve our tag balance.

In other words, at runtime, arguments treated linearly are never duplicated or discarded, exactly as we would hope. We can also immediately prove the following corollary, summarizing everything we know about expressions with unrestricted type:

**Corollary 13** (Unrestricted results). If e is proper,  $\vdash e : \tau$ , and  $\cdot \vdash \tau : \star$ , then there exists some v such that  $e \stackrel{C}{_{D}} \longrightarrow * v$ ,  $\mathfrak{T}(e) \setminus \star \uplus \{C\} \setminus \star = \{D\} \setminus \star$ , and  $\mathfrak{T}(v) \setminus \star = \emptyset$ .

*Proof.* Follows directly from soundness, strong normalization, Lemma 11, and Theorem 12.  $\hfill \Box$ 

# 4.3 Applications

To demonstrate how the above applies to the DFA encoding given in Section 3.2, we first examine the simpler encodings of products and existential types from Section 2.

**Products** Our traces treat pairs exactly as we would expect. Given  $e_1 \xrightarrow[D_1]{D_1} * v_1$  and  $e_2 \xrightarrow[D_2]{D_2} * v_2$ , where  $e_1$  has type  $\tau_1$  and  $e_2$  has type  $\tau_2$ , recall that

$$\begin{array}{ll} (e_1, e_2) & = & (\Lambda \alpha {:} \circ . \ \Lambda \beta {:} \circ . \ \lambda^* x {:} \alpha . \ \lambda^\circ y {:} \beta . \ \Lambda \gamma {:} \circ . \\ & & \lambda^\circ f {:} \alpha \xrightarrow{\circ} \beta \xrightarrow{\circ} \gamma . \ f \ x \ y) \ [\tau_1] \ [\tau_2] \ e_1 \ e_2 \end{array}$$

$$e_2$$

 $e_3$ 

Figure 6. Evaluation of annotated existentials

Clearly, then,

$$(e_1, e_2) \xrightarrow{C_1, C_2, (v_1:(\tau_1:\circ)), (v_2:(\tau_2:\circ))}_{D_1, D_2} \xrightarrow{*} (v_1, v_2)$$

Given our representation of pairs as closures, this is reasonable; though the conscription trace tags  $v_1$  and  $v_2$  with linear types, they will be discharged at these types as soon as  $(v_1, v_2)$  is destructed and, assuming the provided function f is of declared type  $\tau_1 \stackrel{\circ}{\to} \tau_2 \stackrel{\circ}{\to} \gamma$ , reconscripted at their original, untagged types  $\tau_1$  and  $\tau_2$ .

**Existentials** Existential types are a bit more complicated. Given  $e_1 \stackrel{C}{\longrightarrow} * v_1$  where  $e_1$  has type  $\{\alpha \mapsto \tau\}\tau_1, \tau_1$  has kind  $\kappa_1$ , and  $\tau$  has kind  $\kappa$ , let

$$e_{2} = \operatorname{pack} \alpha: \kappa = \tau \operatorname{in} e_{1}: \tau_{1}$$
  
= let  $x = e_{1} \operatorname{in} \Lambda \beta: \circ. \lambda^{\kappa_{1}} f: (\forall \alpha: \kappa. \tau_{1} \xrightarrow{\circ} \beta). f[\tau] x$   
=  $(\lambda^{\circ} x: \{\alpha \mapsto \tau\} \tau_{1}. \Lambda \beta: \circ. \lambda^{\kappa_{1}} f: (\forall \alpha: \kappa. \tau_{1} \xrightarrow{\circ} \beta). f[\tau] x) e_{1}$ 

Letting  $v_2$  refer to the result of fully evaluating  $e_2$ , we have  $e_2 \stackrel{C, (v_1: \{\alpha \mapsto \tau\}\tau_D)}{\longrightarrow} v_2$ ; the details of this reduction can be seen in Figure 6. This seems sensible for the encoding of an existential package, as it reflects exactly that the result of evaluating  $e_1$  has been captured in a function closure, just as with the capture that occurs when applying the product constructor.

Now, to use  $v_2$ , take e' of type  $\tau'$  and let

$$e_3 = \mathbf{unpack} \alpha, x = v_2 \mathbf{in} e' \\ = v_2 [\tau'] (\Lambda \alpha : \kappa. \lambda^{\circ} x : \tau_1. e')$$

The evaluation of this expression, also shown in Figure 6, is a bit lengthier; the interesting thing to note, though, is that, before the existentially wrapped  $v_1$  is passed to the function wrapping e', it is discharged at type  $\{\alpha \mapsto \tau\}\tau_1$  and reconscripted at the type  $\{\alpha \mapsto (\tau:\kappa)\}\tau_1$ . This stricter reconscription stands in direct contrast to what occurs in the elimination of pairs, and, if  $\kappa$  is  $\circ$ , we know that any component of  $v_1$  with type  $(\tau:\kappa)$  must be discharged exactly once in e' for each time it is conscripted, regardless of whether or not  $\tau$  can also be given kind  $\star$ .

Subsequently, of course, it could be reconscripted simply at type  $\tau$ , but thanks to parametricity we know that e' cannot do this on its own—as in the filesystem example, it would need to make use of a function already present in the existential package.

**DFAs** Finally, to show how the arguments made in proving of Lemma 10 are made rigorous by this machinery, consider making two more minor syntactic changes to System  $F^\circ$ : first, rather than writing  $\lambda^{\kappa} x:\tau$ . *e* for functions, write **fn**^{\kappa} name  $x:\tau$ . *e*, and second, rather than the simple L-TAG, add rules for the left and right sides of applications which further add L and R tags to discharged values, along with a final result rule which adds no further tag. It doesn't matter what names we give to functions, but if we choose to use our record labels as function names when possible—and require that all other names are drawn from some set disjoint from said record labels—then our definition of an evaluation reflecting a

$$\begin{split} \tau & ::= \dots \mid (\tau, \dots, \tau) \\ e & ::= \dots \mid (e, \dots, e) \mid \mathsf{let} (x_1, \dots, x_n) = e \; \mathsf{in} \; e \\ v & ::= \dots \mid (v, \dots, v) \\ \end{split} \\ \begin{bmatrix} \mathsf{K}\operatorname{-\mathsf{PROD}} \end{bmatrix} \frac{\Gamma \vdash \tau_1 : \kappa \, \dots \, \Gamma \vdash \tau_n : \kappa}{\Gamma \vdash (\tau_1, \dots, \tau_n) : \kappa} \\ \begin{bmatrix} \mathsf{T}\operatorname{-\mathsf{PROD}} \end{bmatrix} \frac{\Gamma; \Delta_1 \vdash e_1 : \tau_1 \, \dots \, \Gamma; \Delta_n \vdash e_n : \tau_n}{\Delta_1 \sqcup \dots \amalg \Delta_n = \Delta} \\ \frac{\Delta_1 \sqcup \dots \amalg \Delta_n = \Delta}{\Gamma; \Delta \vdash (e_1, \dots, e_n) : (\tau_1, \dots, \tau_n)} \\ \\ \begin{bmatrix} \mathsf{T}\operatorname{-\mathsf{PLET}} \end{bmatrix} \frac{\Gamma; \Delta_2 \mid x_1 : \tau_1, \dots, x_n : \tau_n \supseteq \Gamma', \Delta_2'}{\Gamma; \Delta \vdash \mathsf{let} (x_1, \dots, x_n) = e_1 \; \mathsf{in} \; e_2 : \tau} \end{split}$$

Figure 7. Syntax and static rules for multiplicative products

word becomes very simple: the DFA trace must be visible in D by examining the values tagged with L.

We can now return to the proof given in Section 3 and consider the problem in terms of value traces. After unpacking the existential, we will have the linear initial state conscripted but not yet discharged, and because our entire expression has the unrestricted type Unit, we know from Corollary 13 that it must be discharged before the final result is returned. Parametricity further restricts states to being discharged on the right side of applications. Then, as before, it is simply a matter of following the state as functions are applied to it, noting the labels that can become attached to the value trace; in the end, when some done\_at\_q is applied, our value trace will contain a clear accepting DFA trace.

# 5. Extensions: Towards Practicality

Although System  $F^{\circ}$  as presented so far can serve as an expressive core calculus, to be usable in practice it is necessary to extend it with features suitable for a surface language. This section shows how to add a variety of useful constructs familiar from functional programming. None of these extensions are particularly difficult, though we point out a few places where linearity must be minded.

#### 5.1 Language additions

**Polykinded products** System  $F^{\circ}$  as it stands can already encode both unrestricted and linear (*multiplicative*) products. In practice, however, it is useful to build in support for n-ary tuples. As shown in Figure 7, the typing rules are mostly standard; evaluation rules are completely standard and have been omitted. The kinding rule K-PROD, in combination with subkinding, allows the same syntax

Figure 8. Syntax and rules for the additive products

$$e ::= \dots | \mathbf{fix} f. v$$

$$[\text{E-FIX}] \mathbf{fix} f. v \longrightarrow \{ f \mapsto (\mathbf{fix} f. v) \} v$$

$$[\text{T-FIX}] \frac{\Gamma, f:\tau; \cdot \vdash v : \tau \qquad \Gamma \vdash \tau : \star}{\Gamma; \cdot \vdash \mathbf{fix} f. v : \tau}$$

Figure 9. Syntax and rules for fixpoints.

to be used for both linear and unrestricted tuples; if any component type of a tuple is linear, so is the tuple itself. T-PLET uses an n-ary version of the nondeterministic bind operation (recall T-LAM), allowing unrestricted components of a tuple to be used arbitrarily but binding linear components in the linear context.

Additive products As mentioned in Section 2, we cannot encode linear sums in System F° as presented so far—we have no means of sharing the linear context among different subexpressions, all but one of which will be dropped. Connectives that enable this are referred to as *additive* in the terminology of linear logic; traditionally they are written as  $\tau_1 \oplus \tau_2$  for sums, and  $\tau_1 \& \tau_2$  for *additive products*—lazy linear products eliminated by projection rather than pattern matching. In keeping with our syntactic conventions, however, we introduce n-ary additive products  $\langle \tau_1, \ldots, \tau_n \rangle$ ; sums could be encoded in terms of these products but are also subsumed by the algebraic datatypes discussed below.

Figure 8 gives the new syntax, definitions, and static and dynamic rules for n-ary additive products. Note that additive products necessarily denote suspended computations, as further evaluation progress cannot be made until a branch is chosen. Unlike multiplicative products, additive products always have kind o regardless of the kinds of the component types, as seen in K-ADD.

**Fixpoints** Adding support for fixpoint computation turns out to be remarkably straightforward: recursion implies the ability to repeatedly invoke a function within its own body, which naturally puts us in the unrestricted portion of System  $F^\circ$ . Figure 9 shows the new syntax and rules, where we write **fix** *f*. *v* for the computation that immediately unwinds itself to a value as shown in E-FIX. Since these operational semantics might duplicate *v*, T-FIX requires that fixpoint expressions be closed with respect to the linear context; note that this is compatible with Lemma 4 (3).

*Algebraic datatypes* Figure 10 shows the static rules for general recursive, polymorphic datatypes—as with multiplicative products,

the dynamic rules are standard. We assume a signature  $\Sigma_T$  that maps a datatype constructor T to its (higher) kind  $\overline{\kappa}$ , which describes T's type parameters; T has arity  $k_T$ . Term-level constructors for type T are written  $ctr_i^T$ , and their types are given by the map  $\Sigma_c$ . The type of a term-level constructor is polymorphic over the parameters of the datatype, but it may also expect additional type parameters, which act as through they are existentially bound—this design is similar to that found in Haskell. Importantly, if the result kind of the type constructor T is  $\star$ , then the value parameters to the term-level constructors must also be of kind  $\star$ . This constraint, written as part of the signature well formedness checks in Figure 10, is necessary to ensure that an unrestricted datatype can't hide a linear value and thus allow it to be duplicated or discarded.

The additive nature of datatypes is shown in T-CASE, which uses the same linear context  $\Delta_2$  when checking each branch. The branches themselves are polymorphic linear functions additionally abstracted over any existential type parameters—one of which will be applied to create a result of the appropriate type. As with additive products, each branch will capture the same linear free variables, but only one of the branches will fire at run time.

# 5.2 Syntactic support for linearity

In addition to making it possible for types to express safe protocols over stateful resources, it must also be convenient for programmers to write client code that interacts with these interfaces. For instance, a file copy program, using the interface given in Section 1.2, should look like simple imperative code and not require the programmer to think explicitly about existentially bound linear type variables and the threading of linear filehandles. In the rest of this section we sketch out some possible syntactic support—partly inspired by but simpler than Haskell's type classes—for programming with linearity, which allows us to copy files as follows<sup>8</sup>:

let copy = fix f.  
action ([File] 
$$h_1$$
, [File]  $h_2$ ).  
 $y \leftarrow h_1$ .read;  
if  $y = \text{EOF}$  then return ()  
else  $h_2$ .write  $y$ ;  
 $(h_1, h_2)..f$   
in  
start [File]  $h_1 = \text{open "InFile.txt"}$ ;  
start [File]  $h_2 = \text{open "OutFile.txt"}$ ;  
 $(h_1, h_2)..\text{copy}$ ;  
 $h_1$ .close;  $h_2$ .close

This program would fail to typecheck if either of the calls to close were omitted, if one of the file handles were closed within the body of copy, or if a read or write operation were called after close. In the rest of this section we will show how to translate from this syntax into System  $F^{\circ}$  as described up to this point.

First, we make more explicit the pattern common to the protocol examples we have seen thus far. Let  $\Theta = \alpha_1 \dots \alpha_n$  be a list of type variables representing the states of the protocol.<sup>9</sup> The possible actions of the protocol can be given by a record of operations, each of which causes a state transition from  $\sigma_{in}$  to  $\sigma_{out}$ .

$$Ops(\Theta) = \{ op_1 : \forall \beta_1:\kappa_1.\tau_{in1} \stackrel{\star}{\to} \sigma_{in1} \stackrel{\star}{\to} (\tau_{out_1}, \sigma_{out_1}) \\ \cdots \\ op_m : \forall \overline{\beta_m:\kappa_m}.\tau_{inm} \stackrel{\star}{\to} \sigma_{inm} \stackrel{\star}{\to} (\tau_{out_m}, \sigma_{out_m}) \}$$

Here,  $\sigma_{in_j} \in \Theta$  is the start state of the operation and  $\sigma_{out_j} \in \Theta \cup$ {Unit} is either the end state of the operation or Unit, indicating

<sup>&</sup>lt;sup>8</sup> Technically, to match up with what follows, the return type of write must be (Unit,  $\alpha$ ) rather than just  $\alpha$ ; we also assume that EOF is of type Char.

<sup>&</sup>lt;sup>9</sup> In this section, we write  $\forall \Theta : \circ . \tau$  as short hand for  $\forall \alpha_1 : \circ . . . . \forall \alpha_n : \circ . \tau$ , and similarly for existentials.

$$\begin{split} \bar{\kappa} & ::= \kappa \mid \kappa \Rightarrow \bar{\kappa} \\ \tau & ::= \dots \mid T \tau_{1} \dots \tau_{n} \\ [K-T] & \frac{\sum_{T} \vdash T : \kappa_{1} \Rightarrow \dots \kappa_{k_{T}} \Rightarrow \kappa_{T}}{\Gamma \vdash \tau_{1} : \kappa_{1} \dots \Gamma \vdash \tau_{k_{T}} : \kappa_{k_{T}}} \\ [K-T] & \frac{\sum_{T} \vdash T : \kappa_{1} \Rightarrow \dots \kappa_{k_{T}} \Rightarrow \kappa_{T}}{\Gamma \vdash T \tau_{1} \dots \tau_{k_{T}} : \kappa_{T}} \\ [K-T] & \frac{\sum_{T} \vdash T : \kappa_{1} \Rightarrow \dots \kappa_{k_{T}} \Rightarrow \kappa_{T}}{\Gamma \vdash T \tau_{1} \dots \tau_{k_{T}} : \kappa_{T}} \\ [T-CTR] & \frac{\sum_{c} \vdash ctr_{i}^{T} : \forall \alpha_{1}:\kappa_{1}, \dots, \alpha_{n}:\kappa_{n}.\tau \xrightarrow{*} T \alpha_{1} \dots \alpha_{k_{T}}}{\Gamma; \Delta \vdash ctr_{i}^{T} \tau_{1} \dots \tau_{k_{T}}, \tau_{k_{T}+1}, \dots, \tau_{n} e : T \tau_{1} \dots \tau_{k_{T}}} \\ \frac{\Delta_{1} \Downarrow \Delta_{2} = \Delta \qquad \Gamma; \Delta_{1} \vdash e : T \tau_{1} \dots \tau_{k_{T}} \qquad \alpha_{1} \dots \alpha_{n} \notin \Gamma, \tau}{\Gamma; \Delta \vdash ctr_{i}^{T} : \forall \alpha_{1}:\kappa_{1}, \dots, \alpha_{k_{T}}:\kappa_{k_{T}}, \alpha_{k_{T}+1}:\kappa_{k_{T}+1}, \dots \alpha_{n_{i}}:\kappa_{n_{i}}.\tau_{i} \xrightarrow{*} T \alpha_{1} \dots \alpha_{k_{T}}} \\ \frac{\left(\sum_{c} \vdash ctr_{i}^{T} : \forall \alpha_{1}:\kappa_{1}, \dots, \alpha_{k_{T}}:\kappa_{k_{T}}, \alpha_{k_{T}+1}:\kappa_{k_{T}+1}, \dots \alpha_{n_{i}}:\kappa_{n_{i}}.\tau_{i} \xrightarrow{*} T \alpha_{1} \dots \alpha_{k_{T}}}\right)^{(i \in 1 \dots m)}}{\Gamma; \Delta \vdash case \ e \ of \ ctr_{1}^{T} \cdot v_{1} \mid \dots \mid ctr_{m}^{T} \cdot v_{m} : \tau} \\ \end{split}$$

 $\Sigma_T \vdash \Sigma_c$  holds if and only if for every T such that  $\Sigma_T \vdash T : \kappa_1 \Rightarrow \ldots \kappa_{k_T} \Rightarrow \kappa_T$  and  $ctr_i^T$  such that  $\Sigma_c \vdash ctr_i^T : \forall \alpha_1 : \kappa_1, \ldots, \alpha_{n_i} : \kappa_{n_i} \cdot \tau_i \stackrel{\star}{\to} T \alpha_1 \ldots \alpha_{k_T}$  it is the case that  $\alpha_1 : \kappa_1, \ldots, \alpha_{n_i} : \kappa_{n_i} \vdash \tau_i : \kappa_T$ 

Figure 10. Syntax and static rules for polymorphic, recursive datatypes.  $\Sigma_T$  and  $\Sigma_c$  are global constructor contexts such that  $\Sigma_T \vdash \Sigma_c$ 

that the protocol is complete. We assume that each operation might be polymorphic and, for simplicity, that each takes one input of type  $\tau_{in_j}$  and produces an output of type  $\tau_{out_j}$ , both of which might contain occurrences of variables in  $\Theta$ . Almost all of the protocol examples we have seen can be written in this form.

[

If  $\alpha_{init} \in \Theta$  is the start state, then an initial instance  $e_1$  of the protocol is represented by a value of type  $\exists \Theta : \circ.(\alpha_{init}, \mathsf{Ops}(\Theta))$ . Thanks to  $\alpha_{init}$ , this value has linear kind, so it must be unpacked and the resulting tuple destructed. We suggest a convenient notation to simplify this process:

start [Ops] 
$$h = e_1$$
 in  $e_2 \triangleq$  unpack  $\Theta_h, x = e_1$  in  
let  $(h, ops_h) = x$  in  $e_2$ 

Here,  $\Theta_h$  is a list of fresh type variables, h is bound to the linear state, and  $ops_h$  is the record of operations associated with h. The syntax includes the annotation Ops to help with type checking and to determine what syntactic sugar applies with respect to h in the body  $e_2$ . Note that  $ops_h$  has type  $Ops(\Theta_h)$ , which is unrestricted and hence may be duplicated at will; this invariant will be maintained throughout the interpretation of our syntax.

The intuition behind our scheme is simple: we use the single name h for the linear handle associated with the thread of a protocol—this implicit reuse can never be an issue, since the variable is linear. This handle (which must be typed by a variable in  $\Theta_h$ ) then determines an appropriate record of operations,  $ops_h$ , associated with the protocol. The programmer never mentions  $ops_h$  explicitly; rather it is threaded through the computation automatically, much like a typeclass dictionary.

The basic notation treats the protocol operations as "methods" of the linear handle h with implicit argument  $ops_h$ , binding the result to the variable y and conflating the initial and result handles:

$$y \leftarrow h.op_i [\overline{\tau}] e_1; e_2 \triangleq \text{let} (y,h) = ops_h.op_i [\overline{\tau}] e_1 h \text{ in } e_2$$

If the programmer wishes to write a client function that takes one or more handles, each following its own protocol, the function must be polymorphic over the appropriate *ops* records. We call such functions "actions"; it makes sense to give them a type very similar to the operations of a protocol, since they must take a vector of resource handles each in some protocol state and return a new vector of protocol states.

$$\begin{array}{l} \operatorname{action}\left(\left[\mathsf{Ops}_{1}\right]h_{1},\ldots,\left[\mathsf{Ops}_{n}\right]h_{n}\right)\overline{\beta:\kappa}\left(x:\tau_{in}\right).e^{-}\triangleq\\ \Lambda\Theta_{h_{1}}:\circ.\lambda^{*}ops_{h_{1}}:\mathsf{Ops}_{1}(\Theta_{h_{1}}).\ldots\Lambda\Theta_{h_{n}}:\circ.\lambda^{*}ops_{h_{n}}:\mathsf{Ops}_{1}(\Theta_{h_{n}})\\ \overline{\Lambda\overline{\beta:\kappa}}.\lambda^{*}x:\tau_{in}.\lambda^{\circ}s:(\alpha_{h_{1}},\ldots,\alpha_{h_{n}}).\\ \operatorname{let}\left(h_{1},\ldots,h_{n}\right)=s \text{ in } e \end{array}$$

Here *e* must have type  $(\tau_{out}, (\sigma_{out_1}, \ldots, \sigma_{out_n}))$ , where  $\sigma_{out_j} \in \Theta_{h_j}$ and the type of  $h_j$  must be some  $\alpha_{h_j}$  in  $\Theta_{h_j}$ , representing the state that handle must be in when this function is called. Our syntactic sugar requires that type inference be able to determine the  $\alpha_{h_j}$ 's by inspecting *e*. We expect this to be rather straightforward, since the type is uniquely determined by the first operation invoked on the handle; alternatively we could add a type annotation to  $h_j$ . Inside the body of such an action, we allow a convenient means of packaging the state values to be returned to the caller:

return 
$$e \triangleq (e, (h_1, \ldots, h_n))$$

Taken together, these constraints ensure that the type of an action is a polymorphic operation over a vector of states; *i.e.*, it has the type

$$\forall \Theta_{h_1} : \circ. \operatorname{Ops}_1(\Theta_{h_1}) \xrightarrow{*} \dots \forall \Theta_{h_n} : \circ. \operatorname{Ops}_n(\Theta_{h_n}) \xrightarrow{*} \\ \forall \overline{\beta} : \kappa. \tau_{in} \xrightarrow{*} (\alpha_{h_1}, \dots, \alpha_{h_n}) \xrightarrow{*} (\tau_{out}, (\sigma_{out_1}, \dots, \sigma_{out_n}))$$

Inside the body of an action, the "method call" syntax defined earlier can be used to invoke protocol operations on the handles  $h_j$ . Since an action is just another sort of operation over handles albeit a higher level one—we would like similar syntax for invoking such an operation on a list of handles. Assuming f has the type of an action as given above, we define:

$$y \leftarrow (h_1, \dots, h_n) \dots f[\overline{\tau}] e_1; e_2 \triangleq$$
  
let  $(y, s) = f[\Theta_{h_1}] ops_{h_1} \dots [\Theta_{h_n}] ops_{h_n}[\overline{\tau}] e_1 (h_1, \dots, h_n)$  in  
let  $(h_1, \dots, h_n) = s$  in  $e_2$ 

One can easily define variants of both the above and the earlier protocol operation notation that omit the argument  $e_1$  or the binding for y in the case that one or both of  $\tau_{in}$  and  $\tau_{out}$  are Unit, or that do not rebind the handle variable in cases where the handle is consumed, as with close. Also, if the operation invocation is last in a sequence,  $e_2$  and the corresponding the let binding may be omitted, since the result of the operation is the result of the sequence—this is useful, *e.g.*, for making tail calls.

Of course, there are many ways of using linear types, and the above does not account for all of them. We believe, however, that it does show the feasibility of making System  $F^{\circ}$  protocols palatable to the programmer, and that this provides further evidence that linearity as implemented in System  $F^{\circ}$  deserves to be considered for inclusion in more mainstream functional programming languages.

# 6. Conclusion

We have presented System  $F^{\circ}$ , a simple variant of the linear polymorphic  $\lambda$ -calculus that nevertheless can enforce a rich variety of

protocols. System  $F^{\circ}$  is sound and enjoys parametricity, and we have proved that protocol enforcement is faithful—that is, linear resources are never misused—thanks to these properties. We have demonstrated the applicability of System  $F^{\circ}$  through a variety of examples and by showing how to extend it with features geared towards practical programming with linear types.

# References

- [1] Amal Ahmed, Matthew Fluet, and Greg Morrisett. A step-indexed model of substructural state. In *ICFP '05: Proceedings of the tenth* ACM SIGPLAN international conference on Functional programming, pages 78–91, New York, NY, USA, 2005. ACM.
- [2] Amal Ahmed, Matthew Fluet, and Greg Morrisett. L3: A linear language with locations. *Fundam. Inf.*, 77(4):397–449, 2007.
- [3] Andrew Barber. *Linear Type Theories, Semantics and Action Calculi*. PhD thesis, Edinburgh University, 1997.
- [4] Nick Benton, G. M. Bierman, J. Martin E. Hyland, and Valeria de Paiva. A term calculus for intuitionistic linear logic. In M. Bezem and J. F. Groote, editors, *Proceedings of the International Conference* on Typed Lambda Calculi and Applications, pages 75–90. Springer-Verlag LNCS 664, 1993.
- [5] P. N. Benton. A mixed linear and non-linear logic: proofs, terms and models. In *Proceedings of Computer Science Logic (CSL '94), Kazimierz, Poland.*, pages 121–135. Springer-Verlag, 1995.
- [6] G. M. Bierman, A. M. Pitts, and C. V. Russo. Operational properties of lily, a polymorphic linear lambda calculus with recursion. In *Fourth International Workshop on Higher Order Operational Techniques in Semantics, Montral, volume 41 of Electronic Notes in Theoretical Computer Science.* Elsevier, 2000.
- [7] Gavin M. Bierman. Program equivalence in a linear functional language. Journal of Functional Programming, 10(2), 2000.
- [8] Lars Birkedal, Rasmus Ejlers Møgelberg, and Rasmus Lerchedahl Petersen. Category-theoretic models of Linear Abadi & Plotkin Logic. *Theory and Applications in Categories*, 20(7), 2008.
- [9] Arthur Charguéraud and François Pottier. Functional translation of a calculus of capabilities. In *ICFP '08: Proceeding of the 13th* ACM SIGPLAN international conference on Functional programming, pages 213–224, New York, NY, USA, 2008. ACM.
- [10] Karl Crary, David Walker, and Greg Morrisett. Typed memory management in a calculus of capabilities. In Proc. 26th ACM Symp. on Principles of Programming Languages (POPL), pages 262–275, San Antonio, Texas, January 1999.
- [11] Robert DeLine and Manuel Fähndrich. Enforcing high-level protocols in low-level software. In Proc. of the SIGPLAN Conference on Programming Language Design, pages 59–69, Snowbird, UT, June 2001.
- [12] Manuel Fähndrich, Mark Aiken, Chris Hawblitzel, Orion Hodson, Galen Hunt, James R. Larus, and Steven Levi. Language support for fast and reliable message-based communication in singularity os. *SIGOPS Oper. Syst. Rev.*, 40(4):177–190, 2006.
- [13] Manuel Fähndrich and Robert DeLine. Adoption and focus: Practical linear types for imperative programming. In *Proc. of the SIGPLAN Conference on Programming Language Design*, pages 13–24, Berlin, Germany, June 2002.
- [14] Jean-Yves Girard. Interprétation fonctionnelle et élimination des coupures de l'arith mé tique d'ordre supérieur. Thèse d'état, University of Paris VII, 1972. Summary in J. E. Fenstad, editor, Scandinavian Logic Symposium, pp. 63–92, North-Holland, 1971.
- [15] Jean-Yves Girard. Linear logic. *Theoretical Computer Science*, 50:1– 102, 1987.
- [16] Jean-Yves Girard, Y. Lafont, and P. Taylor. *Proofs and Types*. Cambridge University Press, 1989.
- [17] Michael Hicks, Greg Morrisett, Dan Grossman, and Trevor Jim. Experience with safe manual memory-management in Cyclone. In ISMM

'04: Proceedings of the 4th international symposium on Memory management, pages 73–84, New York, NY, USA, 2004. ACM.

- [18] Kohei Honda, Vasco T. Vasconcelos, and Makoto Kubo. Language primitives and type discipline for structured communication-based programming. In *ESOP98, volume 1381 of LNCS*, pages 122–138. Springer-Verlag, 1998.
- [19] Atsushi Igarashi and Naoki Kobayashi. Resource usage analysis. ACM Trans. Program. Lang. Syst., 27(2):264–313, 2005.
- [20] Oleg Kiselyov and Chung-chieh Shan. Lightweight monadic regions. In Haskell '08: Proceedings of the first ACM SIGPLAN symposium on Haskell, pages 1–12, New York, NY, USA, 2008. ACM.
- [21] Yves Lafont. The linear abstract machine. *Theoretical Computer Science*, 59:157–180, 1988. Some corrections in volume 62 (1988), pp. 327–328.
- [22] Patrick Lincoln and John Mitchell. Operational aspects of linear lambda calculus. In 7th Symposium on Logic in Computer Science, IEEE, pages 235–246. IEEE Computer Society Press, 1992.
- [23] John C. Reynolds. Towards a theory of type structure. In *Programming Symposium*, volume 19 of *Lecture Notes in Computer Science*, pages 408–425. Springer-Verlag, Paris, France, April 1974.
- [24] John C. Reynolds. Types, abstraction, and parametric polymorphism. In R.E.A Mason, editor, *Information Processing*, pages 513–523. Elsevier Science Publishers B.V., 1983.
- [25] Kaku Takeuchi, Kohei Honda, and Makoto Kubo. An interactionbased language and its typing system. In *Proceedings of PARLE'94*, pages 398–413. Springer-Verlag, 1994. Lecture Notes in Computer Science number 817.
- [26] David N. Turner and Philip Wadler. Operational interpretations of linear logic. *Theoretical Computer Science*, 227(1-2):231–248, September 1999.
- [27] Vasco T. Vasconcelos, Simon J. Gay, and António Ravara. Type checking a multithreaded functional language with session types. *Theoreti*cal Computer Science, 368(1–2):64–87, 2006.
- [28] Edsko Vries, Rinus Plasmeijer, and David M. Abrahamson. Uniqueness typing simplified. In *Implementation and Application of Functional Languages: 19th International Workshop, IFL 2007, Freiburg, Germany, September 27-29, 2007. Revised Selected Papers*, pages 201–218, Berlin, Heidelberg, 2008. Springer-Verlag.
- [29] Philip Wadler. Linear types can change the world! In M. Broy and C. Jones, editors, *Progarmming Concepts and Methods*, Sea of Galilee, Israel, April 1990. North Holland. IFIP TC 2 Working Conference.
- [30] Philip Wadler. There's no substitute for linear logic. In 8th International Workshop on the Mathematical Foundations of Programming Semantics, 1992.
- [31] Philip Wadler. A taste of linear logic. In *Mathematical Foundations of Computer Science*, volume 711 of *Lecture Notes in Computer Science*, pages 185–210. Springer-Verlag, 1993.
- [32] David Walker. A type system for expressive security policies. In Proc. 27th ACM Symp. on Principles of Programming Languages (POPL), pages 254–267. ACM Press, Jan 2000.
- [33] David Walker. Advanced Topics in Types and Programming Languages, chapter Substructural Type Systems. MIT Press, 2005.
- [34] Keith Wansbrough and Simon Peyton Jones. Once upon a polymorphic type. In POPL '99: Proceedings of the 26th ACM SIGPLAN-SIGACT symposium on Principles of programming languages, pages 15–28, New York, NY, USA, 1999. ACM.
- [35] Andrew K. Wright and Matthias Felleisen. A syntactic approach to type soundness. *Information and Computation*, 115(1):38–94, 1994. Preliminary version in Rice TR 91-160.
- [36] Dengping Zhu and Hongwei Xi. Safe Programming with Pointers through Stateful Views. In *Proceedings of the 7th International Symposium on Practical Aspects of Declarative Languages*, pages 83–97, Long Beach, CA, January 2005. Springer-Verlag LNCS vol. 3350.