



# A Church-Style Intermediate Language for MLF

Didier Rémy, Boris Yakobowski

► **To cite this version:**

Didier Rémy, Boris Yakobowski. A Church-Style Intermediate Language for MLF. Theoretical Computer Science, Elsevier, 2012, 435, pp.77–105. <<http://dx.doi.org/10.1016/j.tcs.2012.02.026>>. <10.1016/j.tcs.2012.02.026>. <hal-01093719>

**HAL Id: hal-01093719**

**<https://hal.inria.fr/hal-01093719>**

Submitted on 11 Dec 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A Church-Style Intermediate Language for $\text{ML}^F$

Didier Rémy<sup>a</sup>, Boris Yakobowski<sup>b</sup>

<sup>a</sup>INRIA

Rocquencourt - BP 105, 78153 Le Chesnay Cedex

<sup>b</sup>CEA, LIST, Laboratoire Sécurité des Logiciels,

Boîte 94, 91191 Gif-sur-Yvette Cedex, France.

---

## Abstract

$\text{ML}^F$  is a type system that seamlessly merges ML-style implicit but second-class polymorphism with System-F explicit first-class polymorphism. We present  $x\text{ML}^F$ , a Church-style version of  $\text{ML}^F$  with full type information that can easily be maintained during reduction. All parameters of functions are explicitly typed and both type abstraction and type instantiation are explicit. However, type instantiation in  $x\text{ML}^F$  is more general than type application in System F. We equip  $x\text{ML}^F$  with a small-step reduction semantics that allows reduction in any context, and show that this relation is confluent and type preserving. We also show that both subject reduction and progress hold for weak-reduction strategies, including call-by-value with the value-restriction. We exhibit a type preserving encoding of  $\text{ML}^F$  into  $x\text{ML}^F$ , which shows that  $x\text{ML}^F$  can be used as the internal language for  $\text{ML}^F$  after type inference, and also ensures type soundness for the most expressive variant of  $\text{ML}^F$ .

*Keywords:*  $\text{ML}^F$ , System F, Types, Type Generalization, Type Instantiation, Retyping functions, Coercions, Type Soundness, Binders

---

## Introduction

$\text{ML}^F$  (Le Botlan and Rémy, 2003, 2009; Rémy and Yakobowski, 2008) is a type system that seamlessly merges ML-style implicit but second-class polymorphism with System-F explicit first-class polymorphism. This is done by enriching System-F types. Indeed, System F is not well-suited for partial type inference, as illustrated by the following example. Assume that a function, say *choice*, of type  $\forall(\alpha) \alpha \rightarrow \alpha \rightarrow \alpha$ , and the identity function *id*, of type  $\forall(\beta) \beta \rightarrow \beta$ , have been defined. How can the application *choice* to *id* be typed in System F? Should *choice* be applied to the type  $\forall(\beta) \beta \rightarrow \beta$  of the identity, that is itself kept polymorphic? Or should it be applied to the monomorphic type  $\gamma \rightarrow \gamma$ , with the identity being applied to  $\gamma$  (where  $\gamma$  is bound in a type abstraction in front of the application)? Unfortunately, these alternatives have incompatible types, respectively  $(\forall(\alpha) \alpha \rightarrow \alpha) \rightarrow (\forall(\alpha) \alpha \rightarrow \alpha)$  and  $\forall(\gamma) (\gamma \rightarrow \gamma) \rightarrow (\gamma \rightarrow \gamma)$ : none is an instance of the other. Hence, in System F, one is forced to irreversibly choose between one of the two explicitly typed terms.

However, a type inference system cannot choose between the two, as this would sacrifice completeness and be somehow arbitrary. This is why  $\text{ML}^F$  enriches types with instance-bounded polymorphism, which allows to write more expressive types that factor out in a single type

all typechecking alternatives in such cases as the example of *choice*. Now, the type  $\forall(\alpha \geq \forall(\beta) \beta \rightarrow \beta) \alpha \rightarrow \alpha$ , which should be read “ $\alpha \rightarrow \alpha$  where  $\alpha$  is any instance of  $\forall(\beta) \beta \rightarrow \beta$ ”, can be assigned to *choice id*, and the two previous alternatives can be recovered *a posteriori* by choosing different instances for  $\alpha$ .

Currently, the language  $\text{ML}^F$  comes with a Curry-style version,  $i\text{ML}^F$ , where no type information is needed and a type-inference version,  $e\text{ML}^F$ , that requires partial type information (Le Botlan and Rémy, 2009). However,  $e\text{ML}^F$  is not quite in Church style: a large amount of type information is still implicit, and partial type information cannot be easily maintained during reduction. Hence, while  $e\text{ML}^F$  is a good surface language, it is not a good candidate for use as an internal language during the compilation process, where some program transformations, and perhaps some reduction steps, are being performed. This has been a problem for the adoption of  $\text{ML}^F$  in the Haskell community (Peyton Jones, 2003), as the Haskell compilation chain uses an explicitly-typed internal language, especially, but not only, for evidence translation due to the use of qualified types (Jones, 1994).

This is also an obstacle to proving subject reduction, which does not hold in  $e\text{ML}^F$ . In a way, this is unavoidable in a language with non-trivial partial type inference. Indeed, type annotations cannot be completely dropped, but must at least be transformed and reorganized during reduction. Still, one could expect that  $e\text{ML}^F$  is equipped with reduction rules for type annotations. This has actu-

---

URL: <http://gallium.inria.fr/~remy> (Didier Rémy),  
<http://www.yakobowski.org> (Boris Yakobowski)

ally been considered in the original presentation of  $\text{MLF}$ , but only with limited success. The reduction kept track of annotation sites during reduction; this showed, in particular, that no new annotation site needs to be introduced during reduction. Unfortunately, the exact form of annotations could not be maintained during reduction, by lack of an appropriate language to describe their computation. As a result, it has only been shown that some type derivation can be rebuilt after the reduction of a well-typed program, but without exhibiting an algorithm to compute it during reduction.

Independently, Rémy and Jakobowski (2008) have introduced graphic constraints, both to simplify the presentation of  $\text{MLF}$  and to improve its type inference algorithm. This also resulted in a simpler and more expressive definition of  $\text{MLF}$ . Hence, by  $e\text{MLF}$ , we refer to the graphical presentation of  $\text{MLF}$  rather than the original version. Consistently,  $i\text{MLF}$  refers to the graphic Curry’s style version of  $e\text{MLF}$ . We still use the generic name  $\text{MLF}$  when the style of presentation does not matter.

In this article, we present  $x\text{MLF}$ , a Church-style version of  $\text{MLF}$  that contains full type information. In fact, type checking becomes a simple and local verification process—by contrast with type inference in  $e\text{MLF}$ , which is based on unification. In  $x\text{MLF}$ , type abstraction, type instantiation, and all parameters of functions are explicit, as in System F. However, type instantiation is more general and more atomic than type application in System F: we use explicit type instantiation expressions that are proof evidences for the type instance relations.

In addition to the usual  $\beta$ -reduction, we give a series of reduction rules for simplifying type instantiations. These rules are confluent when allowed in any context. Moreover, reduction preserves typings, and is sufficient to reduce all typable expressions to a value when used in either a call-by-value or call-by-name setting. This establishes the soundness of  $\text{MLF}$  for a call-by-name semantics for the first time. Furthermore, we show that  $x\text{MLF}$  is a conservative extension of System F.

The natural definition of  $x\text{MLF}$  is actually more expressive than that of  $\text{MLF}$ . Still, we can restrict type-checking in  $x\text{MLF}$  so that well-typed terms are in closer correspondence with  $\text{MLF}$  terms. This defines a well-behaved subset  $x\text{MLF}_b$  of  $x\text{MLF}$ . Then, all three versions  $i\text{MLF}$ ,  $e\text{MLF}$  and  $x\text{MLF}_b$  have the same expressiveness, and only differ by the amount of type information: terms of  $i\text{MLF}$  contain none, terms of  $e\text{MLF}$  contain some type annotations and no description of type instantiations, while  $x\text{MLF}$  contains all type annotations and a full description of all type instantiations.

A term of  $x\text{MLF}$  can easily be converted into an  $e\text{MLF}$  one by retaining type annotations, but dropping all other type information. The result may in turn be converted into a term of  $i\text{MLF}$  by further dropping all type annotations. Conversely, terms of  $i\text{MLF}$  cannot be automatically translated into terms of  $e\text{MLF}$ , since type inference in  $i\text{MLF}$  is

undecidable—some type annotations are required. Hence, source terms are terms of  $e\text{MLF}$ : type inference can rebuild the type annotations that may be left implicit, or fail if mandatory type annotations have been omitted (or are incorrect). Terms of  $e\text{MLF}$ —for which type inference succeeds—can then be elaborated into terms of  $x\text{MLF}$ .

*Outline.* Perhaps surprisingly, the difficulty in defining an internal language for  $\text{MLF}$  is not reflected in the internal language itself, which remains simple and easy to understand. Rather, the difficulties lie in the translation from  $e\text{MLF}$  to  $x\text{MLF}$ , which is made somewhat complicated by many administrative details. Hence, we present  $x\text{MLF}$  first, and study its meta-theoretical properties independently of  $e\text{MLF}$ . We then describe the elaboration of  $e\text{MLF}$  terms.

More precisely, the paper is organized as follows. We present  $x\text{MLF}$ , its syntax and its static and dynamic semantics in §1. We study its main properties, including type soundness for different evaluations strategies in §2. The elaboration of  $e\text{MLF}$  programs into  $x\text{MLF}$  is described §3. We discuss the expressiveness of  $x\text{MLF}$  in §4 and related and future works in §5.

*Proofs and implementation.* The soundness proof of  $x\text{MLF}$  has been mechanized in the Coq proof assistant<sup>1</sup> and is briefly discussed in Appendix A. Other interesting proofs of §1 and §2 can be found in Appendix B, except for two results, which have already been proved by Manzonetto and Tranquilli (2010). Detailed proofs of §3 can all be found in the dissertation of Jakobowski (2008, Chapters 14 & 15), although for a slightly different—but equivalent—presentation. We do not reproduce them here, as they depend too much on the metatheoretical properties of  $e\text{MLF}$ . The elaboration of  $e\text{MLF}$  into  $x\text{MLF}$  has been implemented in a prototype<sup>2</sup>.

## 1. The calculus

### 1.1. Types, instantiations, terms, and typing environments

All the syntactic definitions of  $x\text{MLF}$  can be found in Figures 1 and 2. We assume given a countable collection of type variables written with letters  $\alpha, \beta, \gamma$ , and  $\delta$ . As usual, types include type variables and arrow types. Other type constructors will be added later—straightforwardly, as the arrow constructor receives no special treatment. Types also include a bottom type  $\perp$  that corresponds to the System-F type  $\forall\alpha.\alpha$ . Finally, a type may also be a form of bounded quantification  $\forall(\alpha \geq \tau) \tau'$ , called *flexible* quantification, that generalizes the  $\forall\alpha.\tau$  form of System F and, intuitively, restricts the variable  $\alpha$  to range only over

<sup>1</sup>The Coq development is available at <http://www.yakobowski.org/publis/2010/xmlf-coq/>. Properties that have been mechanically verified in Coq are marked with the Coq symbol.

<sup>2</sup>Available at <http://gallium.inria.fr/~remy/mlf/proto/>.

$\alpha, \beta, \gamma, \delta$	:	<b>Type variable</b>
$\tau$	::=	<b>Type</b>
		$\alpha$ Type variable
		$\tau \rightarrow \tau$ Arrow type
		$\forall(\alpha \geq \tau) \tau$ Quantification
		$\perp$ Bottom type
$\phi$	::=	<b>Instantiation</b>
		$@\tau$ Bottom
		$!\alpha$ Abstract
		$\forall(\geq \phi)$ Inside
		$\forall(\alpha \geq) \phi$ Under
		$\&$ $\forall$ -elimination
		$\wp$ $\forall$ -introduction
		$\phi; \phi$ Composition
		$\mathbb{1}$ Identity

Figure 1: Grammar of types and instantiations

$a$	::=	<b>Term</b>
		$x$ Variable
		$\lambda(x : \tau) a$ Function
		$a a$ Application
		$\Lambda(\alpha \geq \tau) a$ Type function
		$a \phi$ Instantiation
		$\text{let } x = a \text{ in } a$ Let-binding
$\Gamma$	::=	<b>Environment</b>
		$\emptyset$ Empty
		$\Gamma, \alpha \geq \tau$ Type variable
		$\Gamma, x : \tau$ Term variable

Figure 2: Grammar of terms and typing contexts

instances of  $\tau$ . The variable  $\alpha$  is bound in  $\tau'$  but not in  $\tau$ . We may just write  $\forall(\alpha) \tau'$  when the bound  $\tau$  is  $\perp$ .

In Church-style System F, type instantiation inside terms is simply type application  $a \tau$ . By contrast, in  $x\text{MLF}$ , we use type instantiation  $a \phi$  to detail every intermediate instantiation step, so that it can be checked locally. Intuitively, the *instantiation*  $\phi$  transforms a type  $\tau$  into another type  $\tau'$  that is an instance of  $\tau$ . In a way,  $\phi$  is a witness for the instance relation that holds between  $\tau$  and  $\tau'$ . It is therefore easier to understand instantiations altogether with their static semantics, which will be explained in the next section.

Terms of  $x\text{MLF}$  are those of the  $\lambda$ -calculus enriched with let bindings, with two small differences: type instantiation  $a \phi$  generalizes System-F type application as just described; and type abstractions are extended with an instance bound  $\tau$  and written  $\Lambda(\alpha \geq \tau) a$  where the type variable  $\alpha$  is bound in  $a$ , but not in  $\tau$ . We abbreviate  $\Lambda(\alpha \geq \perp) a$  as  $\Lambda(\alpha) a$ , which simulates the type abstraction  $\Lambda\alpha. a$  of System F.

We write  $\text{ftv}(\tau)$  and  $\text{ftv}(a)$  for the sets of type variables that appear free in  $\tau$  and  $a$ , respectively. We identify types, instantiations, and terms up to the renaming of bound variables. The capture-avoiding substitution of

$$\frac{\alpha \notin \text{dom}(\Gamma) \quad \text{wf}(\Gamma) \quad \text{ftv}(\tau) \subseteq \text{dom}(\Gamma)}{\text{wf}(\Gamma, \alpha \geq \tau)}$$

$$\frac{x \notin \text{dom}(\Gamma) \quad \text{wf}(\Gamma) \quad \text{ftv}(\tau) \subseteq \text{dom}(\Gamma)}{\text{wf}(\Gamma, x : \tau)}$$

Figure 3: Well-formed environments

$$\frac{\text{INST-BOT}}{\Gamma \vdash @\tau : \perp \leq \tau}$$

$$\frac{\text{INST-UNDER} \quad \Gamma, \alpha \geq \tau \vdash \phi : \tau_1 \leq \tau_2}{\Gamma \vdash \forall(\alpha \geq) \phi : \forall(\alpha \geq \tau) \tau_1 \leq \forall(\alpha \geq \tau) \tau_2}$$

$$\frac{\text{INST-ABSTR} \quad \alpha \geq \tau \in \Gamma}{\Gamma \vdash !\alpha : \tau \leq \alpha} \quad \frac{\text{INST-INSIDE} \quad \Gamma \vdash \phi : \tau_1 \leq \tau_2}{\Gamma \vdash \forall(\geq \phi) : \forall(\alpha \geq \tau_1) \tau \leq \forall(\alpha \geq \tau_2) \tau}$$

$$\frac{\text{INST-INTRO} \quad \alpha \notin \text{ftv}(\tau)}{\Gamma \vdash \wp : \tau \leq \forall(\alpha \geq \perp) \tau} \quad \frac{\text{INST-COMP} \quad \Gamma \vdash \phi_1 : \tau_1 \leq \tau_2 \quad \Gamma \vdash \phi_2 : \tau_2 \leq \tau_3}{\Gamma \vdash \phi_1; \phi_2 : \tau_1 \leq \tau_3}$$

$$\frac{\text{INST-ELIM}}{\Gamma \vdash \& : \forall(\alpha \geq \tau) \tau' \leq \tau' \{ \alpha \leftarrow \tau \}} \quad \frac{\text{INST-ID}}{\Gamma \vdash \mathbb{1} : \tau \leq \tau}$$

Figure 4: Type instance

an expression  $s_0$  for a variable  $v$  inside an expression  $s$  is written  $s\{v \leftarrow s_0\}$ .

As usual, type environments assign types to program variables. However, instead of just listing type variables, as is the case in System F, they also assign them a type bound, using the form  $\alpha \geq \tau$ . We write  $\text{dom}(\Gamma)$  for the set of all term variables and type variables that are bound by  $\Gamma$ . We also assume that typing environments are *well-formed*, *i.e.* they do not bind twice the same variable and free type variables appearing in a type of the environment  $\Gamma$  are bound earlier in  $\Gamma$ . Well-formedness rules are given in Figure 3: the empty environment is well-formed; given a well-formed environment  $\Gamma$ , the relations  $x \notin \text{dom}(\Gamma)$ ,  $\alpha \notin \text{dom}(\Gamma)$ , and  $\text{ftv}(\tau) \subseteq \text{dom}(\Gamma)$  must hold to form environments  $\Gamma, x : \tau$  and  $\Gamma, \alpha \geq \tau$ .

### 1.2. Instantiations

Instantiations  $\phi$  are defined in Figure 1. Their typing, described in Figure C.19, are *type instance* judgments of the form  $\Gamma \vdash \phi : \tau \leq \tau'$ , stating that in environment  $\Gamma$ , the instantiation  $\phi$  transforms the type  $\tau$  into the type  $\tau'$ . (For conciseness, the syntax of instantiations uses mathematical symbols  $!$ ,  $\&$ ,  $\wp$ , *etc.* which have no connection at

$$\begin{aligned}
\tau \ (!\alpha) &= \alpha \\
\perp \ (@\tau) &= \tau \\
\tau \ \mathbb{1} &= \tau \\
\tau \ (\phi_1; \phi_2) &= (\tau \ \phi_1) \ \phi_2 \\
\tau \ \wp &= \forall(\alpha \geq \perp) \ \tau \quad \alpha \notin \text{ftv}(\tau) \\
(\forall(\alpha \geq \tau) \ \tau') \ \& &= \tau' \{\alpha \leftarrow \tau\} \\
(\forall(\alpha \geq \tau) \ \tau') \ (\forall(\geq \phi)) &= \forall(\alpha \geq \tau \ \phi) \ \tau' \\
(\forall(\alpha \geq \tau) \ \tau') \ (\forall(\alpha \geq) \ \phi) &= \forall(\alpha \geq \tau) \ (\tau' \ \phi)
\end{aligned}$$

Figure 5: Type instantiation (on types)

all with linear logic.)

The *bottom* instantiation  $\text{@}\tau$  expresses that (any) type  $\tau$  is an instance of the bottom type. The *abstract* instantiation  $\text{!}\alpha$ , which assumes that the hypothesis  $\alpha \geq \tau$  is in the environment, abstracts the bound  $\tau$  of  $\alpha$  as the type variable  $\alpha$ . The *inside* instantiation  $\forall(\geq \phi)$  applies  $\phi$  to the bound  $\tau'$  of a flexible quantification  $\forall(\alpha' \geq \tau') \ \tau$ . Conversely, the *under* instantiation  $\forall(\alpha \geq) \ \phi$  applies  $\phi$  to the type  $\tau$  under the quantification; the type variable  $\alpha$  is bound in  $\phi$  and the environment in the premise of the rule INST-UNDER is increased accordingly. The *quantifier introduction*  $\wp$  introduces a fresh trivial quantification  $\forall(\alpha \geq \perp)$ . Conversely, the *quantifier elimination*  $\&$  eliminates the bound of a type of the form  $\forall(\alpha \geq \tau) \ \tau'$  by substituting  $\tau$  for  $\alpha$  in  $\tau'$ . This amounts to definitely choosing the present bound  $\tau$  for  $\alpha$ , while the bound before the application could be further instantiated by some inside instantiation. The *composition*  $\phi; \phi'$  witnesses the transitivity of type instance, while the *identity* instantiation  $\mathbb{1}$  witnesses reflexivity.

*Example.* Let  $\tau_{\min}$ ,  $\tau_{\text{cmp}}$ , and  $\tau_{\text{and}}$  be the types of the parametric minimum and comparison functions, and of the boolean conjunction:

$$\begin{aligned}
\tau_{\min} &\triangleq \forall(\alpha \geq \perp) \ \alpha \rightarrow \alpha \rightarrow \alpha \\
\tau_{\text{cmp}} &\triangleq \forall(\alpha \geq \perp) \ \alpha \rightarrow \alpha \rightarrow \text{bool} \\
\tau_{\text{and}} &\triangleq \text{bool} \rightarrow \text{bool} \rightarrow \text{bool}
\end{aligned}$$

Let  $\phi$  be the instantiation  $\forall(\geq \text{@bool}); \&$ . Then, both  $\vdash \phi : \tau_{\min} \leq \tau_{\text{and}}$  and  $\vdash \phi : \tau_{\text{cmp}} \leq \tau_{\text{and}}$  hold.

Let  $\tau_K$  be the type  $\forall(\alpha \geq \perp) \ \forall(\beta \geq \perp) \ \alpha \rightarrow \beta \rightarrow \alpha$  (e.g. of the  $\lambda$ -term  $\lambda(x) \lambda(y) x$ ) and  $\phi'$  be the instantiation  $\forall(\alpha \geq) \ (\forall(\geq \text{@}\alpha); \&)$  (the occurrence of  $\alpha$  in the inside instantiation is bound by the under instantiation). Then, the relations  $\vdash \phi' : \tau_K \leq \tau_{\min}$  holds.

*Type application.* As above, we often instantiate a quantification over  $\perp$  and immediately substitute the result. Moreover, this pattern corresponds to the System-F unique instantiation form. Therefore, we define  $\langle \tau \rangle$  as syntactic sugar for  $(\forall(\geq \text{@}\tau); \&)$ . The previous instantiations  $\phi$  and  $\phi'$  can then be abbreviated as  $\langle \text{bool} \rangle$  and  $\forall(\alpha \geq) \langle \alpha \rangle$ .

*Properties of instantiations.* Since instantiations make all steps in the instance relation explicit, their typing is deterministic.

$$\begin{array}{c}
\text{VAR} \\
\frac{x : \tau \in \Gamma}{\Gamma \vdash x : \tau} \\
\text{LET} \\
\frac{\Gamma \vdash a : \tau \quad \Gamma, x : \tau \vdash a' : \tau'}{\Gamma \vdash \text{let } x = a \text{ in } a' : \tau'}
\end{array}$$

$$\text{APP} \\
\frac{\Gamma \vdash a_1 : \tau_2 \rightarrow \tau_1 \quad \Gamma \vdash a_2 : \tau_2}{\Gamma \vdash a_1 a_2 : \tau_1}$$

$$\text{ABS} \\
\frac{\Gamma, x : \tau \vdash a : \tau'}{\Gamma \vdash \lambda(x : \tau) a : \tau \rightarrow \tau'}$$

$$\begin{array}{c}
\text{TABS} \\
\frac{\Gamma, \alpha \geq \tau' \vdash a : \tau \quad \alpha \notin \text{dom}(\Gamma)}{\Gamma \vdash \Lambda(\alpha \geq \tau') a : \forall(\alpha \geq \tau') \ \tau} \\
\text{TAPP} \\
\frac{\Gamma \vdash a : \tau \quad \Gamma \vdash \phi : \tau \leq \tau'}{\Gamma \vdash a \ \phi : \tau'}
\end{array}$$

Figure 6: Typing rules for  $x\text{MLF}$

**Lemma 1.** *If  $\Gamma \vdash \phi : \tau \leq \tau_1$  and  $\Gamma' \vdash \phi : \tau \leq \tau_2$ , then  $\tau_1 = \tau_2$ .* Coq

The use of  $\Gamma'$  instead of  $\Gamma$  may be surprising. However,  $\Gamma$  does not contribute to the instance relation, except in the side condition of rule INST-ABSTR. Hence, the type instance relation defines a partial function, called *type instantiation*<sup>3</sup> that, given an instantiation  $\phi$  and a type  $\tau$ , returns (if it exists) the unique type  $\tau \ \phi$  such that  $\Gamma \vdash \phi : \tau \leq \tau \ \phi$  holds for some  $\Gamma$ . An inductive definition of this function is given in Figure 5. Type instantiation is complete for type instance:

**Lemma 2.** *If  $\Gamma \vdash \phi : \tau \leq \tau'$ , then  $\tau \ \phi = \tau'$ .* Coq

However, the fact that  $\tau \ \phi$  may be defined and equal to  $\tau'$  does not imply that  $\Gamma \vdash \phi : \tau \leq \tau'$  holds for some  $\Gamma$ . Indeed, type instantiation does not check the premise of rule INST-ABSTR. This is intentional, as it avoids parametrizing type instantiation over the type environment. This means that type instantiation is not sound *in general*. This is never a problem, however, since we only use type instantiation originating from well-typed terms for which there always exists some context  $\Gamma$  such that  $\Gamma \vdash \phi : \tau \leq \tau'$ .

We say that types  $\tau$  and  $\tau'$  are equivalent in  $\Gamma$  if there exist  $\phi$  and  $\phi'$  such that  $\Gamma \vdash \phi : \tau \leq \tau'$  and  $\Gamma \vdash \phi' : \tau' \leq \tau$ . Although types of  $x\text{MLF}$  are *syntactically* the same as the types of  $i\text{MLF}$ —the Curry-style version of  $\text{MLF}$  (Le Botlan and Rémy, 2009)—they are richer, because type equivalence in  $x\text{MLF}$  is finer than type equivalence in  $i\text{MLF}$ , as explained in §4.

### 1.3. Typing rules for $x\text{MLF}$

Typing rules are defined in Figure 6. Compared with System F, the novelties are type abstraction and type in-

<sup>3</sup>There should never be any ambiguity between type instantiation  $\tau \ \phi$  and instantiation of expressions  $a \ \phi$ ; moreover, both operations have strong similarities and are closely related.

stantiation, unsurprisingly. The typing of a type abstraction  $\Lambda(\alpha \geq \tau) a$  extends the typing environment with the type variable  $\alpha$  bound by  $\tau$ . The typing of a type instantiation  $a \phi$  resembles the typing of a coercion, as it just requires the instantiation  $\phi$  to transform the type of  $a$  into the type of the result. Of course, it has the full power of the type application rule of System F. For example, the type instantiation  $a \langle \tau \rangle$  has type  $\tau' \{ \alpha \leftarrow \tau \}$  provided the term  $a$  has type  $\forall(\alpha) \tau'$ . As in System F, a well-typed closed term has a unique type and, in fact, a unique typing derivation.

**Lemma 3.** *If  $\Gamma \vdash a : \tau_1$  and  $\Gamma \vdash a : \tau_2$ , then  $\tau_1 = \tau_2$ . Coq*

A let-binding  $\text{let } x = a_1 \text{ in } a_2$  cannot entirely be treated as an abstraction for an immediate application  $(\lambda(x : \tau_1) a_2) a_1$  because the former does not require a type annotation on  $x$  whereas the latter does. This is nothing new, and the same as in System F extended with let-bindings. Notice however that  $\tau_1$ , which is the type of  $a_1$ , is fully determined by  $a_1$  and can be easily synthesized by a typechecker.

*Example.* Let  $\text{id}$  stand for the identity  $\Lambda(\alpha \geq \perp) \lambda(x : \alpha) x$  and  $\tau_{\text{id}}$  be the type  $\forall(\alpha \geq \perp) \alpha \rightarrow \alpha$ . We have  $\vdash \text{id} : \tau_{\text{id}}$ —much as in System F, except that unconstrained universal variables are given the bound  $\perp$ . The function choice mentioned in the introduction may be defined as  $\Lambda(\beta \geq \perp) \lambda(x : \beta) \lambda(y : \beta) x$ . It has type  $\forall(\beta \geq \perp) \beta \rightarrow \beta \rightarrow \beta$ . This is again similar to its typing in System F. We abbreviate this type as  $\tau_{\text{choice}}$ .

The application of choice to  $\text{id}$ , which we refer to below as  $\text{choice\_id}$ , may be defined as  $\Lambda(\beta \geq \tau_{\text{id}}) \text{choice } \langle \beta \rangle (\text{id } !\beta)$  and has type  $\forall(\beta \geq \tau_{\text{id}}) \beta \rightarrow \beta$ . Indeed, its typing derivation ends with:

$$\text{TAPP} \frac{\frac{\Gamma_\beta \vdash \text{choice} : \tau_{\text{choice}} \quad \Gamma_\beta \vdash \text{id} : \tau_{\text{id}}}{\Gamma_\beta \vdash \langle \beta \rangle : \tau_{\text{choice}} \leq \beta \rightarrow \beta \rightarrow \beta} \quad \Gamma_\beta \vdash !\beta : \tau_{\text{id}} \leq \beta \text{ (1)}}{\Gamma_\beta \vdash \text{choice } \langle \beta \rangle : \beta \rightarrow \beta \rightarrow \beta} \quad \Gamma_\beta \vdash \text{id } !\beta : \beta}{\Gamma \vdash \Lambda(\beta \geq \tau_{\text{id}}) \text{choice } \langle \beta \rangle (\text{id } !\beta) : \forall(\beta \geq \tau_{\text{id}}) \beta \rightarrow \beta} \text{APP} \quad \text{ABS}$$

where  $\Gamma_\beta$  is  $\Gamma, \beta \geq \tau_{\text{id}}$  and the key judgment (1), which follows by Rule INST-ABSTR, says that type  $\tau_{\text{id}}$  can be seen as type  $\beta$  whenever  $\beta$  is declared to be an instance of  $\tau_{\text{id}}$ .

The term  $\text{choice\_id}$  may also be given weaker types by type instantiation. For example,  $\text{choice\_id } \&$  has type  $(\forall(\alpha \geq \perp) \alpha \rightarrow \alpha) \rightarrow (\forall(\alpha \geq \perp) \alpha \rightarrow \alpha)$  as in System F, while  $\text{choice\_id } (\& ; \forall(\gamma \geq) (\forall(\geq \langle \gamma \rangle) ; \&))$  has the ML type  $\forall(\gamma \geq \perp) (\gamma \rightarrow \gamma) \rightarrow \gamma \rightarrow \gamma$ . The former expression can be understood directly, by fixing  $\beta$  to its bound  $\tau_{\text{id}}$ . The latter can be understood informally as the introduction of a free type variable  $\gamma$  and then the instantiation of the bound  $\tau_{\text{id}}$  of  $\beta$  to the type  $\gamma \rightarrow \gamma$ , say  $\tau_\gamma$ . Formally, the derivation is a little tedious. Let  $\Gamma$  be the typing environment  $\gamma \geq \perp$ .

$$\begin{array}{ll} (\lambda(x : \tau) a_1) a_2 & \longrightarrow a_1 \{x \leftarrow a_2\} \\ \text{let } x = a_2 \text{ in } a_1 & \longrightarrow a_1 \{x \leftarrow a_2\} \\ a \mathbb{1} & \longrightarrow a \\ a (\phi ; \phi') & \longrightarrow a \phi (\phi') \\ a \& & \longrightarrow \Lambda(\alpha \geq \perp) a \quad \alpha \notin \text{ftv}(a) \\ (\Lambda(\alpha \geq \tau) a) \& & \longrightarrow a \{! \alpha \leftarrow \mathbb{1}\} \{ \alpha \leftarrow \tau \} \\ (\Lambda(\alpha \geq \tau) a) (\forall(\alpha \geq) \phi) & \longrightarrow \Lambda(\alpha \geq \tau) (a \phi) \\ (\Lambda(\alpha \geq \tau) a) (\forall(\geq \phi)) & \longrightarrow \Lambda(\alpha \geq \tau \phi) a \{! \alpha \leftarrow \phi ; ! \alpha\} \\ E[a] & \longrightarrow E[a'] \quad \text{if } a \longrightarrow a' \end{array}$$

Figure 7: Reduction rules

First, we have

$$\begin{array}{llll} \Gamma \vdash @_\gamma & : & \perp & \leq & \gamma & (5) \\ \Gamma \vdash \forall(\geq @_\gamma) & : & \forall(\alpha \geq \perp) \alpha \rightarrow \alpha & \leq & \forall(\alpha \geq \gamma) \alpha \rightarrow \alpha & (6) \\ \Gamma \vdash \& & : & \forall(\alpha \geq \gamma) \alpha \rightarrow \alpha & \leq & \gamma \rightarrow \gamma & (4) \\ \Gamma \vdash \forall(\geq @_\gamma) ; \& & : & \forall(\alpha \geq \perp) \alpha \rightarrow \alpha & \leq & \gamma \rightarrow \gamma & (5) \end{array}$$

(2) is by rule INST-BOT; (3) is by INST-INSIDE and (2); (4) is by INST-ELIM; (5) is by INST-COMP, (3), and (4).

Then,

$$\begin{array}{llll} \Gamma \vdash \langle \gamma \rangle & : & \tau_{\text{id}} & \leq & \gamma \rightarrow \gamma \\ \Gamma \vdash \forall(\geq \langle \gamma \rangle) & : & \forall(\beta \geq \tau_{\text{id}}) \beta \rightarrow \beta & \leq & \forall(\beta \geq \gamma \rightarrow \gamma) \\ \Gamma \vdash \& & : & \forall(\beta \geq \gamma \rightarrow \gamma) \beta \rightarrow \beta & \leq & (\gamma \rightarrow \gamma) \rightarrow \gamma \\ \Gamma \vdash \forall(\geq \langle \gamma \rangle) ; \& & : & \forall(\beta \geq \tau_{\text{id}}) \beta \rightarrow \beta & \leq & (\gamma \rightarrow \gamma) \rightarrow \gamma \end{array}$$

(6) is an abbreviation of (5); (7) is by INST-INSIDE; (8) is by INST-ELIM; (9) is by INST-COMP, (7) and (8). By rule INST-UNDER and (9), we have

$$\vdash \forall(\gamma \geq) (\forall(\geq \langle \gamma \rangle) ; \&) : \forall(\gamma \geq \perp) \forall(\beta \geq \tau_{\text{id}}) \beta \rightarrow \beta \leq \forall(\gamma \geq \perp) (\gamma \rightarrow \gamma) \rightarrow \gamma \rightarrow \gamma \quad (10)$$

Finally, by rule INST-INTRO, (10), and INST-COMP, we have:

$$\text{APP} \vdash \& ; \forall(\gamma \geq) (\forall(\geq \langle \gamma \rangle) ; \&) : \forall(\beta \geq \tau_{\text{id}}) \beta \rightarrow \beta \leq \forall(\gamma \geq \perp) (\gamma \rightarrow \gamma) \rightarrow \gamma$$

As illustrated on this rather simpler example, computing all intermediate steps of a type instantiation is very tedious for a human and usually harder than just checking type instantiation. However,  $x\text{ML}^F$  is only meant to be used as an internal language by a machine.

#### 1.4. Reduction

The semantics of the calculus is given by a small-step reduction semantics. We let reduction occur in any context, including under abstractions. That is, the evaluation contexts are single-hole contexts, given by the grammar:

$$E ::= [\cdot] \mid E \phi \mid \lambda(x : \tau) E \mid \Lambda(\alpha \geq \tau) E \mid E a \mid a E \mid \text{let } x = E \text{ in } a \mid \text{let } x = a \text{ in } E$$

The reduction rules are described in Figure 7. As usual, basic reduction steps contain  $\beta$ -reduction, with the two

variants ( $\beta$ ) and ( $\beta_{\text{let}}$ ). Other basic reduction rules, related to the reduction of type instantiations and called  $\iota$ -steps, are described below. The one-step reduction is closed under the context rule. We write  $\rightarrow_{\beta}$  and  $\rightarrow_{\iota}$  for the two subrelations of  $\rightarrow$  that contain only CONTEXT and  $\beta$ -steps or  $\iota$ -step, respectively. Finally, the reduction is the reflexive and transitive closure  $\rightarrow^*$  of the one-step reduction relation.

*Reduction of type instantiation.* By definition, type instantiation redexes are all of the form  $a \phi$ . The first three rules do not constrain the form of  $a$ . The identity type instantiation is just dropped (Rule  $\iota$ -ID). A type instantiation composition is replaced by the successive corresponding type instantiations (Rule  $\iota$ -SEQ). Rule  $\iota$ -INTRO introduces a new type abstraction in front of  $a$ ; we assume that the bound variable  $\alpha$  is fresh for  $a$ .

The other three rules require the type instantiation to be applied to a type abstraction  $\Lambda(\alpha \geq \tau) a$ . Rule  $\iota$ -UNDER propagates the type instantiation under the bound, on the body  $a$ .

By contrast, Rule  $\iota$ -INSIDE propagates the type instantiation  $\phi$  inside the bound, replacing  $\tau$  by  $\tau \phi$ . However, as the bound of  $\alpha$  has changed, the domain of the type instantiation  $!\alpha$  is no more  $\tau$ , but  $\tau \phi$ . Hence, in order to maintain well-typedness, all the occurrences of the instantiation  $!\alpha$  in  $a$  must be simultaneously replaced by the instantiation  $(\phi; !\alpha)$ . Here, the instantiation  $!\alpha$  is seen as an atomic construct, *i.e.* all occurrences of  $!\alpha$  are substituted, while other occurrences of  $\alpha$  (*i.e.* that are not part of  $!\alpha$ ) are left unchanged. Formally,  $a\{\!\alpha_0 \leftarrow \phi_0\}$  is defined recursively, as described in Figure 8 (abbreviating  $\{\!\alpha_0 \leftarrow \phi_0\}$  by  $\theta$ ). The interesting lines are the two first ones of the second column, as other lines are just lifting the substitution from the leaves to types, type instantiations, and terms in the usual way.

As an example of  $\iota$ -INSIDE, if  $a$  is the term

$$\Lambda(\alpha \geq \tau) \lambda(x : \alpha \rightarrow \alpha) \lambda(y : \perp) y (@(\alpha \rightarrow \alpha)) (z !\alpha)$$

then, the type instantiation  $a (\forall (\geq \phi))$  reduces to:

$$\Lambda(\alpha \geq \tau \phi) \lambda(x : \alpha \rightarrow \alpha) \lambda(y : \perp) y (@(\alpha \rightarrow \alpha)) (z (\phi; !\alpha))$$

Rule  $\iota$ -ELIM eliminates the type abstraction, replacing all the occurrences of  $\alpha$  inside  $a$  by the bound  $\tau$ . All the occurrences of  $!\alpha$  inside  $\tau$  (used to instantiate  $\tau$  into  $\alpha$ ) become vacuous and must be replaced by the identity instantiation. For example, reusing the term  $a$  above,  $a \&$  reduces to

$$\lambda(x : \tau \rightarrow \tau) \lambda(y : \perp) y (@(\tau \rightarrow \tau)) (z \mathbb{1})$$

Finally, notice that type instantiations  $a @\tau$  and  $a !\alpha$  are irreducible.

*Examples of reduction.* Let us reuse the term `choice_id` defined in §1.3 as  $\Lambda(\beta \geq \tau_{\text{id}}) \text{choice } \langle \beta \rangle (\text{id } !\beta)$ . Remember

## Types

$$\tau \theta = \tau$$

## Terms

$$x \theta = x$$

$$(a_1 a_1) \theta = (a_1 \theta) (a_1 \theta)$$

$$(a \phi) \theta = (a \theta) (\phi \theta)$$

$$(\lambda(x : \tau) a) \theta = \lambda(x : \tau \theta) (a \theta)$$

$$(\Lambda(\alpha \geq \tau) a) \theta = \Lambda(\alpha : \tau \theta) (a \theta)$$

## Type instantiations

$$!\alpha \theta = !\alpha \quad \text{if } \alpha \neq \alpha_0$$

$$!\alpha_0 \theta = \phi_0$$

$$(@\tau) \theta = @(\tau \theta)$$

$$(\forall (\geq \phi)) \theta = \forall (\geq \phi \theta)$$

$$(\forall (\alpha \geq) \phi) \theta = \forall (\alpha \geq) (\phi \theta)$$

$$(\phi; \phi') \theta = (\phi \theta); (\phi' \theta)$$

$$\& \theta = \&$$

$$\wp \theta = \wp$$

$$\mathbb{1} \theta = \mathbb{1}$$

Figure 8: Definition of  $a \theta$ , where  $\theta$  is  $\{\!\alpha_0 \leftarrow \phi_0\}$

that  $\langle \tau \rangle$  stands for the System-F type application  $\tau$  and expands to  $(\forall (\geq @\tau); \&)$ . Therefore, the type instantiation choice  $\langle \beta \rangle$  reduces to the term  $\lambda(x : \beta) \lambda(y : \beta) x$  by  $\iota$ -SEQ,  $\iota$ -INSIDE and  $\iota$ -ELIM. Hence, the term `choice_id` reduces by these rules, CONTEXT, and ( $\beta$ ) to the expression  $\Lambda(\beta \geq \tau_{\text{id}}) \lambda(y : \beta) \text{id } !\beta$ .

Below are three specialized versions of `choice_id` (with  $\forall (\alpha) \tau$  and  $\Lambda(\alpha) a$  being abbreviations for  $\forall (\alpha \geq \perp) \tau$  and  $\Lambda(\alpha \geq \perp) a$ ). Here, all type instantiations are eliminated by reduction, but this is of course not always the case in general.

$$\text{choice\_id } (\forall (\geq \langle \text{int} \rangle); \&) : (\text{int} \rightarrow \text{int}) \rightarrow (\text{int} \rightarrow \text{int}) \rightarrow^* \lambda(y : \text{int} \rightarrow \text{int}) \lambda(x : \text{int}) x$$

$$\text{choice\_id } \& : (\forall (\alpha) \alpha \rightarrow \alpha) \rightarrow (\forall (\alpha) \alpha) \rightarrow^* \lambda(y : \forall (\alpha) \alpha \rightarrow \alpha) (\Lambda(\alpha) a)$$

$$\text{choice\_id } (\wp; \forall (\gamma \geq) (\forall (\geq \langle \gamma \rangle); \&)) : \forall (\gamma) (\gamma \rightarrow \gamma) \rightarrow (\gamma \rightarrow \gamma) \rightarrow^* \Lambda(\gamma) \lambda(y : \gamma \rightarrow \gamma) \lambda(x : \gamma) x$$

### 1.5. System F as a subsystem of $x\text{ML}^F$

System F can be seen as a subset of  $x\text{ML}^F$ , using the following syntactic restrictions: all quantifications are of the form  $\forall (\alpha) \tau$  and  $\perp$  is not a valid type anymore (however, as in System F,  $\forall (\alpha) \alpha$  is); all type abstractions are of the form  $\Lambda(\alpha) a$ ; and all type instantiations are of the form  $a \langle \tau \rangle$ . The derived typing rules for  $\Lambda(\alpha) a$  and  $a \langle \tau \rangle$  are exactly the System-F typing rules for type abstraction and type application. Hence, typechecking in this restriction of  $x\text{ML}^F$  corresponds to typechecking in System F. Moreover, the one-step System-F  $\beta$ -reduction  $(\Lambda(\alpha) a) \langle \tau \rangle \rightarrow a\{\alpha \leftarrow \tau\}$  can be performed in  $x\text{ML}^F$  in

three steps:

$$\begin{aligned}
(\Lambda(\alpha) a) \langle \tau \rangle &= (\Lambda(\alpha \geq \perp) a) (\forall (\geq @\tau); \&) & (1) \\
&\longrightarrow (\Lambda(\alpha \geq \perp) a) (\forall (\geq @\tau)) \& & (2) \\
&\longrightarrow (\Lambda(\alpha \geq \perp (\@ \tau)) a \{! \alpha \leftarrow @\tau; ! \alpha\}) \& & (3) \\
&= (\Lambda(\alpha \geq \tau) a) \& & (4) \\
&\longrightarrow a \{! \alpha \leftarrow \mathbb{1}\} \{ \alpha \leftarrow \tau \} & (5) \\
&= a \{ \alpha \leftarrow \tau \} & (6)
\end{aligned}$$

Equality (1) is by definition; step (2) is by  $\iota$ -SEQ; step (3) is by  $\iota$ -INSIDE; step (5) is by  $\iota$ -ELIM; equalities (4) and (6) are by type instantiation and by the assumption that  $a$  is a term of System F thus in which  $! \alpha$  cannot appear.

Conversely, if a term  $a$  is in System F, then its reduction steps in  $x\text{MLF}$  are all of these forms but possibly interleaved. Formally, the Church-Rosser property and the strong normalization lemma stated in §2.2 ensure that any reduction of  $a$  in  $x\text{MLF}$  will eventually terminate with the same normal form, hence with its normal form in System F.

## 2. Properties of reduction

### 2.1. Subject reduction

Reduction in  $x\text{MLF}$ , which can occur in any context, preserves typings. This relies on weakening and substitution lemmas for both instance and typing judgments.

**Lemma 4 (Weakening).** *Let  $\Gamma, \Gamma', \Gamma''$  be a well-formed environment.*

*If  $\Gamma, \Gamma'' \vdash \phi : \tau_1 \leq \tau_2$ , then  $\Gamma, \Gamma', \Gamma'' \vdash \phi : \tau_1 \leq \tau_2$ .*

*If  $\Gamma, \Gamma'' \vdash a : \tau'$ , then  $\Gamma, \Gamma', \Gamma'' \vdash a : \tau'$ .* Coq

**Lemma 5 (Term substitution).**

*If  $\Gamma, x : \tau', \Gamma' \vdash \phi : \tau_1 \leq \tau_2$  then  $\Gamma, \Gamma' \vdash \phi : \tau_1 \leq \tau_2$ .*

*Suppose  $\Gamma \vdash a' : \tau'$ ; if  $\Gamma, x : \tau', \Gamma' \vdash a : \tau$  then  $\Gamma, \Gamma' \vdash a \{x \leftarrow a'\} : \tau$ .* Coq

The next lemma, which expresses that we can substitute an instance bound inside judgments, ensures the correctness of Rule  $\iota$ -ELIM.

**Lemma 6 (Bound substitution).**

*Let  $\vartheta$  and  $\theta$  be respectively the substitutions  $\{\alpha \leftarrow \tau\}$  and  $\{! \alpha \leftarrow \mathbb{1}\} \{ \alpha \leftarrow \tau \}$ .*

*If  $\Gamma, \alpha \geq \tau, \Gamma' \vdash \phi : \tau_1 \leq \tau_2$  then  $\Gamma, \Gamma' \vartheta \vdash \phi \theta : \tau_1 \vartheta \leq \tau_2 \vartheta$ .*

*If  $\Gamma, \alpha \geq \tau, \Gamma' \vdash a : \tau'$  then  $\Gamma, \Gamma' \vartheta \vdash a \theta : \tau' \vartheta$ .* Coq

The result below ensures in turn the correctness of rule  $\iota$ -INSIDE.

**Lemma 7 (Narrowing).** *Assume  $\Gamma \vdash \phi : \tau \leq \tau'$ . Let  $\theta$  be  $\{! \alpha \leftarrow \phi; ! \alpha\}$ .*

*If  $\Gamma, \alpha \geq \tau, \Gamma' \vdash \phi' : \tau_1 \leq \tau_2$  then  $\Gamma, \alpha \geq \tau', \Gamma' \vdash \phi' \theta : \tau_1 \leq \tau_2$ .*

*If  $\Gamma, \alpha \geq \tau, \Gamma' \vdash a : \tau''$  then  $\Gamma, \alpha \geq \tau', \Gamma' \vdash a \theta : \tau''$ .* Coq

Subject reduction is an easy consequence of all these results.

**Theorem 8 (Subject reduction).**

*If  $\Gamma \vdash a : \tau$  and  $a \longrightarrow a'$  then,  $\Gamma \vdash a' : \tau$ .* Coq

### 2.2. Confluence

**Theorem 9.** *The relation  $\longrightarrow_\beta$  is confluent. The relations  $\longrightarrow_\iota$  and  $\longrightarrow$  are confluent on the terms well-typed in some context.*

This result is proved using the standard technique of parallel reductions (Barendregt, 1984). The proof is uninteresting and omitted here; it can be found in (Yakobowski, 2008).

Confluence means that  $\beta$ -reduction and  $\iota$ -reduction are independent. For instance,  $\iota$ -reductions can be performed under  $\lambda$ -abstractions as far as possible while keeping a weak evaluation strategy for  $\beta$ -reduction.

The restriction to well-typed terms for the confluence of  $\iota$ -reduction is due to two things. First, the rule  $\iota$ -INSIDE is not applicable to ill-typed terms in which  $\tau \phi$  cannot be computed, (for example  $(\Lambda(\alpha \geq \text{int}) a) (\forall (\geq \&))$ ). Second,  $\tau \phi$  can sometimes be computed, even though  $\Gamma \vdash \phi : \tau \leq \tau'$  never holds, typically if  $\phi$  is  $! \alpha$  and  $\tau$  is not the bound of  $\alpha$  in  $\Gamma$ . Hence, type errors may be either revealed or silently reduced and perhaps eliminated, depending on the reduction path. As an example, let  $a$  be the term

$$(\Lambda(\alpha \geq \forall(\gamma) \gamma) ((\Lambda(\beta \geq \text{int}) x) (\forall (\geq !\alpha)))) (\forall (\geq \&))$$

It is ill-typed in any context, because  $! \alpha$  coerces a term of type  $\forall(\gamma) \gamma$  into one of type  $\alpha$ , but  $! \alpha$  is here indirectly applied to a term of type  $\text{int}$ . If we reduce the outermost type instantiation first, we are stuck with  $\Lambda(\alpha \geq \perp) ((\Lambda(\beta \geq \text{int}) x) (\forall (\geq \& !\alpha)))$ , which is irreducible since the type instantiation  $\text{int} (\& !\alpha)$  is undefined.

Conversely, if we reduce the innermost type instantiation first, the faulty type instantiation disappears and we obtain the term  $(\Lambda(\alpha \geq \forall(\gamma) \gamma) \Lambda(\beta \geq \alpha) x) (\forall (\geq \&))$ , which further reduces to the normal form  $\Lambda(\alpha \geq \perp) \Lambda(\beta \geq \alpha) x$ .

The fact that ill-typed terms may not be confluent is not new: for instance, this is already the case with  $\eta$ -reduction in System F. We believe this is not a serious issue. In practice, this means that typechecking should be performed before any program simplification, which is usually the case anyway.

### 2.3. Termination of reduction

The termination of reduction has been proved by Manzonetto and Tranquilli (2010).

**Theorem 10.** (Manzonetto-Tranquilli) *The reduction  $\longrightarrow$  is terminating.*

As a corollary of this result and of Theorem 9, we have immediately

**Corollary 11.** *The relation  $\longrightarrow$  is strongly normalizing.*

The proof of Theorem 10 is by translation of  $x\text{MLF}$  into System F, where reductions are known to terminate, and by showing a simulation between reduction in  $x\text{MLF}$



and reduction of the elaborated term in System F. (This is also discussed in §5.1.) As a corollary,  $\rightarrow_{\iota}$  alone is also terminating. The termination of  $\rightarrow$  is useful but not critical, as  $x\text{MLF}$  is meant to be used in a language with general recursion. However, the termination of  $\rightarrow_{\iota}$  is essential for  $x\text{MLF}$  to have a type-erasure semantics.

#### 2.4. Type-erasure semantics

The reduction has been defined so that the type erasure of a reduction sequence in  $x\text{MLF}$  is a reduction sequence in the untyped  $\lambda$ -calculus. Formally, the type erasure of a term  $a$  of  $x\text{MLF}$  is the untyped  $\lambda$ -term  $[a]$  defined inductively by

$$\begin{aligned} [x] &= x & [\text{let } x = a_1 \text{ in } a_2] &= \text{let } x = [a_1] \text{ in } [a_2] \\ [a \phi] &= [a] & [\lambda(x : \tau) a] &= \lambda(x) [a] \\ [a_1 a_2] &= [a_1] [a_2] & [\Lambda(\alpha \geq \tau) a] &= [a] \end{aligned}$$

It is immediate to verify that two terms related by  $\iota$ -reduction have the same type erasure. Moreover, if a term  $a$   $\beta$ -reduces to  $a'$ , then the type erasure of  $a$   $\beta$ -reduces to the type erasure of  $a'$  in one step in the untyped  $\lambda$ -calculus.

**Lemma 12.** *If  $a \rightarrow_{\iota} a'$  then  $[a] = [a']$ . If  $a \rightarrow_{\beta} a'$ , then  $[a] \rightarrow_{\beta} [a']$ .*

The converse direction is also true:

**Lemma 13.** (Manzonetto-Tranquilli) *If  $[a] \rightarrow_{\beta} M$ , then there exist  $a'$  and  $a''$  such that  $a \rightarrow_{\iota}^* a' \rightarrow_{\beta} a''$  and  $[a''] = M$ .*

A proof has been given by Manzonetto and Tranquilli (2010, Appendix B<sup>4</sup>). Combining these two results ensures that  $x\text{MLF}$  has a type-erasure semantics.

#### 2.5. Accommodating weak reduction strategies and constants

In order to show that the calculus may also be used as the core of a programming language, we now introduce constants and we restrict the semantics to a weak evaluation strategy. We then show that subject reduction and progress hold for the main two forms of weak-reduction strategies, namely call-by-value and call-by-name.

We let the letter  $c$  range over constants. Each constant comes with its arity  $|c|$ . The dynamic semantics of constants must be provided by primitive reduction rules, called  $\delta$ -rules. However, these are usually of a certain form. To characterize  $\delta$ -rules (and values), we partition constants into *constructors* and *primitives*, ranged over by letters  $C$  and  $f$ , respectively. The difference between the two lies in their semantics: primitives (such as  $+$ ) are reduced when

fully applied, while constructors (such as `cons`) are irreducible and typically eliminated when passed as argument to primitives.

In order to classify constructed values, we assume given a collection of type constructors  $\kappa$ , together with their arities  $|\kappa|$ . We extend types with constructed types  $\kappa (\tau_1, \dots, \tau_{|\kappa|})$ . We write  $\bar{\alpha}$  for a sequence of variables  $\alpha_1, \dots, \alpha_k$  and  $\forall(\bar{\alpha}) \tau$  for the type  $\forall(\alpha_1) \dots \forall(\alpha_k) \tau$ . The static semantics of constants is given by an initial typing environment  $\Gamma_0$  that assigns to every constant  $c$  a type  $\tau$  of the form  $\forall(\bar{\alpha}) \tau_1 \rightarrow \dots \tau_{|c|} \rightarrow \tau_0$ , where  $\tau_0$  is a constructed type (hence neither bottom, a variable or an arrow type) whenever the constant  $c$  is a constructor.

We distinguish a subset of terms, called *values* and *evaluated*, that are term abstractions, type abstractions, full or partial applications of constructors, or partial applications of primitives. We use an auxiliary letter  $w$  to characterize the arguments of functions, which differ for call-by-value and call-by-name strategies. In values, an application of a constant  $c$  can involve a series of type instantiations, but only evaluated ones and placed before all other arguments. Moreover, the application may only be partial whenever  $c$  is a primitive. Evaluated instantiations  $\theta$  may be quantifier eliminations or either inside or under (general) instantiations. In particular,  $a (@\tau)$  and  $a (!\alpha)$  are *never* values. The grammar for values and evaluated instantiations is as follows:

$$\begin{aligned} v &::= \lambda(x : \tau) a \\ &| \Lambda(\alpha : \tau) a \\ &| C \theta_1 \dots \theta_k w_1 \dots w_n & n \leq |C| \\ &| f \theta_1 \dots \theta_k w_1 \dots w_n & n < |f| \\ \theta &::= \forall(\geq \phi) \mid \forall(\alpha \geq) \phi \mid \& \end{aligned}$$

Importantly, values cannot have type  $\perp$ :

**Lemma 14.** *If  $v$  is a value and  $if \vdash v : \tau$ , then  $\tau$  is not  $\perp$ .*

(Proof p. 25)

Finally, we assume that  $\delta$ -rules are of the form  $f \theta_1 \dots \theta_k w_1 \dots w_{|f|} \rightarrow_f a$  (that is,  $\delta$ -rules may only reduce fully applied primitives).

In addition to this general setting, we make further assumptions to relate the static and dynamic semantics of constants.

**SUBJECT REDUCTION:**  $\delta$ -reduction preserves typings, *i.e.*, for any typing context  $\Gamma$  such that  $\Gamma \vdash a : \tau$  and  $a \rightarrow_f a'$ , the judgment  $\Gamma \vdash a' : \tau$  holds.

**PROGRESS:** Well-typed, full applications of primitives can be reduced, *i.e.*, for any term  $a$  of the form  $f \theta_1 \dots \theta_k w_1 \dots w_{|f|}$  verifying  $\Gamma_0 \vdash a : \tau$ , there exists a term  $a'$  such that  $a \rightarrow_f a'$ .

#### Call-by-value reduction

We now specialize the previous setting to a call-by-value semantics. In this case, arguments of applications in

<sup>4</sup>The indirect proof given in §4 is not correct, since it relies on the subject reduction property for their intermediate System Fc, which unfortunately does not hold.

values are themselves restricted to values, *i.e.*  $w$  is taken equal to  $v$ . Reduction rules of Figure 7 are modified as follows. Rules  $(\beta)$  and  $(\beta_{\text{let}})$  are limited to the substitution of values, that is, to reductions of the form  $(\lambda(x : \tau) a) v \rightarrow a\{x \leftarrow v\}$  and  $\text{let } x = v \text{ in } a \rightarrow a\{x \leftarrow v\}$ . Rules  $\iota\text{-ID}$ ,  $\iota\text{-SEQ}$  and  $\iota\text{-INTRO}$  are also restricted so that they only apply to values (*e.g.*  $a$  is textually replaced by  $v$  in each of these rules). Finally, we restrict rule  $\text{CONTEXT}$  to call-by-value contexts, which are of the form

$$E_v ::= [\cdot] \mid E_v a \mid v E_v \mid E_v \phi \mid \text{let } x = E_v \text{ in } a$$

We write  $\rightarrow_v^*$  the resulting reduction relation. It follows from the above restrictions that the reduction is deterministic. Moreover, since  $\delta$ -reduction preserves typings, by assumption, the relation  $\rightarrow_v^*$  also preserves typings by Theorem 8. Hence, in combination with progress, stated next, the evaluation of well-typed terms “cannot go wrong”.

**Theorem 15 (Progress for call-by-value).**

If  $\Gamma_0 \vdash a : \tau$ , then either  $a$  is a value or  $a \rightarrow_v^* a'$  for some  $a'$ .

(Proof p. 25)

*Call-by-value reduction and the value restriction*

The value-restriction is the standard way of adding side effects in a call-by-value language. We verify that it can be transposed to  $x\text{MLF}$ .

Typically, the *value restriction* amounts to restricting type generalization to non-expansive expressions, that cannot have direct or indirect side effects. Those contain at least value-forms, *i.e.* values and term variables, as well as their type-instantiations. In the case of  $x\text{MLF}$ , which is a target language and not a source one, we obtain a restricted grammar of (potentially) expansive expressions  $a$ , and a subset which is constituted of non-expansive expressions  $u$ .

$$\begin{array}{lcl} a & ::= & u \mid a a \mid \text{let } x = u \text{ in } a \\ u & ::= & x \mid \lambda(x : \tau) a \mid \Lambda(\alpha : \tau) u \mid u \phi \mid \text{let } x = u \text{ in } u \\ & & \mid C \theta_1 \dots \theta_k u_1 \dots u_n \qquad n \leq |C| \\ & & \mid f \theta_1 \dots \theta_k u_1 \dots u_n \qquad n < |f| \end{array}$$

As usual, we restrict let-bound expressions to be non-expansive, since they implicitly contain a type generalization. Hence, a let-bound expression is expansive when its body is expansive—but it remains non-expansive when its body is non-expansive. Notice that, although type instantiations are restricted to non-expansive expressions, this is not a limitation:  $b \phi$  can always be written as  $(\lambda(x : \tau) x \phi) b$ , where  $\tau$  is the type of  $b$ , and similarly for applications of constants to expansive expressions.

Lemma 16, stated below, ensures two things: our restricted grammar has a meaning as a standalone language (as it is stable by reduction); and non-expansive expressions are closed by reduction and are thus harmless in presence of side-effects.

**Lemma 16.** *Expansive and non-expansive expressions are closed by call-by-value reduction.*

As an immediate consequence:

**Corollary 17.** *Subject reduction holds with the value restriction.*

It is then routine work to extend the semantics with a global store to model side effects and verify type soundness for this extension.

*Call-by-name reduction*

In call-by-name reduction semantics, values may contain applications of constants to arbitrary expressions—and not just to values. That is, we take  $a$  for  $w$ . The  $\iota$ -reduction is restricted as for call-by-value, while  $\rightarrow_\beta$  is unchanged. However, evaluation contexts are now  $E_n ::= [\cdot] \mid E_n a \mid E_n \phi$ .

We write  $\rightarrow_n^*$  the resulting reduction relation. As for call-by-value, it is deterministic by construction and preserves typings. Moreover, it may always progress. Hence, call-by-name evaluation of well-typed terms “cannot go wrong”.

**Theorem 18 (Progress for call-by-name).**

If  $\Gamma_0 \vdash a : \tau$ , then either  $a$  is a value or  $a \rightarrow_n^* a'$  for some  $a'$ .

(Proof p. 25)

### 3. Elaboration of graphical $e\text{MLF}$ into $x\text{MLF}$

To verify that, as expected,  $x\text{MLF}$  can be used as an internal language for  $e\text{MLF}$ , we now exhibit a type-preserving type-erasure-preserving translation from  $e\text{MLF}$  to  $x\text{MLF}$ . We use the graphic constraint presentation of  $e\text{MLF}$  (Rémy and Jakobowski, 2008; Jakobowski, 2008) which is more general than the syntactic presentation (Le Botlan and Rémy, 2003, 2009) and also better suited for type inference.

The elaboration of  $e\text{MLF}$  into  $x\text{MLF}$  proceeds in two phases. The first phase is just type inference in  $e\text{MLF}$ , described by Rémy and Jakobowski (2008) and Jakobowski (2008). A source program of  $e\text{MLF}$  is translated into a typing constraint, which can be seen as a decoration of the source program with (1) placeholders for missing types, and (2) type instantiation constraints that relate types (either known or unknown). The constraint is then solved, filling in all unknown types so that all type instantiation constraints become valid. The result of type inference is called a presolution.

The second phase translates a presolution into a term of  $x\text{MLF}$ . The main difficulty is to infer for each instantiation constraint a precise description of the type instantiation steps. Interestingly, this is done by replaying type inference with an instrumented algorithm. More precisely, the instantiation steps are extracted from the proof that the presolution found by type inference is indeed in solved form. It then remains to translate the instrumented presolution, which is represented graphically, into a syntactic form, *i.e.* a term of  $x\text{MLF}$ . This second phase is a form of compilation, which is technically not very deep, but meticulous.

Since the elaboration is based on—and starts with—type inference, it contains many details that require some minimal understanding of  $e\text{MLF}$ . Hence we present an overview of  $e\text{MLF}$ . Sill, other reading might help (Rémy and Jakobowski, 2007, 2008; Jakobowski, 2008). As no other part depends on §3, most details (or even the whole section) can also be skipped in a first reading of the paper.

*Outline.* We first review the graphic constraints type inference framework (§3.1); we then present the main steps of the translation (§3.2); finally, we describe the key steps in details (§3.3-3.5). The elaboration has been implemented in a prototype by Scherer (2010a).

#### 3.1. An overview of graphical $e\text{MLF}$

A full presentation of graphical  $e\text{MLF}$  is out of the scope of this paper. In this section, we only remind the key points about graphic types and associated type instance, which is the basis of the elaboration algorithm. We put more emphasis on the aspects of graphic types that either depart significantly from more traditional syntactic presentation of types, or that play a key role in understanding the elaboration process. Detailed presentations can be found in (Rémy and Jakobowski, 2007, 2008; Jakobowski, 2008).

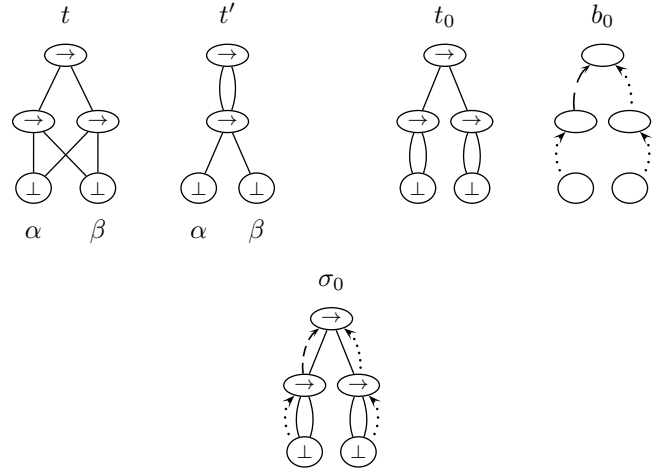


Figure 9: Dags and graphic types

#### 3.1.1. Graphic types

Types of graphical  $e\text{MLF}$  are graphs, designated with letter  $\sigma$ , composed of the superposition of a *term-dag*, representing the structure of the type, and of a *binding tree* encoding polymorphism.

Term-dags are just dag representations of usual tree-like types where all occurrences of the same variable are shared, and inner nodes representing identical subtypes may also be shared. We write  $\sigma(n)$  for the constructor at node  $n$ . Variables are anonymous and represented by the pseudo-constructor  $\perp$ . Term-dag edges are written  $n \circ^i \rightarrow m$ , where  $i$  is an integer that ranges between 1 and the arity of  $\sigma(n)$ ; we also use the notation  $\langle ni \rangle$  to designate  $m$ , the root node being simply noted  $\langle \rangle$ . On pictures, edges are drawn with plain lines, oriented downwards; we leave  $i$  implicit, as outgoing edges are always drawn from left to right.

**Example 1.** The dag  $t$  on Figure 9 represents the first-order type  $(\alpha \rightarrow \beta) \rightarrow (\alpha \rightarrow \beta)$ . The nodes  $\langle 11 \rangle$  and  $\langle 22 \rangle$  are variables (the names  $\alpha$  and  $\beta$  are here to help reading the figure, but formally they are not part of the graphic type). Compared with the tree notation, leaves representing the same variable are merged together; the names of leaves are left anonymous. That is, paths 11 and 21 lead to the same node, which can therefore be designated by  $\langle 11 \rangle$  or  $\langle 21 \rangle$ , indifferently. Similarly, paths 12 and 22 lead to the same node.

The dag structure also allows sharing internal nodes whose subtrees are identical as described by the dag  $t'$  where nodes  $\langle 1 \rangle$  and  $\langle 2 \rangle$  coincide. The dag  $t'$  could be syntactically written as  $(\text{let } \gamma = \alpha \rightarrow \beta \text{ in } \gamma \rightarrow \gamma)$ . In fact, sharing of internal nodes is a key to the efficient implementation of unification algorithms on first-order types. Those typically see  $t'$  as an instance of  $t$ , but not the converse; thus sharing can only be increased, and never lost. However, this refinement of the instance relation needs not be revealed externally, and dag  $t'$  can be displayed as dag  $t$  by splitting (or just reading back) shared internal nodes into separate

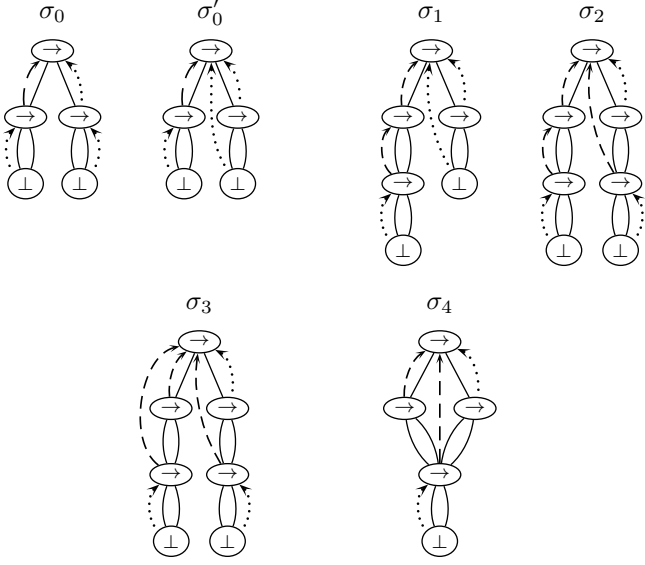


Figure 10: Examples of instance on graphic types

ones.

The second component of graphic types, the binding tree, is an upside-down tree with an edge  $n \succ_{\diamond} m$  leaving from each node  $n$  different from the root, and going to some node  $m$  upper in the term-dag at which  $n$  is bound. Binding edges may be either flexible or rigid, which is represented by labeling the edge with  $\succ_{\geq}$  or  $\succ_{=}$ , respectively. On drawings, these flags are represented by dotted or dashed lines, respectively. We use the flag metavariable  $\diamond$  to range over  $\succ_{\geq}$  and  $\succ_{=}$ .

**Example 2.** Consider the graphic type  $\sigma_0$  of Figure 9. It is the superposition of the first-order term-dag  $t_0$  and the binding tree  $b_0$ . The edge  $\langle 22 \rangle \succ_{\geq} \langle 2 \rangle$  is a flexible binding edge (the rightmost lowermost one), while  $\langle 1 \rangle \succ_{=} \langle \rangle$  is a rigid binding edge (the leftmost uppermost one) and  $\langle 1 \rangle \circ^2 \langle 12 \rangle$  is a structure edge.

Binding edges express polymorphism. They are oriented, and the target of the edge indicates the place where the binding occurs. The node at the source of the edge represents the variable being introduced, while the subtree at that node is the bound of that variable. Binding edges are of two kinds: a *rigid* edge means that polymorphism is required; typically, it is used for the type of an argument that is used polymorphically. By contrast, a *flexible* edge means that polymorphism is available (as with flexible quantification in  $x\text{MLF}$ ) but not required.

**Example 3 (cont.).** The type  $\sigma_0$  of Figures 9 and 10 describes a function  $f$  whose argument must be at least as polymorphic as  $\forall(\alpha) \alpha \rightarrow \alpha$ , and whose result has type  $\forall(\beta) \beta \rightarrow \beta$ , or any instance of it. In other words, the result of an application of  $f$  can be used in place of the

successor function of type  $\text{int} \rightarrow \text{int}$ , but  $f$  cannot be passed the successor function as argument, which is not as polymorphic as required.

The type  $\sigma'_0$  of Figure 10 describes a polymorphic function that, given a type  $\gamma$ , expects an argument of type  $\forall(\alpha) \alpha \rightarrow \alpha$  and returns a value of type  $\gamma \rightarrow \gamma$ . In particular,  $\sigma'_0$  is strictly less polymorphic than  $\sigma_0$ , as in System-F, since  $\gamma \rightarrow \gamma$  is a strict instance of  $\forall(\beta) \beta \rightarrow \beta$ .

Rigid bounds arise from type annotations: the principal type of a term that contains no type annotations (in an environment that contains no types with rigid bounds), uses only flexible bounds. That is, required polymorphism may be offered by type inference, but never requested automatically.

*Classifying nodes.* For the purpose of defining type instance, we distinguish four kinds of nodes according to their position in the binding tree. The kind of each node is used below to determine how they can be transformed during type instantiations. Hence, this classification plays an important role in the translation.

Nodes on which no variable is transitively flexibly bound are called *inert*, as they neither hold nor control polymorphism. They will be discussed in detail further on. All other nodes hold or control some polymorphism and are classified as follows. Nodes whose binding path is flexible up to the root are called *instantiable*: they can be freely instantiated as described in the next section; in  $x\text{MLF}$  these nodes correspond to parts of types that can be transformed by a suitable instantiation expression. Nodes whose binding edge is rigid are called *restricted*, because they cannot be grafted; in  $x\text{MLF}$  they roughly correspond to polymorphic types occurring under some arrow type. Nodes whose binding edge is flexible but whose binding path up to the root contains a rigid edge are called *locked*; they cannot be transformed in any way. In  $x\text{MLF}$ , these nodes roughly correspond to polymorphic types occurring in the bound of quantifiers themselves under some arrow type—they offer polymorphism that is requested and cannot be diminished.

**Example 4 (cont.).** In the type  $\sigma'_0$  of Figure 10, the node  $\langle 2 \rangle$  is inert,  $\langle 21 \rangle$  is instantiable,  $\langle 1 \rangle$  is restricted and  $\langle 11 \rangle$  is locked.

*Type instance.* The instance relation on graphic types, written  $\sqsubseteq$ , can be described as the composition of four atomic operations: *grafting*, *merging*, *raising*, and *weakening*. All four operations are detailed below, and depicted schematically in Figure 11. In the figure, we use the following conventions to constrain the position of nodes in the binding tree: the green (or light gray) node with dotted border is instantiable; blue (or darker gray) nodes with double-line borders are anything but locked; small white nodes are unconstrained.

- $\text{Graft}(\sigma, n)$ , called *grafting*, replaces an instantiable bottom node  $n$  by a closed graph  $\sigma$ . Grafting corresponds

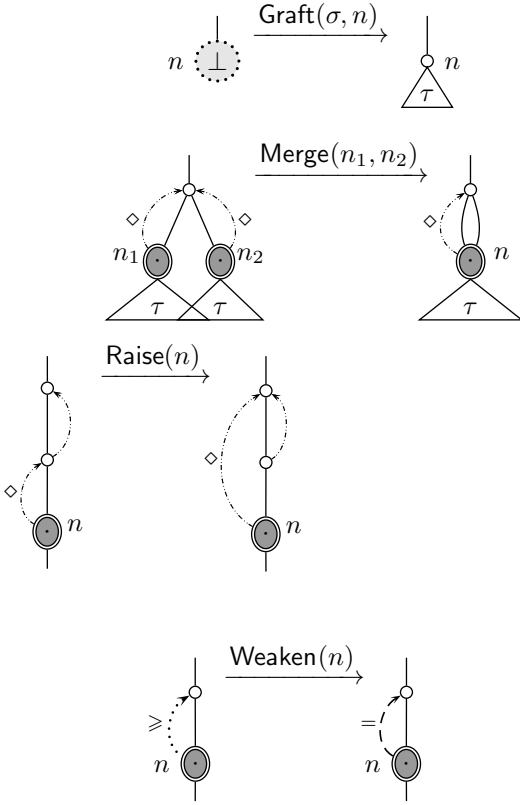


Figure 11: Atomic graphic instance operations

to the INST-BOT rule of  $x\text{MLF}$ . That is,  $\Gamma \vdash @\tau : \perp \leq \tau$  where  $\tau$  is the type describing the graph  $\sigma$ .

- **Merge**( $n_1, n_2$ ), called *merging*, fuses two nodes  $n_1$  and  $n_2$  that are not locked and have identical subgraphs. After merging, the subgraphs will thus be shared and can only be instantiated synchronously. In  $x\text{MLF}$  terms, it replaces two identical quantifications by a unique one, as in

$$\Gamma \vdash \phi : \forall(\alpha \geq \tau) \forall(\beta \geq \tau) \tau' \leq \forall(\alpha \geq \tau) \tau' \{\beta \leftarrow \alpha\}$$

with, for instance,  $\phi$  equal to  $\forall(\alpha \geq) (\forall(\geq !\alpha); \&)$ .

- **Raise**( $n$ ), called *raising*, makes the binding of a node  $n$  that is not locked slide other the binding edge above it. Raising corresponds to a scope extrusion in  $x\text{MLF}$ , as in

$$\Gamma \vdash \phi : \forall(\alpha \geq \forall(\beta \geq \tau) \tau') \tau'' \leq \forall(\beta \geq \tau) \forall(\alpha \geq \tau') \tau''$$

with, for instance,  $\phi$  equal to  $\wp; \forall(\geq @\tau); \forall(\beta \geq) (\forall(\geq \forall(\geq !\beta)); \&)$ .

- **Weaken**( $n$ ), called *weakening*, changes the binding of a flexible node  $n$  that is not locked into a rigid one. This freezes the subgraph under the node, preventing further instance operations on non-inert nodes, and all graftings. When this operation occurs on an instantiable node, it corresponds to the  $x\text{MLF}$  INST-ELIM instantiation:

$$\Gamma \vdash \& : \forall(\alpha \geq \tau) \tau' \leq \tau' \{\alpha \leftarrow \tau\}$$

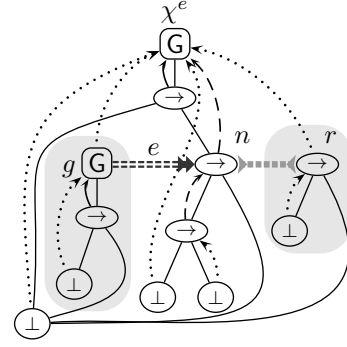


Figure 12: Constraints and expansion

Notice that grafting and merging do not change the bindings of existing nodes, while conversely, raising and weakening only change the bindings of existing nodes.

**Example 5 (cont.).** The type  $\sigma'_0$  of Figure 10 is an instance of  $\sigma_0$  obtained by raising (21). The type  $\sigma_4$  is an instance of  $\sigma_1$ , obtained by grafting then weakening (21) (resulting in  $\sigma_2$ ), raising the node  $\langle 11 \rangle$  (which gives  $\sigma_3$ ), and finally merging  $\langle 11 \rangle$  and  $\langle 21 \rangle$ . Letting  $\sigma$  be the graph corresponding to  $\forall(\alpha) \alpha \rightarrow \alpha$ , we may formally write:

$$\begin{aligned} \sigma_1 &\xrightarrow{\text{Graft}(\sigma, \langle 21 \rangle)} \xrightarrow{\text{Weaken}(\langle 21 \rangle)} \sigma_2 \xrightarrow{\text{Raise}(\langle 11 \rangle)} \\ &\xrightarrow{\text{Merge}(\langle 11 \rangle, \langle 21 \rangle)} \sigma_4 \end{aligned}$$

Hence, the instance  $g_1 \sqsubseteq g_2$  is witnessed by the transformation

$$\text{Graft}(\sigma, \langle 21 \rangle); \text{Weaken}(\langle 21 \rangle); \text{Raise}(\langle 11 \rangle); \text{Merge}(\langle 11 \rangle, \langle 21 \rangle)$$

where “;” is the composition with arguments given in reverse order.

*On the importance of inert nodes.* While inert nodes carry no polymorphism, it is important to treat them especially—so as to allow slightly more instance operations. Intuitively, since these nodes carry no polymorphism, they need not be shared, nor do they need a binding edge. However, it is technically more regular to let every node but the root node be bound to some other node, which we do. Furthermore, we only allow raising, merging or weakening those nodes, not the converse operations; §3.3 will justify why this is technically possible.

### 3.1.2. Type constraints

Type constraints are used to formalize  $\text{MLF}$  typing problems. They generalize graphic types by adding new forms of edges, called constraint edges. These can be either *unification edges*  $\dashv\dashv\dashv$  or *instantiation edges*  $\dashv\dashv\dashv$ . They also generalize let-constraints that have been proposed for type inference in ML by Pottier and Rémy (2005). Instantiation edges are oriented. They relate special nodes, used to represent type schemes and called G-nodes, to regular

nodes. An example of a constraint  $\chi^e$  is shown on Figure 12. The instance on type constraints is exactly as on graphic types—constraint edges are just preserved.

A unification edge is solved when it relates a node to itself (thus, a unification edge forces the nodes it relates to be merged). An instantiation edge  $e$  of the form  $g \dashrightarrow n$  of a constraint  $\chi$  is solved when, informally,  $n$  is an instance of the type scheme represented by  $g$ , or formally, when the expansion of  $e$  in  $\chi$  (defined below) is an instance of  $\chi$ .

A type constraint is solved when all of its constraint edges are solved. A *presolution* of a constraint is one of its solved instances. It still contains all the nodes of the original constraint, many of which may have become irrelevant for describing the resulting type. A *solution* of a constraint is, roughly, a presolution in which such nodes have been removed. We need not define solutions formally since the translation uses presolutions directly.

*Expansion.* In a constraint  $\chi$ , consider an instantiation edge  $e$  defined as  $g \dashrightarrow n$ . We define an *expansion* operation that enforces the constraint represented by  $e$ . The expansion of  $e$  in  $\chi$ , written  $\chi^e$ , is the constraint  $\chi$  extended with both a copy of the type scheme represented by  $g$  and a unification edge between  $n$  and the root  $r$  of the copy. The copy is bound at the same node as  $n$ . Technically, we define the *interior* of  $g$ , written  $\mathcal{I}(g)$  as all the nodes transitively bound to  $g$ . The expansion operation copies all the nodes structurally strictly under  $g$  and in the interior of  $g$ . Intuitively, those nodes are generic at the level of  $g$ . Conversely, the nodes under  $g$  that are not in the interior of  $g$  are not generic at the level of  $g$  and are not copied by the expansion<sup>5</sup> (but are instead shared with the original).

**Example 6.** Let us consider the expansion  $\chi^e$  of Figure 12. The original constraint  $\chi$  can be obtained from  $\chi^e$  by removing the rightmost highlighted nodes, as well as the resulting dangling edges. The interior of  $g$  is composed of the leftmost highlighted nodes. Hence, the copied nodes are  $\langle g1 \rangle$  and  $\langle g11 \rangle$ , but not  $\langle g12 \rangle$ , which is not in  $\mathcal{I}(g)$ . The root of the expansion  $r$  is the copy of  $\langle g1 \rangle$ . It is bound to the bound of  $n$  and connected to  $n$  by an unification edge.

By definition, we say that an instantiation edge  $e$  is *solved* when  $\chi$  is an instance of  $\chi^e$ . This indeed means that the subtype constrained by the instantiation edge is less general than the type scheme at the origin of the edge—as a copy of this scheme can be instantiated into the subtype. We call *instantiation witness* an instance derivation of  $\chi^e \sqsubseteq \chi$  for a solved instantiation edge  $e$ .

<sup>5</sup>Readers familiar with MLF (Rémy and Jakobowski, 2008) may notice a slight change in terminology, as in this work we use the term “expansion” instead of “propagation”, and we solve frontier unification edges on the fly, for conciseness.

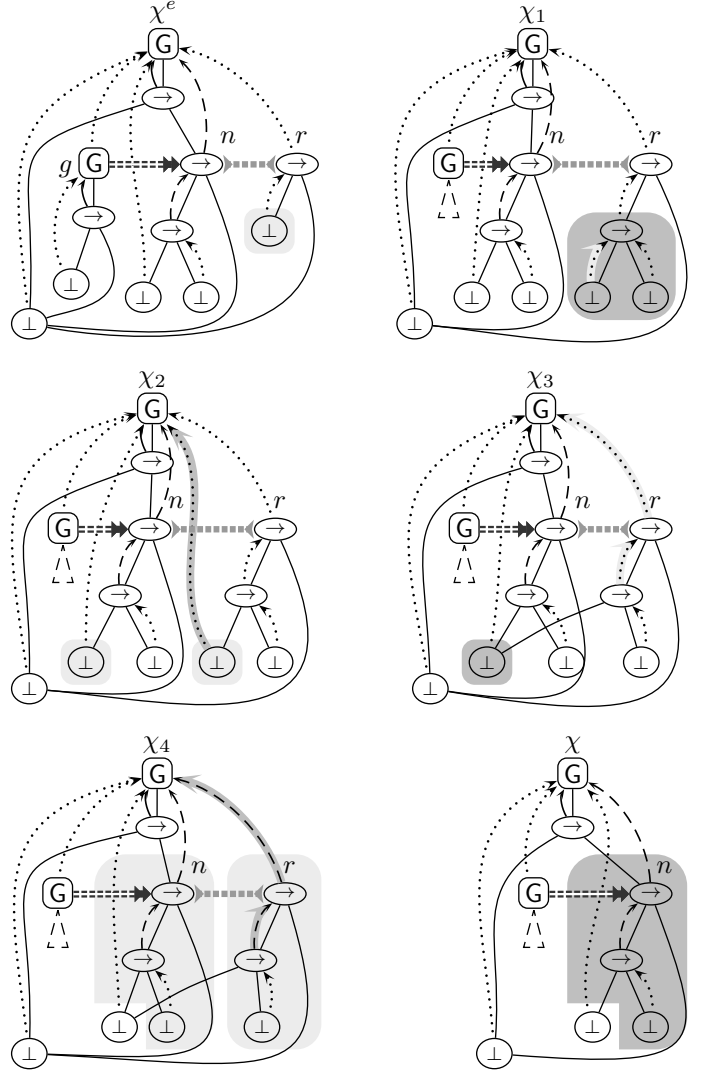


Figure 13: Example of solved instantiation edge

**Example 7 (cont.).** In Figure 12,  $\chi$  is an instance of  $\chi^e$ —hence, the edge  $e$  is solved. This is witnessed by the sequence of transformations given below and depicted in Figure 13.

All nodes below  $g$  are invariant during the transformations and are elided (represented as the  $\triangleleft$  subtree) in all other constraints, for conciseness. Nodes or edges about to change are highlighted in green or in light gray, while those that have just changed are highlighted in red or in dark gray.

Grafting  $\forall(\alpha) \forall(\beta) \alpha \rightarrow \beta$  under  $\langle r1 \rangle$  in  $\chi^e$  leads to  $\chi_1$ ; raising  $\langle r11 \rangle$  twice gives  $\chi_2$ ; mergings nodes  $\langle r11 \rangle$  and  $\langle n11 \rangle$  gives  $\chi_3$ ; weakening node  $\langle r1 \rangle$ , then node  $\langle r \rangle$  leads to  $\chi_4$ ; finally, by merging  $n$  and  $r$ , which is possible as the two subgraphs under them are equal, we end up with exactly  $\chi$ .

Formally, this is the transformation  $\Omega$ :

Graft( $\forall(\alpha) \forall(\beta) \alpha \rightarrow \beta, \langle r1 \rangle$ ); Raise( $\langle r11 \rangle$ ); Raise( $\langle r11 \rangle$ ); Merge( $\langle r11 \rangle, \langle n11 \rangle$ ); Weaken( $\langle r1 \rangle$ ); Weaken( $r$ ); Merge( $r, n$ )

### 3.1.3. From $\lambda$ -terms to typing constraints

Terms of  $e\text{MLF}$  are the partially annotated  $\lambda$ -terms generated by the following grammar:

$$b ::= x \mid \lambda(x) b \mid \lambda(x : \sigma) b \mid b b \mid \text{let } x = b \text{ in } b \mid (b : \sigma)$$

Type inference is performed by translating a source term into a type constraint, solving the constraint into a (principal) presolution, from which a (principal) solution can easily be read.

Type constraints are generated in a compositional manner. Every occurrence of a subexpression  $b$  is associated to a distinct  $G$ -node in the constraint, which we label with  $b$  for readability; however, it should be understood that different occurrences of equal subexpressions are mapped to different nodes. (Formally, occurrences may be identified by their path to the root of the type constraint.) We let  $y$  and  $z$  stand for  $\lambda$ -bound and let-bound variables, respectively. We assume that the source term has been renamed so that every bound variable is distinct from all others.

Constraint generation is described on the top of Figure 14: each case refers to the expression on the left-hand side of the corresponding equality<sup>6</sup> at the bottom of the Figure. The unification edge  $u_y$  in (1) links the node that encodes an occurrence of a  $\lambda$ -bound variable  $y$  to the node  $y$  generated in (4) by the translation of the abstraction binding  $y$ . The instantiation edge  $e_z$  ending in (2) is coming from the  $G$ -node labeled  $b_1$  generated in (3) by the translation of the let expression binding  $z$ . The type of an abstraction  $\lambda(y) b$  is an arrow type whose domain is the type of  $y$  and codomain is an instance of the type of  $b$ , as witnessed by the edge  $e$  (4). The type for an application  $b_1 b_2$  is the codomain of an instance of the type of  $b_1$ , which must itself be an arrow type whose domain is an instance of the type of  $b_2$  (5). The type of a let-expression  $\text{let } x = b_1 \text{ in } b_2$  is just an instance of the type of  $b_2$ : as explained above, the constraints  $b_2$  will contain, for every occurrence of  $x$  in  $b_2$ , one instance edge coming from the type of  $b_1$  and ending at that occurrence. The typing constraint for let-expressions could be optimized to avoid taking an additional instance of  $b_2$ , as done in (Rémy and Jakobowski, 2008; Jakobowski, 2008). The advantage of this unoptimized version, which still preserves the complexity of type inference, is that every subexpression introduces exactly one  $G$ -node; this establishes a one-to-one mapping between subexpressions and  $G$ -nodes that is preserved during constraint resolution ( $G$ -nodes are never merged) and helps define the elaboration after constraint resolution.

**Example 8.** The typing constraint  $\chi$  for the term  $\lambda(x) \lambda(y) x$  is described on the left-hand side of Figure 15. One of its presolutions  $\chi_p$  is drawn on the middle. (We

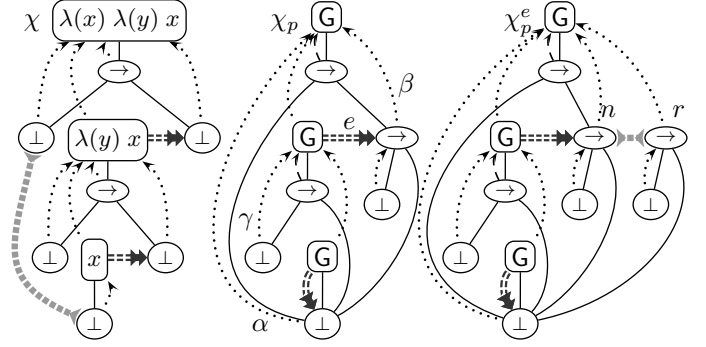


Figure 15: Typing constraints for  $\lambda(x) \lambda(y) x$ .

have dropped the mapping of expressions to  $G$ -nodes for conciseness, and labeled some binding edges that will appear in the  $x\text{MLF}$  translation.) This is not the most general presolution, as some arrow nodes bound at  $G$ -nodes have been made rigid, but an equivalent rigid presolution, as explained in §3.3, that is ready for translation into  $x\text{MLF}$ .

While type inference is out of the scope of this work, we may however easily check that  $\chi_p$  is a presolution, i.e. that both instantiation edges are solved. Consider for example the edge  $e$ . We must verify that  $\chi_p$  is an instance of the expansion  $\chi_p^e$  drawn on the right-hand side, that is, exhibit a sequence of atomic instance operations that transforms  $\chi_p^e$  into  $\chi_p$ . Here, the obvious solution is just to merge the two nodes related by the unification edge, i.e.  $\text{Merge}(n, r)$ .

**Annotated expressions.** The constructions  $\lambda(x : \sigma) b$  and  $(b : \sigma)$  are actually syntactic sugar for  $\lambda(x) \text{let } x = \kappa_\sigma x \text{ in } b$  and  $\kappa_\sigma b$ , respectively<sup>7</sup>, where  $\kappa_\sigma$  is a coercion function that has type  $\forall (\alpha \geq \sigma) \sigma \rightarrow \alpha$  in  $x\text{MLF}$ ; those coercion functions are discussed in more detail in §3.6.

Both constructs are desugared before the translation into constraints. The effect of rebinding  $x$  to  $\kappa_\sigma x$  is to request the parameter  $x$  to be of type  $\sigma$  and simultaneously let all occurrences of  $x$  in  $b$  be typed with possibly different instances of  $\sigma$ . By contrast,  $\lambda(x) b$ , without an annotation, forces the parameter  $x$  and all occurrences of  $x$  in  $b$  to have exactly the same type.

### 3.2. An overview of the translation to $x\text{MLF}$

The elaboration of an  $e\text{MLF}$  term  $b$  to  $x\text{MLF}$  is based on a presolution  $\chi$  of the typing constraint for  $b$ . The translation is based on presolutions rather than solutions, since presolutions still contain the original subconstraints (unlike solutions, which only retain the final type). While typing constraints have principal presolutions, any presolution—not merely the principal one that is returned

<sup>6</sup>The right-hand side is the elaboration of the left-hand side, which will be explained in the next section.

<sup>7</sup>The expression  $\lambda(x) \text{let } x = \kappa_\sigma x \text{ in } b$  is equal to  $\lambda(y) \text{let } x = \kappa_\sigma y \text{ in } b$  whenever  $y$  does not appear free in  $b$ ; using the same variable  $x$  for  $y$  avoids the side condition and so makes the syntactic sugar a purely local transformation.

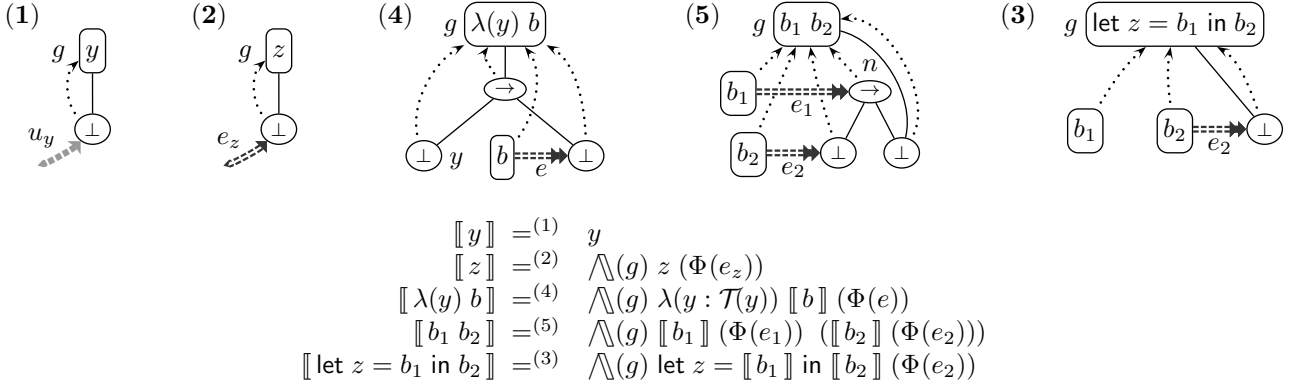


Figure 14: Constraint generation and translation of presolutions

by type inference—can be translated. However, presolutions must be slightly transformed into *rigid* presolutions before translating them, as explained in §3.3—but we may ignore this minor detail for the moment.

Given an original program  $b$  and a (rigid) presolution of the graphic constraint for  $b$ , the translation is inductively defined on the structure of  $b$ , reading auxiliary information on the corresponding nodes in the presolution; we build this way the type of function parameters, type abstractions, and type instantiations. Since presolutions are instances of the original constraint, and type instance preserves both G-nodes and instantiation edges, we can refer to the original nodes and edges in the top of Figure 14 when defining the translation (hence both top and bottom parts of Figure 14 should now be read in parallel to understand the translation). There are two key ingredients:

- For each instantiation edge  $e$  of the form  $g \dashrightarrow n$ , an instantiation  $\Phi(e)$  is inserted to transform the type of the translation of the expression  $b$  corresponding to  $g$  into the type of  $n$ . It can be computed from the proof that  $e$  is solved in  $\chi$ , *i.e.* from the instantiation witness for  $e$ . Details are given in §3.4 and §3.5.
- For each flexible binding edge to a G-node  $n \dashrightarrow g$ , a type abstraction  $\Lambda(\alpha_n \geq \tau_n)$  is inserted in front of the translation of the expression  $b$  corresponding to  $g$ ,  $\tau_n$  being the type of the node  $n$ . Indeed, such an edge corresponds to some polymorphism in  $n$  that must be introduced at the level of  $g$ . We use the notation  $\bigwedge(g)$  to refer to the sequence of all such quantifications at the level of  $g$ , which is a binding prefix of the form  $\Lambda(\alpha_1 \geq \tau_1) \dots \Lambda(\alpha_q \geq \tau_q)$  that will be precisely defined in §3.4.

(Conversely, rigid bindings, which are only useful to make type inference decidable, are inlined during the translation and thus do not give rise to type quantifications.)

The translation is given in Figure 14. We let  $\bigwedge(g)$  and  $\Phi(e)$  abstract for the moment. They will be defined in sections 3.4 and 3.5, respectively.

The translation of a  $\lambda$ -bound variable  $y$  (1) is itself. Indeed, the G-node  $y$  is always monomorphic and there is no polymorphism to introduce; moreover, as the type of  $y$  in the presolution is its only instance, there is no need to add a type instantiation. For all other cases, the translation is of the form  $\bigwedge(g) b'$ ,  $g$  being the G-node for  $b$ . Indeed, generalization is needed in  $\text{ML}^F$  for let-bound expressions (as in ML) and also for applications and abstractions (unlike in ML).

An occurrence of a variable  $z$  (2) bound by some let  $z = b_1$  in  $b_2$  expression is instantiated by  $\Phi(e_z)$  so as to transform the type of  $\llbracket b_1 \rrbracket$  into the type of this occurrence of  $z$ , according to the edge  $e_z$ ; each occurrence of  $z$  in  $\llbracket b_2 \rrbracket$  will potentially pick a different instance. Thus, in the translation of let  $z = b_1$  in  $b_2$  (3), the translation of  $b_1$  is bound to  $z$  uninstantiated (as it suffices to instantiate the occurrences of  $z$ ), while the translation of  $b_2$  is instantiated according to the edge  $e_2$ . In the translation of an abstraction  $\lambda(y) b$  (4), we annotate  $y$  by its type in the presolution (written  $\mathcal{T}(y)$  and defined in §3.4) and coerce  $\llbracket b \rrbracket$  to its type inside the abstraction according to the edge  $e$ . Finally, the translation of an application (5) is the application of the translations, each of which is instantiated according to its constraint edge.

**Example 9.** The presolution  $\chi_p$  in Figure 15 can be translated into the term

$$\Lambda(\alpha) \Lambda(\beta \geq \forall(\delta) \delta \rightarrow \alpha) \lambda(x : \alpha) (\Lambda(\gamma) \lambda(y : \gamma) (x \mathbb{1})) (!\beta)$$

which has type  $\forall(\alpha) \forall(\beta \geq \forall(\delta) \delta \rightarrow \alpha) \alpha \rightarrow \beta$ . Notice the three type quantifications for  $\alpha$ ,  $\beta$ , and  $\gamma$  that correspond to the flexible edges of the same name. The instantiation  $!\beta$  is the translation of  $e$ .

*Type-erasure.* As we will see later,  $\bigwedge(g)$  is only composed of type quantifications, and  $\Phi(e)$  of instantiations. Thus, the translation is type-erasure preserving by construction, which ensures that the semantics of the original and translated terms are the same.

**Theorem 19.** Given a (desugared) term  $b$ , we have  $\llbracket \llbracket b \rrbracket \rrbracket = \llbracket b \rrbracket$ .

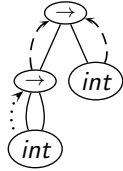


### 3.3. Rigidifying presolutions

All presolutions are not suitable for elaboration into  $x\text{MLF}$ , because rigid and flexible bindings are not treated symmetrically during the translation. Indeed,  $x\text{MLF}$  has flexible quantification, but does not have the rigid form. Rigid quantification is not necessary in  $x\text{MLF}$  because types are fully explicit and rigid nodes can always be explicitly unshared. Unsurprisingly, flexible bindings will be translated to flexible quantification—while rigid bindings will be inlined.

This causes a problem with inert nodes that are flexible but bound under a rigid edge: while they are instantiable in  $e\text{MLF}$  in any context, they would appear in a non-instantiable context in  $x\text{MLF}$  if we translated them as flexible bounds, and there would be no way to instantiate them afterward. One solution is to inline them during the translation, exactly as rigid bounds. However, an even simpler solution is to rigidify them prior to the translation. This is a sound operation in  $e\text{MLF}$ , since inert nodes can always be weakened, and it avoids a special case during the translation.

**Example 10.** For example, the flexible binding edge in the type drawn on the right, which is leaving from the inert node  $\langle 11 \rangle$ , may be weakened in  $e\text{MLF}$ . The two types with or without rigidification are equivalent in  $e\text{MLF}$ . However, they are translated into  $(\forall (\alpha \geq \text{int}) \alpha \rightarrow \alpha) \rightarrow \text{int}$  and  $(\text{int} \rightarrow \text{int}) \rightarrow \text{int}$ , which are not equivalent in  $x\text{MLF}$  (in this case, they are actually incomparable): since type applications are explicit in  $x\text{MLF}$ , a term of the former type must instantiate its argument before applying it, while a term of the latter type can apply its argument directly. This is quite similar to the difference between the two types  $(\forall (\alpha) \text{int} \rightarrow \text{int}) \rightarrow \text{int}$  and  $(\text{int} \rightarrow \text{int}) \rightarrow \text{int}$  in System  $F$ .



For now, let us call *rigidification* the weakening of an inert node. A weakening is in general a strict instance operation in  $e\text{MLF}$ . However, on inert nodes it is a *lossless* one as it right-commutes with all instance operations: a rigidification followed by an instantiation can always be rewritten as an instantiation followed by a rigidification. This means that rigidification will never make typechecking fail when it would not fail without rigidification. Intuitively, when an inert node  $n$  that has been rigidified is unified with another inert node  $m$ , then  $m$  itself can always be rigidified so that unification succeeds, because it is already or can be made inert.<sup>8</sup>

Although inert nodes in non-instantiable contexts are the only nodes that *must* be rigidified, all inert nodes *may* be rigidified. This is easier to implement, but more importantly, it results in simpler and more uniform elaborated terms.

<sup>8</sup> This reasoning can actually be generalized to lowering and splitting of inert nodes, which explains why we only need direct instance operations on such nodes.

For the same reason, we also rigidify flexible existential nodes, even though these are not inert. An existential node is bound to a  $G$ -node but not reachable by structure edges. If it is rigid, it will be inlined by the translation. But no occurrence will be found, so it will be skipped. However, if it is flexible, its translation introduces a (useless) type abstraction over a variable that does not appear in the body of the abstraction but that would still have to be eliminated by some irrelevant type application. Rigidifying flexible existential nodes is always correct and still lossless. Moreover, it avoids useless abstraction and applications in the translated term, as in Example 10.

Since presolutions are instances of the original type constrains (no node and no edge have been lost), we can describe rigidification on the typing constrains of Figure 14. Namely, the following nodes of the presolution are rigidified:

- the node  $\langle g1 \rangle$  in the translation of abstractions (4);
- the node  $n$  in the translation of an application (5);
- the node  $\langle g1 \rangle$  whenever it is bound on  $g$ ;
- any node bound on a  $G$ -node but not reachable from a  $G$ -node by following only structure edges (*i.e.* an existential node).

In the first two cases, rigidification could have been performed during constraint generation since nodes that are rigidified are already inert in the constraint. Conversely, in the two last cases, it is important that the nodes are left flexible *during* type inference when some of the constraints might not have yet been solved, and rigidified only *after* type inference, *i.e.* in presolutions so that rigidification remains a lossless transformation, as argued earlier. Notice that although nodes  $\langle g1 \rangle$  are always bound on  $\langle g \rangle$  in the original constraint, they might be bound above in the presolution, in which case they must not be rigidified—unless they have been merged with other nodes that must be rigidified according to the criteria above.

We call *rigid* a presolution that respects the four conditions above and in which all inert nodes are rigid. We call *rigidification* the transformation of a presolution into a most general, rigid one. The following lemma states the existence of lossless rigid presolutions.

**Lemma 20.** Given a presolution  $\chi_p$  of a constraint  $\chi$ , there exists a rigid presolution  $\chi'_p$  of  $\chi$ , derived from  $\chi_p$  only by rigidifying some nodes, such that the solutions of  $\chi_p$  and  $\chi'_p$  are equivalent up to the weakening of inert nodes.

This result suggests that we could have restricted ourselves to rigid presolutions in the first place, since principal presolutions can be turned into rigid ones in a principal manner. However, rigid presolutions are only useful for the translation of  $e\text{MLF}$  into  $x\text{MLF}$  and useless, if not harmful, for type inference purposes: binding edges can

$$\begin{aligned}
\mathcal{R}_\chi(n) &\triangleq \forall (\mathcal{Q}_\chi(n)) \chi(n) (\mathcal{T}_\chi(\langle n1 \rangle), \dots, \mathcal{T}_\chi(\langle np \rangle)) \\
&\quad \text{where } p \text{ is the arity of } \chi(n) \\
\mathcal{T}_\chi(n) &\triangleq \begin{cases} \mathcal{R}_\chi(n) & \text{if } n \text{ is rigidly bound in } \chi \\ \alpha_n & \text{if } n \text{ is flexibly bound in } \chi \end{cases} \\
\mathcal{Q}_\chi(n) &\triangleq (\alpha_{\langle n_1 \rangle} \geq \mathcal{R}_\chi(n_1) \dots \alpha_{\langle n_k \rangle} \geq \mathcal{R}_\chi(n_k)) \\
&\quad \text{where } n_1, \dots, n_k \text{ are all non } \mathbf{G}\text{-nodes} \\
&\quad \text{flexibly bound to } n \text{ in } \chi, \text{ ordered by } \prec. \\
\mathcal{G}_\chi(g) &\triangleq \forall (\mathcal{Q}_\chi(g)) \mathcal{T}_\chi(\langle g1 \rangle)
\end{aligned}$$

Figure 16: Mapping nodes of  $e\text{MLF}$  to types of  $x\text{MLF}$ .

only be rigidified—without losing solutions—after all the constraint edges under them have been solved. This imposes some synchronization during the constraint resolution. Therefore, we prefer to stay with the more flexible (and simpler) definition of presolutions for  $e\text{MLF}$  and perform rigidification as a first step of the translation into  $x\text{MLF}$ . This way, rigidification needs not be exposed to the user.

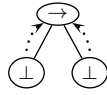
In the remainder of this section, we abstract over a rigid presolution  $\chi$  and an instantiation edge  $e$  of the form  $g \dashrightarrow d$ .

### 3.4. Translating types

*Ordering binders.* In  $e\text{MLF}$ , two binding edges reaching the same node are unordered. It is actually a useful property for type inference not to distinguish between two types that just differ by the order of their quantifiers. However, adjacent quantifiers do not commute in  $x\text{MLF}$ . While they could be explicitly reordered by type instantiation, it is much better to get them in the right order by construction as far as possible, even if reordering of quantifiers remains necessary in some cases, as described below (§3.4, page 18).

The simplest way to order quantifiers is to assume a total ordering  $\prec$  of all the nodes of a constraint  $\chi$ . Of course,  $\prec$  cannot be arbitrary, as it should also ensure the well-scopedness of syntactic types: if  $n \circ \rightarrow n'$  or  $n' \succ \rightarrow n$ , then  $n' \prec n$  must hold.

We choose the leftmost-lowest ordering of nodes for  $\prec$ . That is, if  $n_1, \dots, n_k$  are bound to  $n$ , we first translate the  $n_i$  that is structurally lowest in the type, or leftmost if the  $n_i$  are not ordered by  $\circ \rightarrow$ . This means that the type drawn on the right is always translated as  $\forall (\alpha_1) \forall (\alpha_2) \alpha_1 \rightarrow \alpha_2$ , not as  $\forall (\alpha_2) \forall (\alpha_1) \alpha_1 \rightarrow \alpha_2$ .



*From graphic types to  $x\text{MLF}$  types.* Every node of  $\chi$  can be translated to an  $x\text{MLF}$  type. Moreover, the translation is uniquely determined by the ordering of binders.

We assume that every node  $n$  in  $\chi$  is in bijection with a type variable  $\alpha_n$ . Each non  $\mathbf{G}$ -node  $n$  of  $\chi$  is mapped to a type  $\mathcal{T}_\chi(n)$  of  $x\text{MLF}$  as described in Figure 16. A flexibly

bound node is translated by  $\mathcal{T}_\chi$  as  $\alpha_n$ ; this translation is always used in a context where  $\alpha_n$  is bound. Otherwise,  $n$  is rigidly bound and its type is inlined as  $\mathcal{R}_\chi(n)$  whose definition uses a helper function  $\mathcal{Q}_\chi(n)$  to build a sequence of type quantifications (one for each node flexibly bound to  $n$ ); then  $\mathcal{R}_\chi(n)$  is also used recursively to build the bounds of the type variables in  $\mathcal{Q}(n)$ . When  $\chi$  is clear from context, we omit it for brevity.

**Example 11.** Consider again Figure 12, disregarding the expanded part on the right for now. Let us consider the translation of the node  $\langle n1 \rangle$  (the arrow node under  $n$ ). There is only one node bound on it, the node  $\langle n12 \rangle$ , whose bound is  $\perp$ . Hence,  $\mathcal{T}(\langle n1 \rangle)$  is  $\forall (\alpha_{\langle n12 \rangle} \geq \perp) \alpha_{\langle n11 \rangle} \rightarrow \alpha_{\langle n12 \rangle}$ .

The function  $\mathcal{G}$  is used to translate a  $\mathbf{G}$ -node  $g$ . This is done by introducing the sequence of type quantifications  $\mathcal{Q}(n)$  (representing the type variables generalized at the level of the type scheme that  $g$  stands for), followed by the translation of  $\langle g1 \rangle$ . Notice that some other type quantifications can be introduced when translating  $\langle g1 \rangle$ ; this stands for polymorphism purely local to  $g$ . That is, this polymorphism was already present in  $g$ , has not been instantiated, and needs not be re-introduced. Notice also that, by definition of rigid presolutions,  $\langle g1 \rangle$  cannot be flexibly bound on  $g$ . Hence, the translation is never of the form  $\forall (\dots) \forall (\alpha \geq \tau) \alpha$ .

Finally, we write  $\mathcal{G}(\chi)$  for the translation  $\mathcal{G}(\langle \rangle)$  of the root  $\mathbf{G}$ -node of the whole constraint.

**Example 12 (cont.).** Let us focus on the root of the constraint in Figure 12. The non- $\mathbf{G}$  nodes that are flexibly bound on  $\langle \rangle$  before expansion are  $\langle 11 \rangle$  and  $\langle n11 \rangle$ . As  $n$  is also  $\langle 12 \rangle$ , we have  $\langle 11 \rangle \prec \langle n11 \rangle$ . Thus, the translation  $\mathcal{G}(\langle \rangle)$  of  $\langle \rangle$  is

$$\forall (\alpha_{\langle 11 \rangle} \geq \perp) \forall (\alpha_{\langle n11 \rangle} \geq \perp) \alpha_{\langle 11 \rangle} \rightarrow (\mathcal{T}(\langle n1 \rangle) \rightarrow \alpha_{\langle 11 \rangle})$$

Given all these definitions, we are now able to formally define the notation  $\bigwedge(g)$  used in Figure 14. It is simply  $\Lambda(\mathcal{Q}(g))$ .

*Translating the type of an expansion.* Let  $\chi$  be a constraint containing an instantiation edge  $e$  equal to  $g \dashrightarrow d$ . Let  $\chi'$  be an instance of the expansion  $\chi^e$  of  $e$  in  $\chi$ , such that  $\chi^e \sqsubseteq \chi' \sqsubseteq \chi$ . Let  $r$  be the root of the expanded (i.e. copied) part in  $\chi'$ . In §3.5, we will need to refer to the type under  $r$ , as we will transform this type so that it matches the type under  $d$ . It would be meaningless to translate  $r$  as  $\alpha_{\langle r \rangle}$ , because after any transformation under  $r$ , the translation would still be  $\alpha_{\langle r \rangle}$ . Instead, the correct type is the following: if  $r$  has been created by the expansion, we inline it regardless of its binding flag, and translate it as  $\mathcal{R}_{\chi'}(r)$ . Conversely, if  $r$  is in fact  $d$ , it is translated as  $\mathcal{T}_{\chi'}(d)$  as usual.<sup>9</sup> Formally, the translation  $\mathcal{E}_{\chi'}(r)$  of  $r$  is

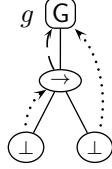
<sup>9</sup>The case  $r$  equal to  $d$  happens either when  $\chi'$  has been instantiated back into  $\chi$  or when  $g$  is degenerate in  $\chi$  and does not hold

defined as

$$\mathcal{E}_{\chi'}(r) \triangleq \begin{cases} \mathcal{T}_{\chi'}(r) & \text{if } r \text{ is } d \\ \mathcal{R}_{\chi'}(r) & \text{otherwise} \end{cases}$$

**Example 13 (cont.).** In Figure 12, the translation of  $r$  is the type  $\forall(\alpha_{\langle r1 \rangle} \geq \perp) \alpha_{\langle r1 \rangle} \rightarrow \alpha_{\langle 11 \rangle}$ , as  $r$  is not part of the initial constraint.

*Type of a G node vs. type of an expansion.* In some cases, the translation of the expansion does not correspond to the translation of  $g$ , regardless of our use of  $\prec$ . This can easily be seen in the example drawn on the right. Here  $\mathcal{G}(g)$  is  $\forall(\beta) \forall(\alpha) \alpha \rightarrow \beta$ , as we start by translating the flexible nodes bound on  $g$ , here  $\langle g12 \rangle$ , before translating  $\langle g1 \rangle$ . However, the expansion of  $g$  has type  $\forall(\alpha) \forall(\beta) \alpha \rightarrow \beta$ : the quantifiers appear in the opposite order.



We believe that this difficulty is actually inherent to elaborating terms for languages with second-order polymorphism, in which second-order polymorphism can be kept local (as here for  $\langle g11 \rangle$ ), or be introduced by generalization (as for  $\langle g12 \rangle$ ). Thankfully, the two translations may differ only by a reordering of quantifiers. In  $x\text{MLF}$ , we can explicitly reorder them using the instantiation

$$\wp; \forall(\geq @\tau_\alpha); \wp; \forall(\geq @\tau_\beta); \forall(\beta \geq) \forall(\alpha \geq) (!\alpha); (!\beta)$$

whose effect is just to commute  $\alpha$  and  $\beta$  in the type  $\forall(\alpha \geq \tau_\alpha) \forall(\beta \geq \tau_\beta) \tau$ .

In the general case, we write  $\Sigma(g)$  for the instantiation that transforms  $\mathcal{G}(g)$  into the translation of its expansion.

### 3.5. From instantiation edges to type instantiations

The main part of the translation is the computation of the type instantiation  $\Phi(e)$  corresponding to an instantiation edge  $e$ . By assumption, the edge is solved; thus  $\chi$  is an instance of the expansion  $\chi^e$  of  $e$  in  $\chi$ . This instantiation can be witnessed by a sequence  $\Omega$  of atomic instance operations. We first build a graphical instantiation  $\Omega$  that will then be translated into a type instantiation in  $x\text{MLF}$ .

*Building  $\Omega$ .* Because  $\Omega$  leaves  $\chi$  unchanged (as otherwise  $\chi^e \sqsubseteq \chi$  would not hold,  $\sqsubseteq$  being antisymmetric), the operations can be rearranged into the following forms (we let  $r$  be the root node of the expansion in  $\chi^e$ ):

- (1)  $\text{Graft}(\sigma, n)$  or  $\text{Weaken}(n)$  with  $n$  in  $\mathcal{I}(r)$ ;
- (2)  $\text{Merge}(n, m)$  with  $n$  and  $m$  in  $\mathcal{I}(r)$ , and  $m \prec n$ ;
- (3)  $\text{Raise}(n)$  with  $n \succ^{\pm} \rightarrow r$ ;
- (4) a sequence  $(\text{Raise}(n))^k; \text{Merge}(n, m)$ , with  $n \in \mathcal{I}(r)$  and  $m \notin \mathcal{I}(r)$ . We write this sequence  $\text{RaiseMerge}(n, m)$  and see it as an atomic operation.

---

polymorphism; see, e.g., the lowermost G-node in Figure 15 in which case both  $r$  and  $d$  are equal to  $\langle g1 \rangle$  in  $\chi^e$ .

An operation  $\text{RaiseMerge}(n, m)$  lets  $n$  leaves the interior of  $r$  and be merged with some node  $m$  of  $\chi$  bound above  $r$ . Conversely, the other operations occur inside the interior of  $r$ . The grouping of operations in  $\text{RaiseMerge}(n, m)$  helps translating the subparts of instantiation witnesses that operate outside of  $\mathcal{I}(r)$ .

Furthermore, since  $\chi$  is a rigid presolution, we may also assume that an operation  $\text{Weaken}(n)$  appears after all the other operations on a node below  $n$  (5). This ensures that  $\Omega$  does not perform any operation under a rigidly bound node, which would not be expressible as an  $x\text{MLF}$  instantiation, as explained in §3.3.

We call *normalized* an instantiation witness that verifies the conditions (1)–(4), and (5) above. Normalized witnesses always exist. A constructive proof of this fact is given by Jakobowski (2008) and it is actually quite easy to establish: performing all instance operations bottom-up, while delaying weakening operations as much as possible, is always possible and results in a normalized witness.

**Example 14.** The constraint edge  $e$  of  $\chi$  in Figure 13 is solved. We recall the witness of  $\chi^e \sqsubseteq \chi$  that we gave in Example 7:

$$\begin{aligned} & \text{Graft}(\forall(\alpha) \forall(\beta) \alpha \rightarrow \beta, \langle r1 \rangle) \\ & \text{Raise}(\langle r11 \rangle); \text{Raise}(\langle r11 \rangle); \text{Merge}(\langle r11 \rangle, \langle n11 \rangle); \\ & \text{Weaken}(\langle r1 \rangle); \text{Weaken}(r); \text{Merge}(r, n) \end{aligned}$$

This transformation is not normalized because node  $\langle r11 \rangle$  is raised twice above the root  $r$ , then merged with  $\langle n11 \rangle$ . We must join those three operations into  $\text{RaiseMerge}(\langle r11 \rangle, \langle n11 \rangle)$ . Similarly, the last operation merges  $n$  and  $r$  and should be replaced by  $\text{RaiseMerge}(r, n)$ . This results in the following normalized derivation:

$$\begin{aligned} & \text{Graft}(\forall(\alpha) \forall(\beta) \alpha \rightarrow \beta, \langle r1 \rangle); \\ & \text{RaiseMerge}(\langle r11 \rangle, \langle n11 \rangle); \\ & \text{Weaken}(\langle r1 \rangle); \text{Weaken}(r); \text{RaiseMerge}(r, n) \end{aligned}$$

Similarly, in Figure 15, we have  $\chi_p^e \sqsubseteq \chi_p$  —as witnessed by  $\text{RaiseMerge}(r, n)$ , which is normalized, hence equal to  $\Omega(e)$ .

*Instantiation contexts.* In order to relate graphic nodes and  $x\text{MLF}$  bounds, we introduce one-hole *instantiation contexts* defined by the following grammar:

$$\mathcal{C} ::= \{ \cdot \} \mid \forall(\geq \mathcal{C}) \mid \forall(\alpha \geq) \mathcal{C}$$

We write  $\mathcal{C}\{\phi\}$  for the replacement of the hole by the instantiation  $\phi$ .

Consider a node  $n$ , and a flexible node  $m$  that is transitively bound to  $n$ . Given our use of  $\prec$  to order nodes, there exists a unique instantiation context  $\mathcal{C}_m^n$  that can be used to descend in front of the quantification corresponding to  $m$  in  $\mathcal{R}(n)$ . For presolutions, in order to avoid  $\alpha$ -conversion-related issues, we build instantiation contexts using variables whose names are based on the nodes they traverse.

Any operation on a node that is transitively bound to the root of an expansion can be expressed using an instantiation context (and a “local” instantiation). Conversely, the operations on rigidly bound or inert-locked nodes cannot. This is unimportant in our case, as normalized witnesses of rigid presolutions only transform nodes transitively flexibly bound to the root of the expansion.

**Example 15.** For example, consider the constraint  $\chi_p$  in Figure 15. The translation  $\mathcal{Q}(\langle \cdot \rangle)$  of the root G-node is

$$\forall (\alpha_{\langle 11 \rangle} \geq \perp) \forall (\alpha_{\langle 12 \rangle} \geq \forall (\alpha_{\langle 121 \rangle} \geq \perp) \alpha_{\langle 121 \rangle} \rightarrow \alpha_{\langle 11 \rangle}) \alpha_{\langle 11 \rangle} \rightarrow \alpha_{\langle 12 \rangle} \Phi_{\xi}(\text{RaiseMerge}(r, m)) = !\alpha_m$$

With the convention above,  $\mathcal{C}_{\langle 11 \rangle}^{\langle \cdot \rangle} = \{\cdot\}$ ,  $\mathcal{C}_{\langle 12 \rangle}^{\langle \cdot \rangle} = \forall (\alpha_{\langle 11 \rangle} \geq) \{\cdot\}$ , and  $\mathcal{C}_{\langle 121 \rangle}^{\langle \cdot \rangle} = \forall (\alpha_{\langle 11 \rangle} \geq) \forall (\geq \{\cdot\})$ .

*Translating normalized derivations into instantiations.* Let us resume the construction of  $\Phi(e)$  by translating a normalized witness  $\Omega$  of  $\chi^e \sqsubseteq \chi$  into a type instantiation in  $x\text{MLF}$ . In fact, we generalize the problem by translating a normalized witness  $\Omega$  of  $\xi \sqsubseteq \chi$  where  $\xi$  is an instance of  $\chi^e$ , i.e. such that  $\chi^e \sqsubseteq \xi \sqsubseteq \chi$ . Inside  $\chi^e$  and  $\xi$ , we let  $r$  be the root of the expansion (inside  $\chi$ ,  $r$  is merged with  $d$ ). We remind that  $\mathcal{E}_{\chi}(r)$  is the translation of  $r$  in the constraint  $\xi$ . By definition, the translation of  $\Omega$ , written  $\Phi_{\xi}(\Omega)$ , must witness the instantiation  $\mathcal{E}_{\xi}(r) \leq \mathcal{E}_{\chi}(r)$ , i.e.

$$\Gamma_d \vdash \Phi_{\xi}(\Omega) : \mathcal{E}_{\xi}(r) \leq \mathcal{E}_{\chi}(r)$$

where  $\Gamma_d$  is the typing context for the node  $d$ .<sup>10</sup> The translation of  $\Omega$  is defined by induction on  $\Omega$  as described in Figure 17. The function  $\Phi_{\xi}$  is overloaded to act on both an instance derivation and a single operation.

The translation of an instance derivation is defined recursively: the translation of an empty derivation is the identity instantiation  $\mathbb{1}$ ; otherwise,  $\Omega$  is of the form  $(\omega; \Omega')$  and we return the composition of the translation of the operation  $\omega$  followed by the translation of the instance derivation  $\Omega'$  applied to the constraint  $\omega(\xi)$ .

The translation of an operation on a rigid node is the identity instantiation  $\mathbb{1}$ , as rigid bounds are inlined. Inert nodes have been weakened into rigid ones and locked nodes cannot be transformed at all. Hence, the remaining and interesting part of the translation is a (single) operation applied to an instantiable node.

The translation of an instance operation on  $r$  (when  $r$  is flexible) is handled especially, as follows.

- The grafting of a type  $\sigma$  is translated to the instantiation  $(@_{\tau})$ , where  $\tau$  is the translation of  $\sigma$  into  $x\text{MLF}$ . (Grafting grafts only closed types, so the constraint in which we translate  $\sigma$  is unimportant.)
- A raise-merge of  $r$  with  $m$  is translated to  $!\alpha_m$ : it must be the last operation of the derivation  $\Omega$ , and  $\alpha_m$  is necessarily bound in the typing environment  $\Gamma_d$ ; hence we may abstract the type of  $r$  under  $\alpha_m$ .

<sup>10</sup>We do not define the typing contexts  $\Gamma_d$  formally, since they are not needed for the translation, but only to state its properties.

### Sequences of operations

$$\begin{aligned} \Phi_{\xi}(\cdot) &= \mathbb{1} \\ \Phi_{\xi}(\omega; \Omega') &= \Phi_{\xi}(\omega); \Phi_{\omega(\xi)}(\Omega') \end{aligned}$$

Operation  $\omega$  on a rigid node  $n$

$$\Phi_{\xi}(\omega) = \mathbb{1}$$

Operation on the (flexible) root  $r$  of the expansion

$$\Phi_{\xi}(\text{Graft}(\sigma, r)) = @(\mathcal{R}(\sigma))$$

$$\Phi_{\xi}(\text{RaiseMerge}(r, m)) = !\alpha_m$$

$$\Phi_{\xi}(\text{Weaken}(r)) = \mathbb{1}$$

Operation on a flexible node different from the root

$$\Phi_{\xi}(\text{Graft}(\sigma, n)) = \mathcal{C}_n^r \{ \forall (\geq @(\mathcal{R}(\sigma))) \}$$

$$\Phi_{\xi}(\text{RaiseMerge}(n, m)) = \mathcal{C}_n^r \{ \forall (\geq !\alpha_m); \& \}$$

$$\Phi_{\xi}(\text{Merge}(n, m)) = \mathcal{C}_n^r \{ \forall (\geq !\alpha_m); \& \}$$

$$\Phi_{\xi}(\text{Weaken}(n)) = \mathcal{C}_n^r \{ \& \}$$

$$\Phi_{\xi}(\text{Raise}(n)) = \mathcal{C}_n^r \{ \&; \forall (\geq @(\mathcal{R}_{\xi}(n))) ;$$

$$\forall (\beta_n \geq) \mathcal{C}_n^m \{ \forall (\geq !\beta_n); \& \} \}$$

$$\text{where } m = \min_{\prec} \{ m \mid n \twoheadrightarrow \twoheadrightarrow \twoheadrightarrow \twoheadrightarrow m \wedge n \prec m \}$$

Figure 17: Translating normalized instance operations

- The weakening of  $r$  is translated to  $\mathbb{1}$ : it must be the next-to-the-last operation in the derivation  $\Omega$ , before the merging of  $r$  with a rigidly bound node, and there is actually nothing to reflect in  $x\text{MLF}$ , as the type of  $r$  itself is unchanged.

In the remaining cases, the operation is applied to an instantiable node  $n$ . Since the derivation is normalized and  $n$  is not rigid,  $n$  must be flexible and transitively bound to  $r$ . Therefore, there exists an instantiation context  $\mathcal{C}_n^r$ , which we call  $\mathcal{C}$ , to reach the bound of  $\alpha_n$  in  $\mathcal{R}_{\xi}(r)$ .

- The grafting of a type  $\sigma$  at  $n$  is translated to  $\mathcal{C} \{ \forall (\geq @(\mathcal{R}(\sigma))) \}$  which transforms the bound  $\perp$  of  $\alpha_n$  into  $\mathcal{R}(\sigma)$ .
- The merging of  $n$  with a node  $m$  is translated to  $\mathcal{C} \{ !\alpha_m \}$ , which first abstracts the bound of  $\alpha_n$  under the name  $\alpha_m$  and immediately eliminates the quantification. (We have assumed  $m \prec n$ , hence  $\alpha_m$  is in scope in the bound of  $n$ .)
- The translation is the same for a raise-merge, but  $\alpha_m$  is bound in the typing environment instead of in  $\mathcal{R}_{\xi}(r)$ .
- The weakening of  $n$  is translated to  $\mathcal{C} \{ \& \}$ , which eliminates the bound of  $n$ .
- Finally, the translation of the raising of  $n$  is more involved, and of the form  $\mathcal{C}_n^r \{ \&; \forall (\geq @(\mathcal{R}_{\xi}(n))) ; \phi \}$ . We first insert a fresh quantification, which will be the one of  $n$  after the raising, inside  $\mathcal{R}_{\xi}(r)$ . The

bound is the current bound of  $n$ , i.e.  $\mathcal{R}_\xi(n)$ . The difficulty consists in finding the node  $m$  in front of which to insert this quantification, so as to respect the ordering between bounds. Notice that the set  $\{m \mid n \multimap \multimap \multimap \multimap m \wedge n \prec m\}$  contains at least the binder of  $n$ , hence its minimum  $m$  is well-defined. Then, the instantiation  $\phi$  equal to  $\forall(\beta_n \geq) \mathcal{C}_n^m \{\forall(\geq !\beta_n); \&\}$  aliases the bound of  $n$  to the quantification just introduced and eliminates the resulting quantification.

The net result of the whole type instantiation is that the type of  $n$  is introduced one level higher than it previously was.

Finally, in order to have a correct instantiation, it remains to reorder quantifiers as described earlier (page 18). Thus we take

$$\Phi(e) = \Sigma(g); \Phi_{\chi^e}(\Omega)$$

**Example 16 (cont.).** *The translation of each step of the normalized witness of Example 14 is:*

Normalized graphic operation	$x\text{MLF}$ translation
Graft( $\forall(\alpha) \forall(\beta) \alpha \rightarrow \beta, \langle r1 \rangle$ )	$\forall(\geq @(\mathcal{R}(\forall(\alpha) \forall(\beta) \alpha \rightarrow \beta)))$
RaiseMerge( $\langle r11 \rangle, \langle n11 \rangle$ )	$\forall(\geq \forall(\geq !\alpha_{\langle n11 \rangle}))$
Weaken( $\langle r1 \rangle$ )	$\&$
Weaken( $\langle r \rangle$ )	$\mathbb{1}$
RaiseMerge( $\langle r \rangle, \langle n \rangle$ )	$!\alpha_n$

Since for the edge  $e$  of  $\chi$  we have  $\Sigma(g) = \mathbb{1}$ , the entire translation of  $e$  is

$$\Phi(e) = \mathbb{1}; \forall(\geq @(\mathcal{R}(\forall(\alpha) \forall(\beta) \alpha \rightarrow \beta))); \forall(\geq \forall(\geq !\alpha_{\langle n11 \rangle})); \&$$

### 3.6. Translating annotated terms

As mentioned in §3.1.3, expressions such as  $(b : \sigma)$  and  $\lambda(y : \sigma) b$  are actually syntactic sugar, for  $\kappa_\sigma b$  and  $\lambda(y) \text{ let } y = \kappa_\sigma y \text{ in } b$ , respectively. The translation  $\mathcal{R}(\kappa_\sigma)$  of the type of the coercion function  $\kappa_\sigma$  in  $x\text{MLF}$  is  $\forall(\alpha \geq \mathcal{R}(\sigma)) \mathcal{R}(\sigma) \rightarrow \alpha$ . Interestingly, coercion functions need not be primitive in  $x\text{MLF}$ —unlike in  $e\text{MLF}$ . Let  $\text{id}_\kappa$  be the expression  $\Lambda(\alpha) \Lambda(\beta \geq \alpha) \lambda(x : \alpha) (x (!\beta))$ . Then, define  $\kappa_\sigma$  as  $\text{id}_\kappa(\mathcal{R}(\sigma))$ . Notice that  $\kappa_\sigma$  behaves as the identity function. Moreover, coercion functions can always be eliminated by strong reduction (as implied by Lemma 13) in the elaboration of the presolution, so that they have no runtime cost.

### 3.7. Soundness of the translation

**Theorem 21.** *Let  $b$  be an  $e\text{MLF}$  term,  $\chi$  a rigid presolution for  $b$ . The translation  $\llbracket b \rrbracket$  of  $\chi$  is well-typed in  $x\text{MLF}$ , of type  $\mathcal{G}(\chi)$ .*

Our translation preserves the type-erasure of programs (Theorem 19). Hence, the soundness of  $x\text{MLF}$  also implies the soundness of  $e\text{MLF}$ —which had previously only been proved for the syntactic versions of  $\text{MLF}$ , but not for the most general, graphical version.

### 3.8. Optimizations

The elaboration is a compilation process, and we have defined it in its simplest form. In practice, some optimizations could be performed during the elaboration process. For instance, raising  $k$  times a node  $n$  (to a position  $n'$ ), is currently done step by step by invoking the atomic  $\text{Raise}(n)$  operation  $k$  times. This could (and should) be translated in a simple step, avoiding intermediate abstractions and applications in  $x\text{MLF}$ . Similarly, contexts could be factored, replacing  $\mathcal{C}_n^r(\phi); \mathcal{C}_n^r(\phi')$  by  $\mathcal{C}_n^r(\phi; \phi')$ . Those optimizations are actually straightforward and significantly simplify elaborated terms—they have been implemented in our prototype (Scherer, 2010b), indeed. Optimizations can also be performed a posteriori, by transforming  $x\text{MLF}$  terms into equivalent ones (with the same type and the same type erasure), as discussed in §5.2.

## 4. Expressiveness of $x\text{MLF}$

The translation of  $e\text{MLF}$  into  $x\text{MLF}$  shows that  $x\text{MLF}$  is at least as expressive as  $e\text{MLF}$ . However, the converse is not true. (This is not entirely surprising: as mentioned in §3.6, coercion functions are primitive in  $e\text{MLF}$ , but not in  $x\text{MLF}$ .) That is, there exist programs of  $x\text{MLF}$  that cannot be typed in  $e\text{MLF}$ . While this is mostly irrelevant when using  $x\text{MLF}$  as an internal language, the question is still interesting from a theoretical point of view, and may help understanding  $\text{MLF}$  independently of any restriction imposed for the purpose of type inference and perhaps suggest other useful extensions.

For the sake of simplicity, we explain the difference between  $x\text{MLF}$  and  $i\text{MLF}$ , the Curry-style version of  $\text{MLF}$  (which has the same expressiveness as  $e\text{MLF}$ , but does not require explicit type annotations in source terms).

### 4.1. A term typable in $x\text{MLF}$ but not in $i\text{MLF}$

Although syntactically identical, the types of  $x\text{MLF}$  and of syntactic  $i\text{MLF}$  differ in their interpretation of quantifications of the form  $\forall(\beta \geq \alpha) \tau$ . Consider, for example, the two types  $\tau_0$  and  $\tau_d$  defined as  $\forall(\alpha \geq \tau) \forall(\beta \geq \alpha) \beta \rightarrow \alpha$  and  $\forall(\alpha \geq \tau) \alpha \rightarrow \alpha$  respectively. In  $i\text{MLF}$ ,  $\beta$  is just an alias for  $\alpha$  and these two types are equivalent. Intuitively, the set of their instances (stripped of toplevel quantifiers) is  $\{\tau' \rightarrow \tau' \mid \tau \leq \tau'\}$ . In  $x\text{MLF}$ , the set of instances of  $\tau_0$  is larger and at least a superset of  $\{\tau'' \rightarrow \tau' \mid \tau \leq \tau' \leq \tau''\}$ , which can be obtained from  $\tau_d$  by all type instantiations of the form  $\forall(\geq \phi); \&; \forall(\geq \phi')$ ;  $\&$  with  $\vdash \phi : \tau \leq \tau'$  and  $\vdash \phi' : \tau' \leq \tau''$ . That is, an instance of  $\tau_0$  can pick for  $\beta$  an instance of the type chosen for  $\alpha$ . This level of generality, possible in  $x\text{MLF}$ , cannot be expressed in  $i\text{MLF}$ .

From this observation, we may easily exhibit an expression  $a$  that is typable in  $x\text{MLF}$  but not in  $i\text{MLF}$ . For readability of the example, we assume primitive products. Let  $a_0$  be the expression

$$\Lambda(\alpha) \Lambda(\beta \geq \alpha) \lambda(x : \alpha) \lambda(y : \beta) (x, \text{choice } \langle \beta \rangle (x (!\beta)) y)$$

$$\begin{aligned}
\llbracket \lambda(x : \tau) a \rrbracket_{\Delta} &= \lambda(x : \exists(\Delta) \tau) \llbracket a \rrbracket_{\Delta} \\
&\triangleq \lambda(x) \text{ let } x = (\exists(\Delta) \tau) x \text{ in } \llbracket a \rrbracket_{\Delta} \\
\llbracket x \rrbracket_{\Delta} &= x \\
\llbracket a_1 a_2 \rrbracket_{\Delta} &= \llbracket a_1 \rrbracket_{\Delta} \llbracket a_2 \rrbracket_{\Delta} \\
\llbracket \Lambda(\alpha \geq \rho) a \rrbracket_{\Delta} &= \llbracket a \rrbracket_{\Delta, \alpha \geq \rho} \\
\llbracket a \phi \rrbracket_{\Delta} &= \llbracket a \rrbracket_{\Delta}
\end{aligned}$$

Figure 18: Translating  $x\text{MLF}$  into  $e\text{MLF}$

of type  $\tau_0 \triangleq \forall(\alpha) \forall(\beta \geq \alpha) \alpha \rightarrow \beta \rightarrow (\alpha \times \beta)$ . Let  $a_1$  and  $a_2$  be defined as

$$\begin{aligned}
a_1 &\triangleq \Lambda(\alpha) \lambda(x : \alpha) x & : & \quad \forall(\alpha) \alpha \rightarrow \alpha \triangleq \tau_1 \\
a_2 &\triangleq \Lambda(\alpha) \lambda(x : \alpha) \lambda(y : \alpha) x : \forall(\alpha) \alpha \rightarrow \alpha \rightarrow \alpha \triangleq \tau_2
\end{aligned}$$

Let  $i$  be 1 or 2 and  $a'_i$  be  $\lambda(x : \tau_i) x \langle \tau_i \rangle x$ . We have  $\vdash a'_i : \tau'_i$ , where  $\tau'_i$  is defined as  $\tau_i \langle \tau_i \rangle$ . If  $f$  has type  $\tau_0$ , then  $f \langle \langle \sigma \rangle; \forall(\geq \phi); \& \rangle$  has type  $\sigma \rightarrow \sigma' \rightarrow (\sigma \times \sigma')$ , for any instantiation  $\phi$  such that  $\phi \vdash \sigma \leq \sigma'$ . Let  $\phi_i$  be  $\langle \tau_i \rangle; \forall(\geq \langle \tau_i \rangle); \&$  and  $\tau''_i$  be  $\tau_i \rightarrow \tau'_i \rightarrow (\tau_i \times \tau'_i)$  and observe that  $\phi_i \vdash \tau_0 \leq \tau''_i$ . Let  $a''_i$  be  $\text{let } (x_i, x'_i) = f \phi_i a_i a'_i \text{ in } x'_i x_i$  and take  $(\lambda(f : \tau_0) (a''_1, a''_2)) a_0$  for  $a$ . The expression  $a$  is well-typed in  $x\text{MLF}$  (and has type  $\tau_1 \times (\tau_2 \rightarrow \tau_2)$ ).

However, the type erasure of  $a$  is ill-typed in  $i\text{MLF}$ , as there is no annotation  $\tau_0$  for the type of the parameter  $f$  that is simultaneously a correct type for  $\llbracket a_0 \rrbracket$  and that can be independently instantiated to  $\tau''_1$  and  $\tau''_2$ —or some other types that allow to simultaneously type both expressions  $a''_1$  and  $a''_2$ . The problem is that, in  $i\text{MLF}$ ,  $\llbracket a_0 \rrbracket$  can only be given a type of the form  $\tau \rightarrow \tau' \rightarrow (\tau \times \tau'')$  or  $\tau' \rightarrow \tau \rightarrow (\tau' \times \tau'')$  with  $\tau \leq \tau' \leq \tau''$ , or of the form  $\forall(\alpha) \alpha \rightarrow \alpha \rightarrow (\alpha \times \alpha)$  (in which both arguments must have identical types), but not simultaneously two such types.

#### 4.2. Restricting $x\text{MLF}$ to match $e\text{MLF}$

The current treatment of variable bounds in  $x\text{MLF}$  is quite natural in a Church-style presentation. Surprisingly, it is also simpler than treating them as in  $e\text{MLF}$ . A restriction  $x\text{MLF}_b$  of  $x\text{MLF}$  without variable bounds that is closed under reduction and in close correspondence with  $i\text{MLF}$  can still be defined a posteriori, by constraining the formation of terms. But the definition is contrived and unnatural, and may not be so appealing in practice (see Appendix C for details). Still, all terms of  $x\text{MLF}_b$  can be translated to  $e\text{MLF}$ .

The translation is actually very similar to that for Church-style System F (Le Botlan and Rémy, 2009) and proceeds by dropping all type abstractions and type applications and translating type annotations of argument of functions. As a result, some type variable may become free in translated types and must be existentially quantified, leading to annotations of the form  $\exists(\Delta) \tau$ . Free variables are kept with their bound in the source. Hence,  $\Delta$  is a list of  $\alpha_i \geq \rho_i$  where  $\rho$  are non-variable types

(see Appendix C). This is a minor difference with System F where all bounds are trivial—and thus need not be tracked. Here, the translation uses an environment to pass this information downward as described in Figure 18. The annotation  $\exists(\Delta) \tau$  stands in  $e\text{MLF}$  for the coercion function of type  $\forall(\Delta) \forall(\alpha = \tau) \forall(\alpha' = \tau) \alpha \rightarrow \alpha'$ , which can easily be translated into some graphic type, as described in (Yakobowski, 2008, Chapter 8).

The restriction to  $x\text{MLF}_b$  prevents the use of variable bounds and therefore of type instantiation between types whose translation into  $e\text{MLF}$  would not be in some instance relationship. This should ensure that the translation of well-typed terms is well-typed, although we have not checked it formally.

Notice that the translation described above annotates all parameters of functions, which is not necessary in  $e\text{MLF}$ . Only parameters of functions that are used polymorphically need to be annotated. A simple optimization is to omit monomorphic type annotations, *i.e.* type annotations of the form  $\exists(\Delta) \tau$  where neither  $\Delta$  nor  $\tau$  contain quantifiers. Still all parameters of functions that have a polymorphic type, whether or not used polymorphically, will be annotated. The image of the translation is then in HML (Leijen, 2008), a strict subset of  $e\text{MLF}$ . Indeed, parameters of functions that are polymorphic may still not be used polymorphically and need not be annotated in  $\text{MLF}$ . However, we do not know whether this can be easily checked during the translation. (In fact, this would amount to detecting and removing useless type-annotations in  $e\text{MLF}$ .)

#### 4.3. Enriching $e\text{MLF}$ to match $x\text{MLF}$ ?

Instead of restricting  $x\text{MLF}$  to match the expressiveness of  $i\text{MLF}$ , a question worth further investigation is whether the treatment of variable bounds could be enhanced in  $i\text{MLF}$  and  $e\text{MLF}$  to match their interpretation in  $x\text{MLF}$  but without compromising type inference. A solution might exist, but it would likely depart from  $e\text{MLF}$ : graphic types have been introduced to simplify the metatheory of the syntactic presentation of  $\text{MLF}$  and one of the simplifications was precisely to disallow variable bounds, which could be written in the syntactic presentation but lead to many complications.

#### 4.4. Comparing $x\text{MLF}$ and $F^\eta$

Type instantiation in  $x\text{MLF}$ , which changes the type of an expression without changing its meaning, can be applied deeply inside a type while it is only superficial in System F. This has some resemblance with retyping functions in  $F^\eta$ , the closure of System F by  $\eta$ -conversion (Mitchell, 1988), which also allows deep type instantiations. However, type instantiations rely on quite different mechanisms in both languages. While it is explicitly expressed in flexible bounds in  $x\text{MLF}$ , it is left implicit and driven by the underlying structure of types in  $F^\eta$ , propagating type instantiation covariantly on the right-hand side of arrow types and contravariantly on their left-hand side.

Both  $F^\eta$  and  $xMLF$  have a little more than System F in common, as our running example `choice id` has both types  $\forall(\alpha) (\alpha \rightarrow \alpha) \rightarrow \alpha \rightarrow \alpha$  and  $(\forall(\alpha) \alpha \rightarrow \alpha) \rightarrow (\forall(\alpha) \alpha \rightarrow \alpha)$ , since the latter can be recovered from the former by type containment, distributing the  $\forall$  over the arrow type constructor.

However,  $F^\eta$  fails on the application, `choice (choice id)`: which is a small variant of `choice id`: this program has type  $\forall(\gamma \geq \forall(\beta \geq \forall(\alpha) \alpha \rightarrow \alpha) \beta \rightarrow \beta) \gamma \rightarrow \gamma$  in  $MLF$ , which admits the three following particular System-F types as instances:

$$\begin{aligned} & (\forall(\alpha) \alpha \rightarrow \alpha) \rightarrow \forall(\alpha) \alpha \rightarrow \alpha \rightarrow (\forall(\alpha) \alpha \rightarrow \alpha) \rightarrow \forall(\alpha) \alpha \rightarrow \alpha \\ & (\forall(\alpha) (\alpha \rightarrow \alpha) \rightarrow \alpha \rightarrow \alpha) \rightarrow \forall(\alpha) (\alpha \rightarrow \alpha) \rightarrow \alpha \rightarrow \alpha \\ & \forall(\alpha) ((\alpha \rightarrow \alpha) \rightarrow \alpha \rightarrow \alpha) \rightarrow (\alpha \rightarrow \alpha) \rightarrow \alpha \rightarrow \alpha \end{aligned}$$

However, `choice (choice id)` does not have any type in  $F^\eta$  of which all these three types are instances.

Conversely, a function of type  $\forall(\beta) (\tau_1 \{ \alpha \leftarrow \tau_2 \} \rightarrow \beta) \rightarrow \beta$  can be seen as one of type  $\forall(\beta) (\forall(\alpha) \tau_1 \rightarrow \beta) \rightarrow \beta$  in  $F^\eta$  by contra-variant type instantiation, which cannot (in general) be expressed  $xMLF$ .

In fact  $xMLF$  and  $F^\eta$  are two rather orthogonal extensions of System F, which could be combined together, as shown in recent work by Cretin and Rémy (2012).

## 5. Discussion

### 5.1. Related works

A strong difference between  $eMLF$  and  $xMLF$  is the use of explicit coercions to trace the derivation of type instantiation judgments. Beside the several papers that describe variants of  $MLF$  and are only indirectly related to this work, most related works are about the use of coercion functions in different ways.

*Elaboration of  $MLF$  into System F.* In a way, the closest work to ours is the elaboration of  $MLF$  into System F, first proposed by Leijen and Löh (2005) to extend  $MLF$  with qualified types and later simplified by Leijen (2007) in the absence of qualified types. Since System F is less expressive than  $MLF$ , an  $MLF$  term  $a$  with a polymorphic type of the form  $\forall(\alpha \geq \tau') \tau$  is elaborated as a function of type  $\forall(\alpha) (\tau'_* \rightarrow \alpha) \rightarrow \tau_*$ , where  $\tau_*$  is a runtime representation of  $\tau$ . The first argument is a *runtime coercion*, which bears strong similarities with our instantiations. However, an important difference is that their coercions are at the level of terms, while our instantiations are at the level of types. In particular, although coercion functions should not change the semantics, this critical result has not been proved, and it is not obvious for a call-by-value language with side effects. In our setting the type-erasure semantics comes for free by construction.

Interestingly, while their translation and ours work on very different inputs—syntactic typing derivations in their case, graphic presolutions in ours—there are strong similarities between the two. The resemblance is even closer

with the improved translation proposed by Leijen (2007), in which rigid bindings are inlined during the translation. Both elaborations use some canonical ordering of quantifiers inside types, with slight differences: while we strive to reduce the number of quantifier reorderings, thus order all the quantifiers, Leijen uses only weaker canonical forms that are sufficient to obtain well-typed terms, but may result in additional reorderings.

*Explicit coercions.* A similar approach has already been used in a language with subtyping and intersection types, proposed as a target for the compilation of bounded polymorphism by Crary (2000). In both cases, coercions are used to make typechecking a trivial process. In our case, they are also exploited to make subject reduction easy—by introducing the language to describe how type instance derivations must be transformed during reduction. We believe that, more generally, the use of explicit coercions is a powerful tool for simplifying subject-reduction proofs. In both approaches, reduction is split into a standard notion of  $\beta$ -reduction and a new form of reduction (which we call  $\iota$ -reduction) that only deals with coercions, preserves type-erasures, and is strongly normalizing. There are also important differences. While both coercion languages have common forms, our coercions intendedly keep the instance-bounded polymorphism form  $\forall(\alpha \geq \tau) \tau'$ . On the opposite, Crary uses the coercions to eliminate the subtype-bounded polymorphism form  $\forall(\alpha \leq \tau) \tau'$ , using intersection types and contravariant arrow coercions instead, which we do not need. Perhaps union types, which Crary (2000) proposes as an extension, could be used to encode away our instance-bounded polymorphism form.

*Harnessing  $MLF$ .* In a recent paper, Manzonetto and Tranquilli (2010) have shown that  $xMLF$  is strongly normalizing by translation into System F, reusing the idea of Leijen and Löh (2005) and their translation of types, recalled above, but starting with  $xMLF$  instead of  $MLF$ . It is unsurprising that the elaboration of  $MLF$  into System F can be decomposed into our elaboration of  $MLF$  into  $xMLF$  followed by a translation of  $xMLF$  into System F. However, the idea of Manzonetto and Tranquilli (2010) is to use the elaboration into System F to prove termination of the reduction in  $xMLF$  in some indirect but simple way, while a direct proof of termination seemed trickier. They show that the elaboration preserves well-typedness and the dynamic semantics via a simulation between the reduction of source terms and target terms. In this process, they also exhibit an intermediate calculus  $F_c$  of term-level retyping functions that mimic our type instantiations. Unfortunately, subject reduction does not hold in  $F_c$  (hence, we can only reuse their direct proof of bisimulation given in Appendix B). Moreover, their intermediate calculus  $F_c$  is tuned to be the target of  $xMLF$ , and cannot express much more. It is actually subsumed by a calculus of erasable coercions  $F_t$  recently proposed by Cretin and Rémy (2012), which contrary to  $F_c$ , enjoys subject reduction. Theorem 10 and Lemma 13

have also been verified by a translation of  $x\text{MLF}$  into  $F_\iota$  (Cretin and Rémy, 2012).

*System with type equalities.* An extension of System F with type equality coercions, called FC or  $\text{FC}_2$  for its revised version (Sulzmann et al., 2007; Weirich et al., 2011) has been proposed to be used as an internal language for Haskell. Type equalities are made explicit through witnesses that have some similarities with our instantiations. System FC has also been designed to be a compiler intermediate language, one of the objectives we have pursued with  $x\text{MLF}$ . However, there are also significant differences: type coercions in  $\text{FC}_2$  are type equality coercions while they are type instantiations in  $\text{MLF}$ . In fact  $\text{FC}_2$  is more related to the system  $F_\iota$  mentioned above, of which  $x\text{MLF}$  is only a particular case. Technically,  $\text{FC}_2$  and  $x\text{MLF}$  seems to be orthogonal extensions of System F which, perhaps, could be combined together. Unfortunately,  $\text{MLF}$ -style polymorphism has been removed from the recent versions of GHC to better accommodate for type inference with GADT. We hope that this is temporary and that both could be eventually recombined.

## 5.2. Future works

The demand for an internal language for  $\text{MLF}$  was first made in the context of using the  $e\text{MLF}$  type system for the Haskell language. We expect  $x\text{MLF}$  to better accommodate qualified types than  $e\text{MLF}$ , since no evidence function should be needed for flexible polymorphism, but it remains to be verified.

While graphical type inference has been designed to keep maximal sharing of types during inference so as to have good practical complexity, our elaboration implementation reads back dags as trees and undoes all the sharing carefully maintained during inference. Even with today’s fast machines, this might be a problem when writing large, automatically generated programs. Hence, it would be worth maintaining the sharing during the translation, perhaps by adding type definitions to  $x\text{MLF}$ .

It was somewhat of a surprise to realize that  $x\text{MLF}$  types are actually more expressive than  $i\text{MLF}$  ones, because of a different interpretation of variable bounds. While the interpretation of  $x\text{MLF}$  seems quite natural in an explicitly typed context, and is in fact similar to the interpretation of subtype bounds in  $F_{<}$ , the  $e\text{MLF}$  interpretation also seemed the obvious choice in the context of type inference. We have left for future work the question of whether the additional power brought by the  $x\text{MLF}$  could be returned back to  $e\text{MLF}$  while retaining type inference. In fact, the problem of choosing the right interpretation for variable bounds reappeared in a recent work by Scherer (2010a) on extending  $\text{MLF}$  to cope with higher-order polymorphism. Indeed, this requires making coexist both implicit and explicit quantifiers, and using the  $x\text{MLF}$  interpretation for explicit quantifiers while retaining the  $\text{MLF}$  more restrictive interpretation for implicit quantifiers.

As noticed in §4.4, type instantiation changes the type of an expression without changing its meaning. It can be performed deeply inside terms, as retyping functions in System  $F^\eta$ . In System  $F^\eta$ , retyping functions can be seen either at the level of terms, as expressions of System F that  $\beta\eta$ -reduce to the identity, or at the level of types as a *type conversion*. In  $x\text{MLF}$ , retyping functions are at the level of types. However, the translation of type instantiations back into coercion functions as done by Manzonetto and Tranquilli (2010) allows one to also see them at the level of terms, bringing  $x\text{MLF}$  and  $F^\eta$  even closer. While the two languages differ in their coercions, they can be combined together as shown in recent work by Cretin and Rémy (2012), allowing a form of abstraction (as in  $x\text{MLF}$ ) over retyping functions (as in  $F^\eta$ ).

Regarding type soundness, it is worth noticing that the proof of subject reduction in  $x\text{MLF}$  does not subsume, but complements, the one in the original presentation of  $\text{MLF}$ . The latter does not explain how to transform type annotations, but shows that annotation sites need not be introduced (but only transformed) during reduction. Because  $x\text{MLF}$  has full type information, it cannot say anything about type information that could be left implicit and inferred. Given a term in  $x\text{MLF}$ , can we rebuild a term in  $i\text{MLF}$  with minimal type annotations? While this should be easy if we require that corresponding subterms have identical types in  $x\text{MLF}$  and  $i\text{MLF}$ , the answer is unclear if we allow subterms to have different types.

The semantics of  $x\text{MLF}$  allows reduction (and elimination) of type instantiations  $a \phi$  through  $\iota$ -reduction but does not allow reduction (and simplification) of instantiations  $\phi$  alone. It would be possible to define a notion of reduction on instantiations  $\phi \rightarrow \phi'$  (such that  $\forall (\geq \phi_1; \phi_2) \rightarrow \forall (\geq \phi_1); \forall (\geq \phi_2)$ , or conversely?) and extend the reduction of terms with a context rule  $a \phi \rightarrow a \phi'$  whenever  $\phi \rightarrow \phi'$ . This might be interesting for more economical representations of type instantiations. However, it is unclear whether there exists an interesting form of reduction that is both Church-Rosser and large enough for optimization purposes. Perhaps, one should rather consider instantiation transformations that preserve observational equivalence; this would leave more freedom in the way one instantiation could be replaced by another.

Less ambitious is to directly generate smaller type instantiations when translating  $e\text{MLF}$  presolutions into  $x\text{MLF}$ , by carefully selecting the instantiation witness to translate—as there usually exist more than one witness for a given instantiation edge. This amounts to using type derivations equivalence in  $e\text{MLF}$  instead of observational equivalence in  $x\text{MLF}$ .

Extending  $\text{MLF}$  with higher-order polymorphism is another ongoing research direction (Herms, 2009; Scherer, 2010a).



## Conclusion

The Church-style version of  $x\text{MLF}$  that was still missing for type-aware compilation and from a theoretical point of view, completes the  $\text{MLF}$  trilogy. The original type-inference version  $e\text{MLF}$ , which requires partial type annotations but does not tell how to track them during reduction, now lies between the Curry-style presentation  $i\text{MLF}$  that ignores all type information and  $x\text{MLF}$  that maintains full type information during reduction.

We have shown that  $x\text{MLF}$  is well-behaved: reduction preserves well-typedness, and the calculus is sound for both call-by-value and call-by-name semantics.

Hence,  $x\text{MLF}$  can be used as an internal language for  $\text{MLF}$ , with either semantics, and also for the many restrictions of  $\text{MLF}$  that have been proposed, including  $\text{HML}$ . Hopefully, this will help the adoption of  $\text{MLF}$  and maintain a powerful form of type inference in modern programming languages that must feature some form of first-class polymorphism.

## Appendix A. Coq formalization

The Coq development is available electronically<sup>11</sup>.

We have proved most of the meta-theoretical results of §2 and §3 using the Coq proof assistant (Coq development team, 2009). In order to deal with alpha-conversion issues—which often represent the most burdensome part of the formalization—we have used the *locally nameless* approach of Aydemir et al. (2008). In this setting, free variables are represented by names, while bound variables are De Bruijn indices. When going through a binder, a term must be *opened* by replacing the bound variable by a free variable. Of course, this variable must be fresh; this is ensured by a cofinite quantification, that allows all names but a given finite set, typically chosen to contain all the free variables of the local typing context.

Given the strong syntactical similarities between  $x\text{MLF}$  and  $F_{<}$ , notably the instance-bounded quantification, we have been able to reuse most of the definitions and results previously established for the examples of (Aydemir et al., 2008). Extending the formalism to add type instantiations was quite natural with a lot of cut-and-paste. We have however found it important to update the tactics<sup>12</sup> contained in the development so that they seamlessly handle the constructs we have added. This way, we have been able to reuse the very high level of automation they provide, which is quite striking in the initial development.

Up-to the use of the locally nameless formalism, our formalization is very faithful to the metatheory of §2. One small difference is that we did not define the operation  $\tau \phi$  as a function, but as a relation. (See below for a justification.) Also, as it is painful to define reduction relations using evaluation contexts, we have inlined rule `CONTEXT` for each context. Finally, characterizing subrelations is also technically heavy, so we have not attempted to formally prove results about call-by-value and call-by-name, but only for  $\longrightarrow$ .

Unfortunately, we have also encountered some difficulties. In particular, defining the operation  $a\{\!|\alpha \leftarrow \phi; \!|\alpha\}$  proved very complicated. To understand why, let us recall rule  $\iota$ -INSIDE:

$$(\Lambda(\alpha \geq \tau) a) (\forall (\geq \phi)) \longrightarrow \Lambda(\alpha \geq \tau \phi) a\{\!|\alpha \leftarrow \phi; \!|\alpha\}$$

The problem lies in the fact that the instantiation  $(\phi; \!|\alpha)$  is not closed in the locally nameless sense when it is substituted instead of  $\!|\alpha$ . That is, the variable  $\alpha$  is not free, but bound in front of  $a\{\!|\alpha \leftarrow \phi; \!|\alpha\}$ . Since bound variables are De Bruijn indices, it is impossible to define the entire operation as a simple recursive operation on  $a$ . Instead,

---

<sup>11</sup>At the url <http://www.yakobowski.org/publis/2010/xmlf-coq/>.

<sup>12</sup>Coq proofs are done using a set of commands, called tactics, which describe in a very high-level way how to build proof terms. The locally nameless examples define some very specialized tactics, that handle *e.g.* the computation of the set of variables against which a variable must be fresh.

we need *e.g.* to shift  $(\phi; !\alpha)$  when crossing a binder. However, this is unsatisfactory, as it requires a considerable amount of new metatheory related to shifting (which the locally nameless approach had been introduced to avoid!). We instead chose to temporarily close  $\phi$  when doing the substitution, by replacing the bound variable  $\alpha$  by a fresh free one. Still (and unsurprisingly), this was not sufficient, as inside the proofs the variable was not “fresh enough”. We thus had to prove that using any fresh free variable, not just the first available one, was equivalent. Those renaming lemmas were quite tedious to prove.

Notice that the exact same problem theoretically occurs when defining the operation  $\tau \phi$ , for the rule  $(\forall (\alpha \geq \tau) \tau')(\forall (\geq \alpha) \phi) = \forall (\alpha \geq \tau) (\tau' \phi)$ . In this case, we did not introduce tedious renaming lemmas, but simply defined  $\tau \phi$  as a relation, instead of as a function.

We tried using the the same solution for  $a\{!\alpha \leftarrow \phi; !\alpha\}$ , which solved some problems related to bound *v.s.* free variables. However, such a solution is only partial. Indeed, when proving progress, we need to give the result term to which a source term reduces to. For rule  $\iota$ -INSIDE, we have to show that both  $\tau \phi$  and the term  $a\{!\alpha \leftarrow \phi; !\alpha\}$  exist. For  $\tau \phi$ , this is easily deduced from the typability of the original term, which requires  $\Gamma \vdash \phi : \tau \leq \tau'$  to hold for some  $\Gamma$ . For  $a\{!\alpha \leftarrow \phi; !\alpha\}$ , this is unfortunately essentially as hard as defining the constructive version of the operation.

## Appendix B. Proofs of §2.5.

### Proof of Lemma 14

Let  $v$  be a value. If it is an abstraction or a type abstraction, the result is immediate. If  $v$  is a partially applied constant, and it is applied to less than its arity, it has either a type of the form  $\forall (\alpha \geq \tau) \tau'$ , or  $\tau \rightarrow \tau'$ . If it is a fully applied constructor, it cannot have type  $\perp$  by hypothesis. ■

### Proof of Theorem 15

The proof is quite standard and proceed by cases on  $a$ . Only the first case is original, but still proceeds without difficulties:

- if  $a$  is  $a' \phi$ , by inversion of typing  $a'$  is typable in the empty environment. If  $a'$  is not a value, it can be further reduced by CONTEXT, and so can  $a$ . Otherwise, we proceed by cases on  $\phi$ :
  - if  $\phi$  is  $\mathbb{1}$ ,  $\wp$  or  $\phi_1; \phi_2$ ,  $a$  can be reduced by rules  $\iota$ -ID,  $\iota$ -INTRO or  $\iota$ -SEQ
  - the case  $\phi = !\alpha$  is impossible in the empty environment;
  - the case  $\phi = @\tau$  is also impossible, as  $a'$  is a value which cannot have type  $\perp$  by Lemma 14.

- in the three last cases,  $a'$  must have type  $\forall (\alpha \geq \tau) \tau'$  for some  $\tau$  and  $\tau'$ . Since it is a value, by inversion of typing it is either a type abstraction of the form  $\Lambda(\alpha \geq \tau) a''$  (and  $a$  can be reduced by  $\iota$ -INSIDE, UNDER or  $\iota$ -ELIM), or it is a partially applied constants, and  $a$  is a value.

- if  $a$  is  $a_1 a_2$ : by inversion of typing,  $a_1$  and  $a_2$  are typable in the empty environment, and  $a_1$  has type  $\tau \rightarrow \tau'$  for some  $\tau$  and  $\tau'$ . If  $a_1$  or  $a_2$  are not values, they can be further reduced, and  $a$  can be further reduced by CONTEXT. Otherwise, since  $a_1$  is a value, of type  $\tau \rightarrow \tau'$ , we proceed by inversion of typing:
  - if  $a_1$  is of the form  $\lambda(x : \tau) a'_1$ ,  $a$  can be reduced by  $(\beta)$ .
  - if  $a_1$  is a partially applied primitive, either  $a$  is a fully applied primitive and it can be reduced by the appropriate  $\delta$  rule, or  $a$  is a value.
  - if  $a_1$  is a partially applied constructor: by hypothesis on the typing of constructors,  $a_1$  is of the form  $C \theta_1 \dots \theta_k v_1 \dots v_n$  with  $n < |C|$  (as a full application would not have an arrow type). Then  $a$  is a value.
- if  $a$  is  $\text{let } x = a_2 \text{ in } a_1$ , by inversion of typing  $a_2$  is typable in the empty environment. If it is not a value, by induction hypothesis it can be reduced. Hence,  $a$  can be reduced by rule CONTEXT. Otherwise,  $a$  can be reduced by rule  $(\beta_{\text{let}})$ .

- variables are not typable in the empty environment;
- constants, abstractions and type abstractions are values; ■

### Proof of Theorem 18

By cases on  $a$ . The cases for variables, constants, abstractions, type abstractions and type applications are the same as for call-by-value.

- If  $a$  is  $a_1 a_2$ : by inversion of typing,  $a_1$  and  $a_2$  are typable in the empty environment, and  $a_1$  has type  $\tau \rightarrow \tau'$  for some  $\tau$  and  $\tau'$ . If  $a_1$  is not a value, by induction hypothesis it can be reduced, and so can  $a$  by rule CONTEXT. Otherwise, by inversion of typing and since  $a_1$  is a value, it is either of the form  $\lambda(x : \tau) a'_1$  (in which case  $a$  can be  $\beta$ -reduced), or a partially applied constant, and the reasoning is the same as for call-by-value.
- If  $a$  is  $\text{let } x = a_2 \text{ in } a_1$ , it can be reduced by rule  $(\beta_{\text{let}})$ . ■

## Appendix C. A restriction of $x\text{MLF}$ without variable bounds

A restriction of  $x\text{MLF}$  without variable bounds that is closed under reduction and in close correspondence with  $e\text{MLF}$  can still be defined a posteriori, by constraining the formation of terms.

The first idea to avoid variable bounds is to restrict the syntax of types and expressions as follows:

$\rho ::=$	$\tau \rightarrow \tau \mid \forall(\alpha \geq \rho) \rho \mid \perp$	Constructed Types
$\tau ::=$	$\alpha \mid \rho \mid \forall(\alpha \geq \rho) \tau$	Types
$a ::=$	$\dots \mid \Lambda(\alpha \geq \rho) a$	Terms
$\Gamma ::=$	$\emptyset \mid \Gamma, \alpha \geq \rho \mid \Gamma, x : \tau$	Environments

The typing rule for type abstraction can be restricted accordingly, replacing  $\tau$  by  $\rho$  in bounds:

$$\frac{\text{TABS} \quad \Gamma, \alpha \geq \rho \vdash a : \tau \quad \alpha \notin \text{ftv}(\Gamma)}{\Gamma \vdash \Lambda(\alpha \geq \rho) a : \forall(\alpha \geq \rho) \tau}$$

There is a slight difficulty however, because new variable bounds could be created during reduction by rule  $\iota$ -INSIDE, turning a bound  $\rho$  into  $\rho \phi$ , which might be a variable. Indeed, assume  $\alpha \geq \rho \vdash \phi : \rho' \leq \alpha$  ( $\phi$  could be  $@\alpha$  if  $\rho'$  is  $\perp$  or of the form  $\phi'; !\alpha$  with  $\alpha \geq \rho \vdash \phi' : \rho' \leq \rho$ ) and consider the reduction sequence:

$$\begin{aligned} & \Lambda(\alpha \geq \rho) (\Lambda(\beta \geq \rho') a) (\forall(\geq \phi); \&) & (1) \\ \rightarrow & \Lambda(\alpha \geq \rho) (\Lambda(\beta \geq \alpha) a\{\beta \leftarrow \phi; !\beta\}) \& & (2) \\ \rightarrow & \Lambda(\alpha \geq \rho) a\{\beta \leftarrow \phi; \mathbb{1}\}\{\beta \leftarrow \alpha\} & (3) \end{aligned}$$

The term (1) is well-formed. However, after one reduction step the bound of  $\beta$  becomes a variable  $\alpha$  and (2) is ill-formed. To prevent this from happening, we may restrict uses of  $\phi$  inside bounds, replacing Rule  $\iota$ -INSIDE by the following variant:

$$\frac{\text{INST-INSIDE} \quad \Gamma \vdash \phi : \rho_1 \leq \rho_2}{\Gamma \vdash \forall(\geq \phi) : \forall(\alpha \geq \rho_1) \tau \leq \forall(\alpha \geq \rho_2) \tau}$$

As expected, this rejects the source term (1) as ill-typed. Unfortunately, this is too restrictive. For instance, it would also reject the application of a polymorphic function. When  $\rho$  and  $\rho'$  are  $\perp$  and  $\phi$  is  $@\alpha$ , (1) is a term of System F, which we must keep!

Notice that the ill-formed term (2) can be further reduced to the term (3), which is well-formed. This suggests another solution to recover type application: making  $\forall(\geq \phi); \&$  a primitive instance operation, say  $\$ \phi$ , and the above reduction sequence atomic, so that one does not see the intermediate ill-formed step.

In summary,  $x\text{MLF}_b$  is defined as follows: first we extend type instantiations with primitive type applications:

$$\phi ::= \dots \mid \$ \phi$$

$$\begin{array}{c} \text{INST-APP} \quad \frac{\Gamma \vdash \phi : \rho \leq \tau}{\Gamma \vdash \$ \phi : \forall(\alpha \geq \rho) \tau_0 \leq \tau_0\{\alpha \leftarrow \tau\}} \quad \text{INST-BOT} \quad \frac{}{\Gamma \vdash @\rho : \perp \leq \rho} \\ \\ \text{INST-UNDER} \quad \frac{\Gamma, \alpha \geq \rho \vdash \phi : \tau_1 \leq \tau_2}{\Gamma \vdash \forall(\alpha \geq \rho) \phi : \forall(\alpha \geq \rho) \tau_1 \leq \forall(\alpha \geq \rho) \tau_2} \\ \\ \text{INST-ABSTR} \quad \frac{\alpha \geq \rho \in \Gamma}{\Gamma \vdash !\alpha : \rho \leq \alpha} \quad \text{INST-INSIDE} \quad \frac{\Gamma \vdash \phi : \rho_1 \leq \rho_2}{\Gamma \vdash \forall(\geq \phi) : \forall(\alpha \geq \rho_1) \tau \leq \forall(\alpha \geq \rho_2) \tau} \\ \\ \text{INST-ELIM} \quad \frac{}{\Gamma \vdash \& : \forall(\alpha \geq \rho) \tau' \leq \tau'\{\alpha \leftarrow \rho\}} \end{array}$$

Figure C.19: Type instance for  $x\text{MLF}_b$

Accordingly, we add the reduction rule

$$(\Lambda(\alpha \geq \rho) a) (\$ \phi) \longrightarrow a\{!\alpha \leftarrow \phi; \mathbb{1}\}\{\alpha \leftarrow \rho \phi\} \quad (\iota\text{-TYPE})$$

and the following case in the recursive definition of type instance:

$$(\forall(\alpha \geq \rho) \tau) (\$ \phi) = \tau\{\alpha \leftarrow \rho \phi\}$$

so that  $\$ \phi$  behaves as its expanded form  $(\forall(\geq \phi); \&)$ . We then restrict the syntax of types and terms as described above, and type instantiation rules as described on Figure C.19 (Rules INST-INTRO, INST-COMP, and INST-ID are omitted as they are left unchanged).

Notice that the intermediate language after the extensions and before the restrictions, say  $x\text{MLF}_b^\dagger$ , is equivalent to  $x\text{MLF}$ : both typing and reduction rules of  $\$ \phi$  are derived; subject reduction hence holds in  $x\text{MLF}_b^\dagger$ .

We show that reduction of  $x\text{MLF}_b^\dagger$  is closed in the  $x\text{MLF}_b$  subset by revisiting the proof of subject reduction for  $x\text{MLF}$ , and checking in each case that the typing derivation rebuilt after reduction is well-formed in  $x\text{MLF}_b$ , having  $\rho$  terms instead of general  $\tau$  terms wherever required by the syntax and the typing rules of  $x\text{MLF}_b$ .

Finally, the target of the translation of  $e\text{MLF}$  into  $x\text{MLF}$ , described in §3, lies in  $x\text{MLF}_b$ . In particular, bounds of variables are  $\rho$ -types  $\mathcal{R}(\cdot)$ . Moreover, the translation of instantiation witnesses described in Figure 17 only applies  $@(\cdot)$  to  $\rho$ -types  $\mathcal{R}(\cdot)$ . Uses of  $!\alpha$  appear either in expressions  $\forall(\geq !\alpha); \&$ , which can be replaced by  $\$(!\alpha)$ , or not under  $\forall(\geq \cdot)$ .

## References

- D. Le Botlan, D. Rémy, MLF: Raising ML to the power of System-F, in: Proceedings of the 8th ACM SIGPLAN International Conference on Functional Programming (ICFP'03), ACM Press, 27–38, URL <http://doi.acm.org/10.1145/944705.944709>, 2003.
- D. Le Botlan, D. Rémy, Recasting MLF, Information and Computation 207 (6) (2009) 726–785, ISSN 0890-5401, URL <http://dx.doi.org/10.1016/j.ic.2008.12.006>.

- D. Rémy, B. Yakobowski, From ML to MLF: Graphic type constraints with efficient type inference, in: Proceedings of the 13th ACM SIGPLAN International Conference on Functional Programming (ICFP '08), ACM Press, 63–74, URL <http://doi.acm.org/10.1145/1411203.1411216>, 2008.
- S. Peyton Jones, Haskell 98 Language and Libraries: The Revised Report, Cambridge University Press, ISBN 0521826144, 2003.
- M. P. Jones, A theory of qualified types, Science of Computer Programming 22 (3) (1994) 231–256, ISSN 0167-6423, URL [http://dx.doi.org/10.1016/0167-6423\(94\)00005-0](http://dx.doi.org/10.1016/0167-6423(94)00005-0).
- G. Manzonetto, P. Tranquilli, Harnessing MLF with the Power of System F, in: Mathematical Foundations of Computer Science, vol. 6281 of LNCS, 525–536, URL <http://perso.ens-lyon.fr/paolo.tranquilli/content/docs/snmlf.pdf>, 2010.
- B. Yakobowski, Graphical types and constraints: second-order polymorphism and inference, Ph.D. thesis, University of Paris 7, URL <http://www.yakobowski.org/phd-dissertation.html>, 2008.
- H. P. Barendregt, The Lambda Calculus: Its Syntax and Semantics, North-Holland, ISBN 0-444-86748-1, 1984.
- D. Rémy, B. Yakobowski, A graphical presentation of MLF types with a linear-time unification algorithm, in: Proceedings of the 2007 ACM SIGPLAN International Workshop on Types in Languages Design and Implementation (TLDI' 07), ACM Press, 27–38, URL <http://doi.acm.org/10.1145/1190315.1190321>, 2007.
- G. Scherer, Extending MLF with Higher-Order Types, Master's thesis, École Normale Supérieure d'Ulm, URL <http://gallium.inria.fr/~remy/mlf/scherer@master2010:mlfomega.pdf>, 2010a.
- F. Pottier, D. Rémy, The Essence of ML Type Inference, in: B. C. Pierce (Ed.), Advanced Topics in Types and Programming Languages, chap. 10, MIT Press, 389–489, URL <http://crystal.inria.fr/attapl/>, 2005.
- G. Scherer, Prototype implementation of MLF, URL <http://gallium.inria.fr/~remy/mlf/mlf-omega/>, 2010b.
- D. Leijen, Flexible types: robust type inference for first-class polymorphism, Tech. Rep. MSR-TR-2008-55, Microsoft Research, URL <ftp://ftp.research.microsoft.com/pub/tr/TR-2008-55.pdf>, 2008.
- J. C. Mitchell, Polymorphic type inference and containment, Information and Computation 2/3 (76) (1988) 211–249.
- J. Cretin, D. Rémy, On the Power of Coercions Abstraction, in: ACM Symposium on Principles of Programming Languages (POPL), Philadelphia, Pennsylvania, URL <http://crystal.inria.fr/~remy/coercions/>, 2012.
- D. Leijen, A. Löb, Qualified types for MLF, in: Proceedings of the 10th International Conference on Functional Programming (ICFP '05), ACM Press, 144–155, URL <http://doi.acm.org/10.1145/1090189.1086385>, 2005.
- D. Leijen, A Type Directed Translation of MLF to System F, in: Proceedings of the 12th International Conference on Functional Programming (ICFP '07), ACM Press, 111–122, URL <http://doi.acm.org/10.1145/1291151.1291169>, 2007.
- K. Crary, Typed compilation of inclusive subtyping, in: Proceedings of the 5th ACM SIGPLAN International Conference on Functional Programming (ICFP '00), ACM Press, ISBN 1-58113-202-6, 68–81, URL <http://doi.acm.org/10.1145/351240.351247>, 2000.
- M. Sulzmann, M. M. T. Chakravarty, S. P. Jones, K. Donnelly, System F with type equality coercions, in: Proceedings of the 2007 ACM SIGPLAN International Workshop on Types in Languages Design and Implementation (TLDI' 07), ACM Press, 53–66, URL <http://doi.acm.org/10.1145/1190315.1190324>, 2007.
- S. Weirich, D. Vytiniotis, S. Peyton Jones, S. Zdancewic, Generative type abstraction and type-level computation, in: Proceedings of the 38th annual ACM SIGPLAN-SIGACT symposium on Principles of programming languages, POPL '11, ACM, New York, NY, USA, ISBN 978-1-4503-0490-0, 227–240, URL <http://doi.acm.org/10.1145/1926385.1926411>, 2011.
- P. Herms, Partial Type Inference with Higher-Order Types, Master's thesis, University of Pisa and INRIA, URL <http://pauillac.inria.fr/~remy/mlf/Herms@master2009:mlf-omega.pdf>, 2009.
- Coq development team, The Coq Proof Assistant Reference Manual, version 8.2, URL <http://coq.inria.fr/refman/>, 2009.
- B. Aydemir, A. Charguéraud, B. C. Pierce, R. Pollack, S. Weirich, Engineering formal metatheory, in: Proceedings of the 35th annual ACM SIGPLAN-SIGACT symposium on Principles Of Programming Languages (POPL '08), ACM Press, ISBN 978-1-59593-689-9, 3–15, URL <http://doi.acm.org/10.1145/1328438.1328443>, 2008.