



Polymorphic Types with Polynomial Sizes

Jean-Louis Colaço, Baptiste Pauget, Marc Pouzet

► To cite this version:

Jean-Louis Colaço, Baptiste Pauget, Marc Pouzet. Polymorphic Types with Polynomial Sizes. 9th ACM SIGPLAN International Workshop on Libraries, Languages and Compilers for Array Programming (ARRAY 2023), Jun 2023, Orlando, United States. pp.36-49, <10.1145/3589246.3595372>. <hal-04491216>

HAL Id: hal-04491216

<https://hal.science/hal-04491216v1>

Submitted on 6 Mar 2024

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.



HAL Authorization

Polymorphic Types with Polynomial Sizes

Jean-Louis Colaço
ANSYS
Toulouse, France
Jean-Louis.Colaco@ansys.com

Baptiste Pauget
ANSYS
Toulouse, France
INRIA
France
Baptiste.Pauget@ansys.com

Marc Pouzet
Ecole normale supérieure – PSL
university
Paris, France
INRIA
Paris, France
Marc.Pouzet@ens.fr

Abstract

This article presents a compile-time analysis for tracking the size of data-structures in a statically typed and strict functional language. This information is valuable for static checking and code generation. Rather than relying on dependent types, we propose a type-system close to that of ML: polymorphism is used to define functions that are generic in types and sizes; both can be inferred. This approach is convenient, in particular for a language used to program critical embedded systems, where sizes are indeed known at compile-time. By using sizes that are multivariate polynomials, we obtain a good compromise between the expressiveness of the size language and its properties (verification, inference).

The article defines a minimal functional language that is sufficient to capture size constraints in types. It presents its dynamic semantics, the type system and inference algorithm. Last, we sketch some practical extensions that matter for a more realistic language.

CCS Concepts: • **Software and its engineering** → **Functional languages; Polymorphism; Recursion; Semantics; Automated static analysis; Embedded software; Software safety; Software usability.**

Keywords: array programming, type systems

ACM Reference Format:

Jean-Louis Colaço, Baptiste Pauget, and Marc Pouzet. 2024. Polymorphic Types with Polynomial Sizes. In *Proceedings of ACM Conference (Conference’17)*. ACM, New York, NY, USA, 22 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 Introduction

We are interested here in the programming, with a high-level language, of certified real-time embedded software submitted to strong safety requirements, such as those found in avionics, railway and automotive (eg, flight control, braking, electrical engine). In this field, the domain-specific programming language SCADE [Colaço et al. 2017], is used for more

than twenty years. It inherits the principles and style of the synchronous language LUSTRE [Halbwachs et al. 1991]. The specific features of these languages are essentially orthogonal to our discussion, which focuses on arrays. However, the constraints imposed by the targeted applications are such that:

- (i) Resources must be statically bounded, both in term of memory and execution time. This ensures that a system can for an arbitrarily long time and meet its deadlines.
- (ii) Programs must be certified by independent authorities. This requires a reference specification, extensive testing, and property checking, both for programs and the tools used to generate code.

To this end, the size any data-structure in SCADE must be known statically. While some functions may depend on size parameters, these sizes get ultimately instantiated at compile-time with a concrete value (e.g., an integer). Moreover, the language and its compiler comply with the highest certification standards for critical software (e.g., DO178C, level A of avionics): the generated code can be used without any further validation that the semantics is preserved.

Modern real-time applications combine complex control code (e.g., hierarchical automata) and intensive computation using arrays (e.g., Kalman filters, neural networks, optimization algorithms). Arrays introduce dynamic accesses to memory that must respect array bounds; otherwise, the risk is, at best a stop of the execution, at worst a silent corruption of the memory. Ensuring access correctness ranges from programmer’s responsibility (e.g., in C) to programmer’s proof obligations (e.g., in SPARK [Barnes 2003]), and include skeptical compilers that generate defensive code in various ways: by throwing exceptions (e.g., OCAML, ADA), by saturating the index value [Gérard et al. 2012] or by adding a default value [Colaço et al. 2017] — two solutions followed in several synchronous compilers (e.g., Heptagon¹, LUSTRE V6² and SCADE³). This results in less efficient generated code and the potential introduction of dead code.⁴

Conference’17, July 2017, Washington, DC, USA

© 2024 Association for Computing Machinery.

This is the author’s version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *Proceedings of ACM Conference (Conference’17)*, <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>.

¹ <https://gitlab.inria.fr/synchrone/heptagon>

² <https://www-verimag.imag.fr/DIST-TOOLS/SYNCHRONE/reactive-toolbox/>

³ <https://www.ansys.com/products/embedded-software/ansys-scade-suite>

⁴ This last point is not without importance: coverage analysis, an activity required for the certification of critical applications, needs justifications for the code that cannot be covered by a test case.

Functional languages have popularized intensional array operations (e.g., map, fold, transpose) [Bird and Wadler 1988] which provide predefined access schemes and thus avoid the explicit manipulation of indexes. Their safety only needs size equalities to be solved instead of inequalities (e.g., bound checking), with algorithms that are simpler and more efficient. data-flow style of SCADE favors this kind of intensional definitions. E.g., the scalar product is written in SCADE:

```
function dot «n» (u, v: int32^n) returns (d: int32)
  d = (fold $+$ «n») (0, (map $*$ «n») (u, v));
```

However, writing complex array operations as the ones found in signal or image processing or AI is rapidly tedious in SCADE for multiple reasons:

- (i) The language is explicitly typed. This leads to long and redundant annotations when size expressions grow or multiple size variables are used. In the above example, the size n appears four times.
- (ii) Current primitives (map, fold, concat, ...) offer a limited expressiveness: only linear relations between sizes are possible. Sampling or filtering are hardly expressed and inefficiently programmed.
- (iii) The SCADE compiler checks sizes at program elaboration (i.e. instantiation), where sizes get constant values. Error detection is thus late and non-modular.
- (iv) All dynamic array accesses are guarded, leading to code with unnecessary conditionals and dead branches.

Contribution. We believe that the possible improvements for the above remarks share a common seed: a type-like knowledge of sizes, available during the whole compilation process, not only after elaboration, would give new perspectives for verification and compilation. In short, the proposed solution is based on the following elements:

- (i) A language of sizes made of multivariate polynomials. It provides a practical compromise between formal manipulations and expressiveness.
- (ii) An ML-like type system that extends polymorphism to sizes. Sizes are generalized at declarations and instantiated as in ML. This allows to handle sizes and types in the unified manner.
- (iii) An inference algorithm. Although incomplete, it allows most sizes to be omitted. Size constraints, that are vanishing polynomials $P[X] = 0$, could be solved with an external solver. However, we propose here a dedicated procedure because using black-box tools is unpractical for certified software.

Although being modest, this extension of the type system deeply affects language properties. First, principal typing (as for dependent types), is lost: some terms may receive multiple (incomparable) types. Second, sizes have a *computational content*, i.e. the language semantics is not type erasable. Both points are challenges for the inference algorithm: it should

not only produce a well-typed term but it must also ensure that the semantics is independent from inference choices.

This presentation is purposely conducted on a toy functional language that contains the minimal constructs to highlight the main issues. The article is organized as follows: [section 2](#) gives a general overview of the proposed contribution. The language and its semantics and typing discipline are defined in [section 3](#). Then, [section 4](#) details type and size inference and establishes their meta-theoretical properties. Practical extensions are sketched in [section 5](#). We briefly present the use of static sizes for both verification and compilation in [section 6](#). We discuss related works in [section 7](#) and conclude in [section 8](#).

An extended version of this type system is implemented in a prototype of compiler for a synchronous language with arrays. A type checker for a simpler, ML-like language with the proposed syntax is available alongside the submission.⁵

2 Overview

For brevity, we introduce a core language \mathcal{L}^η that contains the minimal constructs required to introduce a notion of sizes into types. In particular, it has no primitive notion of arrays: they are considered as functions on a bounded domain. This section gives an informal insight of \mathcal{L}^η , its type system and type inference through simple examples.

\mathcal{L}^η is equipped with a separate language for sizes, namely sizes and expressions are distinct syntactical objects. This size language is made of multivariate polynomials. Sizes (ranged over by η, \dots) and their variables (ranged over by $\iota, \kappa, \delta, \dots$) are handled in a similar way to types (τ, \dots) and type variables ($\alpha, \beta, \gamma, \dots$).

Intensional arrays and size consistency. In numerous programming languages intended for scientific computations such as SISAL [Feo et al. 1990], explicit index manipulations are replaced by operators acting on arrays called *combinators* [Jay and Sekanina 1997; ?]. This style benefits especially to functional and data-flow programming languages by allowing to write array definitions with single expressions. Array combinators provide predefined access patterns that are always correct, avoiding at the same time the need for runtime checks. However, some of these primitives still require array sizes to coincide. To enforce such properties by type checking, sizes need to be expressed in types. The point-wise function application (map), its binary variant (map2) and the array reduction (fold), three operators that are available in SCADE, are given the following type schemes in the proposed language \mathcal{L}^η :

```
val map  :  $\forall \iota. \alpha \cdot \beta. (\alpha \rightarrow \beta) \rightarrow \langle \iota \rangle \rightarrow [\iota] \alpha \rightarrow [\iota] \beta$ 
val fold :  $\forall \iota. \alpha \cdot \beta. (\alpha \rightarrow \beta \rightarrow \alpha) \rightarrow \langle \iota \rangle \rightarrow \alpha \rightarrow [\iota] \beta \rightarrow \alpha$ 
val map2 :  $\forall \iota. \alpha \cdot \beta \cdot \gamma. (\alpha \rightarrow \beta \rightarrow \gamma) \rightarrow \langle \iota \rangle \rightarrow [\iota] \alpha \rightarrow [\iota] \beta \rightarrow [\iota] \gamma$ 
```

⁵<https://gitlab.inria.fr/bpauget/array-2023>

Given a polynomial size η , the type $\langle \eta \rangle$ (read *value of size η*) denotes a *refinement* of the integer base type: $\{x : \text{int} \mid x = \eta\}$. Actually, this is a singleton type. Here, the second argument of each function thus allows to constraint the value of the size variable ι . $[\eta]\tau$ is the type of arrays with length η and elements of type τ . Used as an expression, the syntax $\langle \eta \rangle$ also designates the only value of type $\langle \eta \rangle$, thus the partial application $\text{—map } f \langle 9 \rangle\text{—}$ can only be given an array of length 9. These signatures highlight polymorphism that act both on type variables (α, β, γ) and size variables (ι). Using them, the scalar product is expressed as:

```
let dot =  $\lambda u. \lambda v. \text{fold } (+) \langle \_ \rangle 0 (\text{map2 } (*) \langle \_ \rangle u v)$ 
           [ $\forall \iota. [\iota]\text{int} \rightarrow [\iota]\text{int} \rightarrow \text{int}$ ]
```

In the definition of dot, the size values $\langle _ \rangle$ are omitted for both iterators: they are inferred. The above type scheme, built by the inference, forces the length of input arrays to coincide. The size variable ι cannot be directly constrained: no argument have type $\langle \iota \rangle$. Thus, ι will be deduced from arrays' size. Lets assume the definition of a primitive window defining a sliding window of size κ with step 1:

```
val window :  $\forall \iota. \kappa. \alpha. \langle \kappa \rangle \rightarrow [\iota + \kappa - 1]\alpha \rightarrow [\iota][\kappa]\alpha$ 
```

This function builds a matrix whose rows are slices of length κ of the input array, starting at the element given by the row index. For example $\text{—window } \langle 3 \rangle [0, 1, 2, 3, 4]\text{—}$ produces the matrix $[[0, 1, 2], [1, 2, 3], [2, 3, 4]]$. The size $\iota + \kappa - 1$ encodes the relation between input and output array sizes so that the former is fully read. Filtering data with a kernel K of size κ is a common signal processing operation. It is expressed with a discrete convolution: $(K * I)_i = \sum_{j=0}^{\kappa-1} K_j \cdot I_{i+j}$. This uni-dimensional filter may be defined as:

```
let convolution =  $\lambda k. \lambda i. \text{map } (\text{dot } k) \langle \_ \rangle (\text{window } \langle \_ \rangle i)$ 
                 [ $\forall \iota. \kappa. [\kappa]\text{int} \rightarrow [\iota + \kappa - 1]\text{int} \rightarrow [\iota]\text{int}$ ]
```

Here as well, the inference is able to determine the missing sizes (and types), making the kernel size coincide with slices length. Inference derives the above type scheme for this declaration. Note that, by a change of variables, it is equivalent to $\forall \iota. \kappa. [\kappa]\text{int} \rightarrow [\iota]\text{int} \rightarrow [\iota - \kappa + 1]\text{int}$.

Extensional arrays and bounds propagation. Arrays are not primitive constructs: the type $[\eta]\tau$ is a shortcut for $[\eta] \rightarrow \tau$ where $[\eta]$ (read *index of size η*) is a second refinement of type int denoting positive integers strictly smaller than η : $\{x : \text{int} \mid 0 \leq x < \eta\}$. Although not realistic for compilation, such simplification limits language complexity. Using this refinement, the map2 iterator is expressible in \mathcal{L}^η :

```
let map2 =  $\lambda f. \lambda n. \langle \iota \rangle. \lambda u. \lambda v. \lambda i. [\iota]. f (u i) (v i)$ 
           [ $\forall \iota. \alpha. \beta. \gamma. (\alpha \rightarrow \beta \rightarrow \gamma) \rightarrow \langle \iota \rangle \rightarrow [\iota]\alpha \rightarrow [\iota]\beta \rightarrow [\iota]\gamma$ ]
```

It defines an 'array', i.e. a function with a bounded domain, whose content is obtained by applying f to u and v elements. Their accesses are denoted by the applications $\text{—}(u i) ; (v i)\text{—}$

that respect bounds by construction. The second argument $\text{—}n\text{—}$ of map2 is unused, but the types annotations $\langle \iota \rangle$ and $[\iota]$, where ι is an *anonymous* size variable alike OCAML's ones, forces n to be the size of index i .

The index type only allows to propagate known bounds. Notably, no arithmetic operations are defined for indexes. Values of type $[\eta]$ are obtained by calling functions with a bounded codomain, e.g., the modulo, whose type scheme is $\forall \iota. \text{int} \rightarrow \langle \iota \rangle \rightarrow [\iota]$. Although elementary, this refined type allows to separate bound checking from array accesses.

The ghost size issue. In the previous examples, all unspecified sizes were deducible from the types. However, this is not always so simple. The cst function below defines a constant array with an arbitrary size, thanks to the type annotation $[_]$.

```
let cst =  $\lambda c. \lambda i. [\_]. c$  [ $\forall \iota. \alpha. \alpha \rightarrow [\iota]\alpha$ ]
let even = fold (+)  $\langle \_ \rangle 0$  (cst 2) (Error: Unconstrained size)
```

In even declaration, summing the element of cst 2 without specifying fold's size leads to an undefined value since this size could be chosen arbitrarily. This must be rejected.

Contrary to types in ML like languages, sizes have a *computational content*: they may determine the semantics of expressions. Inference must thus ensure that the semantics of the reconstructed term was already specified in the source. We formalize this property in subsection 4.5: when type inference succeeds, all well typed annotated versions of the source program evaluate to the same result. Our even example becomes unambiguous by adding an argument:

```
let even =  $\lambda n. \text{fold } (+) n 0$  (cst 2) [ $\forall \iota. \langle \iota \rangle \rightarrow \text{int}$ ]
```

3 A typed core functional language with size polymorphism

This section focuses on the type system. It aims at expressing as many relations between sizes as possible while remaining decidable and largely implicit. It is a combination of two widely studied type systems traits: (i) a restricted form of *refinement types*, as defined in Xi and Pfenning [1999], and (ii) the ubiquitous *let-polymorphism* of Milner [1978] extended to sizes.

Such an expressive type system may lead to undecidable type checking and incomplete type inference. Nonetheless, the context of SCADE ensures that all sizes get ultimately known statically: once elaboration is done, size checking becomes trivial but late and non modular. We would like to check most size constraints earlier, during static typing and per declaration, relying on the specialization as a fallback mechanism for the properties that remain unproved. This explains also why the size language is purposely not restricted to decidable (e.g., linear) arithmetic expressions.

For clarity, the tightest possible language \mathcal{L}^η is used: a core ML (λ -calculus with let bindings) augmented with a few

constructs. We propose some extensions that are useful for a realistic language in [section 5](#).

3.1 Syntax and semantics

The syntax of \mathcal{L}^η is summed-up in [Figure 1](#). It is explicitly typed i.e. type annotations are part of expressions. In the subsequent, n denotes an integer. To emphasize on their similarities, sizes, types and their variables are designated by Greek letters whereas the Latin ones will be dedicated to terms and program variables.

Name-spaces, free variables. Because they are syntactically separated, sizes, types and expressions use variables in distinct name-spaces, respectively denoted \mathcal{V}_η , \mathcal{V}_τ and \mathcal{V}_e , allowing for name reuse without masking. However, these sets will be considered disjointed in the formalization. Given a syntactical object o , the set $FV(o) \in \mathcal{V}_\eta \cup \mathcal{V}_\tau \cup \mathcal{V}_e$ of *free variables* is defined by the set of variables that are not bound by one of the following rules: (i) abstraction $\lambda x : \tau. e$ — binds x in e ; (ii) local binding $\text{let } x_V : \tau = e \text{ in } e'$ — binds x in e' and size and type variables V in both τ and e . *Closed* objects are the ones with no free variables.

Sizes and types. The size language is made of multivariate polynomials with integer coefficients: $\mathbb{Z}[\mathcal{V}_\eta]$. The main benefit of this restricted class of arithmetic expressions lies in the existence of a normal form: a weighted sum of products of variables. This allows for symbolic comparison of sizes that are structurally different (e.g., $(\iota - 1)^2 - 1 = \iota * (\iota - 2)$).

Besides functions, the type language contains a single constructor **int**, along with two refinements, as defined by [Xi and Pfenning \[1999\]](#): (i) the type $\langle \eta \rangle$ (read *value of size η*) denotes the singleton $\{\bar{\eta}\}$ and (ii) $[\eta]$ (read *index of size η*) represents the interval $\llbracket 0, \bar{\eta} - 1 \rrbracket$, where $\bar{\eta}$ is the value of η , depending on size variables valuation. In the syntax refinement types, they are respectively expressed as $\{x : \text{int} \mid x = \bar{\eta}\}$ and $\{x : \text{int} \mid 0 \leq x < \bar{\eta}\}$.

Polymorphism. Types, including polymorphism, are explicit in \mathcal{L}^η : variables $-^S x-$ are instantiated with a list S of sizes and types while local bindings $\text{let } x_V : \tau = e \text{ in } e'$ — declare the list⁶ V of size and type variables that are generalized. We shift away from the standard notation $\text{let } x : \forall S. \tau = e \text{ in } e$ — to emphasize on generalized variables' scope: their are bound in both type τ and expression e .

Expressions. Integers occur in two ways: $-n-$ denotes general integer while $-\langle \eta \rangle-$ stands for the only value of type $\langle \eta \rangle$. At this point, no constructs may be given an index type $[\eta]$.

Last, the coercion $-e \triangleright \tau-$ represents an explicit type constraint. Because of refined types, it plays a central role in the definition of the semantics (see [subsection 3.2](#)).

⁶ The use of lists simplifies the association of generalized variables with their instantiations

Arrays. \mathcal{L}^η has no support for array manipulation, neither in types nor in expressions. For typing purposes, arrays are essentially functions on a bounded domain: that is the role of the index refinement. To make examples more intuitive, we will use the notation of Futhark $[\eta]\tau$ [?] as a shorthand for the type $[\eta] \rightarrow \tau$. Allowing sub-typing from functions to arrays seems irrelevant, both for clarity and compilation perspectives. However, the typing issues that arises from refinement types would still occur with a dedicated language support for arrays.

3.2 Semantics

The big-step semantics of $\mathcal{L}^\eta - e \rightsquigarrow v -$ associates to some closed expressions a *value*. As defined in [Figure 1](#), values are either integers or abstractions. The deduction rules are detailed in [Figure 2](#). They are syntax-directed, thus \mathcal{L}^η semantics is deterministic.

Substitutions. Substitutions (ranged over by ρ, \dots) are defined for each syntactical class and variable/element pairs. Their application is written $e\{\rho\}$. Explicit ones are uniformly denoted \cdot/\cdot . Thus, $e\{\eta/\iota\}$ is the substitution in expression e of the occurrences of the size variable ι by the size η , including in sizes and types contained in e . This notation is naturally extended to generalization and instantiation lists, assuming they are compatible, i.e. that each size variable is substituted with a size and likewise for types.

When evaluating let bindings (rule E-LET), each occurrence of the defined variable $-^S x-$ get substituted with its coerced definition in which generalized size and type variables have been instantiated $-(e \triangleright \tau)\{S/V\}$.

Refinements and coercions. The semantics of \mathcal{L}^η is not type-erasable. This obviously transpires in the rule E-SIZE that extracts a value from a size. Moreover, refinements discriminate between values of the same *shape* (or base type) and they must be checked in several places. For instance, the term $-(\lambda x : [\eta]. e) 8-$ should not be reduced further since the argument $-8-$ is not a value of the expected type: $[\eta]$. More generally, for any substitution of a term variable, the substituting value must fulfill the substituted variable refinement. Therefore, the E-APP and E-LET rules insert coercions; hence the need of an explicit coercion construction in expressions.

For integer refinements, coercions check that the refinement is fulfilled (rules C-INT, C-SIZE and C-INDEX). Similarly to λ^H [[Flanagan 2006](#)], function coercions reduce into delayed coercions on argument and result values (rule C-FUN), that will be evaluated upon application. The coercion on argument is actually inserted by the evaluation of the introduced application $-v x-$.

Type independence. Although \mathcal{L}^η semantics is not type-erasable, only sizes have computational content i.e. the observational semantics of an expression does not depend on

$\eta ::=$	ι, κ, δ	Sizes	$V ::= \varepsilon \mid \iota \cdot V \mid \alpha \cdot V$	Generalization
\mid	n	variable	$S ::= \varepsilon \mid \eta \cdot S \mid \tau \cdot S$	Instantiation
\mid	$\eta + \eta$	constant		
\mid	$\eta * \eta$	sum		
		product		
$\tau ::=$	α, β, γ	Types	$e ::=$	Expressions
\mid	$\tau \rightarrow \tau$	variable	\mid	S_x
\mid	int	function	\mid	$e e$
\mid	$\langle \eta \rangle$	integer	\mid	$\lambda x : \tau. e$
\mid	$[\eta]$	singleton	\mid	let d in e
		interval	\mid	$e \triangleright \tau$
			\mid	$\langle \eta \rangle$
			\mid	n
			$d ::=$	$x_V : \tau = e$
				Declarations

Figure 1. Syntax of \mathcal{L}^η . Note that the syntax $\langle \eta \rangle$ is overloaded: it denotes both the singleton type and its unique value. The values used in the semantics are marked with *.

Semantics of closed Expressions	$e \rightsquigarrow v$	Semantics of Coercions	$v \triangleright \tau \rightsquigarrow v'$
E-SIZE $\frac{}{\langle n \rangle \rightsquigarrow n}$	E-APP $\frac{e_1 \rightsquigarrow \lambda x : \tau. e \quad e_2 \triangleright \tau \rightsquigarrow v}{e_1 e_2 \rightsquigarrow v'}$	C-SIZE $\frac{n' = n}{n' \triangleright \langle n \rangle \rightsquigarrow n'}$	C-INDEX $\frac{n' \in [0, n-1]}{n' \triangleright [\eta] \rightsquigarrow n'}$
E-COERCE $\frac{e \rightsquigarrow v \quad v \triangleright \tau \rightsquigarrow v'}{e \triangleright \tau \rightsquigarrow v'}$	E-LET $\frac{e' \{ (e \triangleright \tau) \{ S/V \} / S_x \} \rightsquigarrow v}{\text{let } x_V : \tau = e \text{ in } e' \rightsquigarrow v}$	C-INT $\frac{}{n \triangleright \text{int} \rightsquigarrow n}$	C-FUN $\frac{v = \lambda x : \tau. e}{v \triangleright \tau_1 \rightarrow \tau_2 \rightsquigarrow \lambda x : \tau_1. v x \triangleright \tau_2}$

Figure 2. Semantics of \mathcal{L}^η .

the valuation of its type variables. Changing types (hence possible refinements) only restrict semantics domain.

Definition 3.1 (Observational equivalence). Two closed expressions e_1 and e_2 , are *observationally equivalent* $-e_1 \equiv e_2-$ if and only if, for any closed expressions a_1, \dots, a_k and integers n_1, n_2 :

$$\left\{ \begin{array}{l} e_1 a_1 \dots a_k \rightsquigarrow n_1 \\ e_2 a_1 \dots a_k \rightsquigarrow n_2 \end{array} \right\} \implies n_1 = n_2$$

Expressions for which no such common evaluation environment exist are considered equivalent. Used along with typing assumptions to rule out silly cases, it allows to compare functions on their common domain.

Definition 3.2 (Equality modulo types). Two expressions e_1 and e_2 are *equal modulo types* $-e_1 \approx_\tau e_2-$ if and only if it exists an expression e , free type variables $\bar{\alpha}$ of e and types $\bar{\tau}_1, \bar{\tau}_2$ such that:

$$e_1 = e\{\bar{\tau}_1/\bar{\alpha}\} \wedge e_2 = e\{\bar{\tau}_2/\bar{\alpha}\}$$

Equivalence modulo type is preserved by the semantics:

Theorem 3.3 (Type independence). *Given two closed terms e_1, e_2 such that $e_1 \approx_\tau e_2$, then*

$$\forall v_1 v_2, \left\{ \begin{array}{l} e_1 \rightsquigarrow v_1 \\ e_2 \rightsquigarrow v_2 \end{array} \right\} \implies v_1 \approx_\tau v_2$$

Proof. The above invariant holds across semantics rules thanks to the following observations (detailed in [Appendix A](#)):

- Size substitutions cannot capture types. E-SIZE instances are thus equals and yield the same integer.

- Selecting between C-FUN and other coercion's rules depends only on value shape. \square

Corollary 3.4 (Type observable independence). *Expressions that are equal modulo types are observationally equivalent:*

$$\forall e_1 e_2, e_1 \approx_\tau e_2 \implies e_1 \equiv e_2$$

Proof. Given e_1, e_2 such that $e_1 \approx_\tau e_2$, and closed expressions a_1, \dots, a_k , we immediately have $e_1 a_1 \dots a_k \approx_\tau e_2 a_1 \dots a_k$. Observational equivalence follows because equality modulo types for integers implies equality. \square

3.3 Typing

A type discipline, based on the one of Hindley and Milner [[Hindley 1969](#)], filters the terms for whom a semantics exists, i.e. expressions that may be reduced to a value.

Environment. Expressions are typed in an environment Γ defined as a pair (Γ_v, Γ_e) where $\Gamma_v \subset \mathcal{V}_\eta \cup \mathcal{V}_\tau$ is the set of bound size and type variables and Γ_e is a partial map from term variables to type schemes $-\sigma := \forall V. \tau-$. In the following, terms are supposed to be named so that no clashes occur. The environment is thus unordered: $-\Gamma, x : \forall V. \tau-$ defines variable x , assuming it was not in Γ , and $-\Gamma, V-$ registers V size and type variables into Γ_v , assuming they are unbound too.

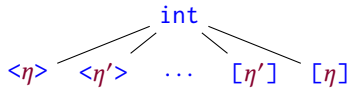
Judgments. The typing judgment $-\Gamma \vdash e : \tau-$ reads ‘in the environment Γ , the expression e has type τ ’. This relation implicitly assumes that τ and e are *well-formed*, i.e. that their free variables are bound in Γ . It is defined alongside the sub-typing relation $-\tau_1 <: \tau_2-$ in [Figure 3](#).

Expressions Typing		$\Gamma \vdash e : \tau$	Sub-typing		$\tau_1 <: \tau_2$
T-VAR	$\frac{\Gamma(x) = \forall V. \tau}{\Gamma \vdash^S x : \tau\{S/V\}}$		S-SIZE	$\frac{}{\langle \eta \rangle <: \text{int}}$	
T-SIZE	$\frac{}{\Gamma \vdash \langle \eta \rangle : \langle \eta \rangle}$		S-INDEX	$\frac{}{[\eta] <: \text{int}}$	
T-INT	$\frac{}{\Gamma \vdash n : \text{int}}$		S-REFL	$\frac{}{\tau <: \tau}$	
T-ABS	$\frac{\Gamma, x : \forall \varepsilon. \tau \vdash e : \tau'}{\Gamma \vdash \lambda x : \tau. e : \tau \rightarrow \tau'}$		S-FUN	$\frac{\tau_2 <: \tau_1 \quad \tau'_1 <: \tau'_2}{\tau_1 \rightarrow \tau'_1 <: \tau_2 \rightarrow \tau'_2}$	
T-APP	$\frac{\Gamma \vdash e_1 : \tau' \rightarrow \tau \quad \Gamma \vdash e_2 : \tau'}{\Gamma \vdash e_1 e_2 : \tau}$		T-LET	$\frac{\Gamma, V \vdash e : \tau \quad \Gamma, x : \forall V. \tau \vdash e' : \tau'}{\Gamma \vdash \text{let } x_V : \tau = e \text{ in } e' : \tau'}$	
T-SUBTYPE	$\frac{\Gamma \vdash e : \tau \quad \tau <: \tau'}{\Gamma \vdash e : \tau'}$		T-COERCE	$\frac{\Gamma \vdash e : \tau}{\Gamma \vdash e \triangleright \tau : \tau}$	

Figure 3. Type system for \mathcal{L}^η , non syntax directed

It is worth mentioning that type equality, used in S-REFL rule among others, requires a formal identity between the sizes appearing in refinements. For instance, types $[\iota]$ and $[2 - \iota]$ are considered different even though they yields equal types when instantiated with $\iota = 1$.

Refinements and sub-typing. General dependent type systems such as DML [Xi and Pfenning 1999] provide rich sub-typing relations based on refinement implication, at the cost of static checking undecidability. For that reason as well as inference perspectives, sub-typing is restricted to inserting or dropping refinements (with respect to the variance). Thus, the relations $[\eta] <: \text{int}$ and $\text{int} \rightarrow \alpha <: [\eta]\alpha$ are valid, whereas the semantically correct relation $[\iota] <: [\iota + 1]$ is invalid. This flat order between refined types, illustrated bellow, is the key restriction to keep type checking decidable: correction only relies on size equalities, instead of general inequalities on polynomials.



Preservation and soundness. This type system enjoys both preservation and soundness: types are preserved by reduction and well-typed terms have a semantics. Formally:

Theorem 3.5 (Preservation).

Given e an expression, τ a type and v a value.

If $\vdash e : \tau$ and $e \rightsquigarrow v$ then $\vdash v : \tau$.

Theorem 3.6 (Soundness).

Given e an expression, type τ such that $\vdash e : \tau$

then there exists a value v such that $e \rightsquigarrow v$.

Proof. The generic construction for big step semantics set up by Dagnino et al. [2020] allows to establish these results from three local properties on the type system and the semantics (see Appendix B): local preservation, \exists -progress and \forall -progress. A key simplification lies in the normalization of typing derivations since our type system is not syntax-directed: the T-SUBTYPE rule may be instantiated anywhere. \square

4 Inference

Although type annotations might be helpful for documentation purposes (e.g., in interfaces), they tend to obfuscate

programs as size expressions get larger. They should be inferred. However, pursuing a full and complete type inference as the HM type discipline enjoys [Hindley 1969] is vain: the size language, that allows non-linear arithmetic expressions, will surely cause unsolvable constraints. Despite this, the size relations that occurs in data-intensive applications are often simple, giving the opportunity to omit most of them. Figure 4 gives implicitly typed definitions of simple linear algebra operation and their inferred type.

One point must be carefully handled: \mathcal{L}^η semantic is specified over *closed typed* terms. Inference must ensure that the semantics of reconstructed terms is fully defined by implicitly typed ones. The *ghost size* issue sketched in section 2 is crucial here: unconstrained size variables should not get defined during reconstruction. As a result, inference must ensure that no unnecessary size relations are introduced.

4.1 Implicitly typed \mathcal{L}^η

As an implicitly typed language, a slight variation of \mathcal{L}^η is used: generalization and instantiation places, i.e. V and S in \mathcal{L}^η syntax, are omitted. Contrary to polymorphism, type annotations are still present in implicitly typed terms, but they might contain size and type variables that are unbound. The inference process builds definitions for polymorphism places, i.e. a list of size and type variables that are generalized or used for instantiation, alongside a substitution of unbound size and type variables. In examples, place-holders ($_$) stand for fresh size or type variables.

4.2 Algorithm

Size equality constraints amount to vanishing polynomials. Unlike types, whose unification is structural, these constraints cannot be solved easily. For that reason, instead of building a substitution on the fly as done by Algorithm \mathcal{W} [Milner 1978], sub-typing constraints and unbound size and type variables are collected by a term traversal and the resulting system is solved at generalization places (i.e. **let**), in the hope of using the simplest constraints to simplify the most complex ones. In the context of sub-typing, similar algorithms were proposed, e.g., by Aiken and Wimmers [1993], where constraints are simplified at generalization points. The constraint collecting algorithm is explained in details in Appendix C.

```

let dot = λu. λv. fold (+) <_> 0 (map2 (*) <_> u v)
let mat_vec = λa. λv. map (vec_vec u) <_> (transpose a)
let mat_mat = λa. λb. map (mat_vec b) <_> a

```

$$\begin{aligned}
& [\forall l. [l] \text{int} \rightarrow [l] \text{int} \rightarrow \text{int}] \\
& [\forall l. \kappa. [l][\kappa] \text{int} \rightarrow [l] \text{int} \rightarrow [\kappa] \text{int}] \\
& [\forall l. \kappa. \delta. [l][\kappa] \text{int} \rightarrow [\kappa][\delta] \text{int} \rightarrow [l][\delta] \text{int}]
\end{aligned}$$

Figure 4. Usual linear algebra primitives defined with iterators

Let bindings introduce generalization: once the definition has been traversed, sub-typing constraints are solved. As for simple ML, the remaining type and size variables that do not appear in the environment are generalized. Moreover, two extra checks are performed on the generalized variables:

1. They should not appear in remaining constraints.
2. They must appear in declaration's type.

The former allows to keep simple polymorphism, while the latter detects unconstrained variables. This last check is crucial since term's semantics depend on sizes: the [section 2](#) gives an example of such ambiguous term:

```
let even = fold (+) <_> 0 (cst 2)  (Error: Unconstrained size)
```

4.3 Principal typing

Before presenting the constraint resolution strategy, let us focus on a thorn in our side: this type system does not enjoy principal types, i.e. some declarations do not have a most general type scheme. Comparison of type schemes is defined by the *subsumption* relation presented by [Jones et al. \[2007\]](#). Informally $\sigma_1 \leq \sigma_2$ if and only if any instance of σ_2 may be obtained by instantiating σ_1 and using sub-typing. σ_1 is then *more general* than σ_2 . It naturally defines a notion of equivalence, that amounts for simple ML types (without sizes), to a renaming of type variables. Because size equality is not structural, this relation widens here: the uni-dimensional convolution defined in [section 2](#) may be given the following type schemes:

```

val convolution : ∀ l. κ. [κ] int → [l] int → [l - κ + 1] int
val convolution : ∀ l. κ. [κ] int → [l + κ - 1] int → [l] int

```

Any instance of the first is an instance of the second, and reciprocally. More importantly, some terms may be given multiple type schemes that have no common generalization; this must be carefully handled by the inference. There are two reasons for this:

1. Polynomial sizes. Allowing more than linear expressions for sizes surely causes constraints with multiple solutions. Given a function `split`, declared below, that transforms a 1-dimensional array into a 2-dimensional one, its application to an 'array' of size 4 raises several possible types, corresponding to different semantics:

```

val split : ∀ l. κ. α. [l * κ] α → [l][κ] α
let mat = split (λi: [4]. 0)

```

$$\begin{cases} [1][4] \text{int} \\ [2][2] \text{int} \\ [4][1] \text{int} \end{cases}$$

In such a situation, the underlying constraint ($l * \kappa - 4 = 0$) will not be solved (see [section 4.4](#)), and inference will fail, asking for more annotations.

2. Sub-typing and simple polymorphism. Unconstrained polymorphism forces refinements to be selected for all occurrences of integer types. The `slope` function below computes the ratio of images' difference over arguments' difference, assuming suitable arithmetic operators defined over integers.

```
let slope = λf. λi. λj. (f i - f j) / (i - j)
```

The subtlety comes from the simultaneous applications of f to i and j : should f 's domain and both argument share the same refinement? Indeed, possible type schemes include:

$$\begin{aligned}
& \forall l. (<l> \rightarrow \text{int}) \rightarrow <l> \rightarrow <l> \rightarrow \text{int} & (\sigma_s^s) \\
& \forall l. \kappa. (\text{int} \rightarrow \text{int}) \rightarrow <l> \rightarrow <\kappa> \rightarrow \text{int} & (\sigma_s^b) \\
& (\text{int} \rightarrow \text{int}) \rightarrow \text{int} \rightarrow \text{int} \rightarrow \text{int} & (\sigma_b^b) \\
& \forall l. \kappa. (\text{int} \rightarrow \text{int}) \rightarrow [l] \rightarrow [\kappa] \rightarrow \text{int} & (\sigma_l^b) \\
& \forall l. ([l] \rightarrow \text{int}) \rightarrow [l] \rightarrow [l] \rightarrow \text{int} & (\sigma_l^l)
\end{aligned}$$

Among them, $\sigma_b^b \leq \sigma_s^s$ and $\sigma_b^b \leq \sigma_l^l$. Others are incompatible pair-wise (denoted $\not\leq$) for multiple reasons: refinements...

1. are incompatible: $\sigma_s^s \not\leq \sigma_l^l$; $\sigma_s^b \not\leq \sigma_l^b$; $\sigma_s^b \not\leq \sigma_l^l$; $\sigma_s^s \not\leq \sigma_b^b$
2. appear covariant and contravariant: $\sigma_s^s \not\leq \sigma_b^b$; $\sigma_l^l \not\leq \sigma_b^b$
3. impose extra size constraints: $\sigma_s^s \not\leq \sigma_s^b$; $\sigma_l^l \not\leq \sigma_l^b$.

Constrained polymorphism. Using simple polymorphism with sub-typing is unusual. The general theory proposed by [Aiken and Wimmers \[1993\]](#) provides constrained types schemes. Shrinking the constraint set at generalization point is then the key to avoid an exponential blow-up of constraints [\[Pottier 2001\]](#). Such systems enjoy the principal types property. In this context, the function `slope` would be given the type scheme:

$$\begin{aligned}
& \forall \alpha. \beta. \gamma. | \alpha <: \text{int} \wedge \beta <: \alpha \wedge \gamma <: \alpha. \\
& (\alpha \rightarrow \text{int}) \rightarrow \beta \rightarrow \gamma \rightarrow \text{int}
\end{aligned}$$

However, modularity would be sacrificed here, by deferring size constraints resolution to monomorphic instantiations. Coupled with the loss of readability of such types, this is the main reason for keeping simple polymorphism.

Inference and semantics. This issue about principal type is all the more crucial because our semantics is not type erasable. Sizes have computational contents in our language. For that reason, inference should ensure that no sizes have been arbitrarily defined. We formalize this in [subsection 4.5](#).


```

let slope = λf:[_] → _ . λi: _ . λj: _ . (f i - f j) / (i - j)      [∀ι. ([ι] → int) → [ι] → [ι] → int]
let slope = λf: _ → _ . λi:[ι]. λj:[ι]. (f i - f j) / (i - j)      [∀ι. ([ι] → int) → [ι] → [ι] → int]
let slope = λf: _ → _ . λi:[_]. λj:[_]. (f i - f j) / (i - j)      [∀ι.κ. (int → int) → [ι] → [κ] → int]

```

Figure 5. Different type annotations in the slope example lead to different type schemes.

4.4 Constraint solving

Solving the constraint system aims at extracting from the set of sub-typing constraints a *most general unifier*, i.e. a necessary substitution of the free variables. This is achieved gradually: (i) types (without refinements) are inferred using structural unification; (ii) necessary refinements of type `int` are selected; (iii) sizes constraints are solved; (iv) refinements are propagated. Similar stratification has been previously used for inference in extended type systems [Knowles and Flanagan 2007; Rondon et al. 2008; Xi and Pfenning 1998]. However these steps are utterly entangled in our proposal: instead of separating phases across multiple passes, types, refinements and sizes get partially defined at each solving point (**let**), allowing an easier handling of polymorphism than it would be possible with disconnected inference passes.

To illustrate our overview of the solving process, three slightly modified version of the slope example used previously are defined in Figure 5: some annotations are added, constraining the refinements at different places.

(i) Types. To begin with, refinements are ignored to build *simple types* that will be made precise in the subsequent phases. By replacing every refinements with `int`, sub-typing relations are turned into equalities. They are solved using structural unification, failing in the usual modalities (e.g., top-level type constructor inequality, cyclic types). It generates a *most general unifier* [Milner 1978]. At that point, Each declaration of slope get type $(\text{int} \rightarrow \text{int}) \rightarrow \text{int} \rightarrow \text{int} \rightarrow \text{int}$.

(ii) Refinements. The integer types previously derived may now be refined: each occurrence of `int` in the substitution are replaced by fresh type variables. Sub-typing constraints are distributed with the S-FUN rule (the usual variance rule), leading to simple constraints containing variables and refined types. The ones of the form $\alpha <: [_]$, $\alpha <: <_>$ and $\text{int} <: \alpha$ define variable α 's refinement while the unsolvable constraints such as $\text{int} <: [_]$ lead to errors that are reported to the programmer. Refinements are not propagated further: this is postponed after size resolution, since adding refinements may introduce constraints between sizes that would otherwise be unrelated. Unconstrained variables get thus substituted with `int`.

In our example, the first definition of the slope function get type $([_] \rightarrow \text{int}) \rightarrow [_] \rightarrow [_] \rightarrow \text{int}$ while two others get $(\text{int} \rightarrow \text{int}) \rightarrow [_] \rightarrow [_] \rightarrow \text{int}$.

(iii) Sizes. Then, sub-typing constraints are distributed again, extracting size equalities, i.e. vanishing polynomials $\eta_1 - \eta_2 = 0$ for the constraints of the form $<\eta_1> <: <\eta_2>$ or

$[\eta_1] <: [\eta_2]$. The resulting polynomial system C^η is solved, by deriving a *most general substitution*:

Definition 4.1 (Most general substitution). Given a constraint set C , a substitution ρ is *most general* if and only if for any substitution ρ' , such that $\vdash C\{\rho'\}$, then there exists a substitution ρ'' such that $\rho' = \rho \circ \rho''$.

For now, we have implemented a simple resolution strategy that eliminates *isolated variables*, i.e. substituting η for ι when a constraint $\iota - \eta = 0$ exists. The resulting substitution is immediately a most general one. This task could be delegated to an external solver, but this is unpractical in the context of safety-critical software for certification purposes⁷. Moreover, this elementary strategy works for most of the size constraints we encountered. Adding some annotations helps for the remaining cases.

Our three versions of slope get respectively types:

1. $([\iota] \rightarrow \text{int}) \rightarrow [\iota] \rightarrow [\iota] \rightarrow \text{int}$
2. $(\text{int} \rightarrow \text{int}) \rightarrow [\iota] \rightarrow [\iota] \rightarrow \text{int}$
3. $(\text{int} \rightarrow \text{int}) \rightarrow [\iota] \rightarrow [\kappa] \rightarrow \text{int}$

In particular, ι and κ sizes are not unified in the last declaration, because sub-typing is at each application.

(iv) Propagation. Last, refinements are propagated. This aims at making types more accurate. Among the type variables introduced during refinement inference, the ones whose lower bound only contains a unique type $[\iota]$ or $<\iota>$, are defined accordingly.

In the third definition, i and j refinements differ: f domain cannot be refined. Conversely, refinements turns out to be equal in the second version: they are propagated.

4.5 Inference properties

The expressiveness of \mathcal{L}^η size language allows no hope for a complete inference. Nevertheless, it is sound and we expect inference to be *non-specializing*, i.e. that it rejects any terms with ambiguous semantics.

Theorem 4.2 (Inference soundness). *Given an implicitly typed expression, if inference succeeds, the reconstructed term is well-typed.*

Proof sketch. The detailed proof is available in Appendix C, alongside a formalization of the inference algorithm. It established the following invariant: for each sub-term and a substitution that solves the constraints gathered by the algorithm, it exists a type derivation for the substituted term, in the same environment. \square

⁷ Even though the compiler is not embedded, it must fulfill the highest certification level, hence the need of a certified solver...

Conjecture 4.3 (Inference non-specialization). *Given an expression e , if inference succeeds and produces a closed term e' , then for any reconstructed term e'' , $e' \equiv e''$.*

This proof has not been fully conducted yet. The main difficulty lies in the handling of let bindings: the constraints set induced by variable instantiation differ. It remains to show that these constraints are stronger in an arbitrary reconstruction than in the inferred one, so that the latter is the most general.

5 Toward a realistic array language

Our simplistic language \mathcal{L}^n provides the basis for an ML-like language where sizes are handled in a same way than types, i.e. with polymorphism. This section gives an insight of some extensions that are necessary for a more realistic array language.

5.1 Locally abstract sizes

OCAML proposed *locally abstract types* that allow to declare types within a scope. These types may not escape their scope, i.e. no substitution of an outer free type variable may capture them. They serve advanced purposes (first-class modules, GADTs, ...) by introducing type variables that can be generalized.

Providing a similar mechanism, for both sizes and types, has a simpler use in our context: local existential quantification — **let size** $\iota = e$ **in** e' . It defines an abstract size ι in e' , using the value of an arbitrary expression e . Such a mechanism is useful to overcome size language limits.

5.2 Polymorphic recursion

Recursive algorithms on array such as the Fast Fourier Transform calls themselves on sub-arrays whose sizes vary at each call. Because sizes are part of types, polymorphic recursion is needed. This extension has been extensively studied [Meertens 1983; Mycroft 1984].

Semantics. While adding recursive declarations requires few changes on the implementation, it impacts deeply the formalization: diverging terms now exist, they must be distinguished from blocked ones in \mathcal{L}^n big-step semantics. Hopefully, the Dagnino et al. [2020] formalization was designed for non-deterministic semantics. By giving a non-deterministic evaluation of fixpoint (either stopping with an error value or reducing further), the preservation and soundness results (subsection 3.3) may be extended.

Inference. As Henglein [1993] showed, polymorphic recursion turns inference into an undecidable problem. We follow the classical approach, by considering fix-points monomorphic unless explicitly generalized at declarations.

An extra check is required to ensure that the actual type of the declaration is indeed as polymorphic as the specified one, as defined in the subsumption relation in subsection 4.3.

This validates *a posteriori* that the recursive occurrences of the introduced variable have been correctly instantiated.

5.3 Explicit coercions

The size language might get too limited, especially when local existential sizes are used. We provide an *explicit coercion* $-e \triangleright \tau-$ to spot some size properties that cannot be checked by the type system: expression's type and τ 's sizes must only have the same structure, allowing for size mismatches.

As mentioned by Jay and Sekanina [1997], coercions may be checked in various ways: at run-time, with defensive code or using alternative formal verification tools. In the context of static sizes, Nielson and Nielson [1988] proposed *binding times*, to ensure that coercions (and local existential sizes) are computable at compiler time.

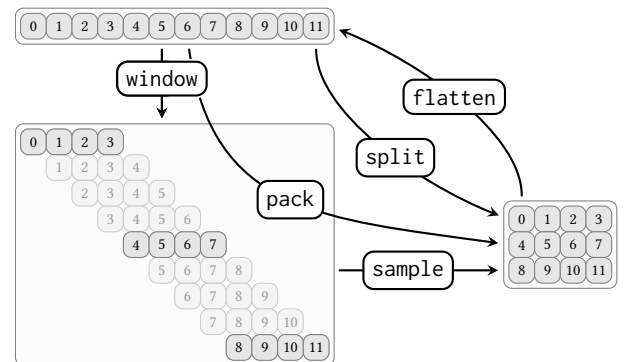
5.4 Language support for arrays

Arrays deserve special language constructs, both for readability and compilation purposes. Besides distinguishing their types from function's one, \mathcal{L}^n should be extended with dedicated syntax for accesses $-e[e']-$ and array definition $-[e, \dots, e]-$. As far as typing is concerned, these constructs amount for type constraint insertions, i.e. $e[e']$ requires sub-expressions to have type $[i]\alpha$ and $[i]$.

To avoid operator overloading, index manipulations are provided as a set of *first order combinators* that transform arrays' shape. They provide a safe way to introduce correct index computation. In addition to the usual transpose, reverse and concat linear primitives⁸, the following ones are added. They are illustrated below.

```
val window :  $\forall i.\kappa.\alpha. \langle \kappa \rangle \rightarrow [i + \kappa - 1]\alpha \rightarrow [i][\kappa]\alpha$ 
val sample :  $\forall i.\kappa.\alpha. \langle \kappa \rangle \rightarrow [i * \kappa - \kappa + 1]\alpha \rightarrow [i]\alpha$ 
val split :  $\forall i.\kappa.\alpha. \langle \kappa \rangle \rightarrow [i * \kappa]\alpha \rightarrow [i][\kappa]\alpha$ 
val flatten :  $\forall i.\kappa.\alpha. \langle \kappa \rangle \rightarrow [i][\kappa]\alpha \rightarrow [i * \kappa]\alpha$ 
```

```
let pack =  $\lambda s. \lambda x. \text{sample } s \text{ (window } \langle \_ \rangle x)$ 
            $[ \forall i.\kappa.\delta.\alpha. \langle \delta \rangle \rightarrow [i * \delta - \delta + \kappa]\alpha \rightarrow [i][\kappa]\alpha ]$ 
```



The window function (see section 2) builds a matrix whose rows are slices of the input array. The sample function extracts one element out of every κ , selecting both ends of the

⁸ With the iterators, these are the available array functions in SCADE.

array. The size of the input array must thus be a multiple of κ plus 1. Composing these functions defines a general sampling operator `pack`. It selects ι slices of size κ that are uniformly distributed and cover both ends of the input array, whose size is obtained by considering a sampling step δ . Note that the size argument of `sample` and `pack` are necessary since the associated variable might not be deduced from array size (because of non linear sizes). Although redundant, the `split` primitive carries a extra information: it defines a bijection between arrays. Defining filters or convolutions requires such building blocks (see [section 2](#)). Here is an example of `pack` application:

$$\text{pack } \langle 2 \rangle (\lambda i: [7]. i) \triangleright [3][3]\text{int} \rightsquigarrow \begin{bmatrix} 0 & 1 & 2 \\ 2 & 3 & 4 \\ 4 & 5 & 6 \end{bmatrix}$$

5.5 Implicit size parameters

Our proposal allows to infer any *size*. However, functions' arguments of type $\langle \eta \rangle$ might not be syntactically omitted: an explicit *size value* expression with an unspecified size $\langle _ \rangle$ must at least be provided.

To make these arguments fully implicit, some syntactical restrictions are needed, so as to determine the places where such unspecified size values should be inserted. For instance, providing n -ary size abstraction $\lambda \langle t_1, \dots, t_n \rangle. e$ and application $e \langle t_1, \dots, t_n \rangle$ without curryfication makes possible to infer missing applications from simple type skeletons.

6 Purposes of sized types

The size information has several usages, both for program verification and compilation. Carrying it through types reveals practical in our ongoing experimentation.

6.1 Verification

As mentioned in [section 2](#), array combinators turn bound checking into size consistency (e.g., `map2` arguments must have the same size). Our type system precisely ensures this property. For decidability purposes, the type system only handles *size equalities*. In particular, it does not ensure that size are positive. In this context, array safety is based on the emptiness of type $[\eta]$ when η is negative or null: none of the language's primitives deliver indexes of negative size.

Thinking of array combinators as pure index computations, as we briefly discuss about below, linear array primitives (`concat`, `reverse`, `window`, ...) as well as index producers (`mapi`, the modulo), are indeed safe, but Pandora's box opens when providing non-linear primitives: the use of a negative step κ in `sample` would allow to build an array of positive (thus nonempty) size from a negatively sized one, hence introducing faulty accesses.

To rule out these hazardous uses, constraints must be added to type schemes. The `split` primitive is restricted to strictly positive steps⁹ with the type:

val `sample` : $\forall \iota. \kappa. \alpha. \langle \kappa \rangle \rightarrow [\iota * \kappa - \kappa + 1] \alpha \rightarrow [\iota] \alpha$ **where** $\kappa > 0$

These constraints are ignored by the inference. They are checked either symbolically or at final instantiation, where they become trivial relations on integer values.¹⁰

6.2 Compilation

Conveying sizes into types is useful for compilation purposes. In particular, the declarative style favors definitions of complex data by pieces that are aggregated (e.g., using array concatenation). To avoid extra memory consumption and data moves, the placement of each part must be carefully chosen. This is for example the role of the *built-in-place* optimization designed by [Gaudiot et al. \[1997\]](#) for SISAL. For arrays, it strongly relies on size information.

Iterator fusion. Complex transformations are expressed by composing extensional primitives that produce intermediate arrays. Fusing these atomic operations is an indispensable compilation pass for functional array languages that target efficient software. In some of them [\[Steuwer et al. 2015; ?\]](#), this is achieved by using a set of rewrite rules with the drawback of requiring new rules for additional primitives or some fallback mechanism.

Other proposals such as OBSIDIAN [\[Claessen et al. 2012\]](#) or DEX [\[Paszke et al. 2021\]](#) rely on the array-function analogy to provide forms that compose arbitrarily. We experiment a similar approach, restricted to array combinators by representing them in a uniform way: functions that map indexes. For instance, reversing an array of size η is described by the function $x \in [\eta] \mapsto \langle \eta - 1 \rangle - x \in [\eta]$, which captures the size. During code generation, these index functions induce computed array accesses that are correct by construction.

Unchecked accesses. Currently in SCADE [\[Colaço et al. 2017\]](#), every dynamic array accesses are guarded, by providing a default value in the event that the index is out of bound. For accesses where bounds are actually met, such as the iterator `mapi` (point wise application with index), this generates dead code and an extra branching. The index refinement allows to decouple array access from bound verification: a value of type $[\eta]$ may be used in several places without any dynamic check that it is indeed within bounds.

7 Discussion and Related works

The definition of a typed functional language with array operations offers several design choices that must be assessed.

⁹ A zero step could make sense here, by accessing the single value of an array (of size 1), but this stricter version enjoys an extra property: it is injective, which gives additional compilation perspectives.

¹⁰ Similar constraints on type variables are already used in SCADE, as shown in the linear algebra examples of [section 7](#)

How expressive are the language and its type system; how difficult and modular are type checking and type reconstruction; what about the verbosity of the code; what is diagnosis like in case of errors?

The motivation of the present work is the extension of the domain-specific functional and synchronous language SCADE [Colaço et al. 2017] that is used for implementing real-time embedded software. SCADE stand out from general purpose functional language by being first-order with a pre-defined set of higher-order operators on arrays, extending a proposal for LUSTRE Maraninchi and Morel [2004]). The expressiveness of the language is purposely limited in order to ensure safety properties (e.g., memory and execution time are bounded and known statically). Moreover, applications in SCADE are almost exclusively developed graphically by connecting blocks in diagrams so that annotating elements (wires and blocks) with explicit types is rapidly cumbersome, in particular when size expressions get larger. Neither type nor size inference are currently available in SCADE.

Hence, we aim at relaxing some constraints of SCADE for writing array operations while keeping the same safety guarantees. The current version of the language limits the manipulation of array through a set of array iterators (e.g., map, fold, and a few others. SCADE does not provide any size or type inference, nor extensional definitions as proposed in the present paper. For example, the matrix product given in Figure 4 is written in SCADE in the following way.

```
-- Scalar product of two vectors: u(n) · v(n)
function dot «n» (u, v: 'T^n)
returns (w: 'T) where 'T numeric
  w = (fold $+$ «n») (0, (map $*$ «n») (u, v));

-- Product of a matrix by a vector: A(m,n) * u(n)
function mat_vec «m, n» (A : 'T^m^n; u : 'T^n)
returns (w: 'T^m) where 'T numeric
  w = (map (dot «n») «m») (transpose (A; 1; 2), u^m);

-- Matrix product: A(m,n) * B(n,p)
function mat_mat «<m, n, p>>» (A : 'T^m^n; B : 'T^n^p)
returns (C: 'T^m^p) where 'T numeric
  C = (map (mat_vec «m, n») «p») (A^p, B);
```

Here, sizes need to be expressed both in types and instantiations of array iterators and functions. The proposition presented in the paper increases significantly what it is possible to express with the current version of SCADE with lighter-weight notations for both definitions and interfaces.

We hope this proposition to be applicable in a wider context than the one of SCADE. Actually, this polymorphism approach deals with features that are not available with SCADE such as higher-order and recursion.

Modularity. Circuit design languages such as LAVA [Bjesse et al. 1998] and WIRED [Axelsson et al. 2005] extensively

use arrays. Because of their target, programs are fully expanded before size checking. This allows using arbitrary (static) expressions in sizes that are evaluated at compile time to ensure correct array use. The same approach is used for HALIDE [Ragan-Kelley et al. 2013], that produces GPU kernels: functions are compiled and optimized once sizes have been given concrete values.

For safety critical embedded software, sizes are also statically fixed, but both checking and compilation gain from modularity, e.g., by allowing easier definition of libraries, or in a view to produce modular code. For error tracking, this type system allows to spot defects of polymorphic definitions before their use in a monomorphic context.

However, since sizes are static, we do not restrict our language to the formally type-provable programs. We provide coercions as a fallback mechanism for remaining checks to be performed after specialization.

A rudimentary system of refined types. Our proposal uses a very restricted form of refinement types, by providing only singleton ($\langle \eta \rangle$) and interval ($[\eta]$) refinements, without sub-typing between them. This is a key for both type checking and inference.

The general theory of dependent types worked out by Xi and Pfenning [1999] allows to express arbitrary predicates in type systems. However, this has a cost: type checking is undecidable in general, mainly because sub-typing amounts to proof obligations of predicate implication. These authors also delineated in [Xi and Pfenning 1998] some restrictions for arrays size checking. Trojahnner and Grelck [2009] extended a similar type system to provide dependently typed rank polymorphism. Both works extract sets of arithmetic constraints that are resolved with an external procedure (SMT solvers). Besides requiring heavy machinery for type checking at the risk of opaque errors, size relations are mainly limited to linear expressions, for the constraint system to be solvable. To lift this reductions, λ^H proposes *hybrid type checking* [Flanagan 2006], a system of refinement types that allows deferring unprovable implications to run time. Our proposal resembles, using static evaluation to eliminate remaining checks.

Inference in extended type systems has been studied in the context of λ^H [Knowles and Flanagan 2007] and LIQUID TYPES [Rondon et al. 2008]. Both approaches are similar to ours: after collecting the set of sub-typing constraints, a most general solution is extracted (using external tools such as SMT solvers).

Explicit proof obligations (coercions). Unless limiting type refinements to simple expressions (linear), some escape mechanism is needed for the terms that cannot be checked by the type systems. In λ^H [Flanagan 2006] and SAC [Tang and Grelck 2012] such coercions are ubiquitous, although implicit: they are systematically inserted at function applications, generating checks that are eliminated at compile time if possible.

In the restricted context of arrays, far less coercions are needed than in the general setting. Thus, we follow the approach of FUTHARK [Henriksen and Elsmann 2021] and DEX [Paszke et al. 2021] (coercions occur at fromOrdinal calls): every possibly failing coercion is explicit. This strengthens the guarantees provided by typing: coercion errors may occur only at explicitly marked points of the program.

Jay and Sekanina [1997] mention several ways to check shape constraints, that apply for our coercions: defensive run-time code generation, static checking with advanced formal methods or partial evaluation at compile time.

A separated size language. Our sizes are syntactically distinguished from general expressions. This matters for cross-compilation (which is common for embedded applications), because sizes are symbolically manipulated at compile-time i.e. on the development host, whereas integer values must be represented on the targeted device with machine (finite) integers that are submitted to overflows. Converting sizes into machine integer is thus non-trivial: then compiler must ensure size's value is actually representable within the concrete type.

A similar size language was proposed by Hughes et al. [1996], so as to bound the size of recursive data-types. Thanks to a distinct language, unbounded integers are extended with their limit ω , to represent data of arbitrary size.

The VEC language of Jay and Sekanina [1997] follows a different path, by enforcing size expressions (a subset of expressions of the language) to be independent of data that come with dedicated typing rules. This results in an analogous restrictions: sizes are static.

Comparison to FUTHARK and DEX. The FUTHARK [Elsmann et al. 2019] and DEX [Paszke et al. 2021] languages shares strong similarities with this work. The founding principle seems similar: most array sizes should be controlled in some inexpensive ways, without trying to fully check arrays, at the risk of limiting language expressiveness. Instead of proving predicates, these type systems keep track of values' properties (bounds), allowing to decouple their verification, either static (argument assumptions) or dynamic (coercions) and their uses (array accesses). This finer control also benefits the compilation by helping to rule out redundant checks.

The main difference between these works lies in the size language. Rather than limiting size to constants or variables, polynomial sizes extend type system expressiveness: neither concat nor reshape may be given a satisfactory type in DEX or FUTHARK. Despite the loss of completeness for inference, we think this extension of the size language useful because it stays easily checkable.

Polymorphism or dependent types? Separating the size language has a direct consequence: notions of scopes, abstractions and applications are needed to express terms that are generic in sizes. Because they are not terms, dependent

types are inadequate here. Following Hughes et al. [1996], we handle similarly sizes and types (let generalization), which is consistent, at least in the context of static sizes.

The dimension type system proposed by Kennedy [1994] also resembles ours: polymorphism over dimensions is considered. Its inference enjoys principal types, since the dimension language is simpler: it is equipped with a single operation, the product, that is both associative and commutative. However, some difficulties still echo to ours: most general types are not unique and polymorphic recursion seems rapidly necessary.

Even FUTHARK, whose type system uses dependent types, shares strong similarities: its normalized form [Henriksen and Elsmann 2021] is obtained by adding let bindings so as to name and scope existentially quantified size variables.

Moreover, size polymorphism may express dynamically sized data-structure. First class polymorphism, as presented by Jones [1997]; Jones et al. [2007], allows data with quantified types (either in size or type, existentially or universally). It relies heavily on the local quantification sketched in section 5. Dynamic arrays could be encoded as a pair $\text{DynArr} : \exists l. \langle l \rangle * [l] \alpha \rightarrow \alpha \text{ dArray}$ where l is existentially quantified, while still locally relying on inference to complete redundant sizes.

8 Conclusion and Perspectives

This article have presented an ML-like type system which adds a size information into types: genericity on sizes is expressed through polymorphism. The size language, made of multivariate polynomials, allows to express a large class of array manipulations, while being easily checked.

Our proposal is not restricted to safety critical languages like SCADE: it may provide the key elements to track array sizes in a functional language, and to highlight the parts that cannot be checked with a simple ML-like type system. This information is valuable for both checking and compilation steps, in particular to reduce the need for defensive code.

Besides assessing the compatibility of this type system with temporal constructs of synchronous languages, the efficient compilation of the proposed array manipulations (array iterators, recursion on size) in the context of safety critical software will be studied next. In particular, targeting devices that are used for intensive computation such as GPUs remains an open question for such applications.

References

- Aiken, A. and Wimmers, E. L. (1993). Type inclusion constraints and type inference. In Williams, J., editor, *Proceedings of the conference on Functional programming languages and computer architecture, FPCA 1993, Copenhagen, Denmark, June 9-11, 1993*, pages 31–41. ACM.
- Axelsson, E., Claessen, K., and Sheeran, M. (2005). Wired: Wire-aware circuit design. In Borriore, D. and Paul, W. J., editors, *Correct Hardware Design and Verification Methods, 13th IFIP WG 10.5 Advanced Research Working Conference, CHARME 2005, Saarbrücken, Germany, October 3-6, 2005, Proceedings*, volume 3725 of *Lecture Notes in Computer Science*,

- pages 5–19. Springer.
- Barnes, J. G. P. (2003). *High Integrity Software - The SPARK Approach to Safety and Security*. Addison-Wesley.
- Bird, R. S. and Wadler, P. (1988). *Introduction to functional programming*. Prentice Hall International series in computer science. Prentice Hall.
- Bjesse, P., Claessen, K., Sheeran, M., and Singh, S. (1998). Lava: Hardware design in haskell. In Felleisen, M., Hudak, P., and Queinnec, C., editors, *Proceedings of the third ACM SIGPLAN International Conference on Functional Programming (ICFP '98)*, Baltimore, Maryland, USA, September 27–29, 1998, pages 174–184. ACM.
- Claessen, K., Sheeran, M., and Svensson, J. (2012). Expressive array constructs in an embedded GPU kernel programming language. In Acar, U. A. and Costa, V. S., editors, *Proceedings of the POPL 2012 Workshop on Declarative Aspects of Multicore Programming, DAMP 2012*, Philadelphia, PA, USA, Saturday, January 28, 2012, pages 21–30. ACM.
- Colaco, J., Pagano, B., and Pouzet, M. (2017). SCADE 6: A formal language for embedded critical software development (invited paper). In Mallet, F., Zhang, M., and Madelaine, E., editors, *11th International Symposium on Theoretical Aspects of Software Engineering, TASE 2017*, Sophia Antipolis, France, September 13–15, 2017, pages 1–11. IEEE Computer Society.
- Dagnino, F., Bono, V., Zucca, E., and Dezani-Ciancaglini, M. (2020). Soundness conditions for big-step semantics. In Müller, P., editor, *Programming Languages and Systems - 29th European Symposium on Programming, ESOP 2020, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2020, Dublin, Ireland, April 25–30, 2020*, *Proceedings*, volume 12075 of *Lecture Notes in Computer Science*, pages 169–196. Springer.
- Elsman, M., Henriksen, T., and Serup, N. G. W. (2019). Data-parallel flattening by expansion. In Gibbons, J., editor, *Proceedings of the 6th ACM SIGPLAN International Workshop on Libraries, Languages and Compilers for Array Programming, ARRAY@PLDI 2019*, Phoenix, AZ, USA, June 22, 2019, pages 14–24. ACM.
- Feo, J., Cann, D. C., and Oldehoeft, R. R. (1990). A report on the Sisal language project. *J. Parallel Distributed Comput.*, 10(4):349–366.
- Flanagan, C. (2006). Hybrid type checking. In Morrisett, J. G. and Jones, S. L. P., editors, *Proceedings of the 33rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2006*, Charleston, South Carolina, USA, January 11–13, 2006, pages 245–256. ACM.
- Gaudiot, J.-L., Bohm, W., Najjar, W., DeBoni, T., Feo, J., and Miller, P. (1997). The Sisal model of functional programming and its implementation. In *Parallel Algorithms/Architecture Synthesis*, *Proceedings*. Second Aizu International Symposium, pages 112–123. IEEE.
- Gérard, L., Guatto, A., Pasteur, C., and Pouzet, M. (2012). A modular memory optimization for synchronous data-flow languages: application to arrays in a lustre compiler. In Wilhelm, R., Falk, H., and Yi, W., editors, *SIGPLAN/SIGBED Conference on Languages, Compilers and Tools for Embedded Systems 2012, LCTES '12*, Beijing, China - June 12 - 13, 2012, pages 51–60. ACM.
- Halbwachs, N., Caspi, P., Raymond, P., and Pilaud, D. (1991). The synchronous data flow programming language LUSTRE. *Proc. IEEE*, 79(9):1305–1320.
- Henglein, F. (1993). Type inference with polymorphic recursion. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 15(2):253–289.
- Henriksen, T. and Elsmann, M. (2021). Towards size-dependent types for array programming. In Low, T. M. and Gibbons, J., editors, *ARRAY 2021: Proceedings of the 7th ACM SIGPLAN International Workshop on Libraries, Languages and Compilers for Array Programming*, Virtual Event, Canada, 21 June, 2021, pages 1–14. ACM.
- Hindley, R. (1969). The principal type-scheme of an object in combinatory logic. *Transactions of the american mathematical society*, 146:29–60.
- Hughes, J., Pareto, L., and Sabry, A. (1996). Proving the correctness of reactive systems using sized types. In Boehm, H. and Jr., G. L. S., editors, *Conference Record of POPL '96: The 23rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Papers Presented at the Symposium*, St. Petersburg Beach, Florida, USA, January 21–24, 1996, pages 410–423. ACM Press.
- Jay, C. B. and Sekanina, M. (1997). Shape checking of array programs. In *Computing: the Australasian Theory Seminar, Proceedings*, volume 19 of *Australian Computer Science Communications*, pages 113–121. University of Technology, Sydney, Australia.
- Jones, M. P. (1997). First-class polymorphism with type inference. In Lee, P., Henglein, F., and Jones, N. D., editors, *Conference Record of POPL '97: The 24th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Papers Presented at the Symposium*, Paris, France, 15–17 January 1997, pages 483–496. ACM Press.
- Jones, S. L. P., Vytiniotis, D., Weirich, S., and Shields, M. (2007). Practical type inference for arbitrary-rank types. *Journal of functional programming*, 17(1):1–82.
- Kennedy, A. (1994). Dimension types. In Sannella, D., editor, *Programming Languages and Systems - ESOP '94, 5th European Symposium on Programming*, Edinburgh, UK, April 11–13, 1994, *Proceedings*, volume 788 of *Lecture Notes in Computer Science*, pages 348–362. Springer.
- Knowles, K. L. and Flanagan, C. (2007). Type reconstruction for general refinement types. In Nicola, R. D., editor, *Programming Languages and Systems, 16th European Symposium on Programming, ESOP 2007, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2007, Braga, Portugal, March 24 - April 1, 2007*, *Proceedings*, volume 4421 of *Lecture Notes in Computer Science*, pages 505–519. Springer.
- Maraninchi, F. and Morel, L. (2004). Arrays and contracts for the specification and analysis of regular systems. In *4th International Conference on Application of Concurrency to System Design (ACSD 2004)*, 16–18 June 2004, Hamilton, Canada, pages 57–66. IEEE Computer Society.
- Meertens, L. G. L. T. (1983). Incremental polymorphic type checking in B. In Wright, J. R., Landweber, L., Demers, A. J., and Teitelbaum, T., editors, *Conference Record of the Tenth Annual ACM Symposium on Principles of Programming Languages*, Austin, Texas, USA, January 1983, pages 265–275. ACM Press.
- Milner, R. (1978). A theory of type polymorphism in programming. *Journal of computer and system sciences*, 17(3):348–375.
- Mycroft, A. (1984). Polymorphic type schemes and recursive definitions. In Paul, M. and Robinet, B. J., editors, *International Symposium on Programming, 6th Colloquium*, Toulouse, France, April 17–19, 1984, *Proceedings*, volume 167 of *Lecture Notes in Computer Science*, pages 217–228. Springer.
- Nielson, H. R. and Nielson, F. (1988). Automatic binding time analysis for a typed lambda-calculus. *Science of computer programming*, 10(1):139–176.
- Paszke, A., Johnson, D. D., Duvenaud, D., Vytiniotis, D., Radul, A., Johnson, M. J., Ragan-Kelley, J., and Maclaurin, D. (2021). Getting to the point: index sets and parallelism-preserving autodiff for pointful array programming. *Proceedings of the ACM on Programming Languages*, 5(ICFP):1–29.
- Pottier, F. (2001). Simplifying subtyping constraints: A theory. *Information and computation*, 170(2):153–183.
- Ragan-Kelley, J., Barnes, C., Adams, A., Paris, S., Durand, F., and Amarasinghe, S. P. (2013). Halide: a language and compiler for optimizing parallelism, locality, and recomputation in image processing pipelines. In Boehm, H. and Flanagan, C., editors, *ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '13*, Seattle, WA, USA, June 16–19, 2013, pages 519–530. ACM.
- Rondon, P. M., Kawaguchi, M., and Jhala, R. (2008). Liquid types. In Gupta, R. and Amarasinghe, S. P., editors, *Proceedings of the ACM SIGPLAN 2008 Conference on Programming Language Design and Implementation*, Tucson, AZ, USA, June 7–13, 2008, pages 159–169. ACM.
- Steuer, M., Fensch, C., Lindley, S., and Dubach, C. (2015). Generating performance portable code using rewrite rules: from high-level functional expressions to high-performance OpenCL code. In Fisher, K. and Reppey, J. H., editors, *Proceedings of the 20th ACM SIGPLAN International Conference on Functional Programming, ICFP 2015*, Vancouver, BC, Canada, September 1–3, 2015, pages 205–217. ACM.

- Tang, F. and Grelck, C. (2012). User-defined shape constraints in Sac. In Hinze, R., editor, *24th International Symposium on Implementation and Application of Functional Languages (IFL'12)*, Oxford, UK. University of Oxford.
- Trojahnner, K. and Grelck, C. (2009). Dependently typed array programs don't go wrong. *J. Log. Algebraic Methods Program.*, 78(7):643–664.
- Xi, H. and Pfenning, F. (1998). Eliminating array bound checking through dependent types. In Davidson, J. W., Cooper, K. D., and Berman, A. M., editors, *Proceedings of the ACM SIGPLAN '98 Conference on Programming Language Design and Implementation (PLDI)*, Montreal, Canada, June 17-19, 1998, pages 249–257. ACM.
- Xi, H. and Pfenning, F. (1999). Dependent types in practical programming. In Appel, A. W. and Aiken, A., editors, *POPL '99, Proceedings of the 26th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, San Antonio, TX, USA, January 20-22, 1999, pages 214–227. ACM.

A Type independence for the semantics

We recall here the formalization given in [subsection 3.2](#).

Definition 3.1 (Observational equivalence). Two closed expressions e_1 and e_2 , are *observationally equivalent* $-e_1 \equiv e_2-$ if and only if, for any closed expressions a_1, \dots, a_k and integers n_1, n_2 :

$$\begin{cases} e_1 a_1 \dots a_k \rightsquigarrow n_1 \\ e_2 a_1 \dots a_k \rightsquigarrow n_2 \end{cases} \implies n_1 = n_2$$

Definition 3.2 (Equality modulo types). Two expressions e_1 and e_2 are *equal modulo types* $-e_1 \approx_\tau e_2-$ if and only if it exists an expression e , free type variables $\bar{\alpha}$ of e and types $\bar{\tau}_1, \bar{\tau}_2$ such that:

$$e_1 = e\{\bar{\tau}_1/\bar{\alpha}\} \wedge e_2 = e\{\bar{\tau}_2/\bar{\alpha}\}$$

Lemma A.1. *Preservation by substitution* Given expressions e_1, e_2 , term variable x and expressions e'_1, e'_2 , then

$$\begin{cases} e_1 \approx_\tau e_2 \\ e'_1 \approx_\tau e'_2 \end{cases} \implies e_1\{e'_1/x\} \approx_\tau e_2\{e'_2/x\}$$

Theorem 3.3 (Type independence). *Given two closed terms e_1, e_2 such that $e_1 \approx_\tau e_2$, then*

$$\forall v_1 v_2, \begin{cases} e_1 \rightsquigarrow v_1 \\ e_2 \rightsquigarrow v_2 \end{cases} \implies v_1 \approx_\tau v_2$$

Proof. Let e_1, e_2 and v_1, v_2 such that $e_1 \approx_\tau e_2$, $e_1 \rightsquigarrow v_1$ and $e_2 \rightsquigarrow v_2$. Let examine the rules of expression semantics to check that the above invariant holds.

E-SIZE Since types only may differ, $e_1 \approx_\tau e_2 \implies e_1 = \langle n \rangle = e_2 \implies v_1 = n = v_2 \implies v_1 \approx_\tau v_2$

E-APP Applying induction hypothesis on the two first premises, and the substitution lemma for the last one gives the equality modulo type of resulting values.

E-COERCE *idem*

E-LET We have: $e_1 = \text{let } x_{\mathbf{V}_1} : \tau_1 = e_1 \text{ in } e'_1$ and $e_2 = \text{let } x_{\mathbf{V}_2} : \tau_2 = e_2 \text{ in } e'_2$. It follows that $e'_1 \approx_\tau e'_2$ and $e_1 \triangleright \tau_1 \approx_\tau e_2 \triangleright \tau_2$. By the substitution lemma, the premise yields equal modulo type values.

It remains to prove that coercions of equal modulo types values and possibly different types yields equal modulo types values. This is established by examining value's shape:

n For both coercions, one of C-SIZE, C-INDEX or C-INT applies. Each of them produces the same result: the original value. Thus results are equal modulo types.

$\lambda x : \tau. e$ Coercions are reduced with C-FUN rule, thus both types have the form $\cdot \rightarrow \cdot$. The produced abstractions are immediately equal modulo types. \square

Remark. This result could be extended to expressions whose type polymorphism differ, i.e. in which generalization lists \mathbf{V} do not contains the same type variables.

B Properties of type system

In presence of recursion, the correction proof of a type system toward a big step semantics cannot be derived from usual progress and type preservation properties of the reduction, as done for small step semantics, because blocked and diverging terms are undistinguishable. However, a general analysis of big step semantics for soundness conditions [Dagnino et al. 2020] provides a few similar properties to check in order to verify soundness. Because of non-deterministic semantics, two kinds of corrections are distinguished:

- *Soundness-must*: None of possible reduction is blocked
- *Soundness-may*: At least one of possible reduction is not blocked

A mechanized derivation of these global properties (soundness-must in our setting) from local ones was proposed by Dagnino et al. [2020]. They follow from usual properties on the type system that we detail first.

B.1 Normalization and preliminary lemmas

Definition B.1 (Normalized typing derivation). A typing derivation $\Gamma \vdash e : \tau$ is *normalized* if the instances of rule T-SUBTYPE appear in one of the following position:

Bottom most rule	First premise of applications
$\text{T-SUBTYPE} \frac{T \quad S}{\Gamma \vdash e : \tau}$	$\text{T-SUBTYPE} \frac{T_1 \quad S}{\Gamma \vdash e_1 : \tau'_2 \rightarrow \tau'} \quad \text{T-APP} \frac{T_2}{\Gamma \vdash e_2 : \tau''}$
Coercions	First premise of declarations
$\text{T-SUBTYPE} \frac{T \quad S}{\Gamma \vdash e : \tau'} \quad \text{T-COERCE} \frac{\Gamma \vdash e : \tau'}{\Gamma \vdash e \triangleright \tau' : \tau}$	$\text{T-SUBTYPE} \frac{T \quad S}{\Gamma, V \vdash e : \tau} \quad \text{T-LET} \frac{T_2}{\Gamma, x : \forall V. \tau \vdash e' : \tau'}$

Definition B.2. Given τ, τ_1, τ_2 types and S_1, S_2 derivations of the sub-typing relations $\tau_1 <: \tau$ and $\tau <: \tau_2$, $S_1 \bowtie_\tau S_2$ is the following sub-typing derivation of $\tau_1 <: \tau_2$, regarding τ shape:

- $\tau = \text{int}$ — Then S_2 is REFL $\frac{}{\text{int} <: \text{int}}$ and $\tau = \tau_2$. The derivation $S_1 \bowtie_\tau S_2$ is defined as S_1 .
- $\tau = \langle \eta \rangle$ — Then S_1 is REFL $\frac{}{\langle \eta \rangle <: \langle \eta \rangle}$ and $\tau = \tau_1$. The derivation $S_1 \bowtie_\tau S_2$ is defined as S_2 .
- $\tau = [\eta]$ — Then S_1 is REFL $\frac{}{[\eta] <: [\eta]}$ and $\tau = \tau_1$. The derivation $S_1 \bowtie_\tau S_2$ is defined as S_2 .
- $\tau = \tau^d \rightarrow \tau^c$ — Then $\tau_1 = \tau_1^d \rightarrow \tau_1^c$, S_1 is FUN $\frac{S_1^d \quad S_1^c}{\tau_1^d <: \tau_1^c}$, $\tau_2 = \tau_2^d \rightarrow \tau_2^c$, S_2 is FUN $\frac{S_2^d \quad S_2^c}{\tau_2^d <: \tau_2^c}$.

$$\text{The derivation } S_1 \bowtie_\tau S_2 \text{ is defined as FUN } \frac{\frac{S_1^d \bowtie_{\tau^d} S_2^d}{\tau_1^d <: \tau_1^c} \quad \frac{S_1^c \bowtie_{\tau^c} S_2^c}{\tau_1^c <: \tau_1^c}}{\tau_1^d \rightarrow \tau_1^c <: \tau_2^d \rightarrow \tau_2^c}.$$

Corollary B.3. The sub-typing relation $<:$ is transitive.

Theorem B.4 (Normalization of typing derivation). *If there exists typing derivation, there exists a normalized typing derivation.*

Proof. Typing derivations are normalized using the following rewrite rules:

$\text{SUBTYPE} \frac{\text{SUBTYPE} \frac{P}{\Gamma \vdash e : \tau_2} \quad \text{SUBTYPE} \frac{S_2}{\tau_2 <: \tau_1}}{\Gamma \vdash e : \tau_1} \quad \text{SUBTYPE} \frac{\Gamma \vdash e : \tau_1}{\Gamma \vdash e : \tau} \quad \text{SUBTYPE} \frac{\Gamma \vdash e : \tau}{\Gamma \vdash e : \tau}$ \downarrow $\text{SUBTYPE} \frac{P}{\Gamma \vdash e : \tau_2} \quad \text{SUBTYPE} \frac{S_1 \bowtie_{\tau_1} S_2}{\tau_2 <: \tau} \quad \text{SUBTYPE} \frac{\Gamma \vdash e : \tau_2}{\Gamma \vdash e : \tau}$	$\text{T-SUBTYPE} \frac{\text{T-APP} \frac{T_1}{\Gamma \vdash e_1 : \tau'_2 \rightarrow \tau'} \quad \text{T-SUBTYPE} \frac{T_2 \quad S_2}{\Gamma \vdash e_2 : \tau_2 \quad \tau_2 <: \tau'_2}}{\Gamma \vdash e_1 e_2 : \tau} \quad \text{T-SUBTYPE} \frac{\Gamma \vdash e_1 e_2 : \tau}{\Gamma \vdash e_1 e_2 : \tau}$ \downarrow $\text{T-SUBTYPE} \frac{T_1}{\Gamma \vdash e_1 : \tau'_2 \rightarrow \tau'} \quad \text{S-FUN} \frac{S_2 \quad S_1}{\tau'_2 <: \tau'_2 \quad \tau'_2 <: \tau} \quad \text{T-APP} \frac{\Gamma \vdash e_1 : \tau_2 \rightarrow \tau \quad \Gamma \vdash e_2 : \tau_2}{\Gamma \vdash e_1 e_2 : \tau}$
$\text{T-SUBTYPE} \frac{T \quad S}{\Gamma, x : \tau \vdash e : \tau'_1 \quad \tau'_1 <: \tau_1} \quad \text{T-ABS} \frac{\Gamma, x : \tau \vdash e : \tau_1}{\Gamma \vdash \lambda x : \tau. e : \tau \rightarrow \tau_1}$ \downarrow $\text{T-ABS} \frac{T}{\Gamma \vdash \lambda x : \tau. e : \tau \rightarrow \tau'_1} \quad \text{S-REFL} \frac{S}{\tau <: \tau} \quad \text{S-FUN} \frac{\tau \rightarrow \tau'_1 <: \tau \rightarrow \tau_1}{\Gamma \vdash \lambda x : \tau. e : \tau \rightarrow \tau_1}$	$\text{T-LET} \frac{T_1 \quad \text{T-SUBTYPE} \frac{T_2 \quad S}{\Gamma, x : \forall V. \tau \vdash e' : \tau'_1 \quad \tau'_1 <: \tau_1}}{\Gamma \vdash \text{let } x_{\forall V} : \tau = e \text{ in } e' : \tau_1}$ \downarrow $\text{T-LET} \frac{T_1 \quad T_2 \quad S}{\Gamma \vdash \text{let } x_{\forall V} : \tau = e \text{ in } e' : \tau_1}$

This rewrite system terminates since instances of the rule T-SUBTYPE are either pushed downward or introduced in a normalized position. Moreover, any non-normalized positions for T-SUBTYPE rule reduces with one of the above rewrite rules. \square

Lemma B.5 (Inversion). *Let Γ, e, τ such that $\Gamma \vdash e : \tau$*

1. If $e = {}^S x$, then $\Gamma(x) = \forall V. \tau'$ and $\tau' \{S/V\} <: \tau$
2. If $e = e_1 e_2$, then there exists τ' such that $\Gamma, x : \forall \mathcal{E}. \tau \vdash e_2 : \tau' \rightarrow \tau$ and $\Gamma, x : \forall \mathcal{E}. \tau \vdash e_1 : \tau'$
3. If $e = \lambda x : \tau. e'$, then $\tau = \tau_1 \rightarrow \tau_2$ et $\Gamma, x : \forall \mathcal{E}. \tau_1 \vdash e' : \tau_2$
4. If $e = n$, then $\tau = \text{int}$
5. If $e = \langle \eta \rangle$, then $\langle \eta \rangle <: \tau$
6. If $e = e \triangleright \tau'$, then $\tau' <: \tau$ and $\Gamma \vdash e : \tau'$
7. If $e = \text{let } x_V : \tau' = e_1 \text{ in } e_2$, then $\Gamma, V \vdash e_1 : \tau'$ and $\Gamma, x : \forall V. \tau' \vdash e_2 : \tau$

Proof. Case-based reasoning on expression shape and analysis of applicable rules. \square

Lemma B.6 (Substitution). *The type of an expression is preserved by well-typed substitution:*

$$\left. \begin{array}{l} \Gamma, x : \forall V. \tau' \vdash e : \tau \\ \Gamma, V \vdash e' : \tau'' \\ \tau' <: \tau'' \end{array} \right\} \Rightarrow \Gamma \vdash e\{e' \{S/V\} / {}^S x\} : \tau$$

Proof. A finite derivation tree of $\Gamma \vdash e\{e'/x\} : \tau$ is obtained by substituting in a finite derivation of $\Gamma, x : \forall V. \tau' \vdash e : \tau$ every occurrence of rule VAR (in finite number) by a derivation of $\Gamma \vdash e' : \tau''$ (also finite) and an instance of SUBTYPE rule. Cares is required with environment that are distinct for each occurrence (they might be extended). This lemma is valid because we considered only name resolved terms such that no clashes occur. Thus, typing is preserved by extension of the environment since added variables are fresh and may not mask existing ones. \square

Lemma B.7 (Canonical form). *Given a type τ , $\{\tau\}$ denotes $\{v \in \mathcal{V} \mid v : \tau\}$. Then,*

$$\begin{aligned} \{\text{int}\} &= \{n \mid n \in \mathbb{Z}\} \\ \{\cdot \rightarrow \cdot\} &= \{\lambda \cdot : \cdot\} \end{aligned}$$

Proof. Case-based analysis on value's shape. \square

B.2 Soundness sufficient conditions

To prove soundness and preservation of a type system toward a big-step semantics, Dagnino et al. [2020] proposed a general reduction to three local properties. These sufficient conditions only require rule examination, while the induction is conducted by the generic construction. To help presenting them, we borrow the proposed syntax of inline format for instances of rules:

$$(e_1 \rightsquigarrow v_1, \dots, e_n \rightsquigarrow v_n, e_{n+1} \rightsquigarrow v_{n+1}, e) \stackrel{\text{def}}{=} \frac{e_1 \rightsquigarrow v_1 \quad \dots \quad e_n \rightsquigarrow v_n \quad e_{n+1} \rightsquigarrow v_{n+1}}{e \rightsquigarrow v_{n+1}}$$

The $e_1 \rightsquigarrow v_1, \dots, e_n \rightsquigarrow v_n$ are rule's *premises* and $e_{n+1} \rightsquigarrow v_{n+1}$ is the *continuation*, that produces the result of rule instance. For rules with no natural continuation, a trivial one is inserted: $v_{n+1} \rightsquigarrow v_{n+1}$.

Lemma B.8 ((S1) Local Preservation). *For any instance $(e_1 \rightsquigarrow v_1, \dots, e_n \rightsquigarrow v_n, e_{n+1} \rightsquigarrow v_{n+1}, e)$, such that $\vdash e : \tau$, there exists $\tau_1, \dots, \tau_{n+1}$ with $\tau_{n+1} = \tau$ such that :*

$$\forall k \in \llbracket 1, n+1 \rrbracket, (\forall h \in \llbracket 1, k-1 \rrbracket, \vdash v_h : \tau_h) \Rightarrow \vdash e_k : \tau_k$$

Proof. Case-base reasoning on instances de semantics rules, using extensively the inversion lemma:

E-SIZE $e = \langle n \rangle$ – There is only the continuation $n \rightsquigarrow n$ that verifies $\vdash n : \text{int}$

E-APP $e = e_1 e_2$ – Lets assumes there exists τ such that $\vdash e : \tau$.

By inversion lemma (2), there exists τ' such that $\vdash e_1 : \tau \rightarrow \tau'$ and $\vdash e_2 : \tau$ Lets find the right type for each premises:

– $e_1 \rightsquigarrow \lambda x : \tau. e$ – Immediately, $\vdash e_1 : \tau \rightarrow \tau'$

– $e_2 \triangleright \tau \rightsquigarrow v$ – Immediately, $\vdash e_2 \triangleright \tau : \tau$

– $e\{v/x\} \rightsquigarrow v'$ – By hypothesis, $\vdash \lambda x : \tau. e : \tau \rightarrow \tau'$. The inversion lemma (3) gives $x : \forall \mathcal{E}. \tau \vdash e : \tau'$ allowing to conclude with substitution lemma: $\vdash e\{v/x\} : \tau'$

E-COERCE $e = e' \triangleright \tau'$ – By inversion lemma (7), there exists τ' such that $\vdash e : \tau''$ and $\tau' <: \tau$.

– $e \rightsquigarrow v$ – Immediately, $\vdash e : \tau''$

Expression Inference		$\Gamma \vdash e : \tau \vdash \mathcal{U}$
I-VAR	$\frac{\Gamma(x) = \forall V. \tau \quad S = \text{Fresh}(V) \quad \Gamma \vdash S \vdash \mathcal{U}}{\Gamma \vdash \lambda x : \tau \{S/V\} \vdash \mathcal{U} \cdot \{l \mapsto S\}}$	
I-SIZE	$\frac{\Gamma \vdash \eta \vdash \mathcal{U}}{\Gamma \vdash \langle \eta \rangle : \langle \eta \rangle \vdash \mathcal{U}}$	
I-INT	$\frac{}{\Gamma \vdash n : \text{int} \vdash \{\}}$	
I-COERCE	$\frac{\Gamma \vdash (e : \tau) \vdash \mathcal{U}}{\Gamma \vdash e \triangleright \tau : \tau \vdash \mathcal{U} \cdot \mathcal{U}_\tau}$	
I-ABS	$\frac{\Gamma \vdash \tau \vdash \mathcal{U}_\tau \quad \Gamma, x : \forall \varepsilon. \tau \vdash e : \tau' \vdash \mathcal{U}}{\Gamma \vdash \lambda x : \tau. e : \tau \rightarrow \tau' \vdash \mathcal{U}_\tau \cdot \mathcal{U}}$	
I-APP	$\frac{\Gamma \vdash e' : \tau' \vdash \mathcal{U}' \quad \Gamma \vdash (e : \tau' \rightarrow \alpha) \vdash \mathcal{U}}{\Gamma \vdash e e' : \alpha \vdash \mathcal{U} \cdot \mathcal{U}'}$	
I-LET	$\frac{\Gamma \vdash d \Rightarrow x : \sigma \vdash \mathcal{U} \quad \Gamma, x : \sigma \vdash e' : \tau' \vdash \mathcal{U}'}{\Gamma \vdash \text{let } d \text{ in } e : \tau' \vdash \mathcal{U} \cdot \mathcal{U}'}$	
Declaration Inference		$\Gamma \vdash d \Rightarrow x : \sigma \vdash \mathcal{U}$
I-DECL	$\frac{\Gamma \vdash (e : \tau) \vdash (V, C, \pi, \rho) \quad \rho' = \text{Solve}(V, C) \quad \check{V}' = \text{Gen}(V, C, \rho')}{\Gamma \vdash x_l : \tau = e \Rightarrow x : \forall V'. \tau\{\rho'\} \vdash (\emptyset, C\{\rho'\}, \pi \cdot \{l \mapsto V'\}, \rho \circ \rho')}$	
Constraint Insertion		$\Gamma \vdash (e : \tau) \vdash \mathcal{U}$
I-CSTR	$\frac{\Gamma \vdash e : \tau' \vdash \mathcal{U} \quad \Gamma \vdash \tau \vdash \mathcal{U}_\tau}{\Gamma \vdash (e : \tau) \vdash \mathcal{U} \cdot \mathcal{U}_\tau \cdot \{\tau' <: \tau\}}$	

Figure 6. The inference algorithm. The function $\vdash \Gamma \vdash S \vdash \mathcal{U}$ registers free size and type variables (the ones that are unbound in Γ) of instantiation list S . The function $\vdash \Gamma \vdash (e : \tau) \vdash \mathcal{U}$ where expression and type are bracketed combines expression inference, type's free variables registering and sub-typing constraint insertion.

- $v \triangleright \tau' \leadsto v'$ – Immediately, $\vdash v \triangleright \tau' : \tau'$
- E-LET $e = \text{let } x_V : \tau' = e_1 \text{ in } e_2$ – Combined used of inversion lemma and substitution one.

□

Lemma B.9 ((S2) \exists -progress). *For any $e \notin \mathcal{V}$, if there exists τ such that $\vdash e : \tau$, then there exists a rule instance of the form $(j_1, \dots, j_n, j_{n+1}, e)$*

Proof. Trivial case-based reasoning on expression shape. □

Lemma B.10 ((S3) \forall -progress). *For any rule instance $(e_1 \leadsto v_1, \dots, e_n \leadsto v_n, e_{n+1} \leadsto v_{n+1}, e)$, assuming there exists τ such that $\vdash e : \tau$, then, for any $k \in \llbracket 1, n+1 \rrbracket$,*

$$\begin{aligned} & \text{assuming for any } h < k, e_h \leadsto v_h \text{ and } e_k \leadsto v, \\ & \text{then there exists a rule instance } (j'_1, \dots, j'_{n'}, j'_{n'+1}, e') \text{ such that} \\ & \forall h < k, j'_h = j_h, e' = e, \text{ et } j_k = e'' \leadsto v \end{aligned}$$

Intuitively, it amounts to check that the evaluation of sub-expressions gives results that fulfill their use (the expected form of value).

Proof. Case-based analysis on semantics rule instances.

- E-APP – By typing, (rule T-APP), $\vdash e_1 : \tau_1 \rightarrow \tau_2$. Thus the canonical form lemma gives $e_1 \leadsto \lambda x : \tau. e$. The rule E-APP can be instantiated. The continuation may be freely instantiated.
- Other rules do not constrain the result of evaluation of their sub-expressions, thus fulfilling the property.

□

Theorem B.11. *As proved in [Dagnino et al. 2020] the three properties allow to deduce Theorem 3.5 and Theorem 3.6:*

$$\begin{aligned} (S1) & \implies (\text{Type preservation}) \\ (S1) + (S2) + (S3) & \implies (\text{Type soundness}) \end{aligned}$$

C Inference properties

C.1 Algorithm

We formalize here the inference algorithm sketched in section 4.

Constraint collecting. Our algorithm builds a *unifier* $\mathcal{U} = (V, C, \pi, \rho)$ by traversing expressions bottom-up (from leaves that are variable occurrences, size values and constants to the top-level declaration). It collects (i) $V \subset \mathcal{V}_\eta \cup \mathcal{V}_\tau$ a set of free size and type variables ; (ii) C a set of sub-typing constraints ; (iii) π the definition of encountered polymorphism labels and (iv) ρ the substitution of some size and type variables, supposed acyclic. Substitutions' domains, free variables and variables that appear in a generalization list are supposed disjoint.

The empty unifier is $\{\}$. The union of unifiers is denoted $\mathcal{U}_1 \cdot \mathcal{U}_2$. It requires the substitutions' domains and free variable sets to be disjoint. This property holds during inference at the condition that size and type variables appear only in one place in traversed term. Singleton unifiers are unambiguously denoted $\{\alpha\}$ (a single free type variable), $\{\tau <: \tau'\}$ (a single sub-typing constraint), $\{l \mapsto V\}$ or $\{l \mapsto S\}$ (a single polymorphism label definition). Last, $\mathcal{U} \setminus V$ denotes the unifier \mathcal{U} where the size and type variables V have been removed from its free variables.

Inference is made of a few mutually recursive functions defined in Figure 6, that use environments as introduced in ?? . To present them, their outputs are underlined. The main *expression inference* function $\Gamma \vdash e : \underline{\tau} \vdash \underline{\mathcal{U}}$ collects the constraints and builds expression's type. It is accompanied with a *registering* function $\Gamma \vdash \mathcal{S} \vdash \underline{\mathcal{U}}$ that returns a unifier containing the size and type variables of \mathcal{S} that are unbounded in the environment¹¹. For convenience, the *constraint insertion* function $\Gamma \vdash (e : \tau) \vdash \underline{\mathcal{U}}$, with bracketed expression and type combines expression inference and sub-typing constraint insertion, thus only producing a unifier. Last, the handling of declarations is set apart: the *declaration inference* function $\Gamma \vdash d \Rightarrow \underline{x} : \underline{\sigma} \vdash \underline{\mathcal{U}}$ builds the type scheme of introduced variables as well as declaration unifier.

Variables introduction. Variable are immediately instantiated with fresh sizes and types (rule VAR and auxiliary function *Fresh*). The definition of this instantiation label $-l-$ is registered in the unifier as well as the generated sizes and types. When handling abstractions (rule ABS), the free variables of the type τ are registered and the environment is extended with the monomorphic introduced variable. For applications (rule APP), a fresh type variable α is picked. It allows constructing expression type without solving any type constraints, namely that the type of e in rule APP should be a function, which is enforced by adding a constraint.

Polymorphism. Let bindings introduce generalization (rule I-DECL): once expression has been traversed, function *Solve* turns as many unifier's sub-typing constraints as possible into a substitution of its free variables (see subsection 4.4). Function *Gen* extracts the free size and type variables, i.e. the ones that are not substituted, and checks that they do not appear in any remaining constraints so that they might be generalized. The resulting unifier is built by composing substitutions and registering generalization label's definition.

Reconstruction. The reconstruction of top level terms $\Gamma \vdash e \mapsto e'$ is defined with the unique rule

$$\text{I-TOP} \frac{\Gamma \vdash e : \tau \vdash (\emptyset, \emptyset, \pi, \rho)}{\Gamma \vdash e \mapsto e\{\pi\}\{\rho\}}$$

It requires all constraints to be solved, and no free variables to remain.¹² The resulting definition of polymorphism markers is applied $-\{\pi\}-$, then the substitution $-\{\rho\}-$. Note that the variables used in π for instantiation might get substituted, hence the order.

Definition C.1 (Reconstruction). Given an environment Γ , expressions e and e' , e' is a *reconstruction* of e , denoted $\Gamma \vdash e \leftarrow e'$ if and only if:

$$\exists \pi \rho \tau, \begin{cases} e' = e\{\pi\}\{\rho\} \\ \Gamma \vdash e' : \tau \end{cases}$$

Theorem C.2 (Inference soundness). *Inference produces well-typed terms, i.e. given expressions e and e' ,*

$$\Gamma \vdash e \mapsto e' \implies \Gamma \vdash e \leftarrow e'$$

Definition C.3. Relation & rule instance substitution Given a n-ary relation $\mathcal{R}(x_1, \dots, x_n)$ and a substitution ρ , the substituted relation is defined as:

$$\mathcal{R}(x_1, \dots, x_n)\{\rho\} \stackrel{\text{def}}{=} \mathcal{R}(x_1\{\rho\}, \dots, x_n\{\rho\})$$

Similarly, given we define substituted rule instances as:

$$\text{RULE} \frac{p_1 \quad \dots \quad p_k}{c}\{\rho\} \stackrel{\text{def}}{=} \text{RULE} \frac{p_1\{\rho\} \quad \dots \quad p_k\{\rho\}}{c\{\rho\}}$$

Proof. By construction, inference produces a type τ , polymorphism definitions π and a substitution ρ such that $e' = e\{\pi\}\{\rho\}$. It remains to show that $\Gamma \vdash e : \tau$, by proving the following invariant: given an environment Γ an expression e , a type τ , size and type variables V , constraints C , polymorphism definitions π and a substitution ρ such that $\Gamma \vdash e : \tau \vdash (V, C, \pi, \rho)$, then

$$\forall \rho', \Gamma \vdash C\{\rho'\} \implies (\Gamma \vdash e\{\pi\}\{\rho\} : \tau)\{\rho'\}$$

Intuitively, this amounts to show that any substitution that solves the remaining constraints leads to a well-type term, i.e. constraint collection captures all the necessary type relations.

$$\left. \Gamma \vdash e : \tau \vdash (V, C, \rho, \pi) \right\} \implies \text{T} \frac{\dots}{\Gamma \vdash e\{\pi\}\{\rho\} : \tau} \{\rho'\}$$

¹¹ At this point, binding size or type variable is impossible in \mathcal{L}^n , but extensions (polymorphic recursion) will allow it.

¹² An extra declaration might be added to introduce a constraint solving point, i.e. **let** $x : _ = e$ **in** x

Because inference is syntax directed, we proceed by induction on expressions: given the resulting unifier and a substitution that solves the remaining constraints, we build a correct typing derivation for e :

n – Inference has the following form:

$$\text{I-INT} \frac{}{\Gamma \vdash n : \text{int} \vdash (\emptyset, \emptyset, \varepsilon, \varepsilon)}$$

Given ρ' a substitution that solves \emptyset , the following typing derivation is correct, since $n\{\pi\}\{\rho\} = n$:

$$\text{T-INT} \frac{}{\Gamma \vdash n : \text{int}} \{\rho'\}$$

$^l x$ – Inference has the following form:

$$\text{I-VAR} \frac{\Gamma(x) = \forall V. \tau \quad S = \text{Fresh}(V)}{\Gamma \vdash ^l x : \tau\{S/V\} \vdash (S, \emptyset, \{l \mapsto S\}, \varepsilon)}$$

Given ρ' a substitution that solves \emptyset , the following typing derivation is correct, since $^l x\{\pi\}\{\rho\} = ^s x$:

$$\text{T-VAR} \frac{\Gamma(x) = \forall V. \tau}{\Gamma \vdash ^s x : \tau\{S/V\}} \{\rho'\}$$

$\langle \eta \rangle$ – *idem*

$e \triangleright \tau'$ – Inference has the following form:

$$\begin{array}{c} \text{I} \frac{\dots}{\Gamma \vdash e_1 : \tau' \vdash (V, C, \pi, \rho)} \quad \frac{}{\Gamma \vdash \tau' \vdash (V_\tau, \emptyset, \varepsilon, \varepsilon)} \\ \text{I-CSTR} \frac{}{\Gamma \vdash (e : \tau') \vdash (V \cdot \mathcal{V}_\tau, C \cdot \{\tau' <: \tau\}, \pi, \rho)} \\ \text{I-COERCE} \frac{}{\Gamma \vdash e \triangleright \tau' : \tau \vdash (V \cdot \mathcal{V}_\tau, C \cdot \{\tau' <: \tau\}, \pi, \rho)} \end{array}$$

Given ρ' a substitution that solves $C \cdot \{\tau' <: \tau\}$ then $\vdash C\{\rho'\}$, the induction hypothesis allows to construct T a typing derivation of $e\{\pi\}\{\rho\}$. Moreover, because $\tau'\{\rho'\} <: \tau\{\rho'\}$, it exists a sub-typing derivation S such that the following derivation is valid:

$$\begin{array}{c} \text{T} \frac{\dots}{\Gamma \vdash e\{\pi\}\{\rho\} : \tau'} \{\rho'\} \\ \text{T-COERCE} \frac{}{\Gamma \vdash e\{\pi\}\{\rho\} \triangleright \tau' : \tau'} \{\rho'\} \quad \text{S} \frac{\dots}{\tau' <: \tau} \{\rho'\} \\ \text{T-SUBTYPE} \frac{}{\Gamma \vdash e\{\pi\}\{\rho\} \triangleright \tau' : \tau} \{\rho'\} \end{array}$$

$e_1 e_2$ – Inference has the following form:

$$\begin{array}{c} \text{I-2} \frac{\dots}{\Gamma \vdash e_2 : \tau_2 \vdash (V_2, C_2, \pi_2, \rho_2)} \quad \text{I-1} \frac{\dots}{\Gamma \vdash e_1 : \tau_1 \vdash (V_1, C_1, \pi_1, \rho_1)} \quad \frac{}{\Gamma \vdash \tau_2 \rightarrow \alpha \vdash (\{\alpha\}, \emptyset, \varepsilon, \varepsilon)} \\ \text{I-CSTR} \frac{}{\Gamma \vdash (e_1 : \tau_2 \rightarrow \alpha) \vdash (V_1 \cdot \{\alpha\}, C_1 \cdot \{\tau_1 <: \tau_2 \rightarrow \alpha\}, \pi_1, \rho_1)} \\ \text{I-APP} \frac{}{\Gamma \vdash e_1 e_2 : \alpha \vdash (V_1 \cdot V_2 \cdot \{\alpha\}, C_1 \cdot C_2 \cdot \{\tau_1 <: \tau_2 \rightarrow \alpha\}, \pi_1 \cdot \pi_2, \rho_1 \cdot \rho_2)} \end{array}$$

Given ρ' a substitution that solves C_1, C_2 and $\{\tau_1 <: \tau_2 \rightarrow \alpha\}$, using induction hypothesis on I-1 and I-2 defining T-1 and T-2 as above and a sub-typing derivation for $\{\tau_1 <: \tau_2 \rightarrow \alpha\}$, the following typing derivation of $e\{\pi\}\{\rho'\} = e_1\{\pi_1\}\{\rho_1\} e_2\{\pi_2\}\{\rho_2\}$ is correct:

$$\begin{array}{c} \text{T-1} \frac{\dots}{\Gamma \vdash e_1\{\pi_1\}\{\rho_1\} : \tau_1} \{\rho'\} \quad \text{S} \frac{\dots}{\tau_1 <: \tau_2 \rightarrow \alpha} \{\rho'\} \\ \text{T-SUBTYPE} \frac{}{\Gamma \vdash e_1\{\pi_1\}\{\rho_1\} : \tau_2 \rightarrow \alpha} \{\rho'\} \quad \text{T-2} \frac{\dots}{\Gamma \vdash e_2\{\pi_2\}\{\rho_2\} : \tau_2} \{\rho'\} \\ \text{T-APP} \frac{}{\Gamma \vdash e_1\{\pi_1\}\{\rho_1\} e_2\{\pi_2\}\{\rho_2\} : \alpha} \{\rho'\} \end{array}$$

$\lambda x:\tau. e'$ – Inference has the following form:

$$\text{I-Abs} \frac{\frac{\Gamma \vdash \tau \vdash (V_\tau, \emptyset, \varepsilon, \varepsilon)}{\Gamma \vdash \lambda x:\tau. e' : \tau \rightarrow \tau' \vdash (V_\tau \cdot V, C, \pi, \rho)} \quad \text{I} \frac{\dots}{\Gamma, x:\forall \varepsilon. \tau \vdash e' : \tau' \vdash (V, C, \pi, \rho)}}{\Gamma \vdash \lambda x:\tau. e' : \tau \rightarrow \tau' \vdash (V_\tau \cdot V, C, \pi, \rho)}$$

Given ρ' a substitution that solves C , using induction hypothesis on I (defining T), the following typing derivation of $e\{\pi\}\{\rho\}$ is correct since τ variables are registered in collected free variables, hence get defined by ρ' :

$$\text{T-Abs} \frac{\text{T} \frac{\dots}{\Gamma, x:\forall \varepsilon. \tau \vdash e'\{\pi_1\}\{\rho_1\} : \tau'} \{\rho'\}}{\Gamma \vdash \lambda x:\tau. e\{\pi_1\}\{\rho_1\}' : \tau \rightarrow \tau' \{\rho'\}}$$

let $x_V:\tau' = e_1$ **in** e_2 – Inference has the following form:

$$\text{I-LET} \frac{\text{I-1} \frac{\frac{\Gamma \vdash e_1 : \tau \vdash (V_1, C_1, \pi_1, \rho_1)}{\Gamma \vdash (e_1 : \tau_1) \vdash (V_1 \cdot V_\tau, C_1 \cdot \{\tau <: \tau_1\}, \pi_1, \rho_1)} \quad \frac{\Gamma \vdash \tau_1 \vdash (V_\tau, \emptyset, \varepsilon, \varepsilon)}{V, \rho = \text{Solve}(V_1 \cdot V_\tau, C_1 \cdot \{\tau <: \tau_1\})} \quad \text{I-DECL} \frac{\Gamma \vdash x_l:\tau_1 = e_1 \Rightarrow x : \forall V. \tau_1\{\rho\} \vdash (\emptyset, C_1\{\rho\} \cdot \{\tau\{\rho\} <: \tau_1\{\rho\}\}, \pi_1 \cdot \{l \rightarrow V\}, \rho_1 \circ \rho)}{\Gamma \vdash \text{let } x_l:\tau_1 = e_1 \text{ in } e_2 : \tau_2 \vdash (V_2, C_2, \pi_2, \rho_2)} \quad \text{I-2} \frac{\dots}{\Gamma, x:\forall V. \tau_1\{\rho\} \vdash e_2 : \tau_2 \vdash (V_2, C_2, \pi_2, \rho_2)}}{\Gamma \vdash \text{let } x_l:\tau_1 = e_1 \text{ in } e_2 : \tau_2 \vdash (V_2, C_1\{\rho\} \cdot \{\tau\{\rho\} <: \tau_1\{\rho\}\} \cdot C_2, \pi_1 \cdot \{l \rightarrow V\} \cdot \pi_2, \rho_1 \circ \rho_2)}$$

Here, Solve combines the Solve and Gen function of the algorithms. It defines a substitution ρ of variables $V_1 \cdot V_\tau \setminus V$. The constraint set C_1 is substituted and transmitted in the resulting unifier. Thus, given ρ' a substitution that solves C , $\rho \circ \rho'$ solves the constraints $C_1 \cdot \{\tau <: \tau_1\}$. The reconstructed term has the following shape: **let** $x_V:\tau_1\{\rho_1 \circ \rho\} = e_1\{\pi_1\}\{\rho_1 \circ \rho\}$ **in** $e_2\{\pi_2\}\{\rho_2\}$. Its type can be derived with:

$$\text{T-LET} \frac{\text{T-1} \frac{\frac{\Gamma, V \vdash e_1\{\pi_1\}\{\rho_1\} : \tau\{\rho_1\}}{\Gamma, V \vdash e_1\{\pi_1\}\{\rho_1\} : \tau_1\{\rho_1\}} \quad \text{S} \frac{\dots}{\tau\{\rho_1\} <: \tau_1\{\rho_1\}} \{\rho \circ \rho'\}}{\Gamma, V \vdash e_1\{\pi_1\}\{\rho_1\} : \tau_1\{\rho_1\}} \quad \text{T-2} \frac{\dots}{\Gamma, x:\forall V. \tau \vdash e_2\{\pi_2\}\{\rho_2\} : \tau_2\{\rho_2\}} \{\rho'\}}{\Gamma \vdash \text{let } x_V:\tau_1\{\rho_1 \circ \rho\} = e_1\{\pi_1\}\{\rho_1 \circ \rho\} \text{ in } e_2\{\pi_2\}\{\rho_2\} : \tau_2} \{\rho'\}$$

Once established, this invariant allows an easy deriving of inference soundness: for top-level terms, the constraint set must be empty, hence choosing the empty substitution yields a valid type derivation. \square

Conjecture C.4 (Inference non-specialization). *Given an expression e and two reconstructed terms e_1, e_2 where e_1 is the result of the inference, then they have equivalent observable semantics. Formally:*

$$\begin{cases} \Gamma \vdash e \rightarrow e_1 \\ \Gamma \vdash e \leftarrow e_2 \end{cases} \implies e_1 \equiv e_2$$

This property states that inference rejects any terms that might be given multiple semantics, hence that the implicitly typed language is deterministic.

The proof has not been fully conducted yet. The main difficulty lies in the handling of diverging environments: let bindings introduce variables that can be given multiple type schemes. At instantiation places, these variables induce different types leading to different constraint sets. We must then prove that the size constraints extracted by the inference are indeed fulfilled by the arbitrary reconstruction. This would allow to deduce that the substitution built by the inference builds a *most general size unifier*.

Let define a characterization of type schemes built by the inference.

Definition C.5 (Subsumption). Given two type schemes $\sigma_1 := \forall V_1. \tau_1$ and $\sigma_2 := \forall V_2. \tau_2$, the *subsumption* relation $\sigma_1 \preccurlyeq \sigma_2$ holds if and only if any instance of the second is a sub-type of an instance of the first. Formally, one of the two equivalent formulation must hold:

$$\begin{aligned} & \forall S_2, \exists S_1, \tau_1\{S_1/V_1\} <: \tau_2\{S_2/V_2\} \\ & \exists S, \tau_1\{S/V_1\} <: \tau_2 \text{ (where } FV(S) \in V_2) \end{aligned}$$

Because our type systems does not have principal types, the types produced by the inference cannot be the most polymorphic one, i.e. such that is subsume to any other valid type. This property must be refined for sizes:

Definition C.6 (Size subsumption). Given two type schemes σ_1 and σ_2 , the *size subsumption* relation $\sigma_1 \leqslant_\eta \sigma_2$ holds if and only if:

$$\exists \sigma, \bar{\alpha}, \bar{\tau}. \begin{cases} \sigma_1 \leqslant \sigma\{\overline{\text{int}}/\bar{\alpha}\} \\ \sigma\{\bar{\tau}/\bar{\alpha}\} \leqslant \sigma_2 \\ \bar{\tau} <: \overline{\text{int}} \end{cases}$$

The intuition is: if $\sigma_1 \leqslant_\eta \sigma_2$, the second type scheme is obtained by adding some refinements (in both positive and negative positions) in place of some `int` in σ_1 . In a term, this would guarantee that the semantics is independent of those extra refinements (since a semantics exists without). Thus showing that inference builds such minimal term for the size subsumption relation would help establishing our conjecture.