Heapless Functional Programming

```
Ellis Kesterton<sup>[0009-0006-9369-8672]</sup> and Edwin Brady<sup>[0000-0002-9734-367X]</sup>

University of St Andrews, UK {erk4,ecb10}@st-andrews.ac.uk
```

Abstract. Typed functional programming languages like Haskell and OCaml make heavy use of the heap at run-time. This makes them largely unsuitable for systems programming, where resources are limited and programs are often expected to run on bare-metal. This paper demonstrates how a (slightly restricted) high-level, pure, functional language can be compiled to machine code which does not use the heap at all. Despite usually requiring a heap at run-time, features such as higher-order functions, polymorphism and typeclasses are all supported by the surface language. This is made possible through partial evaluation [12]: by carefully reducing the program at compile-time, we can eliminate these highlevel features entirely, resulting in a residual program which is trivial to compile to stack-based machine code. This paper describes the operation of this partial evaluator, justifies its design, and introduces a novel type system which guarantees that the partial evaluator will always succeed in removing all heap-using features.

1 Introduction

At the core of most functional languages is a set of high-level features including higher-order functions, polymorphism and algebraic data types. Most of the time these features are a boon: they allow the programmer to write concise, maintainable, readable code. However, all of these features typically require a heap at run-time. This is problematic if we want to use functional languages for systems programming, where it is common to target environments with limited resources and real time constraints. Consider the following Haskell program:

```
main :: I0 ()
main = do
    n :: Int <- readInt
    m :: Int <- readInt
    printInt (n + m)</pre>
```

While this program may appear simple, it implicitly makes use of a wide variety of high-level (heap-using) features. Alone, using do-notation requires type-classes (to resolve the Monad instance), higher-order functions (for the implicit calls to >>=) and polymorphism (also for >>=)! It is clear that every non-trivial

Haskell program will use the heap in some capacity – features such as higherorder functions are too ubiquitous to avoid. Does this mean that high-level functional programming is entirely reliant on the heap? This paper demonstrates that this is not the case.

Reconsider our example from a computational perspective: it is unclear why a program which adds two integers should need a heap at all! In this case, features such as typeclasses and higher-order functions are only abstractions for the programmer; they are not adding any real computational value. We show that most of the time we can actually use these abstractions for free. Through careful partial evaluation at compile-time, it is possible to transform a large class of high-level functional programs into equivalent low-level programs which do not require the heap at all. We make use of the type system to ensure that input programs always fall into this transformable class, guaranteeing that partial evaluation will always succeed at removing high-level features. Compiling a desugared version of the earlier example with our prototype implementation produces identical machine code to the following Rust program¹.

```
fn main() {
    let n = readInt();
    let m = readInt();
    printInt(n + m);
}
```

2 Contributions

The goal of this paper is to describe a method for designing and implementing high-level functional languages which do not use the heap at run-time. We believe this goal is worthwhile as not every low-level environment can feasibly support the high amount of heap allocations typically required by a functional language — by only using the stack, it becomes possible to use functional programming in these restricted environments. We hope that this work will help facilitate the design of functional languages which are suitable for systems programming on low-resource devices such as microcontrollers.

We focus on the design and implementation of a core language \mathcal{C} , which aims to be an elaboration target for high-level strict functional languages. Specifically, we make the following contributions:

- We define a strict, polymorphic, higher-order core language \mathcal{C} which serves as an elaboration target for surface languages. \mathcal{C} has a novel type system which guarantees that it can be compiled to stack-based code.
- We describe a pipeline which compiles \mathcal{C} programs down to machine code. We primarily focus on how partial evaluation can be used to eliminate all of \mathcal{C} 's high-level features (higher-order functions, polymorphism) at compile-time.

We are assuming that readInt and printInt are externally defined primitives in both code snippets.

- We informally describe how one could build a practical surface language on top of \mathcal{C} . We cover how common features such as typeclasses (incl. monads), implicit polymorphism and data types can be implemented.

We have implemented a prototype for the approach described in this paper, whose source code is publicly available².

3 Core Language

3.1 Outline

This section will define the core language C and describe how it can be compiled to stack-based machine code. The compilation process from C to machine code is broken down into several distinct steps (occurring in the following order):

- 1. Partial Evaluation
- 2. Uncurrying
- 3. Lambda Lifting
- 4. Code Generation

At the core of our approach lies a partial evaluator. By carefully performing reductions on the input program at compile-time, we can completely eliminate high-level features such as higher-order functions and polymorphism. By eliminating these features at compile-time, we no longer need to deal with the problem of representing them at run-time. The result of partial evaluation is therefore a program written in a very restricted subset of $\mathcal C$, greatly simplifying the rest of the compilation process. So that we can guarantee that partial evaluation will eliminate all occurences of these high-level features, we slightly restrict $\mathcal C$ programs through the type system. These non-standard typing rules will be justified in detail throughout the remainder of this section.

After partial evaluation, we follow up with two auxiliary source-to-source transformations: uncurrying and lambda lifting. These steps remove curried and nested functions respectively, and are implemented using well-studied algorithms. Unlike standard approaches to uncurrying [9], the type system and partial evaluator guarantee that the input is in a restricted form where function definitions are maximally η -expanded and there are no partial applications, which greatly simplifies the process. Our approach to lambda lifting is entirely standard, and ultimately the choice of algorithm is irrelevant. Quadratic-time algorithms have been specified in detail in existing literature [22].

Finally, the result of lambda lifting is compiled to machine code. At this point in the pipeline the language is extremely restricted, with only first-order top-level function definitions remaining. Compilation is therefore trivial using standard techniques for first-order languages — our sample implementation uses a straightforward conversion to Rust [14] to complete the compilation process.

² Download link: https://github.com/eayus/heapless

3.2 Syntax

The core language \mathcal{C} is an extension of System F_{ω} with let-bindings (both regular and recursive), products, base types, and primitive operations. We assume the obvious strict semantics. The full syntax is given in Figure 1. In general, we use a single colon for annotating expressions with their type (e.g. $e:\tau$), and we use a double colon for annotating types with their kind (e.g. $\tau:\kappa$).

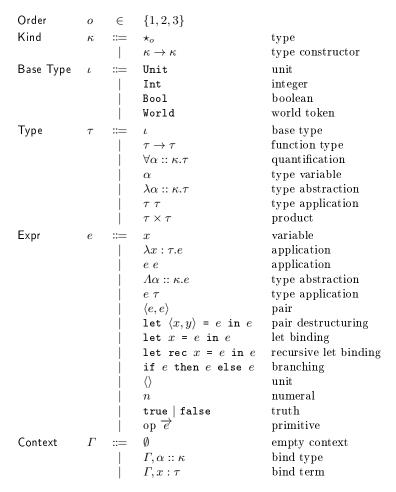


Fig. 1. Syntax of the core language C.

The reason for choosing System F_{ω} as the base for our core language is that it is expressive enough to support all of the high-level features we are interested in via trivial elaboration. Notably, the type-level evaluation facilities of System F_{ω} can be used to implement typeclasses, discussed in detail in Section 6.3.

As a further syntactic convention, throughout the paper we will often give examples in Haskell-like syntax rather than \mathcal{C} . This is to aid readability, as programming directly in a core language is very terse. We assume the natural elaboration from the Haskell-like syntax to \mathcal{C} , which mostly amounts to inferring things like explicit type applications and type binders; inferring whether a definition is bound with let or let rec; and ASCII rather than unicode syntax (i.e. Type1 rather than \star_1).

3.3 Type System

The kinding and typing rules for C are given in Figures 2 and 3 respectively. The typing rules for primitive operations (including constants and arithmetic/logic operations) are omitted for brevity, assuming the usual types.

The most unusual part of the type system is the kind \star_o , indexed by an order o. The full justification for this system is given in Section 4.5, but we give an intuitive overview here. The basic premise is that we want to differentiate first-order and higher-order terms. Types are categorised by their kind by their kind as follows.

- \star_1 Basic types, such as Int, Bool, etc.
- *2 First-order function types, modulo currying. Int -> Int -> Bool is included in this kind, but (Int -> Int) -> Int is not.
- \star_3 Everything else, including higher-order function types and types using universal quantification.

Intuitively, the kinds form a hierarchy where kinds with a higher number are more permissive in the types they allow. This hierarchy is realised through the sub-typing rule for kinds KIND-SUB, meaning each kind is included in the one above.

Most of the rules are straightforward, however some of the rules require some special attention. In general, types with kind \star_3 should be *contagious*, meaning that if they are used as a component of a composite type, then the composite type should also have kind \star_3 . For example, a product type where one of the elements is a higher order function (of kind \star_3) must also be considered higher order, and hence be given the kind \star_3 . This is to ensure that higher-order functions and polymorphism are still removed, even when they are used as part of another type (such as a product). Formally, we represent this in the typing rules as taking the maximum of two orders (see max in KIND-PROD and KIND-ARROW). The typing rule TYPE-LET-REC makes use of the kinding information attached to types to enforce that recursive functions are not higher-order and do not use universal quantification in their type – the justification behind this rule is described in detail in Section 4.4.

3.4 Handling IO

IO functions are an important part of any practical language. In general, impurity does not mix well with program manipulation and partial evaluation —

$$\frac{\text{KIND-BASE}}{\Gamma \vdash \iota :: \star_{1}} \qquad \frac{\prod\limits_{\Gamma \vdash \alpha}^{\text{KIND-FORALL}} \Gamma, \alpha :: \kappa \vdash \tau :: \star_{i}}{\Gamma \vdash (\forall \alpha :: \kappa.\tau) :: \star_{3}} \qquad \frac{\prod\limits_{\Gamma \vdash \alpha}^{\text{KIND-VAR}} \Gamma}{\Gamma \vdash \alpha :: \kappa}$$

$$\frac{\prod\limits_{\Gamma \vdash \alpha}^{\text{KIND-ARROW}} \Gamma \vdash \tau_{1} :: \star_{i}}{\Gamma \vdash \tau_{1} :: \star_{i}} \qquad \frac{\Gamma \vdash \tau_{2} :: \star_{j}}{\Gamma \vdash \tau_{1} :: \star_{i}} \qquad \frac{K \mid \text{ND-ARROW}}{K \mid \text{ND-ABS}}$$

$$\frac{\Gamma, \alpha :: \kappa_{1} \vdash \tau :: \kappa_{2}}{\Gamma \vdash (\lambda \alpha :: \kappa_{1}.\tau) :: \kappa_{1} \to \kappa_{2}} \qquad \frac{\prod\limits_{\Gamma \vdash \alpha}^{\text{KIND-APP}} \Gamma \vdash \tau_{1} :: \kappa_{1} \to \kappa_{2}}{\Gamma \vdash \tau_{1}\tau_{2} :: \kappa_{2}}$$

$$\frac{\prod\limits_{\Gamma \vdash \alpha}^{\text{KIND-PROD}} \Gamma \vdash \tau_{1} :: \star_{i}}{\Gamma \vdash \tau_{1} :: \star_{i}} \qquad \frac{\Gamma \vdash \tau_{2} :: \star_{i}}{\Gamma \vdash \tau :: \star_{i}} \qquad i < j}{\Gamma \vdash \tau :: \star_{j}}$$

Fig. 2. Kinding rules for C.

even something as simple as inlining a let-binding can become problematic if the bound expression has a side-effect. The most popular solution in modern functional programming is the IO monad, but it isn't trivial to represent at run-time without a heap. Instead, we thread a linear World token throughout the program, similar to the approach used by Clean [15]. At run-time the World token can be treated as the unit type. In Section 6.4 we show how to reintroduce the IO monad in the surface language userspace.

Note that the typing rules presented in this paper do not enforce the linearity of the world token – this is done intentionally to keep the presentation of the type system concise. However, real implementations should enforce this property properly, as we do in our sample implementation. We use an approach which supports linearity "on the arrow" using multiplicities, almost identical to the approach used in Linear Haskell [4].

4 Partial Evaluator

4.1 Rules

We characterise the partial evaluator by a set of source-to-source reduction rules on \mathcal{C} . The notation $t \leadsto u$ means that subterms matching the pattern t should be replaced with u during partial evaluation. Unless specified otherwise, there are no restrictions on where these reductions can take place, meaning that we reduce under binders by default. Occasionally, rules may have an attached side-condition which restricts when the rule should be applied – this means the rules

³ In this context, the distinction between linear and unique types is irrelevant, as only one World will ever be present.

$$\begin{array}{c} \begin{array}{c} \text{TYPE-VAR} \\ (x:\tau) \in \varGamma} \\ \hline \varGamma \vdash x:\tau \end{array} & \begin{array}{c} \Gamma \text{YPE-ABS} \\ \varGamma \vdash x:\tau \end{array} & \begin{array}{c} \Gamma \text{YPE-APP} \\ \varGamma \vdash x:\tau \end{array} & \begin{array}{c} \Gamma \vdash e:\tau_2 \\ \hline \varGamma \vdash e:\tau_1 \rightarrow \tau_2 \end{array} & \begin{array}{c} \Gamma \vdash e:\tau_1 \\ \hline \varGamma \vdash e:\tau_1 \rightarrow \tau_2 \end{array} & \begin{array}{c} \Gamma \vdash e:\tau_1 \\ \hline \varGamma \vdash e:\tau_2 \rightarrow \tau_2 \end{array} \\ \hline \end{array} \\ \begin{array}{c} \text{TYPE-FORALL} \\ \hline \varGamma \vdash (\varLambda\alpha::\kappa \vdash e:\tau) \\ \hline \varGamma \vdash (\varLambda\alpha::\kappa \vdash e:\tau) \end{array} & \begin{array}{c} \Gamma \text{YPE-INST} \\ \hline \varGamma \vdash e:\forall\alpha::\kappa \vdash \tau_1 & \varGamma \vdash \tau_2::\kappa \\ \hline \varGamma \vdash e:\tau_1 & \varGamma \vdash e:\tau_2 \end{array} \\ \hline \begin{array}{c} \Gamma \vdash e:\forall\alpha::\kappa \vdash \tau_1 & \varGamma \vdash \tau_2::\kappa \\ \hline \varGamma \vdash e:\tau_1 & \varGamma \vdash e:\tau_2 \end{array} \\ \hline \Gamma \vdash (\exists \vdash x_1) & \Gamma \vdash e:\tau_2 \end{array} & \begin{array}{c} \Gamma \text{YPE-LET-PROD} \\ \hline \varGamma \vdash e:\tau_1 & \Gamma \vdash x \rightarrow \tau_2 \end{array} & \begin{array}{c} \Gamma \text{YPE-LET-PROD} \\ \hline \varGamma \vdash e:\tau_1 & \Gamma \vdash x \rightarrow \tau_2 \end{array} & \begin{array}{c} \Gamma \vdash e:\tau_1 \rightarrow \tau_2 \end{array} \\ \hline \Gamma \vdash (\exists \vdash x,\tau_1) & \Gamma \vdash x \rightarrow \tau_2 \end{array} & \begin{array}{c} \Gamma \vdash e:\tau_1 \rightarrow \tau_2 \end{array} \\ \hline \Gamma \vdash (\exists \vdash x,\tau_1) & \Gamma \vdash x \rightarrow \tau_2 \end{array} & \begin{array}{c} \Gamma \vdash e:\tau_1 \rightarrow \tau_2 \end{array} \\ \hline \Gamma \vdash (\exists \vdash x,\tau_1) & \Gamma \vdash x \rightarrow \tau_2 \end{array} & \begin{array}{c} \Gamma \vdash x \rightarrow \tau_2 \rightarrow \tau_3 \\ \hline \Gamma \vdash (\exists \vdash x,\tau_1) & \Gamma \vdash x \rightarrow \tau_2 \end{array} \\ \hline \begin{array}{c} \Gamma \vdash x \rightarrow \tau_1 \rightarrow \tau_2 \rightarrow \tau_2 \\ \hline \Gamma \vdash (\exists \vdash x,\tau_1) & \Gamma \vdash x \rightarrow \tau_2 \end{array} & \begin{array}{c} \Gamma \vdash x \rightarrow \tau_1 \rightarrow \tau_2 \rightarrow \tau_2 \\ \hline \Gamma \vdash (\exists \vdash x,\tau_1) & \Gamma \vdash x \rightarrow \tau_2 \rightarrow \tau_3 \end{array} \\ \hline \begin{array}{c} \Gamma \vdash x \rightarrow \tau_1 \rightarrow \tau_2 \rightarrow \tau_2 \rightarrow \tau_3 \\ \hline \Gamma \vdash (\exists \vdash x,\tau_1) & \Gamma \vdash x \rightarrow \tau_2 \rightarrow \tau_3 \end{array} \\ \hline \begin{array}{c} \Gamma \vdash x \rightarrow \tau_1 \rightarrow \tau_2 \rightarrow \tau_2 \rightarrow \tau_3 \rightarrow$$

Fig. 3. Typing rules for C.

may not form a true rewriting system [10], but we borrow the notation and terminology for conciseness.

Figures 4 and 5 detail the partial evaluator's reduction rules for types and expressions respectively.

$$(\lambda \alpha :: \kappa.\tau_1) \ \tau_2 \leadsto \tau_1[\alpha \mapsto \tau_2] \tag{1}$$

Fig. 4. Reduction rules for partially evaluating types.

4.2 Elimination

Ultimately, the partial evaluator removes high-level features by applying the appropriate elimination rule. For example, higher-order functions are eliminated through β -reduction (β); polymorphic functions are eliminated through monomorphisation (mono). Almost all of the other rules are there to facilitate these eliminators by rewriting expressions into a form which is more amenable to further reduction. The most straightforward example of a "facilitating rule" is let-inlining

```
e_1 \leadsto \lambda x.e_1 \ x \quad \text{if } e_1: \tau_1 \to \tau_2, \ e_1 \text{ is not a $\lambda$ abstraction} \qquad (\eta) (\lambda x.e_1)e_2 \leadsto e_1[x \mapsto e_2] \qquad (\beta) (if e_1 then e_2 else e_3) e_4 \leadsto \text{if } e_1 then e_2 e_4 else e_3 e_4 \qquad (\textit{if-distr}) let x = e_1 in e_2 \leadsto e_2[x \mapsto e_1] if e_1 is higher-order (let-inline) (let x = e_1 in e_2) e_3 \leadsto \text{let } x = e_1 in e_2 e_3 \qquad (\textit{let-distr}) (let e_1 = e_2 = e_3 \qquad (e_1 = e_2)) (let e_2 = e_3 =
```

Fig. 5. Reduction rules for partially evaluating expressions.

(let-inline). Consider the following expression (assume g is a free variable of the appropriate type):

```
let twice = \lambda f : Int \rightarrow Int. \ \lambda x : Int. \ f \ (f \ x) in twice \ g \ n
```

At the moment there are no possible β -reductions, yet a higher-order function remains. We must therefore inline the higher-order function bound by let so that β -reduction can take place. After running the above program through the partial evaluator we are left with a first-order residual program:

The same idea applies to polymorphic functions which quantify over type variables – they must be inlined and monomorphised. This process is closely related to approaches which specialize higher-order functions to their functional arguments [6]. A deeper comparison with specialisation is given in Section 7.2.

4.3 Distribution

We must be careful that other constructs in the core language do not interfere with the elimination of high-level features such as higher-order functions and polymorphism. For constructs which are around at run-time but also allow higher-order functions as a sub-expression, we must introduce an extra rule to distribute function applications inside the construct.

If-Expressions If-expressions allow any type of data in their branches as long as both branches have the same type. We therefore should consider what happens if both branches contain higher-order functions, like in the following expression:

(if
$$e_1$$
 then e_2 else e_3) ($\lambda x.e_4$)

With just β -reduction, there is no way to make progress, yet higher-order functions still remain (e_2 and e_3) in the residual program! One might try adding extra reduction rules for if:

```
if true then e_t else e_f \leadsto e_t if false then e_t else e_f \leadsto e_f
```

However, if e_1 is a free variable (whose value may not be known until runtime), then we are still stuck. Instead, we add the rule *if-distr* which distributes applications into the branches of the if. This allows us to progress with reductions even if we do not know the value of e_1 ahead of time. After applying this rule, we are left with the following residual expression:

if
$$e_1$$
 then e_2 $(\lambda x.e_4)$ else e_3 $(\lambda x.e_4)$

Let Bindings Our rule for inlining let-bindings covers the case when the bound expression is a higher-order function, but what about when the body of the let is a higher-order function? The following expression shows that there is a similar problem as with if-expressions:

(let
$$n = 0$$
 in e_1) $(\lambda x.e_2)$

The partial evaluator will not inline the let as it does not bind a higher-order term. We are once again stuck, with no way to β -reduce the higher-order function bound in the body of the let. The solution is analogous to the one for if-expressions: introduce a rule which distributes applications inside the let's body. The rules let-distr, let-rec-distr and let-prod-distr implement this behaviour for the various types of let-bindings.

4.4 Recursion

Practical functional programming languages almost always support general recursion, usually in the form of recursive definitions. General recursion is a particularly interesting feature as it makes the language turing-complete, and means that the language is no longer strongly-normalising. This can cause problems for a partial evaluator, as we can no longer blindly reduce expressions without risking non-termination (this is true even if the input program is total).

In general, any rule which indiscriminantly unfolds recursion (i.e. let rec bindings) will be problematic. However, if we do not unfold recursion, then we run into problems when recursion is used in combination with higher-order functions. If we define a function which is both higher-order and recursive, then it will remain in the residual program (which should be first-order). Unlike non-recursive higher-order functions, we cannot inline and reduce the definition due to the let rec blocking further evaluation.

To resolve this, we introduce the following simple restriction: recursive functions must not be higher-order or use universal quantification in their type. This is realised through the kind system, where the typing rule TYPE-LET-REC requires that the type of recursive definitions must have kind \star_2 . This ensures that recursive definitions are first-order (modulo currying) in the residual program. The reason why we use the kind system to enforce this rather than a syntactic check is explained in Section 4.5.

With this set up, no problematic reduction rules are needed for unfolding recursion, as we can only bind first-order functions with let rec (which are easy to compile). However, this restrction does seem impractical: plenty of useful functions are both higher-order and recursive! Fortunately, we can easily rewrite the majority of these functions in an equivalent way which is accepted by the type system. More precisely, we can try to extract out the recursive part of the function into a local definition which is first-order. As an example, suppose we wanted to write the following function:

```
nTimes :: Int -> (Int -> Int) -> (Int -> Int)
nTimes 0 f x = x
nTimes n f x = f (nTimes (n - 1) f x)
```

In its current form, the function is both higher order and recursive, and so would not be permitted by our type system. However, we can refactor it into the following equivalent definition which would be permitted:

```
nTimes :: Int -> (Int -> Int) -> (Int -> Int)
nTimes n f = go n
  where
    go :: Int -> Int -> Int
    go 0 x = x
    go m x = f (go (m - 1) x)
```

The local definition go is recursive but not higher-order, while the exposed outer definition nTimes is higher-order but not recursive. To achieve this, the function go captures the function f, but does not have it as a parameter – when used, the function nTimes will be inlined, and the function f also inlined into the body of go. The result is a first-order program.

It is possible to do this refactoring for the majority of higher-order functions. Precisely, as long as the functional argument (in this case, $f:a \rightarrow b$) remains constant throughout all recursive calls, the refactoring is possible and straightforward⁴.

4.5 Polymorphism

Another trademark feature of modern functional programming languages is *polymorphism*. So that one function can operate on many different types of data, most

⁴ A practical compiler might choose to perform this refactoring automatically.

functional languages use a uniform representation for all data — a pointer to a heap allocated cell. Polymorphic functions can then indiscriminantly treat all inputs in the same way, regardless of type. Since this is not an option without a heap, we use a variant of the other predominant method: monomorphisation [19]. This is realised through the mono rule in the partial evaluator.

Recursive Functions Like higher-order functions, we are relying on polymorphism being eliminated at compile-time by the partial evaluator. However, once again, recursive function definitions can cause issues when used in conjunction with polymorphism. As already discussed, recursive functions cannot be unfolded due to issues with termination. Unfortunately, this means that any function that is both recursive and polymorphic will remain in the residual program, as the partial evaluator cannot monomorphise without inlining first. Similar to higher-order functions, the typing rule TYPE-LET-REC prevents this from happening as the kind of any type which uses universal quantification will be \star_3 .

To implement standard functions which are both polymorphic and recursive we can employ a similar trick to the one for higher-order functions: abstract out the recursive part into an inner auxiliary definition. This is always possible provided recursive calls are made with the same type variables to what the function was given originally (i.e. no *polymorphic recursion*). For example, we can improve our earlier example nTimes by making it polymorphic:

```
nTimes :: forall (a :: Type1). Int -> (a -> a) -> (a -> a)
nTimes n f = go n
where
   go :: Int -> a -> a
   go 0 x = x
   go m x = f (go (m - 1) x)
```

Note how, even though the function go uses type variables, because it does not use universal quantification, it is still considered a first-order function (i.e. it has kind \star_2). This is because the type variable **a** is bound with kind \star_1 .

Kind System In addition to the direct interaction between polymorphism and recursion, there is a more subtle nuance. When the typing rule for let rec was introduced, you may have wondered why a kind system is used to ensure the function is first-order (modulo currying), as opposed to a syntactic condition. There are a few reasons for this.

Primarily, the problem with a syntactic check relates to polymorphic type variables. It raises the question – should a type variable be considered a first-order or higher-order type? Consider a recursive function which takes an argument whose type is a type variable⁵:

⁵ Note that in its current form, the function choose is invalid as it does not give a kind for the type variable a. This is done intentionally to demonstrate the necessity of the kind system.

```
choose :: forall a. Int -> a -> a -> a
choose = go
  where
    go :: Int -> a -> a -> a
    go 0 x y = x
    go n x y = go (n - 1) y x
```

Depending on whether the variable a is later instantiated with a basic type or function type, the type of the monomorphised function may be either first-order (modulo currying) or higher-order. Because of this, we cannot treat all type variables as being basic types: doing so would allow the programmer to circumvent the restrictions on higher-order functions by hiding functions types behind a type variable. By separating type variables based on whether they represent functions or not, we can reason about whether a polymorphic function definition is truly first-order for every possible type instantiation. This idea of encoding information in kinds rather than types to handle polymorphism has been explored previously [8].

4.6 Semantic Preservation

An obvious consequence of introducing general recursion is that input programs are no longer guaranteed to always terminate. Not only does this cause issues with the termination of the partial evaluator itself, but it can also pose problems related to the preservation of program semantics. In general, we would expect the program resulting from partial evaluation to have the same semantics as the input program. However, since general recursion introduces the implicit side effect of non-termination, this is not always the case. Consider the following program:

```
let rec f : Int \rightarrow Int = f in 0
```

Before any transformation, this program would not terminate as we we would get stuck trying to evaluate f (recall we are following strict semantics). Applying the partial evaluator would produce the following residual program:

```
let rec f : Int \rightarrow Int = \lambda x.fx in 0
```

The program semantics have changed! Now, this program will terminate with a value of 0. The only difference between the two programs is that the right-hand side of f has been η -expanded by the partial evaluator. This means that now, we would only run into an infinite loop if f was applied to an argument, which is not the case.

The root of the problem is that non-termination is an untracked side effect, making the language impure. We propose two ways of handling this problem:

- Modify the language so it is pure and total by introducing a termination checker
- Only provide the guarantee of semantics preservation for programs that terminate.

For the purpose of this paper we choose the latter, though a practical implementation may prefer to use a termination checker so that the compiler is more predictable. If the programmer wishes to intentionally perform an infinite loop, then this is still possible by using an IO primitive.

4.7 Further Optimisation

In addition to eliminating high-level constructs, partial evaluation can also be used as a general optimisation method. The partial evaluator presented in this paper uses only the minimal rules necessary to eliminate high-level constructs, but a practical implementation may wish to add additional rules for optimisation purposes.

A simple example is to add rules which implement *constant folding*, a process which reduces primitive operations at compile-time if all operands are known. For example, 3*(2+2) should be reduced to 12 at compile-time. Integrating this optimisation into the partial evaluator could remove the need for a separate pass later on in the compilation process.

The inlining of non-essential definitions (i.e. those which are first-order) could also be easily implemented as an optimisation. Currently, the partial evaluator never inlines anything which isn't necessary (with the idea of avoiding code duplication), but if the bound expression is small then inlining may yield a performance benefit.

5 Implementation

5.1 Partial Evaluator

The most straightforward way to implement the partial evaluator is to continuously traverse the input expression, identifying and reducing expressions which match a reduction rule. For β -redexes, reduction will involve substitution, which is another traversal of the expression. It is well-known that this method of implementing normalisers/partial evaluators is inefficient due to the repeated term traversals, even if optimisations such as De-Bruijn indices [3] are used to speed up substitution.

A modern approach to implementing normalisers is normalisation by evaluation (NbE), where expressions are first evaluated into semantic values and then later reified back into syntactic terms [1,2]. Expensive substitutions are avoided entirely, with environments and variable lookups being used instead. NbE has been widely employed in the implementation of type checkers for dependently typed languages such as Idris and Agda. Since normalisers and partial evaluators are fundamentally very similar, it is relatively straightforward to adapt NbE

algorithms to implement our partial evaluator efficiently – this is the approach we use in our sample implementation.

There are a few ways in which our implementation differs from a conventional NbE algorithm. First, the definition of our normal and neutral terms is more lenient than usual, as redexes like let-bindings should be permitted in the result (provided the let-binding does not bind a higher-order function of course). Second, because we are working in a language based on System F_{ω} , we must perform both term-level and type-level reductions. This requires two distinct evaluation environments. Third, core language terms must carry thorough typing annotations as we are compiling to a typed language backend – the algorithm must ensure that these typing annotations are preserved during evaluation. The complete algorithm is verbose and relatively uninteresting, and so we refer the reader to our sample implementation for the details.

5.2 Uncurrying

The goal of this step in the compilation pipeline is to remove curried functions by transforming them to ones which operate on tuples instead. Due to the η rule in our partial evaluator, functions will always be maximally η -expanded and maximally applied. This makes this process extremely trivial, and we implement this as two simple syntactic transformations:

$$ft_1t_2 \dots t_n \leadsto f(t_1, t_2, \dots t_n)$$
$$\lambda x_1.\lambda x_2.\dots \lambda x_n.t \leadsto \lambda(x_1, x_2, \dots x_n).t$$

Fig. 6. Reduction rules for uncurrying.

5.3 Lambda Lifting

At this stage the program may still have local function definitions. We remove these using the standard method of lambda lifting. This has been discussed in detail in existing literature [7,22], and so we do not repeat the details here.

5.4 Backend

At this point in the compilation pipeline, we have a relatively simple first-order language. We offload the remainder of the compilation process to a backend – there are many possible choices, but in our sample implementation we generate Rust code⁶. This would be an unusual choice for a production language as the

⁶ Note that we do not depend on any code in the standard library which uses the heap.

Rust compiler will also perform type checking, which is unnecessary and inefficient. However, it is well suited to a prototype implementation as Rust supports first class tuples and nested let-bindings, of which both are utilized in \mathcal{C} . Compared to an alternative backend like LLVM [18], this simplifies the compilation process significantly as we do not need extra passes to handle these features.

6 Elaboration

Rather than continuously extend our partial evaluator with new features, we use a common architectural pattern: keep the core language small, but have a more expressive surface language which *elaborates* down to the core language. This section details several common features which can be implemented this way which cumulatively form a Haskell-like language.

6.1 Implicit Polymorphism

An easy extension to our core language which significantly improves practical usability is *implicit polymorphism*, where type applications and abstractions may be omitted and left to be inferred by the compiler. This is a well understood idea and so we do not discuss this in detail. One concrete approach is to use a modification of the Hindley-Milner algorithm [20] which produces type annotated terms (elaboration).

6.2 Algebraic Data Types

Non-Recursive Non-recursive algebraic datatypes are trivial to implement via a simple sum-of-products representation. Record and product types can be represented by nested product types, and sum types may be implemented using tagged unions.

Recursive Implementing recursive algebraic datatypes is significantly more challenging. Clearly, it is not possible to represent these types at runtime as they have no fixed size and we do not have access to a heap⁷. However, like many other features such as higher-order functions, we can permit recursive datatypes in the surface language provided we can guarantee that they will be eliminated at compile time. A remarkably simple way of achieving this is to use Church-encoding to represent recursive datatypes as functions (à la System F), transformed during the elaboration step. We already know how to deal with higher-order functions in the core language, and so this can be implemented with little complexity. For example, after elaboration, the standard definition of a linked list would be represented by the following type:

$$List \equiv \lambda(A :: \star_3). \forall X. X \rightarrow (A \rightarrow X \rightarrow X) \rightarrow X$$

⁷ The lack of a heap means that we cannot use the conventional linked representation.

There are a few important technical details. First, in the surface language any recursive algebraic data type must have kind \star_3 , as it will eventually be elaborated into a polymorphic higher-order function. Second, elimination of recursive algebraic datatypes through recursive functions should be disallowed, and instead a primitive fold operator should be introduced. The fold operator can be elaborated down into a simple application of the church-encoded function.

This idea of elaborating data types into their church encoding so that they will be eliminated at compile time is essentially a form of *deforestation* [24]. Our approach is much more simplistic than many existing approaches to deforestation, mainly because we do not allow arbitrary recursive definitions (instead exposing a primitive fold).

6.3 Typeclasses

Typeclasses are another common feature in surface-level languages which we can support via elaboration. Typeclass constraints can be viewed as implicit function arguments, where the typeclass's functions are stored as a record – this technique is commonly referred to as "dictionary passing". For example, suppose we had a Number class:

```
class Number (a :: Type1) where
  add :: a -> a -> a
  sub :: a -> a -> a
```

This could be elaborated to the following record type:

```
\texttt{Number} \equiv \lambda(a :: \star_1).(a \to a \to a) \times (a \to a \to a)
```

As with all features introduced via elaboration, we must be careful to ensure that the elaborated code is always valid with respect to the core language's typing rules. Since a typeclass constraint may eventually be elaborated into a higher-order type, then it means that in the surface language types using typeclasses should live in \star_3 . To write a recursive function which uses typeclasses, we must apply our classic transformation of factoring out the recursive part of the function into a local definition.

6.4 Monadic Example

We now culminate these features into a simple high-level example program. Suppose we want to write a program which reads in two numbers from the command line, n and k, and then prints out k for n repetitions. Since this program heavily interacts with IO, it would make sense to write our program within the IO monad. Fortunately, despite not being a primitive in the language, we can still use the IO monad by building it in userspace:

```
newtype IO (a :: Type1) = IO { runIO :: World -> (World, a) }
```

The IO type has trivial instances for Functor, Applicative and Monad. There is some nuance in choosing the right kind for these typeclasses – we discuss this in more detail in Section 8.1, but for now we take a pragmatic approach and assume that all monads have kind $\star_1 \to \star_2$. See the appendix for the complete example code, including the definitions for the typeclasses and their instances.

```
instance Functor IO where ...
instance Applicative IO where ...
instance Monad IO where ...
```

Primitive functions should be wrapped in a format which is compatible with the IO monad.

```
readInt :: IO Int
readInt = IO $ \w -> primReadInt w

printInt :: Int -> IO ()
printInt n = IO $ \w -> (primPrintInt n w, ())
```

We can define the familiar replicateM, despite it being both higher-order and recursive:

```
replicateM_ :: forall m a. Monad m => Int -> m a -> m ()
replicateM_ n x = go n
  where
    go :: Int -> m ()
    go 0 = pure ()
    go m = do
        x
        go (m - 1)
```

Finally, we can write our program in idiomatic functional style, with no worries that it will be using the heap at runtime. Note that up until now, all of the code that we have written would likely be part of a standard library, and so the programmer would directly skip to this last step. Our program looks identical to regular Haskell code, which we believe is a large benefit of our approach.

```
printK :: IO ()
printK = do
   n <- readInt
   k <- readInt
   replicateM_ n (printInt k)</pre>
```

7 Related Work

7.1 Multi-Stage Programming

Partial evaluators and multi-stage languages have always been closely related. While the former lets the compiler automatically manipulate the program [11],

the latter gives full control to the programmer and allows them to choose which reductions will happen at compile-time [23]. Kovács's recent work [17] has shown how a heterogenous multi-stage language can be used to generate programs which contain no closures, a goal shared by this paper. We believe there are advantages to both Kovács's and our approaches, inherited from the fundamental differences between multi-stage programming and partial evaluation respectively – our approach has less syntactic overhead as we do not require extra staging annotations, but the programmer no longer has full control over the reductions that happen at compile-time.

7.2 Higher-Order Removal

The removal of higher-order functions by translating to a first-order representation has been very well-studied [5,6,21]. The usual motivation is to simplify the compilation process as first-order programs are much simpler to compile to machine-code. Unfortunately, existing methods for removing higher-order functions are not applicable to our use-case of heapless programming, as in the process of removing higher-order functions they introduce other language constructs which also require the heap at run-time.

Defunctionalisation Defunctionalisation observes that there are only a finite amount of lambda abstractions in a program, and so to represent a function we can enumerate all of the possibilities encoded in an algebraic data type. If a lambda abstraction captures variables from the surrounding scope, parameters should be added to the corresponding ADT constructor representing the captured data.

For example, if a program contained two lambda abstractions:

$$\Gamma \vdash \lambda x : \text{Int.} x + x$$
 (2)

$$\Gamma, y : \text{Int} \vdash \lambda x : \text{Int}.x + y$$
 (3)

Then we encode functions by the following ADT:

The problem with this method for our use-case occurs when one function captures another function from the surrouding scope. The captured function must be added as a parameter to the data type, making the data type recursive. This would require a heap at runtime. More recent work such as *lambda set specialisation* [5] does not change anything in this regard.

Closure Conversion The predominant method for compiling higher-order functions is closure conversion, where functions are represented by a function pointer and values for any captured variables [21]. Since different functions (with the same type) may capture different amounts of variables, closures are variably sized, and so the usual approach is to allocate them on the heap behind a uniformly sized pointer.

While heap-allocated closures are obviously not an option without a heap, some low-level languages like Rust and C++ support stack-allocated closures instead. The basic idea is to allocate closures on the stack when they are created, and then pass them by reference so that they have a uniform representation. There are a few undesirable complications introduced with this method:

- When one closure captures another, the first closure is stored by reference in the second. This means that the second closure must not outlive the first, or we will access a dangling reference. Rust solves this with its signature "lifetimes", but we believe such an approach would be quite unwieldy given the ubiquity of closures in functional programming.
- Since different closures of the same type may be different sizes, every closure is assumed to have a unique (mutually incompatible) type. For example, an if-expression which returns different closures in different branches will be rejected by the type checker, since the types of branches will be different too. The approach in this paper does not suffer from this problem due to distribution rules (Section 4.3).

Specialisation Our approach is closely related to *specialisation*, which generates new definitions for higher-order functions every time they are used with a new higher-order argument [6]. Unlike existing approaches to specialisation, our type system ensures complete removal of higher-order functions, rather than a "best effort" approach which aims to remove as much as possible. Specialising methods are better at avoiding duplication of code, as they do not require inlining of every higher-order definition. This difference becomes important when the same function is called with the same higher-order argument in many different places: our approach will inline the definition many times, while specialisation will produce a single top-level definition which is re-used. While this may seem like a significant disadvantage, in practice most non-recursive higher-order functions are only a few operations long, and so are prime candidates for inlining anyway.

On the other hand, there are other areas in which our approach appears to improve over specialising approaches. For example, our method can eliminate higher-order functions which appear in data structures such as products, which is vital for implementing typeclasses in a surface language (Section 6.3). We also cover how to handle other language features that may usually use the heap, such as the IO monad. We believe our approach is simpler to implement than most specialising algorithms, as it is a simple modification of normalisation rather than an entirely bespoke algorithm.

8 Future Work

8.1 Order Polymorphism

In the current formalisation of the core language, a polymorphic term must fix a specific kind for the type variables it quantifies over (e.g. \star_1 , \star_2 etc.). This can be undesirable as often we want to use the same function or datatype for many different kinds. For example, consider the reader monad:

 $Reader \equiv \lambda R.\lambda A.R \rightarrow A$

What kind should we give this type? There are many different possibilities: $\star_1 \to \star_1 \to \star_2$ and $\star_2 \to \star_2 \to \star_3$ are just two examples. It is undesirable to repeat ourselves and write multiple versions of the same type that only differ by kind, and so it would be better if the core language provided a way to write a single polymorphic definition.

Looking at the typing rules for the \rightarrow type constructor, the most general kind for this would be polymorphic over *orders* (i.e. the number that the kind \star is indexed by). We propose a system very similar to the universe level polymorphism [16] we see in dependently type languages like Agda. Constructs such as \uparrow (increment order) and \sqcup (maximum of two orders) could be exposed in the syntax for orders, along with type constructors which quantify orders. The kind for the reader monad could then be specified as $\forall ij.\star_i \to \star_j \to \star_{(\uparrow i) \sqcup j}$, which is maximally general. The details of such a system are yet to be worked out, however given the stark similarity to universe levels, we believe that such a system is feasible.

8.2 Heap Effect

So far we have discussed the design of a language which avoids using the heap entirely, so that it can be ran in environments where a heap would be inefficient or impractical. However, real-world applications are unlikely to be so clear-cut, and sometimes we may actually have (partial) access to the heap. To better support this use-case, we could allow finer grained control over what parts of the program use the heap.

Rather than banning heap usage across the entire program, a monad or algebraic effect [13] could be used to track heap usage in the type system. Functions which are disallowed in our core language (i.e. recursive and higher-order definitions) could be allowed on the condition that their type indicates that they use the heap. A handler for the heap effect could be implemented in different ways depending on the run-time environment, potentially after some stack-based initialisation code. For example, an operating system might first query the system memory map, identify a free location, set up a heap, and finally handle the heap effect so that subsequent code can use the heap. This is only possible in a language which supports modularity over where the heap is used.

9 Conclusion

We have introduced a novel approach for designing pure functional languages which do not require the heap at runtime. We have shown how, supported by an appropriate type system, partial evaluation can be used to entirely eliminate the high-level features of a language at compile time. We hope that this work will help facilitate more research into functional languages suitable for systems programming.

Acknowledgements

We would like to thank the reviewers of this paper for their time and their invaluable feedback.

References

- 1. Abel, A., Chapman, J.: Normalization by Evaluation in the Delay Monad: A Case Study for Coinduction via Copatterns and Sized Types. arXiv preprint arXiv:1406.2059 (2014)
- Berger, U., Eberl, M., Schwichtenberg, H.: Normalization by evaluation. Prospects for Hardware Foundations: ESPRIT Working Group 8533 NADA—New Hardware Design Methods Survey Chapters pp. 117–137 (1998)
- 3. Berghofer, S., Urban, C.: A Head-to-Head Comparison of de Bruijn Indices and Names. Electronic Notes in Theoretical Computer Science 174(5), 53-67 (2007)
- 4. Bernardy, J.P., Boespflug, M., Newton, R.R., Peyton Jones, S., Spiwack, A.: Linear Haskell: Practical Linearity in a Higher-Order Polymorphic Language. Proceedings of the ACM on Programming Languages 2(POPL), 1–29 (2017)
- Brandon, W., Driscoll, B., Dai, F., Berkow, W., Milano, M.: Better Defunctionalization through Lambda Set Specialization. Proceedings of the ACM on Programming Languages 7(PLDI), 977-1000 (2023)
- Chin, W.N., Darlington, J.: A Higher-Order Removal Method. Lisp and Symbolic Computation 9, 287–322 (1996)
- 7. Danvy, O., Schultz, U.P.: Lambda-Lifting in Quadratic Time. In: International Symposium on Functional and Logic Programming. pp. 134–151. Springer (2002)
- 8. Downen, P., Ariola, Z.M., Peyton Jones, S., Eisenberg, R.A.: Kinds are Calling Conventions. Proceedings of the ACM on Programming Languages 4(ICFP), 1–29 (2020)
- 9. Hannan, J., Hicks, P.: Higher-Order UnCurrying. In: Proceedings of the 25th ACM SIGPLAN-SIGACT symposium on Principles of programming languages. pp. 1–11 (1998)
- Jay, C.B.: Long Bn Normal Forms and Confluence. Tech. rep., Tech. Rep. LFCS-91-183 (and its revised version) (1991)
- 11. Jones, N.D.: An Introduction to Partial Evaluation. ACM Computing Surveys (CSUR) 28(3), 480–503 (1996)
- 12. Jones, N.D., Gomard, C.K., Sestoft, P.: Partial Evaluation and Automatic Program Generation. Peter Sestoft (1993)

- 13. Karachalias, G., Koprivec, F., Pretnar, M., Schrijvers, T.: Efficient Compilation of Algebraic Effect Handlers. Proceedings of the ACM on Programming Languages 5(OOPSLA), 1–28 (2021)
- 14. Klabnik, S., Nichols, C.: The Rust Programming Language. No Starch Press (2023)
- Koopman, P., Plasmeijer, R., van Eekelen, M., Smetsers, S.: Functional Programming in Clean (2002)
- Kovács, A.: Generalized Universe Hierarchies and First-Class Universe Levels. arXiv preprint arXiv:2103.00223 (2021)
- 17. Kovács, A.: Closure-Free Functional Programming in a Two-Level Type Theory. Proceedings of the ACM on Programming Languages 8(ICFP), 659-692 (2024)
- 18. Lattner, C., Adve, V.: LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation. In: International symposium on code generation and optimization, 2004. CGO 2004. pp. 75–86. IEEE (2004)
- Lutze, M., Schuster, P., Brachthäuser, J.I.: The Simple Essence of Monomorphization. Proceedings of the ACM on Programming Languages 9(OOPSLA1), 1015– 1041 (2025)
- Milner, R.: A Theory of Type Polymorphism in Programming. Journal of computer and system sciences 17(3), 348–375 (1978)
- Minamide, Y., Morrisett, G., Harper, R.: Typed Closure Conversion. In: Proceedings of the 23rd ACM SIGPLAN-SIGACT symposium on principles of programming languages. pp. 271–283 (1996)
- Morazán, M.T., Schultz, U.P.: Optimal Lambda Lifting in Quadratic Time. In: Implementation and Application of Functional Languages: 19th International Workshop, IFL 2007, Freiburg, Germany, September 27-29, 2007. Revised Selected Papers 19. pp. 37-56. Springer (2008)
- Taha, W.: A Gentle Introduction to Multi-stage Programming. In: Domain-Specific Program Generation: International Seminar, Dagstuhl Castle, Germany, March 23-28, 2003. Revised Papers. pp. 30-50. Springer (2004)
- 24. Wadler, P.: Deforestation: Transforming Programs to Eliminate Trees. In: European Symposium on Programming. pp. 344-358. Springer (1988)

Appendix

```
newtype IO (a :: Type1) = IO { runIO :: World -> (World, a) }
class Functor (f :: Type1 -> Type2) where
  fmap :: (a -> b) -> f a -> f b

class Functor f => Applicative (f :: Type1 -> Type2) where
  pure :: a -> f a
  (<*>) :: f (a -> b) -> f a -> f b

class Applicative f => Monad (f :: Type1 -> Type2) where
  (>>=) :: f a -> (a -> f b) -> f b

instance Functor IO where
  fmap f (IO g) = IO $ \w -> let (w', x) = g w in (w', f x)
```

```
instance Applicative IO where
  pure x = IO  w \rightarrow (w, x)
  IO f <*> IO x = IO \$ \ \ w \rightarrow
    let (w', f') = f w
        (W'', X') = X W'
    in (w'', f' x')
instance Monad IO where
  (IO g) >>= f = IO \ \w -> let (w', x) = g w in runIO (f x) w'
readInt :: IO Int
readInt = IO $ \w -> primReadInt w
printInt :: Int -> IO ()
printInt n = IO $ \w -> (primPrintInt n w, ())
replicateM_ :: forall m a. Monad m => Int -> m a -> m ()
replicateM_n x = go n
 where
    go :: Int -> m ()
    go 0 = pure ()
    go m = do
     Х
      go (m - 1)
printK :: IO ()
printK = do
 n <- readInt
 k <- readInt
  replicateM_ n (printInt k)
main :: World -> World
main w = let (w', ()) = runIO printK w in w'
```