

A Graph-Based Higher-Order Intermediate Representation

Roland Leißa Marcel Köster Sebastian Hack

Department of Computer Science, Saarland University
(leissa, koester, hack}@cs.uni-saarland.de



Abstract

Many modern programming languages support both imperative and functional idioms. However, state-of-the-art imperative intermediate representations (IRs) cannot natively represent crucial functional concepts (like higher-order functions). On the other hand, functional IRs employ an explicit scope nesting, which is cumbersome to maintain across certain transformations.

In this paper we present Thorin: a higher-order, functional IR based on continuation-passing style that abandons explicit scope nesting in favor of a dependency graph. This makes Thorin an attractive IR for both imperative as well as functional languages. Furthermore, we present a novel program transformation to eliminate the overhead caused by higher-order functions. The main component of this transformation is *lambda mangling*: an important transformation primitive in Thorin. We demonstrate that lambda mangling subsumes many classic program transformations like tail-recursion elimination, loop unrolling or (partial) inlining. In our experiments we show that higher-order programs translated with Thorin are consistently as fast as C programs.

1. Introduction

Nowadays, many (existing and novel) programming languages (C++11, Java 8, Scala, Go, Rust, etc.) support imperative *and* functional programming. Essentially, a function f may be passed to another *higher-order function* (HOF) g while f may contain free variables. The common way to implement a function value is a closure: A record that contains the pointer to a function and value bindings of its free variables. For example, the C++ code in Figure 1a results into the (stylized) code in Figure 1b. This transformation is called *closure-conversion* [3, chap. 10]. However, implementing closures straightforwardly can incur a significant performance penalty. Ideally, the code in the example is compiled into a simple loop.

The IRs used in compilers for imperative languages (LLVM, Java Bytecode, etc.) are too low-level to represent free variables directly. Therefore, these compilers implement closure-conversion already in the front-end. This has several drawbacks:

1. The implementation is language-specific and can hardly be reused in another front-end.

2. The IR code is significantly bloated. For every function abstraction a new struct is created. For example, the LLVM code for the simple example consists of over 600 lines.
3. Finally, it is inelegant and inefficient to lower constructs that have to be restored later on. LLVM, for example, uses a combination of carefully coordinated analyses and transformations to eliminate closures: Inline the call to the closure's function pointer to be able to SSA-construct (`mem2reg`) the closure struct and finally dissolve the struct to scalar values. This strategy, for instance, fails to optimize recursive HOFs, like the one in the example above: After optimizations, the example in Figure 1a is still more than 250 lines of LLVM bitcode long.

IRs for functional languages [e.g. 3, 17] represent free variables naturally. Closure conversion is much less involved on these IRs because they contain closures as first-class citizens. Therefore, the compiler does not need to re-discover closures from their implementation as in imperative IRs.

However, one particular point makes the use of existing functional IRs unattractive: They rely on scope nesting to bind the use of a variable to its definition. Scopes are necessary to disambiguate a variable use when the variable name is defined more than once in the program. Although functional programs share with SSA the property that every use has exactly one reaching definition, they allow reusing the same variable name. Figure 2a shows an imperative program, (b) the SSA form of it in a stylized imperative IR, (c) the corresponding version in a stylized functional IR that uses continuation-passing style (CPS). Because the scope nesting is a feature of the *syntax* of the IR, all program transformations must maintain a correct scope nesting, which involves, among other things, proper renaming of the variables to avoid name capturing. This makes control-flow restructuring program transformations cumbersome to implement.

In this paper, we present the IR Thorin that blends concepts of imperative and functional IRs and is equally well-suited to represent imperative as well as functional programs. Thorin represents functions as first-class citizens and exploits the long-known correspondence between SSA and CPS [16] to *uniformly* represent control flow by continuations: jumps to basic blocks, function calls, generators, exceptions, etc. In that respect, Thorin is similar to an IR of a compiler for functional languages.

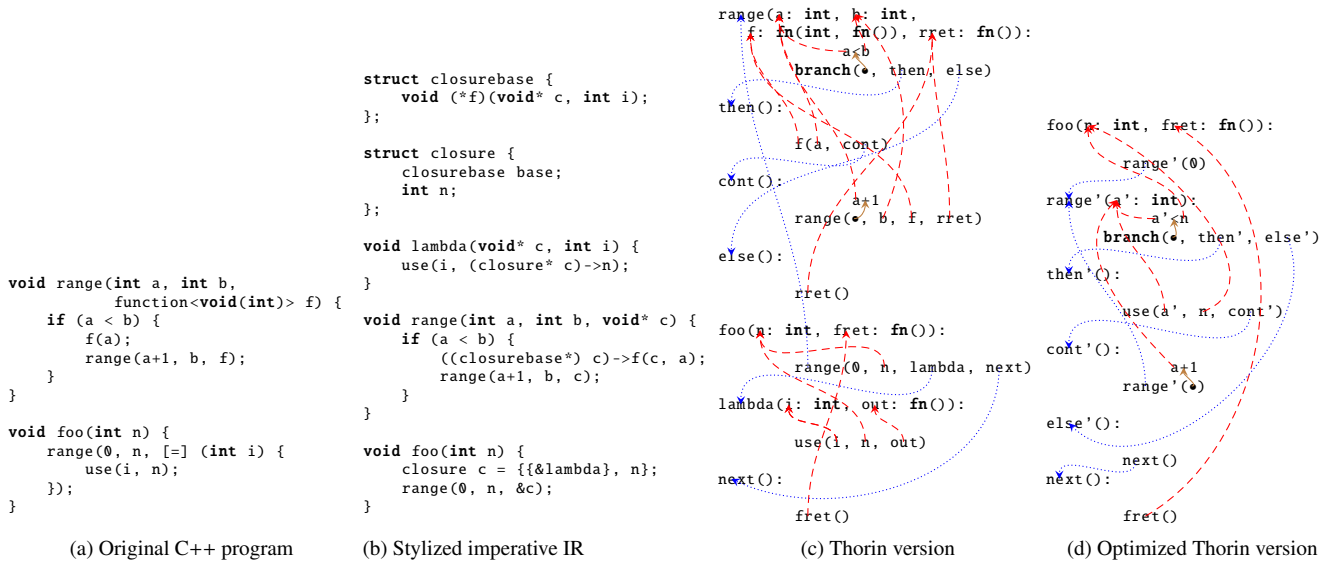


Figure 1. The higher-order function range and a call site within foo.

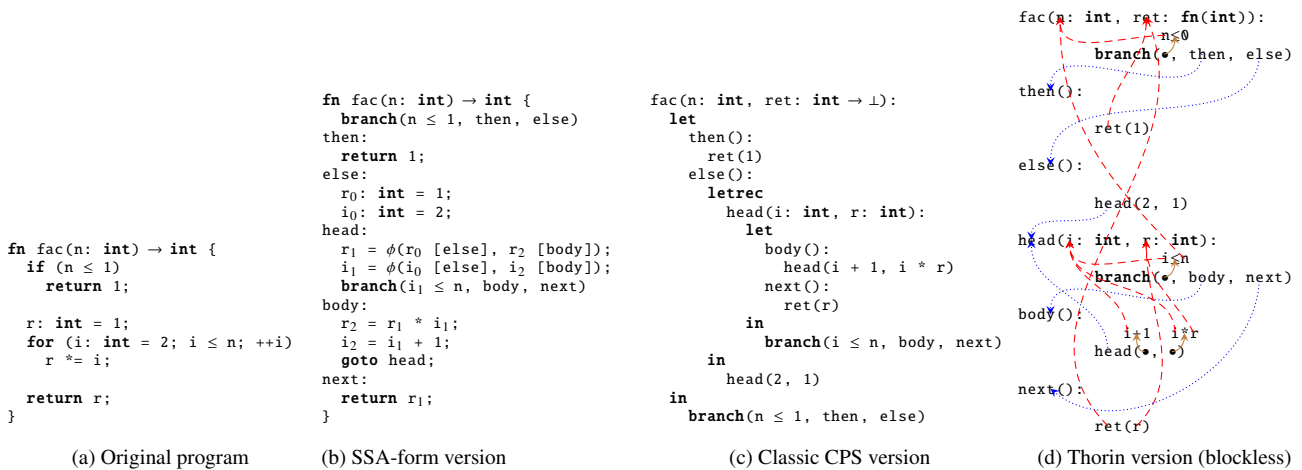


Figure 2. Factorial in different versions. In Thorin, each definition is a node. Every use is an edge to this node. Names are solely present for readability. They have no semantic meaning.

However, like low-level imperative IRs but unlike functional IRs, Thorin has no scope nesting. Instead, Thorin is graph-based like progressive imperative IRs (e.g. HotSpot [7, 22] or Firm [5]). That is, Thorin does not use named variables. In Thorin, each value is a node in a graph and every reference to a value is an edge to this node. Therefore, Thorin does not require explicit scope nesting. Figure 2d shows the Thorin graph for the example: A Thorin program consists of a set of functions. A function introduces parameters and in turn solely consists of a call. As Thorin uses CPS, the call does not return and, hence, function types (denoted by $\text{fn}(T)$) do not have a return type. The callee and the arguments to that call reference other definitions: functions (dotted edges), function parameters (dashed edges) or primops (solid edges). A primitive operation (primop) is a simple operation which references other definitions to produce a new value. Note that

in this paper we nevertheless use names in Thorin programs to make the presentation more accessible for humans. Names have no meaning otherwise.

In summary, this paper makes the following contributions:

- We discuss Thorin formally by presenting its syntax as well as its static and dynamic semantics (Section 3).
- We present *lambda mangling* (Section 4), a novel program transformation that partially in- and outlines functions. Partial inlining is the crucial step in closure elimination. We show that lambda mangling is substantially simpler to implement on Thorin than on existing IRs because Thorin does not have explicit scope nesting. Furthermore, we show how CPS helps to reduce other important control flow transformations, such as tail-recursion elimination, loop unrolling, and loop peeling, to lambda mangling.

<pre> a(x: int, ret: int → ⊥) → ⊥ let b() → ⊥ let c() → ⊥ ret(x) in c() in b() </pre> <p>(a) Classic CPS version</p>	<pre> a(x: int, ret: int → ⊥) → ⊥ let b() → ⊥ c() z() → ⊥ c() c() → ⊥ ret(x) in branch(x = 0, b, z) </pre> <p>(b) Classic CPS with new branch z</p>	<pre> a(x: int, ret: fn(int)): b() b(): c() c(): ret(x) </pre> <p>(c) Thorin version</p>	<pre> a(x: int, ret: fn(int)): branch(x = 0, b, z) b(): c() z(): c() c(): ret(x) </pre> <p>(d) Thorin version with new branch z</p>
--	---	--	---

Figure 3. This example illustrates Thorin’s blockless representation.

<table border="0" style="width: 100%;"> <tr> <td style="padding-right: 10px;">type</td> <td style="padding-right: 10px;">$t ::=$</td> <td style="padding-right: 10px;">int bool</td> <td>(int/bool type)</td> </tr> <tr> <td></td> <td></td> <td></td> <td>fn(\bar{t})</td> </tr> <tr> <td></td> <td></td> <td></td> <td>(function type)</td> </tr> <tr> <td>program</td> <td>$p ::=$</td> <td>$\mathcal{L} \rightarrow \text{fn}(\bar{t}) : b$</td> <td></td> </tr> <tr> <td>body</td> <td>$b ::=$</td> <td>$e(\bar{e})$</td> <td></td> </tr> <tr> <td>expression</td> <td>$e ::=$</td> <td>$\boxtimes(\bar{e})$</td> <td>(primop)</td> </tr> <tr> <td></td> <td></td> <td>ℓ</td> <td>(abstraction)</td> </tr> <tr> <td></td> <td></td> <td>$\ell_i \quad i \in \mathbb{N}$</td> <td>(parameter)</td> </tr> <tr> <td>value</td> <td>$v ::=$</td> <td>$c \in \mathcal{V}$</td> <td>(literal)</td> </tr> <tr> <td></td> <td></td> <td>$\langle \ell, \sigma \rangle$</td> <td>(closure)</td> </tr> </table> <hr/> <p>Body: $p, \sigma, \ell \Rightarrow \sigma', \ell'$</p> <p>E-Body $\frac{\sigma, e_0 \Downarrow \langle \ell', \sigma' \rangle \quad \sigma, e_1 \Downarrow v_1 \quad \dots \quad \sigma, e_n \Downarrow v_n}{p, \sigma, \ell \Rightarrow \sigma'[\ell \mapsto v_1, \dots, v_n], \ell'}$</p> <p>Expression: $\sigma, e \Downarrow v$</p> <p>E-Primop $\frac{\sigma, e_1 \Downarrow v_1 \quad \dots \quad \sigma, e_n \Downarrow v_n}{\sigma, \boxtimes(\bar{e}) \Downarrow \boxtimes(\bar{v})}$ E-Abs $\frac{}{\sigma, \ell \Downarrow \langle \ell, \sigma \rangle}$</p> <p>E-Param $\frac{\sigma[\ell] = v_1, \dots, v_n \quad 1 \leq i \leq n}{\sigma, \ell_i \Downarrow v_i}$</p>	type	$t ::=$	int bool	(int/bool type)				fn(\bar{t})				(function type)	program	$p ::=$	$\mathcal{L} \rightarrow \text{fn}(\bar{t}) : b$		body	$b ::=$	$e(\bar{e})$		expression	$e ::=$	$\boxtimes(\bar{e})$	(primop)			ℓ	(abstraction)			$\ell_i \quad i \in \mathbb{N}$	(parameter)	value	$v ::=$	$c \in \mathcal{V}$	(literal)			$\langle \ell, \sigma \rangle$	(closure)	<p>Program: $\vdash p$</p> <p>T-Prg $\frac{p := [\ell_1 \mapsto _ : b_1, \dots, \ell_n \mapsto _ : b_n] \quad p \vdash b_1 \quad \dots \quad p \vdash b_n}{\vdash p}$</p> <p>Body: $p \vdash b$</p> <p>T-Body $\frac{p \vdash e_0 : \text{fn}(\bar{t}) \quad p \vdash e_1 : t_1 \quad \dots \quad p \vdash e_n : t_n}{p \vdash e_0(\bar{e})}$</p> <p>Expression: $p \vdash e : t$</p> <p>T-Primop $\frac{t = \text{check}_{\boxtimes}(\bar{e})}{p \vdash \boxtimes(\bar{e}) : t}$ T-Abs $\frac{p[\ell] = \text{fn}(\bar{t}) : _}{p \vdash \ell : \text{fn}(\bar{t})}$</p> <p>T-Param $\frac{p[\ell] = \text{fn}(t_1, \dots, t_n) : _ \quad 1 \leq i \leq n}{p \vdash \ell_i : t_i}$</p> <hr/> <p>L-Param $\frac{\ell \neq \ell' \quad p[\ell'] = _ : b_{\ell'} \quad \exists i : \ell_i \leq b_{\ell'}}{p \vdash \ell \text{ live } \ell'}$ $p \vdash b \text{ live } b'$</p> <p>L-Abs $\frac{\ell \neq \ell' \quad p \vdash \ell \text{ live } \ell'' \quad p[\ell'] = _ : b_{\ell'} \quad \ell'' \leq b_{\ell'}}{p \vdash \ell \text{ live } \ell'}$</p>
type	$t ::=$	int bool	(int/bool type)																																						
			fn(\bar{t})																																						
			(function type)																																						
program	$p ::=$	$\mathcal{L} \rightarrow \text{fn}(\bar{t}) : b$																																							
body	$b ::=$	$e(\bar{e})$																																							
expression	$e ::=$	$\boxtimes(\bar{e})$	(primop)																																						
		ℓ	(abstraction)																																						
		$\ell_i \quad i \in \mathbb{N}$	(parameter)																																						
value	$v ::=$	$c \in \mathcal{V}$	(literal)																																						
		$\langle \ell, \sigma \rangle$	(closure)																																						

Figure 4. Syntax, Semantics, Typing and Liveness in Thorin

function which expects two arguments: An integer and a first-order function which expects a boolean. As functions do not return (because of CPS), they have no return type. Note that function types are at least of first order. Parameters can be of zeroth order or higher.

A Thorin *program* consists of a map which maps a label $\ell \in \mathcal{L}$ to a function $\text{fn}(\bar{t}) : b$. Note that we use the label ℓ only in the formalism of Thorin to refer to the node that contains the function $\text{fn}(\bar{t}) : b$. In Thorin’s implementation, labels do not exist. In addition, the textual representation in the examples introduces parameter names. A function consists of its type or signature $\text{fn}(\bar{t})$ and its *body* $b = e_0(\bar{e})$ which is a call to e_0 with arguments \bar{e} . The examples use a function *branch* of type $\text{fn}(\text{bool}, \text{fn}(), \text{fn}())$ for conditional branches. An *expression* can be

1. a *primop* $\boxtimes(\bar{e})$, where \boxtimes is an operator (e.g. $+$, $-$, \dots),
2. the i^{th} *parameter* ℓ_i of a function ℓ or
3. the abstraction of a function ℓ .

Typing. Rule T-Prg checks the body of each function in the program. Rule T-Body gathers all argument types and checks whether each matches the type of the callee. Note that the body itself does not possess a type as the call never returns.

A primop only expects specific input types and in turn has a specific output type. The function $\text{check}_{\boxtimes}(\bar{e})$ in T-Primop handles these primop-specific rules. The type of a function is resolved by looking up its signature in the program map (T-Abs). A reference to a parameter projects the corresponding type from the function’s signature (T-Param).

Semantics. Expression evaluation is defined by a big-step semantics. Memory operations or side-effects are explicitly tracked by a functional store [27, 28]. The environment $\sigma : \mathcal{L} \rightarrow \mathcal{V} \times \dots \times \mathcal{V}$ maps a function label to the *values* of its parameters. Each *value* is either some literal c of a set \mathcal{V} (here: integers and booleans) or a *closure*. This is a pair consisting of a function label and environment at the moment the function is captured (E-Abs). E-Primop evaluates all arguments to values and determines a new value (operator \oplus denotes the arithmetic operation as opposed to the syntactic terminal \boxtimes). The expression ℓ_i extracts the i^{th} parameter value from ℓ ’s parameter list (E-Param).

Bodies are evaluated in small step: We evaluate the arguments \bar{e} to values. The calling expression e_0 must evaluate to a closure ℓ' to which the program steps to. The new state σ' is formed by replacing the old parameter values of ℓ' by the newly evaluated arguments.

The Scope: Identifying Function Dependencies. As outlined in previous sections, functions in a Thorin program are not explicitly nested. For example, in [Figure 2d](#) the functions then and next *directly* depend on fac as both functions reference fac’s parameter ret. Likewise, head *directly* depends on fac because head uses fac’s parameter n. In [Figure 2d](#) function else depends on fac albeit else’s body does not directly use any of fac’s parameters. However, else invokes head, which directly depends on fac. Therefore, function else *indirectly depends* on fac and, hence, is also *implicitly* nested in fac. For many analyses and transformations (including lambda mangling, as presented in [Section 4](#)), we need to know all direct and indirect dependences from the view of a function. We obtain this information by applying a liveness analysis. It holds $p \vdash \ell$ live ℓ' if function ℓ is live in ℓ' . L-Param finds direct dependencies whereas L-Abs transitively determines indirect dependencies. We call all functions which are live from the view of an another function ℓ_e the *scope of ℓ_e* : $\text{scope}_p(\ell_e) := \{\ell \mid p \vdash \ell_e \text{ live } \ell\}$. We call ℓ_e the *entry function* of that scope.

4. Lambda Mangling

In this section, we present Thorin’s main transformation primitive: *lambda mangling*—a combination of lambda lifting [15] and dropping [9]. We demonstrate how many traditional compiler optimizations can be implemented with this transformation and how to deal with simple as well as mutual recursion.

When reading [Figure 5a](#) from left to right, we see how function g is *lambda-lifted* out of f by introducing a new parameter ret' which eliminates g ’s free variable ret . The new lifted version is called g' . Likewise—when reading [Figure 5a](#) from right to left—we see how function g' is *lambda-dropped* into the body of f by eliminating g' ’s parameter ret' and introducing a free variable ret .

[Figure 5b](#) depicts the same procedure for a Thorin program. Since Thorin programs are blockless, we do not have to move g/g' out of/into f . The scope analysis ([Section 3](#)) will identify that g is a function nested in f (in the left box) whereas g' and f will be discovered as two independent functions (in the right box).

4.1 Combining Lambda Lifting and Dropping

The local function $\text{pow}(a, b)$ in [Figure 6a](#) computes a^b . The function has two call sites: in calcx and calcy . As both callers pass 3 for parameter b to pow , we want to specialize pow . For this reason, we drop pow ’s parameter b and substitute each occurrence of b with 3. In doing so, we also apply local optimizations (constant propagation, common subexpression elimination, etc.). Thus, the check $b = 0$ and related blocks are eliminated in the new function pow_d ([Figure 6b](#)). Moreover, we update the call sites in calcx and calcy to call the new function pow_d instead.

$$\begin{array}{c}
\text{Body:} \\
\text{M-Body} \frac{M, e_0 \triangleright e'_0 \quad \dots \quad M, e_n \triangleright e'_n \quad \boxed{M, b \triangleright b'}}{M, e_0(e_1, \dots, e_n) \triangleright e'_0(e'_1, \dots, e'_n)} \\
\\
\text{Expression:} \\
\boxed{M, e \triangleright e'} \\
\\
\text{M-Mapped} \frac{M[e] = e'}{M, e \triangleright e'} \quad \text{M-Primop} \frac{\boxtimes(e_1, \dots, e_n) \notin \text{dom}(M) \quad M, e_1 \triangleright e'_1 \quad \dots \quad M, e_n \triangleright e'_n}{M, \boxtimes(e_1, \dots, e_n) \triangleright \boxtimes(e'_1, \dots, e'_n)} \\
\\
\text{M-Abs} \frac{\ell \notin \text{dom}(M)}{M, \ell \triangleright \ell} \quad \text{M-Param} \frac{\ell_i \notin \text{dom}(M) \quad M, \ell \triangleright \ell'}{M, \ell_i \triangleright \ell'_i}
\end{array}$$

Figure 7. Term rewriting for lambda mangling

```

function mangle(p, ℓe, t, M)
  foreach ℓ ∈ scopep(ℓe) \ ℓe do M[ℓ] ← new label
  foreach ℓ ∈ scopep(ℓe) \ ℓe do
    fn(̄i) : b ← p[ℓ]           ◀ get ℓ's signature and body
    ℓ' = M[ℓ]                 ◀ get ℓ's associated new label
    p[ℓ'] ← fn(̄i) : b'       ◀ insert new ℓ' where M, b ▷ b'
  end
  ℓ'e ← new label             ◀ now deal with entry: create new label
  -: be ← p[ℓe]           ◀ get ℓe's body
  p[ℓ'e] ← t : b'e         ◀ insert new entry where M, be ▷ b'e
  return ℓ'e                 ◀ return entry to new mangled region
end

```

Algorithm 1. Lambda mangling expects the program p to work on, the entry label ℓ_e of the scope to mangle, the new signature t and a map M which maps ℓ_e ’s signature to t .

The scope of function pow in [Figure 6a](#) uses f ’s parameter ret as free variable. Assume we would like to eliminate pow ’s dependency on f . To this end, we want to lift pow by introducing a new parameter ret_1 ([Figure 6c](#)).

If we want to apply both transformations, we could either first drop and then lift, or first lift and then drop ([Figure 6d](#)). Or more elegantly: apply both transformations simultaneously. Both, dropping and lifting, extend the program p by a clone of ℓ ’s scope. In this clone all expressions are rewritten with respect to a map M ([Figure 7](#)).

Reconsider the example in [Figure 6](#). We substitute b with 3 and keep the parameter a called a_d in the dropped version. We specify this mapping as follows: $M := \{a \mapsto a_d, b \mapsto 3\}$. We substitute ret with pow_1 ’s new parameter ret_1 and keep parameters a and b called a_1 and b_1 in the lifted version: $M := \{a \mapsto a_1, b \mapsto b_1, \text{ret} \mapsto \text{ret}_1\}$. We can simultaneously drop b with 3 and lift pow ’s free variable ret to a new parameter ret_m while obeying the mapping $M := \{a \mapsto a_m, b \mapsto 3, \text{ret} \mapsto \text{ret}_m\}$.

[Algorithm 1](#) performs this reconstruction for the whole scope of an entry label ℓ_e , with the new signature t_e (in the example $\text{fn}(\text{int}, \text{fn}(\text{int}))$) and a corresponding mapping M . The new region is added to the program p .

4.2 Recursion

As all function calls in CPS occur in tail position, it is tempting to think that in a CPS program only tail-recursion happens. This is not the case as can be seen in [Figure 8](#). The function else passes cont to fac . Reconsider Thorin

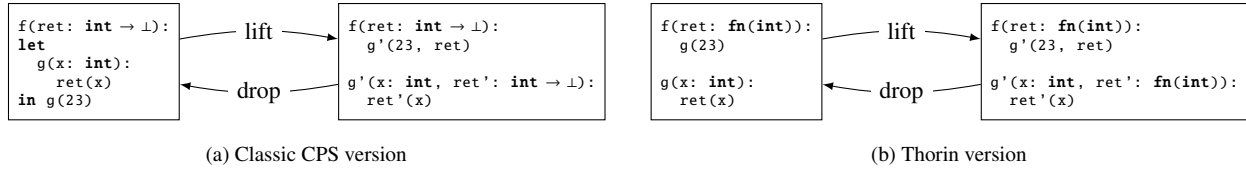


Figure 5. Lambda lifting/dropping

```

f(x: int, y: int, ret: fn(int)): f(x: int, y: int, ret: fn(int)): pow_l(a_l: int, b_l: int,
  branch(..., calcx, calcy)      branch(..., calcx, calcy)      ret_l: fn(int)):
pow(a: int, b: int):             pow_d(a_d: int):             then():
  branch(b = 0, then, else)       head(0, a_d)                       ret_l(1)
then():                           else():                             head(0, a_l)
  ret(1)                           head(i: int, r: int):             head(i: int, r: int):
else():                             branch(i < b_1, body, next)         branch(i < b_1, body, next)
  head(0, a)                         body():                             body():
head(i: int, r: int):             head(i+1, r*a_d)                   head(i+1, r*a_l)
  branch(i < b, body, next)         next():                             next():
body():                             ret(r)                               ret_l(r)
  head(i+1, r*a)                    next():                             head(i+1, r*a_m)
next():                             ret(r)                               body():
  ret(r)                             calcx():                             head(i+1, r*a_m)
calcx():                             pow_d(x)                             next():
  pow(x, 3)                           calcy():                             ret_m(r)
calcy():                             pow_l(y, 3, ret)                   ret_m(r)
  pow(y, 3)                           pow_l(x, 3, ret)                   next():
                                     pow_m(y, 3, ret)                   ret_m(r)
                                     pow_m(x, ret)
                                     calcy():
                                     pow_m(y, ret)

```

(a) The nested pow computes a^b . (b) Dropped pow.d computes $a.d^3$. (c) Lifted pow.l doesn't use free variables. (d) Dropped and lifted pow.m.

Figure 6. The function pow in (a) is always invoked with 3 as second argument. Hence, we want to specialize pow. Also, we would like to get rid of the free variable ret. Lambda dropping performs the specialization, lambda lifting eliminates pow's dependency on f and lambda mangling applies both transformations simultaneously.

```

fn fac(n: int) -> int {
  if (n <= 1)
    return 1;
  return n * fac(n-1);
}

fac(n: int, ret: fn(int)):
  branch(n <= 1, then, else)
then():
  ret(1)
else():
  fac(n-1, cont)
cont(res: int):
  ret(n*res)

```

(a) Original program (b) Thorin version

Figure 8. Naïve, recursive implementation of factorial

semantics (Figure 4): cont is a function abstraction and is evaluated to a closure. This closure captures the state of this iteration. For this reason, fac itself is not tail-recursive albeit the recursive call of fac occurs in else's tail position.

Compare this to the tail-recursive implementation in Figure 9a. It consists of a base case fac and a helper function help which implements the loop. The recursive call in help passes h_ret to help. Parameter h_ret is not a function abstraction but a parameter that is just passed around. Therefore, no closure is captured.

Concluding, it is important to distinguish tail-recursion and recursive tail-calls. Therefore, standard text book definitions for tail-recursion [e.g., 14, chap. 17.4] do not apply to CPS programs. Thus, we introduce our own nomenclature that is based on the following higher-order call graph:

1. Create a node for each Thorin function
2. For each function ℓ , draw an edge to all functions which occur in ℓ 's body.

A strongly connected component (SCC) in this graph forms a recursion. We classify:

recursive call: Any call within an SCC.

simple recursive call: A recursive call within the scope of the callee.

mutual recursive call: All other recursive calls.

static parameter: A parameter which does not change its value within an SCC.

first-order recursive call: A recursive call which only uses static parameters as higher-order arguments.

loop: An SCC formed by recursive calls that are only of zeroth order.

Example. Reconsider (Figure 9a). The call of help in then2 is simple recursive since then2 belong's to help's scope. The call of help in else is *not* recursive as it is not part of the SCC formed by help and then2. The parameter h_ret is static in that SCC. This makes the call of help in then2 first-order recursive. The call of fac in else is not first-order recursive as the call passes the higher-order argument cont (Figure 8b). Note how the definition of first-order recursion reflects non-CPS tail-recursion.

4.2.1 Mangling Simple Recursion

Let us drop help into fac by dropping help's parameters h_n and h_ret with fac's parameters n and f_ret. However, our algorithm would not touch the recursive call within help's scope because we do not put the entry $\ell_e = \text{help}$ into the map M . We would obtain a dropped version of help, say help_d, which still calls the original function help. Just substituting ℓ_e with the new mangled function would create an ill-typed call site: help_d(i+1, r*i, h_n, h_ret). The recursive call help(i+1, r*i, h_n, h_ret) uses as third

argument `h.n` and fourth argument `h.ret`. These are `help`'s static parameters. For this reason, we replace this call with `help_d(i+1, r*i)` (Figure 9b). Moreover, we replace the call site `help(1, 2, n, f.ret)` in `else` nested inside `fac` with `help_d(1, 2)` as holds:

```
help(1, 2, n, f.ret) = help(1, 2, h.n, h.ret) = help_d(1, 2) .
```

Note that the resulting program is identical to the iterative implementation (Figure 2d). In other words, with lambda mangling we performed *tail-recursion elimination* by transforming `help` to a loop `help_d`.

When we drop all parameters of the resulting function `help_d` with 1 and 2, we perform *loop peeling* (see Figure 9c). In this case we cannot substitute the recursive call `help_d(i+1, i*r)` by `help_p()` as the arguments are not static parameters.

Based on the program in Figure 9b, we can also drop all parameters of `help_d` with `i+1` and `r*i`. This performs *loop unrolling* (Figure 9d). For the same reason as above, we cannot substitute the recursive call `help_d(i+1, i*r)` by calling the new function `help_u()`, either.

In general, we substitute calls to the entry function by calls to the mangled one if all dropped parameters are static. In the case of simple recursion, we can take this insight into account by checking for this pattern when rewriting the body in Algorithm 1.

4.2.2 Mangling Mutual Recursion

In Figure 10a functions `iseven` and `isodd` invoke each other in a mutual recursive way. There exists only a sole user of this construct: Function `foo` calls `iseven`. For performance reasons, we would like to drop `iseven` and `isodd` into `foo` by substituting `eret/oret` with `foo`'s parameter `ret`. Currently, our mangling algorithm does not support this transformation: When dropping `iseven` to `iseven'` we do not know that we are going to analogously drop `isodd`, too, and thus, cannot substitute the recursive call `isodd(ei-1, eret)` with `isodd'(ei'-1)`. Likewise, we cannot replace the call `iseven(oi-1, oret)` with `iseven'(oi'-1)` when dropping `isodd`.

However, we know beforehand the mapping $M_e := \{ei \mapsto ei', eret \mapsto ret\}$ when dropping `iseven` to `iseven'` and the mapping $M_o := \{oi \mapsto oi', oret \mapsto ret\}$ when dropping `isodd` to `isodd'`. When considering this during mangling, we can directly substitute any call of the form `iseven(X, eret)` with `iseven'(X)` and any call of the form `isodd(Y, oret)` with `isodd'(Y)` as both `eret` and `oret` are static parameters in the SCC formed by `iseven`, `ethen`, `isodd` and `othern`. Finally, the program in Figure 10b emerges. Note that the recursive calls now form a loop. Thus, we performed *mutual tail-recursion elimination*.

5. Code Generation

In order to translate Thorin programs to a lower-level program representation like machine code or an SSA-based rep-

```
foo(i: int, ret: fn(bool)):      foo(i: int, ret: fn(bool)):
  iseven(i, ret)                  iseven'(i)

iseven(ei: int, eret: fn(bool)): iseven'(ei': int):
  branch(ei>0, ethen, eelse)     branch(ei'>0, ethen', eelse')
ethen():                          ethen'():
  isodd(ei-1, eret)              isodd'(ei'-1)
eelse():                          eelse'():
  eret(true)                     ret(true)

isodd(oi: int, oret: fn(bool)):  isodd'(oi': int):
  branch(oi>0, othen, oelse)     branch(oi'>0, ethen, oelse)
othern():                          othern'():
  iseven(oi-1, oret)             iseven'(oi'-1)
oelse():                          oelse'():
  oret(false)                    ret(false)
```

(a) Functions `iseven` and `isodd` are first-order recursive. (b) The optimized version consists of a loop.

Figure 10. Lambda mangling to eliminate mutual tail-recursion formed by `iseven`, `ethen`, `isodd` and `othern`.

resentation (like LLVM) we classify Thorin functions in the following way:

basic block (BB)-like function: A first-order function.

returning function: A second-order function with exactly one first-order parameter.

top-level function: A function whose scope does not contain free variables.

bad function: A function that is neither BB-like, nor top-level and returning.

control flow form: A scope that does not use bad functions.

Example. Reconsider Figure 2: All basic blocks in the SSA-form version are BB-like functions in Thorin. The function `fac` is returning and top-level. As `fac`'s scope does not contain any bad functions the scope is in CFF.

5.1 Converting Thorin to SSA Form

It is straightforward to generate code from a CFF program. For example, we can use Kelsey's algorithm [16] to translate the program to SSA form:

- All returning functions become ordinary functions in the SSA-form program. This first-order parameter acts as "return".
- All BB-like functions become basic blocks. For each parameter we introduce a ϕ -function. The corresponding arguments of the BB's predecessors determine the ϕ -function's arguments.
- Calls to returning functions become "normal" calls. The parameter of the call's continuation forms the result value in the SSA-form program.

Example. Using this algorithm we obtain Figure 2b from Figure 2d.

5.2 Beyond Control Flow Form

The remaining top-level functions that are not in CFF, can be translated with standard code generation techniques for higher-order CPS programs [1, 3, 26]. An alternative is to eliminate bad functions. Algorithm 2 sketches how to keep

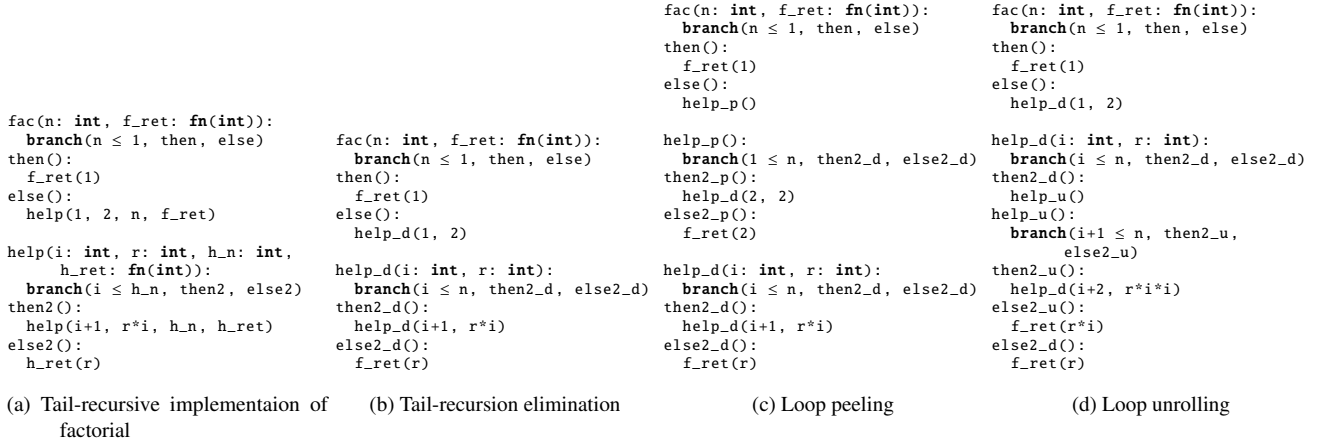


Figure 9. Tail-recursion elimination, loop peeling and loop unrolling by using lambda mangling

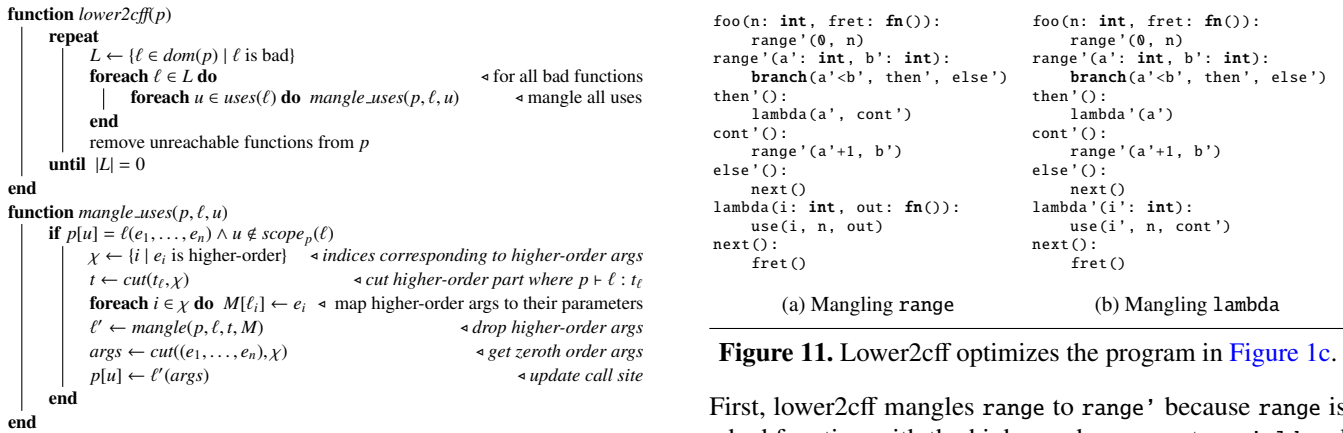


Figure 11. Lower2cff optimizes the program in Figure 1c.

Algorithm 2. This algorithm eliminates bad functions in program p .

specializing higher-order arguments to bad functions. This routine will only terminate if all bad functions eventually disappear. In general, it is undecidable whether this will be the case for an arbitrary program p : Basically, each lambda dropping performs a function specialization. This in turn partially evaluates the program. Due to the halting problem, we cannot statically prove that all bad functions will eventually be lowered.

However, the presented lowering algorithm always reduces non-recursive calls of bad functions. Each specialized call removes one use of the bad function. This property still holds when mangling first-order recursive functions because the mangled version does not reference the original function anymore (see Section 4.2). We call such programs *CFF-convertible*. When mangling recursive but not first-order recursive functions, those references usually remain. In this case, specializing the call does not decrease the number of uses of the bad function.

Example. The program in Figure 1c uses the HOF range to iterate over an interval $[a, b]$. We now apply lower2cff.

First, lower2cff mangles range to range' because range is a bad function with the higher-order parameters yield and ret. Note how mangling eliminates the static parameters (Figure 11a). The returning function lambda with the higher-order parameter out is bad since it is nested inside of f due to its dependency on f's parameter n. For this reason, lower2cff mangles lambda and the program in Figure 11b emerges which is in CFF and the algorithm terminates. Using lambda mangling, we also inline lambda' and specialize range''s parameter b' to n to finally obtain the program in Figure 1d.

Enhancements. The presented lowering algorithm always uses lambda dropping. In certain cases it is worthwhile to use lambda lifting instead as lifting will not increase the code size. However, lambda lifting is only reasonable in the case of returning, non-top-level functions which only contain free variables of zeroth order (lambda in Figure 11b, for example). Then, the lifted function will be a returning, top-level function. Otherwise, when lifting free variables of order one or higher, the resulting function is not a returning function anymore.

Moreover, lower2cff always specializes all higher-order arguments. However, in order to use Kelsey's algorithm it suffices to have top-level, returning functions. This means, that it is fine to keep one first-order, potentially non-static parameter (the "return") of top-level bad functions. The property of first-order recursion must then only hold for the

	C	Impala	Rust	GHC
aobench	1.220	1.357	n/a	22.540
fannkuch-redux	27.137	28.070	n/a	34.670
fasta	2.313	1.517	n/a	1.443
mandelbrot	2.143	2.113	n/a	2.013
meteor-contest	0.047	0.043	0.050	0.327
n-body	5.497	6.130	5.163	6.867
pidigits	0.710	0.763	4.940	0.903
regex	6.477	6.470	18.020	7.720
reverse-complement	1.090	1.220	n/a	1.300
spectral-norm	4.423	4.480	n/a	19.347

Table 1. Median execution times of eleven runs in seconds

remaining parameters. This essentially allows to transform a greater class of programs to CFF.

Summary. Compilation works as follows:

1. Identify all bad functions.
2. Use `lower2cff` to eliminate all non-recursive and first-order recursive uses. For such uses `lower2cff` will terminate. If all uses are of these kinds, the original bad function will become unreachable and can be removed.
3. Translate all functions which are not bad with [Kelsey’s](#) algorithm.
4. Translate remaining functions with conservative code generation techniques.

6. Implementation and Evaluation

6.1 Performance

We implemented a language called *Impala*²—a dialect of Rust³. Besides semantic analyses the Impala compiler does not perform any further analyses or transformations on the abstract syntax tree (AST). As Rust directly compiles to LLVM, the compiler implements many complicated analyses and transformations before translating to LLVM as pointed out in [Section 1](#). The Impala compiler directly translates the AST—including higher-order, nested and/or polymorphic functions—to Thorin⁴. Thorin then performs its own set of optimizations including the presented `lower2cff` phase to eliminate closures before translating the Thorin program with [Kelsey’s](#) algorithm to LLVM.

In order to evaluate the effectiveness of our approach, we ported *The Computer Benchmark Game*⁵ to Impala. We elided some programs which focus on measuring API or runtime overhead. Additionally, we ported `aobench`⁶. [Table 1](#) reports the median execution times of eleven runs for C, Impala, Rust (if available) and GHC in seconds. Our test ran on an Intel® Ivy Bridge Core™ i7-3770K CPU. We consider the C implementations as baseline. We included Rust in our measurements as Impala has similar syntax. We also included Haskell versions of the benchmarks in order to see how well

²<https://github.com/AnyDSL/impala>

³<http://www.rust-lang.org>

⁴<https://github.com/AnyDSL/thorin>

⁵<http://benchmarksgame.alioth.debian.org/>

⁶<https://code.google.com/p/aobench>

	SLoC	Volume	Difficulty	Effort
<code>CloneFunction.cpp</code>	359	21298	107	2269693
<code>CodeExtractor.cpp</code>	523	34599	124	4287983
<code>InlineFunction.cpp</code>	526	32288	109	3511359
<code>LoopUnroll.cpp</code>	279	15393	80	1229623
total	1687	120421	207	24968883
<code>mangle.cpp</code>	132	6636	75	496757

Table 2. Source lines of code (SLoC) and Halstead numbers for LLVM’s C++ implementations compared to Thorin’s mangle implementation.

GHC’s Core IR performs. From the original benchmark suite we selected programs which were neither hand-vectorized nor hand-parallelized. We used `clang 3.4.2`, `rustc 0.11` and `GHC 7.8.3`. The exact compile flags and benchmarks are apparent in our benchmark suite.⁷

Although we used the C versions of the benchmarks as a template for the Impala version, it is important to note that Impala does not offer C-style `for`-loops. Instead, we use HOFs to write appropriate generators and rely on Thorin’s `lower2cff` phase. For example, in order to iterate over an interval we use a HOF range as presented in [Section 5.2](#).

The performance of the Impala programs is mostly on a par with the C implementations except for `fasta` where Impala is about 1.5 times faster than C. We are still investigating this slowdown for C. Rust’s performance stands or falls by the quality of the used libraries. GHC is for the most part roughly on a par with C/Impala. However, some benchmarks—in particular `aobench`—run significantly slower.

6.2 Engineering Effort

In order to estimate the engineering effort to develop code transformations we compare the Halstead metric [[13](#)] of Thorin’s lambda mangling implementation versus LLVM 3.4.2 ([Table 2](#)).⁸ On the one hand, LLVM is a full featured compiler suite which makes these metrics biased towards Thorin. On the other hand, lambda mangling is much more versatile: For LLVM we did not include any source code to eliminate tail-recursion ([Section 4.2](#)). Furthermore, LLVM completely lacks functionality for *partial* inlining, *partial* outlining or to represent HOFs at all. Still, LLVM’s pendants are in total roughly 2.8 times more difficult to implement while it takes about 50 times more time to program.

7. Discussion and Related Work

Graph-Based IRs and CPS. Graph-based IRs for imperative languages [like [5](#), [7](#), [22](#)] are not designed for functional aspects and do not natively support HOFs and closures as Thorin does. On the other hand, higher-order CPS-based IRs used in compilers for functional languages use explicit scope nesting and an AST rather than a graph in contrast

⁷<https://github.com/AnyDSL/benchmarks-impala>

⁸SLoC were generated using David A. Wheeler’s ‘SLOCCount’; Halstead numbers were computed with `c3ms` [[12](#)].

to Thorin [1, 3, 26]. We discussed that graphs instead of names and implicit scope nesting simplify program transformations (Section 2 and 4). An alternative to CPS is A-normal form (ANF) [11] or a monadic language [21]. These languages do not use CPS but are in the spirit of CPS by binding each new temporary to a new name. However, as Kennedy [17] points out, such languages introduce problems not present in a faithful CPS presentation: For example, ANF needs a re-normalization phase after β -reduction.

SSA Form. SSA form was invented by Rosen et al. [24] and became popular after Cytron et al. [8] described an efficient algorithm for computing SSA form. As already outlined, SSA form is a restricted form of CPS. In fact, Impala exploits an SSA construction algorithm [6] to directly construct a Thorin program from the AST without further analyses like the computation of a dominance tree. The scope of a top-level function in CFF is akin to an SSA-form program because Kelsey’s algorithm is only of syntactic nature (Section 5.1). However, in contrast to classic SSA programs Thorin programs can be higher-order. Even classic imperative languages would profit from a higher-order IR. It is often important to annotate code regions (for example, to mark a loop for parallelization). LLVM uses metadata to annotate code. This introduces many subtle problems. For instance, transformations must pay attention to not mistakenly destroy metadata designated for a different pass. HOFs solve this problem in a clean and type-safe way by wrapping the code region to be annotated in a HOF. Köster et al. [18] use this technique to annotate code regions for SIMD vectorization and GPU execution.

Lambda Lifting and Dropping. Lambda lifting was invented by Johnsson [15]. His algorithm uses currying in order to abstract free variables:

$$\lambda x. \dots y \dots \Rightarrow \lambda y. \lambda x. \dots$$

Danvy and Schultz [9] observe that we can directly append y to the parameter list in the case of first-order programs with n -ary functions:

$$\lambda(x). \dots y \dots \Rightarrow \lambda(x, y). \dots$$

They invented lambda dropping as reverse transformation to Johnsson’s algorithm. When dealing with HOFs, their algorithm uses Johnsson’s currying approach since in general a compiler cannot rewrite call sites of a function passed as argument to another function. As we want to eliminate closures, currying is not an option for us. After all, a curried function is not a returning function anymore (Section 5). For this reason, lambda mangling does not use currying to introduce or eliminate parameters. It rather produces one new generalized and/or specialized function with an updated signature which is at that point not connected to the rest of the program. It is in the responsibility of other passes to orchestrate mangling in a reasonable way and to connect a mangled function to the rest of the program properly. Due to

Thorin’s blockless representation we can fuse both algorithms into one simple recursive reconstruction algorithm which does not need the block floating/sinking pass of the original algorithms.

Static Argument Transformation. We are not the first to note that an argument/parameter pair of a recursive call is superfluous if the argument is just the parameter. The *static argument transformation* [10] identifies such recursive calls and eliminates them in the following way:

$$\begin{array}{l}
 f: (a, b) \\
 \dots \text{ use}(a) \dots \\
 f(a, x)
 \end{array}
 \Rightarrow
 \begin{array}{l}
 f: (a, b) \\
 \text{letrec } f': (b) \\
 \dots \text{ use}(a) \dots \\
 f'(x) \\
 \text{in } f'(b)
 \end{array}$$

Note that f itself is not recursive anymore and, hence, is a potential candidate for inlining. For mutual recursive functions, we must find out which parameters do not change their values within an SCC. For this reason, we rather speak of static *parameters* than of *arguments*. However, applying the static argument transformation on a mutual recursive function still leaves the function recursive whereby further inlining becomes problematic. Our algorithm on the other hand considers all functions within one SCC simultaneously and, thus, can eliminate static parameters (Section 4.2).

Super- β inlining. Super- β inlining enables inlining of a closure call c with a function literal f . This transformation is only valid if (1) all applications of c are closures over f and (2) the environment at c is always equivalent to the one where the closure is captured. A control flow analysis (CFA) [19, 25] addresses the first condition. Δ -CFA [20] also takes the environment where a closure is captured into account. This allows aggressive inlining and, thus, aggressive closure elimination. The focus in our work is different. The presented lower2cff algorithm specializes (in contrast to inlining) HOFs even more aggressively. Consider the function range from Figure 1. Suppose, there are two call-sites of range which pass different function literals g_1 and g_2 to range’s parameter f :

```
range(..., g1, ...)
range(..., g2, ...)
```

A Δ -CFA discovers that inlining the closure call of f within range’s scope is not possible since condition (1) is violated. Thorin’s lower2cff on the other hand, specializes each call-site and, hence, gets rid of the HOF range (Section 5.2).

Other Closure-Elimination Strategies. The SML/NJ compiler [3] allocates garbage-collected closures on the heap. An additional phase tries to find closures whose lifetime can be statically proven to be nested in a simple way. These closures can be allocated on the stack instead. Additionally, the compiler lambda-lifts functions that are not passed as argument to other functions. Then, these so-called *known* functions (as opposed to *escaping* functions) do not require a closure at all. An additional η -split phase splits functions that are used in a known as well as in an escaping context. Finally, SML/NJ’s

inliner tries to decrease *escaping* contexts. Other compilers like Scala⁹ also rely on inlining to eliminate explicit closures.

In summary, state-of-the-art functional compilers rely on *inline heuristics* to eliminate closures. Thorin *guarantees* to eliminate all closures in programs that are CFF-convertible. As future work we want to design a type system that captures this property in order to inform the programmer whether his higher-order program will need explicit closures.

8. Conclusions

We presented the functional, CPS-based IR Thorin. Thorin is equally suited to represent imperative as well as functional programs. Its novelty is that it does not use scope nesting to associate the use of a value with its definition, but is graph-based. We built on this property to devise *lambda mangling*, which generalizes partial in- and outlining. Lambda mangling serves as a building block for other important control-flow transformations, especially closure elimination. Finally, we presented a novel closure elimination algorithm on the basis of lambda mangling.

We evaluate Thorin experimentally using our research language Impala, a dialect of Rust, on *The Computer Benchmark Game*, a program suite to benchmark performance across several languages (and their implementations). Impala uses HOFs and closures to represent loop constructs in a generic way. Nevertheless, Thorin’s closure elimination algorithm successfully removes all overhead of closure allocation on this benchmark set and produces programs that match the performance of C.

Acknowledgments

This work is supported by the Federal Ministry of Education and Research (BMBF) as part of the ECOUSS project as well as by the Intel Visual Computing Institute Saarbrücken.

References

- [1] N. Adams, D. Kranz, R. Kelsey, J. Rees, P. Hudak, and J. Philbin. Orbit: an optimizing compiler for scheme. In *Symposium on Compiler Construction*, SIGPLAN, 1986.
- [2] A. W. Appel. SSA is functional programming. *SIGPLAN Not.*, 1998.
- [3] A. W. Appel. *Compiling with Continuations (corr. version)*. Cambridge University Press, 2006.
- [4] H. Barendregt. *The Lambda Calculus: Its Syntax and Semantics*. Elsevier Science, 1985.
- [5] M. Braun, S. Buchwald, and A. Zwinkau. Firm—a graph-based intermediate representation. Technical report, 2011.
- [6] M. Braun, S. Buchwald, S. Hack, R. Leiða, C. Mallon, and A. Zwinkau. Simple and efficient construction of static single assignment form. In *CC*, 2013.
- [7] C. N. Click and Jr. *Combining Analyses, Combining Optimizations*. PhD thesis, Rice University, 1995.
- [8] R. Cytron, J. Ferrante, B. K. Rosen, M. N. Wegman, and F. K. Zadeck. Efficiently computing static single assignment form and the control dependence graph. *TOPLAS*, 1991.
- [9] O. Danvy and U. P. Schultz. Lambda-dropping: Transforming recursive equations into programs with block structure, 2001. This is an extended version of Danvy’s and Schultz’ original paper of the same title which appeared at PEPM’97.
- [10] A. L. de M. Santos. *Compilation by Transformation in Non-Strict Functional Languages*. PhD thesis, University of Glasgow, 1995.
- [11] C. Flanagan, A. Sabry, B. F. Duba, and M. Felleisen. The essence of compiling with continuations. In *PLDI*, 1993.
- [12] C. H. Gonzalez and B. B. Fraguera. A generic algorithm template for divide-and-conquer in multicore systems. In *HPCCC*, 2010.
- [13] M. H. Halstead. *Elements of Software Science (Operating and Programming Systems Series)*. Elsevier Science Inc., 1977.
- [14] M. R. Hansen and H. Rischel. *Introduction to Programming using SML*. Addison Wesley Longman, 1999.
- [15] T. Johnsson. Lambda lifting: Transforming programs to recursive equations. In *Functional programming languages and computer architecture*, 1985.
- [16] R. A. Kelsey. A correspondence between continuation passing style and static single assignment form. *SIGPLAN Not.*, 1995.
- [17] A. Kennedy. Compiling with continuations, continued. In *ICFP*, 2007.
- [18] M. Köster, R. Leiða, S. Hack, R. Membarth, and P. Slusallek. Code refinement of stencil codes. *PPL*, 24, 2014.
- [19] J. Midtgaard. Control-flow analysis of functional programs. *ACM Computing Surveys*, 2012.
- [20] M. Might and O. Shivers. Environment analysis via Δ -CFA. In *POPL*, 2006.
- [21] E. Moggi. Notions of computation and monads. *IC*, 1991.
- [22] M. Paleczny, C. Vick, and C. Click. The Java HotSpot™ server compiler. In *JVM*, 2001.
- [23] S. Peyton Jones and S. Marlow. Secrets of the glasgow Haskell compiler inliner. *Journal of Functional Programming*, 2002.
- [24] B. K. Rosen, M. N. Wegman, and F. K. Zadeck. Global value numbers and redundant computations. *POPL*, 1988.
- [25] O. Shivers. *Control-Flow Analysis of Higher-Order Languages*. PhD thesis, Carnegie Mellon University, 1991.
- [26] G. L. Steele, Jr. Rabbit: A compiler for scheme. Technical report, 1978.
- [27] B. Steensgaard. Sparse functional stores for imperative programs. In *ACM SIGPLAN Workshop on Intermediate Representations*, pages 62–70, 1995.
- [28] C. Strachey. Fundamental concepts in programming languages. *Higher-Order and Symbolic Computation*, 2000.

⁹<http://magarciaepfl.github.io/scala>