

Catching fire: Undefined behaviour for equational reasoning of lazy languages with strictness annotations

Subtitle

Anonymous Author(s)

Abstract

We propose a new style of semantics for non-strict languages that provide access to many optimizations that are currently unsound for lazy languages. Lazy languages, while providing equational reasoning, are unable to reason effectively about strictness annotations. On the other hand, most strict languages are impure, and thus unable to provide equational reasoning due to the presence of side effects. Research on optimizing non-strict languages has focused on optimizing away non-strictness with demands, and efficient lowering of strictness onto modern hardware. We take a different approach: we propose to use undefined behaviour judiciously while defining the semantics of our language to enable many key optimizations that are possible in the strict world, that are not possible in the non-strict world. We model divergent computations as undefined behaviour, giving the compiler far more freedom to optimize and reorder a mix of strict and non-strict code.

1 Introduction

Lazy languages such as Haskell offer many benefits, chief of which is the ability to reason equationally about programs. For performance, considerations, the prototypical lazy language, Haskell, provides “bang patterns” to mark strictness. Herein lies the trouble; these “bang patterns” do not permit equational reasoning for the user. Worse, they form a sequence point in the otherwise pure semantics, taking away most of the optimization freedom from the compiler. We propose to remedy the situation by making judicious use of undefined behaviour semantics. While undefined behaviour is often reviled in discourse, it is in fact a major reason why C and C++ can be optimized as well as they can; defining certain rare or impossible conditions as undefined behaviour provides the compiler a great deal of freedom by being able to ignore these cases. We choose to mark divergence as undefined behaviour. This solves the correctness problem in reordering strict computations; The only thing our compiler (Lizzy) will reason about is performance when reordering strict computations. Furthermore, this is not too insane a proposal; After all, the C standard also defines side-effect free infinite loops as undefined behaviour. Next, we show a variety of examples where this simple extension to the language semantics permits far easier reasoning about mixing

strictness and non-strictness. We then put this to the test, by extending MLIR with a new non-strict IR, called `lz`. Note that the undefined behaviour semantics are critical in order to reason effectively about the mix of strict and non strict behaviour we wish to model within MLIR. Finally, to validate our approach, we consider a series of benchmarks where our programming language, Lizzy outperforms similar programs in Haskell compiled by the Glasgow Haskell Compiler.

Our contributions are:

- A new dialect (`lz`) for MLIR to encode these semantics, and provide modular non-strictness.
- A translator from STG to `lz` that exhibits the ability to eliminate laziness and leverage MLIR’s loop optimisation framework.
- A rewrite driven strategy for worker/wrapper that interleaves worker/wrapper analysis along with worker/wrapper rewriting, by phrasing the worker/wrapper problem as an outlining/inlining problem.

2 Stuff that doesn’t work

- circular references: SSA by default cannot handle circular references. We need to use the “relaxed SSA” that allows graphs. It’s unclear how well this adapts.
- A garbage collector. The usual story here, we don’t really make a fair comparison.
- Any realistic program exhibiting “difficult” laziness. I feel what would be interesting would be to implement, say, GRIN, along with a sophisticated pointer analysis using the MLIR analysis-as-rewrite story here.
- A proof of correctness of the naive demand analysis I wrote.

3 Stuff that does work

- An encoding of laziness that works for the simple stuff I’ve tried
- A simple version of “demand analysis” that eliminates first order laziness.
- Some toy examples where we generate better IR than GHC because we can expose the “real work” to LLVM/MLIR in a way that they understand, by eliminating the noise that is laziness.

4 What I feel would make a good paper/thesis

- Our story of performing worker/wrapper by outlining/inlining is intriguing. I don't think it's strong. So, to remedy that, we can:
- (1) Make our laziness encoding bulletproof, connect to CBPV for example, handle circular references, etc.
- (2) Make a strong case that MLIR is an interesting target for this breed of language. So, perform unification in MLIR, typeclass resolution in MLIR, perform some analysis/optimisation. This way, we show an end-to-end prototype of something that sorta works, instead of one thing that sorta works (?)

TODO: fill in semantics here.

5 lz : Its Syntax, Type system, semantics

5.1 A high level overview of MLIR

5.2 A high level overview of lz

We extend the basic std dialect modestly, by adding primitives that can describe non-strictness. We first introduce a primitive known as `lz.thunk`, which is responsible for creating non-strict thunks. To represent lazy function application, we introduce a new operation called `lz.ap`, which is the lazy sibling of `std.call`. This creates a thunk, which when evaluated, invokes the function with the given arguments. To force a thunk, we introduce `lz.force`, which receives a thunk as input and returns the forced value of the thunk as output. In effect, we have an encoding of administrative normal form as (ANF) an MLIR dialect.

To represent algebraic data types, which enable reasoning of control flow and are the primary mode of abstraction in most functional languages, as well as newer imperative languages such as Rust, we introduce the `lz.construct` and case operators, which are the introduction and elimination forms of structured data in the dialect.

TODO: this is broken, because it doesn't scope properly. We need to describe `lz.thunk` as building a thunk in memory/on the heap, and the value stored in the register file is the pointer to this thunk.

5.3 Purifying `lz.force`

Let us for a moment assume that divergence is not undefined behaviour. Now consider the code in ???. The variable called `%loopv` is unused. However, because divergence `%loopv` is divergent, it is incorrect to eliminate the variable definition of `%loopv`, for the divergence of `loop` occurs when `%loopt` is forced; That is, when the instruction `%loopv = lz.force(%loopv)` is executed. This implies that the instruction `lz.force` is side-effecting.

If we wished to regain purity of `lz.force` and leverage the power of SSA, then we must be allowed to eliminate a call to `lz.force` if the value is unused. This is exactly the same

requirement as being able to equationally reason with bangs in a `let` binding. In that use-case, we wished for undefined behaviour semantics for ease-of-reasoning. Here, we wish for undefined behaviour semantics for ease-of-optimization.

```
1 // loop = loop
2 lz.func @loop () -> !lz.adtBox> lz.return(// main =
    let !x = loop in 42lz.func @main() ->
    !lz.adt<@Int> : !lz.thunk<lz.adt<@Box>>:
    !lz.adt<@Box> lz.return(
```

Listing 1. The SSA encoding of a lazy program with divergence in the `lz` dialect. Note the unused divergent variable `%loopv`

```
1 // loop = loop
2 lz.func @loop () -> !lz.adtBox> lz.return(// main =
    let !x = loop in 42lz.func @main() ->
    !lz.adt<@Int> // // : !lz.thunk<lz.adt<@Box>>//
    // : !lz.adt<@Box> lz.return(
```

Listing 2. Under the assumption that divergence is UB we can eliminate the call to `lz.force`. This regains SSA semantics.

5.4 Purifying `lz.construct`

Consider the program in ???. For us to eliminate the call to `lz.construct`, we need to be sure that memory allocation has no side effects. In our case, since we assume that allocations are not visible to the user in terms of their implementation details (eg. The address of the pointer of the block of memory), we can safely remove the value `%unused = lz.construct(@Box)` as `%unused` is unused.

On the flipside, consider a program where we want to copy the instruction `lz.construct`. If we were not guaranteed the existence of a garbage collector, then this would not be pure; we would have created a chunk of memory that is never freed. **TODO: Is there a good example where we want to copy a value for optimization? Something like, copy an SSA value both into a loop and outside or something?**

```
1 lz.func() -> !hask.adtBox> lz.ret(
```

Listing 3. An unused called to `lz.construct`

6 Worker wrapper by local rewrites: A first stab

In this section, we consider a small example program, which we will explain how to perform a classical transform (the worker-wrapper transform) purely by using local rewrites. Traditionally, such a transformation is performed by first performing a demand analysis whose results are used to drive the rewrite.

We explore how this transformation can be achieved without the need for demand analysis by performing local rewrites. This provides both a simpler manifestation of the algorithm,

```

221 toplevel := func | global
222 func := <fn-name> formal-param-list -> <ret-type> <region>
223 region := list [bb]
224 bb := <bb-name> formal-param-list : list inst; terminator-inst
225 inst := retval "=" <op-name> arg-list
226 terminator-inst := std.return(<name>) | br <name> | condbr <name> <bb1> <param-list> <bb2> <param-list>

```

Figure 1. The std dialect syntax

$$\begin{array}{c}
\frac{(\%vid = \text{std.constant } const, R) \xrightarrow{\text{asgn}} R[\%vid \mapsto const]}{(\%vid = f(x_1, x_2, \dots, x_n), R) \xrightarrow{\text{asgn}} R[\%vid \mapsto y]} \quad \frac{R[x_1] = v_1 \quad R[x_2] = v_2 \quad \dots \quad R[x_n] = v_n \quad \llbracket f \rrbracket(v_1, v_2, \dots, v_n) = y}{(\%vid = f(x_1, x_2, \dots, x_n), R) \xrightarrow{\text{asgn}} R[\%vid \mapsto y]} \\
\\
\frac{P[pc] = \text{condbr } c \text{ } ^{bb1}(xs) \text{ } ^{bb2}(ys) \quad P[^{bb1}] = ^{bb1}(ls) \quad R[c] \neq 0}{(pc, R) \xrightarrow{\text{ctrl}} (pc', R[ls \mapsto R[xs]])} \quad \frac{P[pc] = \text{br } c \text{ } l(v_1) \text{ } r(v_2) \quad P[r] = pc' \quad R[c] = 0}{(pc, R) \xrightarrow{\text{ctrl}} (pc', R[r.input \mapsto R[v_2]])} \\
\\
\frac{P[pc] = \text{br } n(v) \quad P[n] = pc'}{(pc, R) \xrightarrow{\text{ctrl}} (pc', R[n.input \mapsto R[v]])}
\end{array}$$

Figure 2. The operational semantics of the std dialect

```

250 thunk(x): T -> Thunk<T>
251 ap(f, v1, v2, ..., vn): ((T1,...Tn) -> R) x T1 x ... Tn -> Thunk<R>
252 force(x): Thunk<T> -> T
253 construct(ConsName, v1, ..., vn): Symbol x T1 x T2 ... x Tn -> ADT<T>
254 case(x: T) of ... : ADT<T> -> R

```

Figure 3. Syntax of lz extensions

as well as a potentially faster implementation, as the MLIR infrastructure is capable of performing rewrites in parallel.

```

261 1 SimpleInt f(Thunk<SimpleInt> i) {
262 2     SimpleInt icons = force(i);
263 3     int ihash = casedefault/icons);
264 4     if (ihash <= 0) {
265 5         return SimpleInt(42);
266 6     } else {
267 7         int prev = ihash - 1;
268 8         SimpleInt siprev = SimpleInt(prev);
269 9         Thunk<SimpleInt> siprev_t = thunkify(siprev);
270 10        SimpleInt f_prev_v = apStrict(f, siprev_t);
271 11        return f_prev_v;
272 12    }
273 13 }
274 14

```

```

15 int main() {
16     printf("%d\n", f(thunkify(SimpleInt(1))).v);
17 }

```

Listing 4. initial source code

```

1 SimpleInt f2(SimpleInt icons) {
2     int ihash = casedefault/icons);
3     if (ihash <= 0) {
4         return SimpleInt(42);
5     } else {
6         int prev = ihash - 1;
7         SimpleInt siprev = SimpleInt(prev);
8         Thunk<SimpleInt> siprev_t = thunkify(siprev);
9         SimpleInt f_prev_v = f(siprev_t);
10        return f_prev_v;
11    }

```

331	$P[pc] =$	$P[pc] = \%vid = lz.ap(\%fref, \%x1, \dots, \%xn)$	386
332	$\%vid = lz.construct(@ConsName, \%x1, \dots, \%xn)$	$P[\%f1] = ref(@fn)$	387
333	$P[\%x1] = v_1, \dots, P[\%xn] = v_n$	$P[\%x1] = v_1, \dots, P[\%xn] = v_n$	388
334	\hline	\hline	389
335	$R[\%vid \mapsto constructor(@ConsName, v_1, \dots, v_n)]$	$R[\%vid \mapsto thunk(@fn, v_1, \dots, v_n)]$	390
336			391
337			392
338	$P[pc] = \%vid = lz.case(\%x, @Cons_1, r_1, \dots, @Cons_N, r_N)$		393
339	$@Cons_i = @XCons$	$p[pc] = \%vid = lz.force(\%thnk)$	394
340	$P[\%x] = construct(@XCons, v_1, \dots, v_n)$	$p[\%thnk] = thunk(@fn, v_1, \dots, v_n)$	395
341	$r_i = \{ ^{entry}(arg1, \dots, argn): \dots \}$	$\llbracket @fn \rrbracket (v_1, v_2, \dots, v_n) = y$	396
342	\hline	\hline	397
343	$(pc, R) \xrightarrow{ctrl} (P[r_i], R[arg_i \mapsto v_i])$	$r[\%vid \mapsto y]$	398
344			399

Figure 4. Operational semantics of lz extensions

```

348 12 }
349 13 SimpleInt f(Thunk<SimpleInt> i) {
350 14     SimpleInt icons = force(i);
351 15     SimpleInt ret = f2(Icons);
352 16     return ret;
353 17 }
354 Listing 5. Step 1: outlining everything after the initial
355 lz.forceinto f2
356
357 1 SimpleInt f2(SimpleInt icons) {
358 2     int ihash = casedefault(Icons);
359 3     if (ihash <= 0) {
360 4         return SimpleInt(42);
361 5     } else {
362 6         int prev = ihash - 1;
363 7         SimpleInt siprev = SimpleInt(prev);
364 8         Thunk<SimpleInt> siprev_t = thunkify(siprev);
365 9         SimpleInt f_prev_v = f(siprev_t);
366 10        Thunk<SimpleInt> f_prev_v = f2(siprev);
367 11        return f_prev_v;
368 12    }
369 13 }
370 14
371 15 SimpleInt f(Thunk<SimpleInt> i) {
372 16     SimpleInt icons = force(i);
373 17     SimpleInt ret = f2(Icons);
374 18     return ret;
375 19 }

```

Listing 6. Step 2: Removing the un-necessary lazy recursive call

```

379 1 SimpleInt f3(int icons) {
380 2     if (ihash <= 0) {
381 3         return SimpleInt(42);
382 4     } else {
383 5         int prev = ihash - 1;
384 6         SimpleInt siprev = SimpleInt(prev);
385

```

```

7     Thunk<SimpleInt> f_prev_v = f2(siprev);
8     return f_prev_v;
9 }
10 }
11
12 SimpleInt f2(SimpleInt icons) {
13     int ihash = casedefault(Icons);
14     SimpleInt ret = f3(ihash);
15
16     return ret;
17 }
18 ...

```

Listing 7. Step 3: outline everything after casedefault into f3

```

1 SimpleInt f3(int icons) {
2     if (ihash <= 0) {
3         return SimpleInt(42);
4     } else {
5         int prev = ihash - 1;
6         SimpleInt siprev = SimpleInt(prev);
7         Thunk<SimpleInt> f_prev_v = f2(siprev);
8         Thunk<SimpleInt> f_prev_v = f3(prev);
9         return f_prev_v;
10    }
11 }
12 ...

```

Listing 8. replace call f2(SimpleInt(prev)) to f3(prev)

```

1 SimpleInt f3(int icons) {
2     if (ihash <= 0) {
3         return SimpleInt(42);
4     } else {
5         int prev = ihash - 1;
6         Thunk<SimpleInt> f_prev_v = f3(prev);
7         return f_prev_v;
8     }

```

```

441 9 }
442 10
443 11 SimpleInt f2(SimpleInt icons) {
444 12     int ihash = casedefault.icons;
445 13     SimpleInt ret = f3(ihash);
446 14     return ret;
447 15 }
448 16 SimpleInt f(Thunk<SimpleInt> i) {
449 17     SimpleInt icons = force(i);
450 18     SimpleInt ret = f2.icons;
451 19     return ret;
452 20 }

```

Listing 9. taking stock

At this point, we have a clean program that has been worker/wrappered. We can see that the call chain looks as follows:

$$f \xrightarrow{\text{force}} f_2 \xrightarrow{\text{unwrap}} f_3 \rightarrow f_3 \dots \xrightarrow{\text{wrap}} f_2 \xrightarrow{\text{thunk}} f$$

Of course, much remains to be discussed: what about sum types? what about non-tail recursion? Extensions to the same idea will handle the above problems, while retaining the pleasing simplicity of this outlining/inlining paradigm.

7 Worker wrapper by local rewrites: Non-tail-calls

```

468 1 SimpleInt g(Thunk<SimpleInt> i) {
469 2     SimpleInt icons = force(i);
470 3     int ihash = casedefault.icons;
471 4     if (ihash <= 0) {
472 5         return SimpleInt(42);
473 6     } else {
474 7         int prev = ihash - 1;
475 8         SimpleInt siprev = SimpleInt(prev);
476 9         Thunk<SimpleInt> siprev_t = thunkify(siprev);
477 10        SimpleInt g_prev_v = g(siprev_t);
478 11        int g_prev_v_hash = casedefault.g_prev_v;
479 12        int rethash = g_prev_v_hash + 2;
480 13        SimpleInt ret = SimpleInt(rethash);
481 14        return ret;
482 15    }
483 16 }

```

Listing 10. Step 0: the initial program

```

485 1 SimpleInt g2(SimpleInt i) {
486 2     int ihash = casedefault.icons;
487 3     if (ihash <= 0) {
488 4         return SimpleInt(42);
489 5     } else {
490 6         int prev = ihash - 1;
491 7         SimpleInt siprev = SimpleInt(prev);
492 8         Thunk<SimpleInt> siprev_t = thunkify(siprev);
493 9
494 10
495 11

```

```

9     SimpleInt g_prev_v = g(siprev_t);
10    int g_prev_v_hash = casedefault.g_prev_v;
11    int rethash = g_prev_v_hash + 2;
12    SimpleInt ret = SimpleInt(rethash);
13    return ret;
14 }
15 }
16 SimpleInt g(Thunk<SimpleInt> i) {
17     SimpleInt icons = force(i);
18     g2.icons;
19 }

```

Listing 11. Step 1: outline everything after force

```

1 SimpleInt g2(SimpleInt i) {
2     int ihash = casedefault.icons;
3     if (ihash <= 0) {
4         return SimpleInt(42);
5     } else {
6         int prev = ihash - 1;
7         SimpleInt siprev = SimpleInt(prev);
8         Thunk<SimpleInt> siprev_t = thunkify(siprev);
9         SimpleInt g_prev_v = g(siprev_t);
10        int g_prev_v_hash = casedefault.g_prev_v;
11        int rethash = g_prev_v_hash + 2;
12        SimpleInt ret = SimpleInt(rethash);
13        return ret;
14    }
15 }
16 SimpleInt g(Thunk<SimpleInt> i) {
17     SimpleInt icons = force(i);
18     g2.icons;
19 }

```

Listing 12. Step 2: outline everything after force

```

1 SimpleInt g2(SimpleInt i) {
2     int ihash = casedefault.icons;
3     if (ihash <= 0) {
4         return SimpleInt(42);
5     } else {
6         int prev = ihash - 1;
7         SimpleInt siprev = SimpleInt(prev);
8         Thunk<SimpleInt> siprev_t = thunkify(siprev) ;
9         SimpleInt g_prev_v = g(siprev_t);
10        SimpleInt g_prev_v = g2(siprev);
11        int g_prev_v_hash = casedefault.g_prev_v;
12        int rethash = g_prev_v_hash + 2;
13        SimpleInt ret = SimpleInt(rethash);
14        return ret;
15    }
16 }
17 SimpleInt g(Thunk<SimpleInt> i) {
18     SimpleInt icons = force(i);
19     g2.icons;
20 }

```


551 20 }

Listing 13. Step 2: replace recursive call to g

```

552
553
554 1 SimpleInt g3(int i) {
555 2   if (ihash <= 0) {
556 3     return SimpleInt(42);
557 4   } else {
558 5     int prev = ihash - 1;
559 6     SimpleInt siprev = SimpleInt(prev);
560 7     g_prev_v = g2(siprev);
561 8     int g_prev_v_hash = casedefault(g_prev_v);
562 9     int rethash = g_prev_v_hash + 2;
563 10    SimpleInt ret = SimpleInt(rethash);
564 11    return ret;
565 12  }
566 13 }
567 14
568 15 SimpleInt g2(SimpleInt i) {
569 16   int ihash = casedefault(icons);
570 17   return g3(ihash);
571 18 }

```

Listing 14. Step 3: outline block after casedefault

```

572
573
574 1 SimpleInt g3(int i) {
575 2   if (ihash <= 0) {
576 3     return SimpleInt(42);
577 4   } else {
578 5     int prev = ihash - 1;
579 6     SimpleInt siprev = SimpleInt(prev);
580 7     g_prev_v = g2(siprev);
581 8     g_prev_v = g3(prev);
582 9     int g_prev_v_hash = casedefault(g_prev_v);
583 10    int rethash = g_prev_v_hash + 2;
584 11    SimpleInt ret = SimpleInt(rethash);
585 12    return ret;
586 13  }
587 14 }

```

Listing 15. Step 4: replace g2 with g3

```

588
589
590 1 SimpleInt g3(int ihash) {
591 2   int out;
592 3   if (ihash <= 0) {
593 4     return SimpleInt(42);
594 5     out = 42;
595 6   } else {
596 7     int prev = ihash - 1;
597 8     SimpleInt g_prev_v = g3(prev);
598 9     int g_prev_v_hash = casedefault(g_prev_v);
599 10    int rethash = g_prev_v_hash + 2;
600 11    SimpleInt ret = SimpleInt(rethash);
601 12    return ret;
602 13    out = rethash;
603 14  }
604 15  return SimpleInt(out);
605

```

16 }

Listing 16. Step 5: float out SimpleInt

```

1
2   int out;
3   if (ihash <= 0) {
4     out = 42;
5   } else {
6     int prev = ihash - 1;
7     SimpleInt g_prev_v = g3(siprev);
8     int g_prev_v_hash = casedefault(g_prev_v);
9     int rethash = g_prev_v_hash + 2;
10    out = rethash;
11  }
12  return SimpleInt(out);
13 }
14 SimpleInt g2(SimpleInt i) {
15   int ihash = casedefault(icons);
16   return g3(ihash);
17 }
18
19 SimpleInt g(Thunk<SimpleInt> i) {
20   SimpleInt icons = force(i);
21   g2(icons);
22 }

```

Listing 17. Step 5: taking stock**8 Worker wrapper by local rewrites: Sum types**

The transformation of floating out common patterns works equally well for branches as it does for sum types, since a case analysis on a sum type is the same as a branch on the tag.

9 Worker wrapper by local rewrites: Sum types + non-tail-calls**10 Demand analysis by rewrites: The full algorithm****11 Evaluation**

As a baseline, we compare our performance on nofib test-suite, which is the test suite that GHC is tested with and performance improvements to GHC are reported on. We only consider a representative subset of programs from nofib ; The full test suite also extensively tests Haskell's semantics of parallelism, software transactional memory, and other features that are orthogonal to the problem of representing and optimizing non strictness.

11.1 A toy example: eliminating wrapper overhead

```

1 module {
2 // fact 0 = 1
3 // fact n = n*fact(n-1)
4 func @f (%i : !lz.thunk<i64>) -> i64 {
5   %ival = lz.force(%i):i64
6   %out = lz.caseint %ival
7   [0 -> {
8     ^entry:

```

6

```

661 %x = ap(%f, %v1, ..., %vn): !lz.thunk<T>
662 %y = force(%x): T
663 -----
664 %y = %f(%v1, ..., %vn): T

```

(a) force of a known function application: remove laziness

```

665 func f(%x1, ... %txi: thunk<T>, ..., %xn) {
666   %xi = force(%txi)
667   %zi = ... ; %tzi = thunkify(%zi)
668   f(%y1, ..., %tzi, ..., %yn)
669 }
670 -----
671 func f_strict_i(%x1, ... %xi: T, ..., %xn) {
672   %zi = ...
673   f_strict_i(%y1, ..., %zi, ..., %yn)
674 }
675 func f(%x1, ... %txi: thunk<T>, ..., %xn) {
676   %xi = force(%txi)
677   %zi = ...
678   f_strict_i(%y1, %y2, ..., %zi, ..., %yn)
679 }

```

(c) outlining recursive call that is immediately forced

```

680
681
682
683
684
685
686
687
688
689 %x = constructor(@Constructor, %v1, ..., %vm)
690 %y = case %x
691   [C1 -> {^entry1(%z11, ..., %z1n1): ... }]
692   [...]
693   [Ci -> {^entry1(%zi1, ..., %zini): ... }]
694 -----
695 Inline ^entryi with %zik = %vk

```

(e) Case of known constructor: remove indirection

```

696
697
698
699
700
701
702
703
704
705
706
707
708 -----
709 BAR
710

```

(g) Outline pattern matching branches on a function input

```

716 %x = thunkify(%v) : !lz.thunk<T>
717 %y = force(%x): T
718 -----
719 %y = %x

```

(b) force of a thunk: remove laziness

```

720
721
722 func f(%x1, ... %wrapxi: @ADT, ..., %xn) {
723   %xi = case (%wrapxi) [Wrapper -> ^ entry(%v) { return %v }]
724 }
725 -----
726 func f_work_i(%x1, ... %xi: T, ..., %xn) {
727   %zi = ...
728   f_strict_i(%y1, %y2, ..., %zi, ..., %yn)
729 }
730 func f(%x1, ... %txi: thunk<T>, ..., %xn) {
731   %xi = case (%wrapxi) [Wrapper -> ^ entry(%v) { return %v }]
732   %zi = ...
733   f_work_i(%y1, ..., %zi, ..., %yn)
734 }

```

(d) outlining of recursion of a monovariant wrapper

```

735
736
737 data C = MkC(V)
738 @f(%inc: C)
739   case inc
740   [C inv ->
741     %inv = extract(@MKC, %inc) : V
742     %w = ... : V; %wc = construct(@MKC, w) : C
743     %rec = apEager(@f, wc)]
744 -----
745 @frec(%inv: V)
746   %inv = extract(@MKC, %inc) : V
747   %w = ... : V; %wc = construct(@MKC, w) : C
748   %rec = apEager(@f, wc)
749 @f(%inc: C)
750   %inv = extract(%inc) : V
751   apEager(@frec, inv)

```

(f) Outline recursive call of constructor that is immediately unwrapped

```

752 data C = MKC(V)
753 @f(...)
754 ...
755 %out = constructor(C, %w)
756 lz.return (%out)
757 -----
758 data C = MKC(V)
759 @finner(...)
760 ...
761 lz.return %w
762
763 @f(...)
764 %w = call %f(...)
765 %out = constructor(C, %w)
766 lz.return (%out)

```

(h) Outline return of constructor**Figure 5.** Local rewrites performed to eliminate laziness (2)

```

771 9      %c1 = constant 1 : i64
772 10      lz.return %c1 : i64
773 11  }]
774 12  [@default -> {
775 13      ^entry:
776 14      %c1 = constant 1 : i64
777 15      %idec = subi %ival, %c1 : i64
778 16      %idecthnk = lz.thunkify(%idec : i64)
779 17      %f = constant @f : (!lz.thunk<i64>) -> i64
780 18      %f_idec_thnk = lz.ap(%f: (!lz.thunk<i64>) -> i64,
781 19          %idecthnk)
782 20      %f_idec_val = lz.force(%f_idec_thnk) : i64
783 21      %prod = muli %ival, %f_idec_val : i64
784 22      lz.return %prod : i64
785 23  }]
786 24  return %out : i64
787 25  }
788 26
789 27 func @c8 () -> i64 {
790 28     %v = std.constant 8 : i64
791 29     return %v: i64
792 30 }
793 31
794 32 func @main() -> i64 {
795 33     %f = constant @f : (!lz.thunk<i64>) -> i64
796 34     %c8f = constant @c8 : () -> i64
797 35     %c8t = lz.ap(%c8f : () -> i64)
798 36     %outt = lz.ap(%f: (!lz.thunk<i64>) -> i64, %c8t)
799 37     %out = lz.force(%outt) : i64
800 38     return %out : i64
801 39 }
802 40 }
803
804 Listing 18. "Original program with wrapper and laziness
805 overhead"
806
807     After worker/wrapper, we get the optimized program:
808 1 module {
809 2     func @c8() -> i64 {
810 3         %c8_i64 = constant 8 : i64
811 4         return %c8_i64 : i64
812 5     }
813 6     func @main() -> i64 {
814 7         %c7_i64 = constant 7 : i64
815 8         %f = constant @frec_force_outline : (i64) -> i64
816 9         %c8_i64 = constant 8 : i64
817 10        %0 = "lz.apEager"(%f, %c7_i64) : ((i64) -> i64,
818 11            i64) -> i64
819 12        %1 = muli %0, %c8_i64 : i64
820 13        return %1 : i64
821 14    }
822 15    func @frec_force_outline(%arg0: i64) -> i64 {
823 16        %c1_i64 = constant 1 : i64
824 17        %f = constant @frec_force_outline : (i64) -> i64
825 18        %0 = "lz.caseint"(%arg0) ( {

```

```

18      lz.return %c1_i64 : i64
19  }, {
20      %1 = subi %arg0, %c1_i64 : i64
21      %2 = "lz.apEager"(%f, %1) : ((i64) -> i64, i64)
22      -> i64
23      %3 = muli %arg0, %2 : i64
24      lz.return %3 : i64
25  }) {alt0 = 0 : i64, alt1 = @default} : (i64) ->
26  i64
27  return %0 : i64
28  }

```

Listing 19. "factorial that has been worker/wrapper'd"

We will explain this transformation step by step. Finally, when lowered to LLVM, LLVM's strong loop analyses is able to compute a closed form for `frec_force_outline`, ending in a constant time computation.

11.2 A toy example: eliminating laziness and exposing loop optimisation

```

1 module {
2     // sum up all values in the buffer
3     func @sum(%buffert: !lz.thunk<memref<?xi64>>) -> i64
4     {
5         %buffer = lz.force(%buffert) : memref<?xi64>
6         %c0 = constant 0 : index
7         %N = dim %buffer, %c0 : memref<?xi64>
8         %sum_0 = constant 0 : i64
9         %sum = affine.for %i = 0 to %N step 1
10        iter_args(%sum_iter = %sum_0) -> (i64) {
11            %t = affine.load %buffer[%i] : memref<?xi64>
12            %sum_next = std.addi %sum_iter, %t : i64
13            affine.yield %sum_next : i64
14        }
15        return %sum : i64
16    }
17    // create a sequence [0..upper_bound)
18    func @seq(%upper_bound: i64) -> memref<?xi64> {
19        %upper_bound_ix = std.index_cast %upper_bound : i64
20        to index
21        %buf = alloc(%upper_bound_ix) : memref<?xi64>
22        affine.for %i = 0 to %upper_bound_ix step 1 {
23            %ival = std.index_cast %i : index to i64
24            affine.store %ival, %buf[%i] : memref<?xi64>
25        }
26        return %buf : memref<?xi64>
27    }
28    func private @printInt(%i: i64)
29
30    // computes sum of numbers upto 1023?
31    func @main () -> i64 {
32        %seqf = constant @seq : (i64) -> memref<?xi64>
33        %size = std.constant 1024 : i64
34        %seqt = lz.ap(%seqf: (i64) -> memref<?xi64>, %size)
35    }

```

sid: Explain how this happens

840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880


```
881 35
882 36 %sumf = constant @sum : (!lz.thunk<memref<?xi64>>)
883     -> i64
884 37 %outt = lz.ap(%sumf : (!lz.thunk<memref<?xi64>>) ->
885     i64, %seqt)
886 38 %outv = lz.force(%outt): i64
887 39 call @printInt(%outv) : (i64) -> ()
888 40 %zero = constant 0 : i64
889 41 return %zero : i64
890 42 }
891 43 }
```

Listing 20. "example of laziness with tensor computations"

```
892
893 1 func @main() -> i64 {
894 2   %c0_i64 = constant 0 : i64
895 3   %0 = alloc() : memref<1xi64>
896 4   %1 = affine.for %arg0 = 0 to 1024 iter_args(%arg1 =
897     %c0_i64) -> (i64) {
898 5     %2 = index_cast %arg0 : index to i64
899 6     affine.store %2, %0[0] : memref<1xi64>
900 7     %3 = affine.load %0[0] : memref<1xi64>
901 8     %4 = addi %arg1, %3 : i64
902 9     affine.yield %4 : i64
903 10  }
904 11 call @printInt(%1) : (i64) -> ()
905 12 return %c0_i64 : i64
906 13 }
```

Listing 21. "final program after eliminating laziness and loop fusion"

```
907
908 1 define i64 @main() local_unnamed_addr !dbg !35 {
909 2   tail call void @printInt(i64 523776), !dbg !36
910 3   ret i64 0, !dbg !38
911 4 }
```

Listing 22. "final LLVM program after lowering"

12 Related Work

12.1 Demand analysis

[10] "Projections for demand analysis" performs demand analysis computation by modelling demands as projections on the semantic domain. We use their broad framework, adapted to our setting. [3], [2] extend the demand analysis with finer grained information, such as Call Arity analysis.

12.2 GHC's intermediate representation

The Glasgow Haskell compiler [4] is a mature optimizing compiler for Haskell, whose plugin infrastructure we hook into to develop a test set for lz, and whose demand analysis we use as a benchmark to measure against.

12.3 Intel Haskell research compiler IR

The Intel labs Haskell Research compiler [8] models Core in a strict ANF language. The compiler performs demand analysis and abstract simplification on ANF. The demand analysis is performed using traditional abstract interpretation techniques. Later, this demand information is used to interpret the program and perform abstract simplification. (TODO: figure out details of this step).

It then compiles to an intermediate representation called MIL, which is a loosely typed CFG based, SSA-lite intermediate representation. They represent laziness as heap values, and manipulate the heap. Thus, their representation and usage of lazy values reasons with memory, instead of providing value semantics. MIL also does not have a notion of nested regions. Therefore, MIL extends the traditional control flow controls with more finer-grained information, called as cut and interproc.

12.4 Mixing dataflow analysis and transformation

This is based on the theory of compositional dataflow analysis and transformations, [6] which proves that we can interleave dataflow analyses and transformations safely.

Hoopl [9] is the dataflow analysis and transformation library within GHC. However, GHC does not use this for performing dataflow analysis over Core, as the Hoopl library is designed to work with a CFG based intermediate representation, while GHC Core is a typed functional, expression oriented intermediate representation. While Hoopl is used for certain simplifications within C- in GHC, it is not used extensively due to poor performance characteristics exhibited by the implementation. (TODO: bench Hoopl on contrived examples)

12.5 Alternative encodings of laziness

GRIN [1] is an alternative intermediate representation for lazy and strict functional programming languages which explicitly represents heap manipulation and case analysis. However, it is not SSA based.

12.6 Theoretical justification for our encoding

Our encoding is based on call-by-push-value [7], which breaks down call-by-value and call-by-name paradigms into simple primitives. In our dialect lz, we expose these primitives to enable a clean mixture of call by name and call by value. This enjoys the many properties that are proven to hold for call-by-push-value (TODO: which ones?)

12.7 Technology substrate

MLIR [5] is a compiler framework for building reusable and extensible compiler infrastructure. MLIR aims to address software fragmentation, improve compilation for heterogeneous hardware, significantly reduce the cost of building domain specific compilers, and aid in connecting existing compilers together.

13 Conclusion

We provide a baseline implementation of non-strictness for the MLIR framework which with very minimal extensions, allows us to fully express non strict semantics while being compatible with the rest of MLIR's strict dialects. As an experience report, we explore MLIR's strengths and weaknesses at representing non-strictness within an SSA based framework. We also provide a baseline implementation of demand analysis, as inspired by Wadler and Hugh's "projections for demand analysis". We evaluate our example programs against the nofib benchmark test suite. We also provide provocative examples where the potent combination of MLIR's loop optimization infrastructure along with our modest extensions allows us to beat native haskell's performance.

14 TODO

14.1 Performance benchmarking

14.2 Correct encoding for mutual recursion

14.3 Correct encoding for partially applied functions/closures

We currently have a lz.lambda but it's basically untested, because we didn't have many things that used closures. We need to make sure this works properly, and encodes data correctly.

14.3.1 Lambdas in GRIN. GRIN is also a low level IR, so it's useful to recall how this is encoded within GRIN. To quite the GRIN thesis:

sid: MLIR can work harder, but it does not. By the time we get to LLVM, we have a closed form

“ end. Using hbcc we get well optimised code in a low level functional style, comparable to for example the Core language [PJ96] used by the Glasgow Haskell compiler. The code is lambda lifted, i.e., has only super combinators, and most high level Haskell constructions, like overloading, have been transformed away. ”

14.3.2 Implementing lambda lifting. Recall that optimal lambda lifting is $O(n^2)$. I don’t know the algorithm; I’d have to study it before I implement it.

14.4 Stretch: Separate dialect for pattern matching

14.5 Stretch: rewrite based demand analysis

14.6 Stretch: rewrite based unification

We introduce a dialect for performing unification.

14.7 Stretch: rewrite based tabled typeclass unification

14.8 Stretch: As static as possible garbage collection

References

- [1] Urban Boquist and Thomas Johnsson. 1996. The GRIN project: A highly optimising back end for lazy functional languages. In Symposium on Implementation and Application of Functional Languages. Springer, 58–84.
- [2] Joachim Breitner. 2014. Call arity. In International Symposium on Trends in Functional Programming. Springer, 34–50.
- [3] Sebastian Graf. [n. d.]. Call Arity vs. Demand Analysis. ([n. d.]).
- [4] SL Peyton Jones, Cordy Hall, Kevin Hammond, Will Partain, and Philip Wadler. 1993. The Glasgow Haskell compiler: a technical overview. In Proc. UK Joint Framework for Information Technology (JFIT) Technical Conference, Vol. 93.
- [5] Chris Lattner, Jacques Pienaar, Mehdi Amini, Uday Bondhugula, River Riddle, Albert Cohen, Tatiana Shpeisman, Andy Davis, Nicolas Vasilache, and Oleksandr Zinenko. 2020. MLIR: A Compiler Infrastructure for the End of Moore’s Law. arXiv preprint arXiv:2002.11054 (2020).
- [6] Sorin Lerner, David Grove, and Craig Chambers. 2002. Composing dataflow analyses and transformations. ACM SIGPLAN Notices 37, 1 (2002), 270–282.
- [7] Paul Blain Levy. 2012. Call-by-push-value: A Functional/imperative Synthesis. Vol. 2. Springer Science & Business Media.
- [8] Hai Liu, Neal Glew, Leaf Petersen, and Todd A Anderson. 2013. The Intel labs Haskell research compiler. In Proceedings of the 2013 ACM SIGPLAN symposium on Haskell. 105–116.
- [9] Norman Ramsey, Joao Dias, and Simon Peyton Jones. 2010. Hoopl: a modular, reusable library for dataflow analysis and transformation. ACM Sigplan Notices 45, 11 (2010), 121–134.
- [10] Philip Wadler and R John M Hughes. 1987. Projections for strictness analysis. In Conference on Functional Programming Languages and Computer Architecture. Springer, 385–407.

A Appendix

Text of appendix