# Deforestation and Program Fusion: Applying Equational Transforms to Automatically Simplify Programs

`github.com/ryanorendorff/xxxxxxxxxxxx`

Ryan Orendorff, PhD?[1]

May 2017

[1]Department of Bioengineering
University of California, Berkeley
University of California, San Francisco

## Table of Contents

## Common way to process a list: map and fold!

As an example, say we want to square all the elements in a list and then sum the result.

$$process :: [Int] \rightarrow Int$$
$$process \ xs = sum \circ map \ sq \ \$ \ xs$$

Where we have defined the functions as follows.

$$map \ \_ \ [] \qquad = []$$
$$map \ f \ (x : xs) = f \ x : map \ f \ xs$$
$$sq \ x = x * x$$

## Common way to process a list: map and fold!

As an example, say we want to square all the elements in a list and then sum the result.

$$process :: [Int] \to Int$$
$$process\ xs = sum \circ map\ sq\ \$\ xs$$

Where we have defined the functions as follows.

$$map\ \_\ [] \qquad = []$$
$$map\ f\ (x : xs) = f\ x : map\ f\ xs$$
$$sq\ x = x * x$$


$$foldr\ \_\ b\ [] \qquad = b$$
$$foldr\ f\ b\ (a : as) = foldr\ f\ (f\ a\ b)\ as$$
$$sum = foldr\ (+)\ 0$$

**How fast is** *process***?**

So now that we have our process function, how fast does it run?

$process :: [Int] \rightarrow Int$

$process\ xs = sum \circ map\ sq\ \$\ xs$

Let's try to process a million elements with our *process* and *process*′, which uses the standard Prelude *sum* and *map*.

$process\ [0 .. 1,000,000];\ process'\ [0 .. 1,000,000]$

4

**How fast is** *process***?**

So now that we have our process function, how fast does it run?

$$process :: [Int] \rightarrow Int$$
$$process \; xs = sum \circ map \; sq \; \$ \; xs$$

Let's try to process a million elements with our *process* and *process'*, which uses the standard Prelude *sum* and *map*.

$$process \; [0 \ldots 1,000,000]; process' \; [0 \ldots 1,000,000]$$

| Function | Time (ms) | Memory (MB) |
|----------|-----------|-------------|
| *process* | 220.0 | 265.26 |
| *process'* | 5.1 | 80.00 |

**How fast is *process*?**

So now that we have our process function, how fast does it run?

$$process :: [Int] \rightarrow Int$$
$$process \ xs = sum \circ map \ sq \ \$ \ xs$$

Let's try to process a million elements with our *process* and *process'*, which uses the standard Prelude *sum* and *map*.

$$process \ [0 \mathinner{\ldotp\ldotp} 1,000,000]; process' \ [0 \mathinner{\ldotp\ldotp} 1,000,000]$$

| Function | Time (ms) | Memory (MB) |
|----------|-----------|-------------|
| *process* | 220.0 | 265.26 |
| *process'* | 5.1 | 80.00 |

*How does the Prelude do so much better with the same functions?*

4

## Let's try to optimize with an accumulator

What if we try to optimize *process* using an accumulator?

$process :: [Int] \rightarrow Int$
$process\ xs = sum \circ map\ sq\ \$\ xs$

## Let's try to optimize with an accumulator

What if we try to optimize *process* using an accumulator?

$$process :: [Int] \rightarrow Int$$
$$process\ xs = sum \circ map\ sq\ \$\ xs$$

This means we make a tail-recursive "loop".

$$processFused :: [Int] \rightarrow Int$$
$$processFused = loop\ 0$$
$$\textbf{where}$$
$$loop\ n\ [\ ] \qquad = n$$
$$loop\ n\ (x:xs) = loop\ (n + x * x)\ xs$$

## Let's try to optimize with an accumulator

What if we try to optimize *process* using an accumulator?

$$process :: [Int] \rightarrow Int$$
$$process\ xs = sum \circ map\ sq\ \$\ xs$$

This means we make a tail-recursive "loop".

$$processFused :: [Int] \rightarrow Int$$
$$processFused = loop\ 0$$
> **where**
> $$loop\ n\ [\ ] \qquad = n$$
> $$loop\ n\ (x:xs) = loop\ (n + x * x)\ xs$$

This mirrors the version of the program one would write
imperatively.

## How well does our manually fused version do?

$$process\ xs = sum \circ map\ sq\ \$\ xs$$
$$process'\ xs = Prelude.sum \circ Prelude.map\ sq\ \$\ xs$$
$$processFused :: [Int] \rightarrow Int$$
$$processFused = loop\ 0$$
**where**
$$loop\ n\ [\ ] \qquad = n$$
$$loop\ n\ (x:xs) = loop\ (n + x * x)\ xs$$

## How well does our manually fused version do?

$$process \; xs = sum \circ map \; sq \; \$ \; xs$$
$$process' \; xs = Prelude.sum \circ Prelude.map \; sq \; \$ \; xs$$
$$processFused :: [Int] \rightarrow Int$$
$$processFused = loop \; 0$$
**where**
$$loop \; n \; [\,] \qquad = n$$
$$loop \; n \; (x : xs) = loop \; (n + x * x) \; xs$$

| Function | Time (ms) | Memory (MB) |
|---|---|---|
| $process$ | 220.0 | 265.26 |
| $process'$ | 5.1 | 80.00 |
| $processFused$ | 4.7 | 80.00 |
| process.c | 2.6 | $8 \times 10^{-5}$ |

## GHC generated the loop version automatically

Our manual version $processFused$.

$$processFused :: [Int] \rightarrow Int$$
$$processFused = loop\ 0$$
$$\quad \textbf{where}$$
$$\quad\quad loop\ n\ [\ ] \quad\quad = n$$
$$\quad\quad loop\ n\ (x : xs) = loop\ (n + x * x)\ xs$$

and when we compile the Prelude defined $process'$, GHC produces

$$process' :: [Int] \rightarrow Int$$
$$process'\ xs = loop'\ xs\ 0$$
$$loop' :: [Int] \rightarrow Int \rightarrow Int$$
$$loop'\ [\ ] \quad\quad n = n$$
$$loop'\ (x : xs)\ n = loop'\ xs\ (n + x * x)$$

**We want simpe/readible programs that are also performant**

Our original version was simple to understand, but not fast.

$$process\ xs = sum \circ map\ sq\ \$\ xs$$

Our accumulator one is fast but not simple.

$$processFused :: [Int] \rightarrow Int$$
$$processFused = loop\ 0$$
$$\quad \textbf{where}$$
$$\quad\quad loop\ n\ [\,] \quad\quad = n$$
$$\quad\quad loop\ n\ (x : xs) = loop\ (n + x * x)\ xs$$

*How can we leverage the compile to write simple code that is fast?*

## Table of Contents

## The GHC Compilation Pipeline converts Haskell into an intermediate language and then bytecode
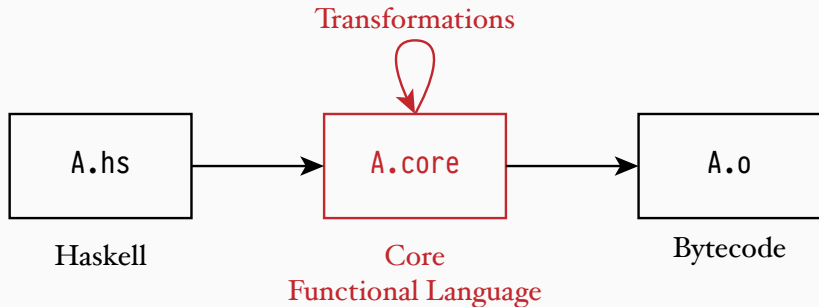
When GHC compiles a Haskell program, it converts the code into an intermediate language called "Core", which is then (eventually) turned into byte code.

Transformations

# The GHC Compilation Pipeline converts Haskell into an intermediate language and then bytecode

When GHC compiles a Haskell program, it converts the code into an intermediate language called "Core", which is then (eventually) turned into byte code.



Transformations

| A.hs | A.core | A.o |
| --- | --- | --- |

Haskell

Core
Functional Language

Bytecode

## GHC performs several program transformations on Core to optimize the code

When GHC is given a Core program, it performs several types of transformations on the program.

- Inlining functions

## GHC performs several program transformations on Core to optimize the code

When GHC is given a Core program, it performs several types of transformations on the program.

- Inlining functions
- Removing redundant lambdas

## GHC performs several program transformations on Core to optimize the code

When GHC is given a Core program, it performs several types of transformations on the program.

- Inlining functions
- Removing redundant lambdas
- Simplifying constant expressions $((x + 8) - 1)$

## GHC performs several program transformations on Core to optimize the code

When GHC is given a Core program, it performs several types of transformations on the program.

- Inlining functions
- Removing redundant lambdas
- Simplifying constant expressions $((x + 8) - 1)$
- Combining type casts

# GHC performs several program transformations on Core to optimize the code

When GHC is given a Core program, it performs several types of transformations on the program.

- Inlining functions
- Removing redundant lambdas
- Simplifying constant expressions $((x + 8) - 1)$
- Combining type casts
- *Applying rewrite rules*

## GHC performs several program transformations on Core to optimize the code

When GHC is given a Core program, it performs several types of transformations on the program.

- Inlining functions

- Removing redundant lambdas

- Simplifying constant expressions $((x + 8) - 1)$

- Combining type casts

- *Applying rewrite rules*

- . . .

## Rewrite Rules allow us to say two expressions are equivalent

Rewrite rules allow us to replace terms in the program with equivalent terms.

$$\{-\# \text{ RULES "name" forall x. id x} = \text{x } \#-\}$$

## Rewrite Rules allow us to say two expressions are equivalent

Rewrite rules allow us to replace terms in the program with equivalent terms.

$$\{-\# \text{ RULES "name" forall x. id x} = \text{x } \#-\}$$

"Any time we see the term $id\ x$, replace it with $x$".

## Rules have some restrictions

Rewrite rules have some gotchas.

- Rules doesn't prevent you from doing something silly

    {-# RULES "id5" forall x. id x = 5 #-}

## Rules have some restrictions

Rewrite rules have some gotchas.

- Rules doesn't prevent you from doing something silly

    {-# RULES "id5" forall x. id x = 5 #-}

- The left hand side is only substituted for the right, not the other way around.

    {-# RULES "id" forall x. id x = x #-}
    $x \nRightarrow id\ x$

## Rules have some restrictions

Rewrite rules have some gotchas.

- Rules doesn't prevent you from doing something silly

  {-# RULES "id5" forall x. id x = 5 #-}

- The left hand side is only substituted for the right, not the other way around.

  {-# RULES "id" forall x. id x = x #-}
  $x \nRightarrow id\ x$

- You can make the compiler go into an infinite loop.

  {-# RULES "fxy" forall x y. f x y = f y x #-}

## Rules have some restrictions

Rewrite rules have some gotchas.

- Rules doesn't prevent you from doing something silly

  {-# RULES "id5" forall x. id x = 5 #-}

- The left hand side is only substituted for the right, not the other way around.

  {-# RULES "id" forall x. id x = x #-}
  $x \nRightarrow id\ x$

- You can make the compiler go into an infinite loop.

  {-# RULES "fxy" forall x y. f x y = f y x #-}

- If multiple rules are possible, GHC will randomly choose one.

## We can combine maps to traverse a list once

Let us introduce the following rule about maps.

$\{-\#$ RULES "$map_f/map_f$" forall f g xs. $map_f$ f ($map_f$ g xs) = $map_f$

**We can combine maps to traverse a list once**

Let us introduce the following rule about maps.

{-# RULES "map$_f$/map$_f$" forall f g xs. map$_f$ f (map$_f$ g xs) = map$_f$

$mapTest0 :: [Int] \rightarrow [Int]$
$mapTest0\ xs = map\ (+1)\ (map\ (*2)\ xs)$

$mapTest1 :: [Int] \rightarrow [Int]$
$mapTest1\ xs = map_f\ (+1)\ (map_f\ (*2)\ xs)$

## Our map fusion performs (a bit) better!

We can test our functions on a million elements

$mapTest0\ xs = map\ (+1)\ (map\ (*2)\ xs)$
$mapTest0\ [1\mathbin{..}1,000,000]$

$mapTest1\ xs = map_f\ (+1)\ (map_f\ (*2)\ xs)$
$mapTest1\ [1\mathbin{..}1,000,000]$

and find we get a bit better time and space performance.

| Function | Time (ms) | Memory (MB) |
|----------|-----------|-------------|
| $mapTest0$ | 26.4 | 256.00 |
| $mapTest1$ | 17.6 | 184.00 |

## Through rules, GHC performs fusion

Some of the rules work together to perform *fusion*: to combine terms in such a way as to pass over a data structure once.

In our *process* function, we create an intermediate list

$$process\ xs = sum \circ map\ sq\ \$\ xs$$

whereas our "fused" form did not make any intermediate structure, and used an accumulator instead.

$$processFused :: [Int] \rightarrow Int$$
$$processFused = loop\ 0$$
$$\quad \textbf{where}$$
$$\quad\quad loop\ n\ [\,] \quad\quad = n$$
$$\quad\quad loop\ n\ (x:xs) = loop\ (n + x * x)\ xs$$

## Table of Contents

## *foldr* / *build* **fusion is used to simplify list computations**

GHC accomplishes fusion with two functions: foldr and build.

GHC accomplishes fusion with two functions: foldr and build.

*foldr* combines the elements of a list

$$foldr \_ b \, [\,] \qquad = b$$
$$foldr \, f \, b \, (a : as) = foldr \, f \, (f \, a \, b) \, as$$

## $foldr$ / $build$ **fusion is used to simplify list computations**

GHC accomplishes fusion with two functions: foldr and build.

$foldr$ combines the elements of a list

$$foldr \ \_ \ b \ [\,] \qquad = b$$
$$foldr \ f \ b \ (a : as) = foldr \ f \ (f \ a \ b) \ as$$

while $build$ builds up a list from a generating function.

$$build :: (forall \ b \circ (a \rightarrow b \rightarrow b) \rightarrow b \rightarrow b) \rightarrow [\,a\,]$$
$$build \ g = g \ (:) \ [\,]$$

$$build \ l \equiv [1, 2, 3]$$
$$\quad \textbf{where}$$
$$\quad\quad l \ cons \ nil = 1 \ `cons` \ (2 \ `cons` \ (3 \ `cons` \ nil))$$

## The $foldr/build$ rule removes intermediate fold/build pairs

To remove intermediate data structures (those created by $build$),
we eliminate $foldr/build$ pairs with a rule.

$\{-\#$ RULES $"build/foldr"$ forall f z (g :: forall b. (a -¿ b -¿ b) -¿

$foldr\ (+)\ 0\ (build\ l) \equiv l\ (+)\ 0 \equiv 1 + (2 + (3 + 0))$
   **where**
     $l\ cons\ nil = 1\ `cons`\ (2\ `cons`\ (3\ `cons`\ nil))$

## We need a few extra rules to convert maps into fold/builds

To convert our definition of maps into a fold/build pair, we need the following helper function.

$$map_{fb} :: (b \to l \to l) \to (a \to b) \to a \to l \to l$$
$$map_{fb} \; cons \; f = \lambda x \; ys \to f \; x \; `cons` \; ys$$

## We need a few extra rules to convert maps into fold/builds

To convert our definition of maps into a fold/build pair, we need the following helper function.

$$map_{fb} :: (b \rightarrow l \rightarrow l) \rightarrow (a \rightarrow b) \rightarrow a \rightarrow l \rightarrow l$$
$$map_{fb} \; cons \; f = \lambda x \; ys \rightarrow f \; x \; `cons` \; ys$$

With that, we have all we to convert map into build/fold.

$\{-\# \text{ RULES } "map" \text{ forall f xs. } map \text{ f xs} = build \; (\text{ŋ -¿ } foldr \; (map$

## We need a few extra rules to convert maps into fold/builds

To convert our definition of maps into a fold/build pair, we need the following helper function.

$$map_{fb} :: (b \to l \to l) \to (a \to b) \to a \to l \to l$$
$$map_{fb} \ cons \ f = \lambda x \ ys \to f \ x \ `cons` \ ys$$

With that, we have all we to convert map into build/fold.

{-# RULES "$map$" forall f xs. $map$ f xs $= build$ (ŋ -ç $foldr \ (map$

We also provide a way to combine sequential $map_{fb}$ functions.

{-# RULES "mapFB" forall c f g. mapFB (mapFB c f) g $=$ mapF

## This is where the example is expanded completely

Will show how the map sum example is converted into the loop we saw.

Will need

$$foldr :: (a \to b \to b) \to b \to [\,a\,] \to b$$
$$foldr\ f\ z = go$$
$$\quad \textbf{where}$$
$$\quad\quad go\ [\,] = z$$
$$\quad\quad go\ (y : ys) = y\ `f`\ go\ ys$$
{-# INLINE [0] foldr2 #-}

## Table of Contents

## Table of Contents

## We can make *process* even faster with $Data.Vector$

The $Data.Vector$ package uses stream fusion and many other rewrite rules behind the scenes in order to optimize array based computations.

$$process\ xs = sum \circ map\ sq\ \$\ xs$$

## We can make *process* even faster with $Data.Vector$

The $Data.Vector$ package uses stream fusion and many other rewrite rules behind the scenes in order to optimize array based computations.

$$process\ xs = sum \circ map\ sq\ \$\ xs$$

The vector version looks very similar.

**import** *qualified Data.Vector as V*

$$processVec\ n = V.sum\ \$\ V.map\ sq\ \$\ V.enumFromTo\ 1\ (n :: Int)$$

# We can make $process$ **even faster with** $Data.Vector$

But has incredible performance!

| Function | Time (ms) | Memory (MB) |
|---|---|---|
| $process$ | 220.0 | 265.26 |
| $process'$ | 5.1 | 80.00 |
| $processFused$ | 4.7 | 80.00 |
| process.c | 2.6 | $8 \times 10^{-5}$ |
| processVec | 0.7 | $16 \times 10^{-5}$ |

While we wrote this in our program

$$processVec\ n = V.sum\ \$\ V.map\ sq\ \$\ V.enumFromTo\ 1\ (n :: Int)$$

GHC ended up generating the following Core code.

```
val_ $ s $ wfoldlM'_loop
val_ $ s $ wfoldlM'_loop =
  λsc sc1 →
    case tagToEnum # (<= #sc 100000000#) of _ {
      False → sc1;
      True → val_ $ s $ wfoldlM'_loop (+#sc 1#) (+#sc1 (*#sc
    }
```

simplify this core, ignore unboxing

# Repa: A numerical Haskell Library using Fusion

Repa also uses fusion in order to handle array operations.

# Data Parallel Haskell: Nested Data Parallelism made easy

$processDPH :: [:Int:] \rightarrow Int$
$processDPH = sumDPH \circ mapDPH \ sq \ \$ \ xs$

Does dispatch by MPI.