
Rewriting Haskell Strings



Duncan Coutts

Don Stewart

Roman Leshchinskiy

Programming Tools Group
Oxford University

Programming Languages and Systems
University of New South Wales



SAPLING 2006

IO IN HASKELL

Lazy Haskell IO can be beautiful:

```
return · foldl' k 5381 · map toLower · filter isAlpha =<< readFile f  
where k h c = h * 33 + ord c
```

IO IN HASKELL

Lazy Haskell IO can be beautiful:

```
return · foldl' k 5381 · map toLower · filter isAlpha =<< readFile f  
  where k h c = h * 33 + ord c
```

But naive code can be slow:

```
$ time ./a.out  
bf805325  
./a.out 23.04s total
```

AND ENTIRELY UNLIKE C

An equivalent naive C implementation:

```
int c;
long h = 5381;
FILE *fp = fopen(f, "r");
while ((c = fgetc(fp)) != EOF)
    if (isalpha(c))
        h = h * 33 + tolower(c);
fclose(fp);
return h;
```

AND ENTIRELY UNLIKE C

An equivalent naive C implementation:

```
int c;
long h = 5381;
FILE *fp = fopen(f, "r");
while ((c = fgetc(fp)) != EOF)
    if (isalpha(c))
        h = h * 33 + tolower(c);
fclose(fp);
return h;
```

```
$ time ./a.out
bf805325
./a.out 3.93s total
```

LET'S FIX IT!

Use ByteStrings with stream fusion

```
import Data.ByteString.Lazy  
return · foldl' k 5381 · map toLower · filter isAlpha ==<< readFile f  
where k h c = h * 33 + ord c
```

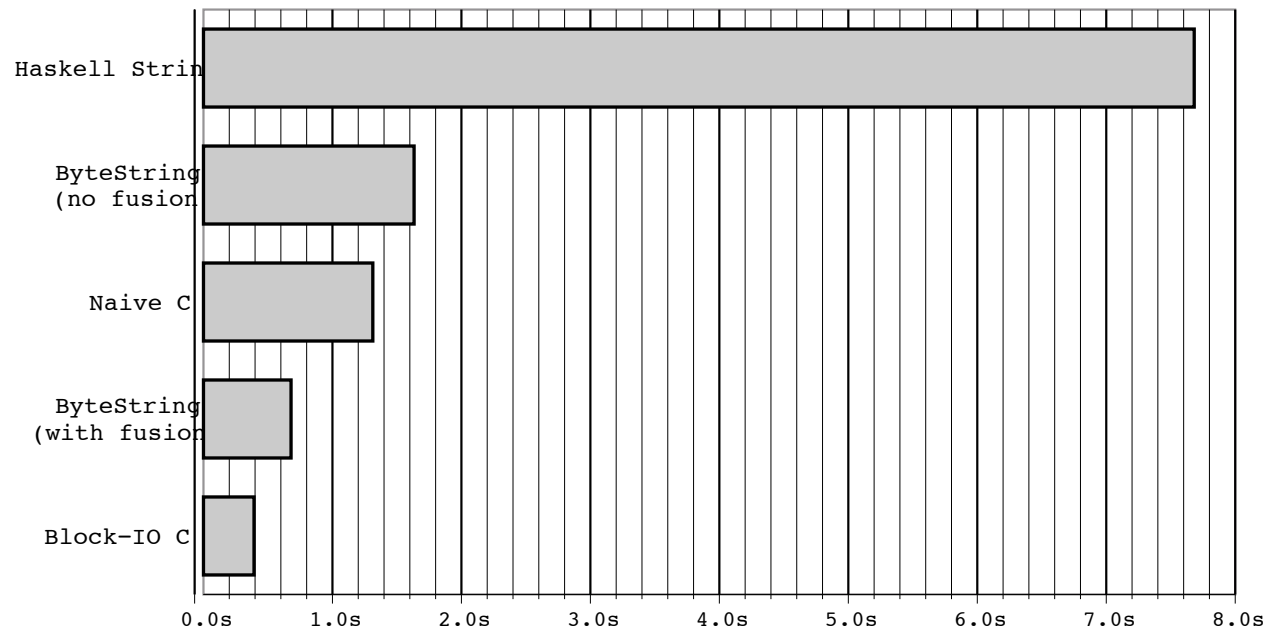
LET'S FIX IT!

Use ByteStrings with stream fusion

```
import Data.ByteString.Lazy  
return · foldl' k 5381 · map toLower · filter isAlpha ==<< readFile f  
  where k h c = h * 33 + ord c
```

```
$ time ./a.out  
bf805325  
./a.out 2.04s total
```

COMPARING STRINGS, BYTESTRINGS AND C



You have to rewrite the C version to be block oriented to win.

REPRESENTING STRINGS

The Lovely, Delicate [Char]

```
data [] a = [] | a : [a]  
type String = [Char]
```

- Strings are represented as linked lists of pointers to chars.
- Every node requires an indirection, and another to retrieve character
- Bad for the cache, bad for gcc, bad for performance.

REPRESENTING STRINGS

The Lovely, Delicate [Char]

```
data [] a = [] | a : [a]
type String = [Char]
```

- Strings are represented as linked lists of pointers to chars.
- Every node requires an indirection, and another to retrieve character
- Bad for the cache, bad for gcc, bad for performance.

The Ruthless, Cut-throat ByteString

```
data ByteString = BS !(ForeignPtr Word8) !Int !Int
```

Listen to your hardware!

Changing the type is a big win, but we can play better games yet ...

BOTTLENECK: TEMPORARY DATA STRUCTURES

```
print . minimum . filter ('isPrefixOf' x) . map toUpper
```

BOTTLENECK: TEMPORARY DATA STRUCTURES

```
print . minimum . filter ('isPrefixOf' x) . map toUpper
```

And this is compiled to:

```
case map toUpper s           of s'   ->
case filter ('isPrefixOf' x) s' of s'' ->
case minimum s''            of s'''  -> print s'''
```

Lots of temporary arrays are being allocated and discarded!

THE BIG IDEA: DEFORESTATION AND FUSION

Ideally, we'd like to make just one pass over the input data

We'd write:

```
map f . map g
```

THE BIG IDEA: DEFORESTATION AND FUSION

Ideally, we'd like to make just one pass over the input data

We'd write:

```
map f . map g
```

and the compiler would emit:

```
map (f . g)
```

THE BIG IDEA: DEFORESTATION AND FUSION

Ideally, we'd like to make just one pass over the input data

We'd write:

```
map f . map g
```

and the compiler would emit:

```
map (f . g)
```

We can teach the compiler how to do this:

$$\langle \text{map/map fusion} \rangle \quad \text{map } f \cdot \text{map } g \mapsto \text{map } (f \cdot g)$$

FUSION SYSTEMS

A variety of general purpose fusion systems exist:

$$\langle \mathbf{foldr/build\ fusion} \rangle \forall g k z .$$
$$\mathit{foldr} k z (\mathit{build} g) \mapsto g k z$$

FUSION SYSTEMS

A variety of general purpose fusion systems exist:

$\langle \mathbf{foldr/build\ fusion} \rangle \forall g k z .$
 $foldr\ k\ z\ (build\ g) \mapsto g\ k\ z$

$\langle \mathbf{destroy/unfoldr\ fusion} \rangle \forall g f e .$
 $destroy\ g\ (unfoldr\ f\ e) \mapsto g\ f\ e$

FUSION SYSTEMS

A variety of general purpose fusion systems exist:

$\langle \mathbf{foldr/build\ fusion} \rangle \forall g k z .$
 $foldr\ k\ z\ (build\ g) \mapsto g\ k\ z$

$\langle \mathbf{destroy/unfoldr\ fusion} \rangle \forall g f e .$
 $destroy\ g\ (unfoldr\ f\ e) \mapsto g\ f\ e$

$\langle \mathbf{function\ array\ fusion} \rangle \forall f g s t .$
 $loop\ f\ s \cdot fst \cdot loop\ g\ t \mapsto loop\ (fuse\ f\ g)\ (s, t)$

FUSION SYSTEMS

A variety of general purpose fusion systems exist:

$\langle \text{foldr/build fusion} \rangle \forall g k z .$
 $\text{foldr } k z (\text{build } g) \mapsto g k z$

$\langle \text{destroy/unfoldr fusion} \rangle \forall g f e .$
 $\text{destroy } g (\text{unfoldr } f e) \mapsto g f e$

$\langle \text{function array fusion} \rangle \forall f g s t .$
 $\text{loop } f s \cdot \text{fst} \cdot \text{loop } g t \mapsto \text{loop } (\text{fuse } f g) (s, t)$

But none support all of: a wide range of operations, flat non-inductive structures, *and* compile to fast code in GHC.

STREAM FUSION

To support *array* access patterns in $\mathcal{O}(1)$, we need to factor out the three phases:

- read the array into a stream of elements
- process the stream elements
- write the resulting stream into a new array

THE STREAM TYPE

We need an efficient stream abstraction:

```
data Stream =  $\exists s$ . Stream (s  $\rightarrow$  Step s) s Int
data Step s = Done
            | Yield Word8 s
            | Skip s
```

An existentially wrapped seed, and stepper function, generating one of 3 results for each element of the stream:

- *Done* Finished processing the stream
- *Yield* Produce a transformed element
- *Skip* Filter this element out

Also carry a hint about the resulting size, to guide reallocations.

BUILDING A STREAM FROM AN ARRAY

Arrays to Streams:

$read \quad :: \text{ByteString} \rightarrow \text{Stream}$

$read\ s = \text{Stream}\ next\ 0\ n$

where

$n \quad \quad \quad = \text{length}\ s$

$next\ i \mid i < n \quad = \text{Yield}\ (\text{index}\ s\ i)\ (i + 1)$

$\mid \text{otherwise} = \text{Done}$

Writing streams back out to array is also fairly straight forward.

REMOVING INTERMEDIATE STREAMS

Pairs of read and write are just the identity function on streams, so they can be removed, yielding our fusion rule:

$$\langle \text{read/write fusion} \rangle \quad \text{read} \cdot \text{write} \mapsto \text{id}$$

- Whenever we see these pairs in the user's code, it is safe for the compiler to remove them.
- The library author supplies rewrite rules specific to the library, extending the compiler's optimisation range.

STREAM TRANSFORMERS

What about actually transforming the data?

Build a stream, transform it, and write it back:

$$\begin{aligned} \text{map} &:: (\text{Word8} \rightarrow \text{Word8}) \rightarrow \text{ByteString} \rightarrow \text{ByteString} \\ \text{map } f &= \text{write} \cdot \text{mapS } f \cdot \text{read} \end{aligned}$$

STREAM TRANSFORMERS

What about actually transforming the data?

Build a stream, transform it, and write it back:

$$\begin{aligned} \text{map} &:: (\text{Word8} \rightarrow \text{Word8}) \rightarrow \text{ByteString} \rightarrow \text{ByteString} \\ \text{map } f &= \text{write} \cdot \text{mapS } f \cdot \text{read} \end{aligned}$$

And applying f to each element:

$$\begin{aligned} \text{mapS} &:: (\text{Word8} \rightarrow \text{Word8}) \rightarrow \text{Stream} \rightarrow \text{Stream} \\ \text{mapS } f (\text{Stream next } s \ n) &= \text{Stream next}' s \ n \end{aligned}$$

where

$$\text{next}' s = \mathbf{case} \ \text{next } s \ \mathbf{of}$$
$$\text{Done} \quad \rightarrow \text{Done}$$
$$\text{Yield } x \ s' \rightarrow \text{Yield } (f \ x) \ s'$$
$$\text{Skip } s' \quad \rightarrow \text{Skip } s'$$

FUSING MAP

$$\mathit{map} f \cdot \mathit{map} g$$

FUSING MAP

$map\ f \cdot map\ g$

$= write \cdot mapS\ f \cdot read \cdot \{inline\ map\ \times 2\}$
 $write \cdot mapS\ g \cdot read$

FUSING MAP

$map\ f \cdot map\ g$

$= write \cdot mapS\ f \cdot read \cdot \{inline\ map\ \times 2\}$
 $write \cdot mapS\ g \cdot read$

$= write \cdot mapS\ f \cdot mapS\ g \cdot read \quad \{read/write\ fusion\}$

FUSING MAP

$map\ f \cdot map\ g$

$= write \cdot mapS\ f \cdot read \cdot \{inline\ map\ \times 2\}$
 $write \cdot mapS\ g \cdot read$

$= write \cdot mapS\ f \cdot mapS\ g \cdot read \quad \{read/write\ fusion\}$

$= write \cdot mapS\ (f \cdot g) \cdot read \quad \{map/map\ fusion\}$

FILTER

$filterS :: (Word8 \rightarrow Bool) \rightarrow Stream \rightarrow Stream$

$filterS\ p\ (Stream\ next\ s\ n) = Stream\ next'\ s\ n$

where

$next'\ s = \mathbf{case\ next\ s\ of}$

$Done \quad \quad \quad \rightarrow Done$

$Yield\ x\ s' \mid p\ x \quad \rightarrow Yield\ x\ s'$

$\quad \quad \quad \mid otherwise \rightarrow Skip\ s'$

$Skip\ s' \quad \quad \quad \rightarrow Skip\ s'$

FOLD

$foldlS' :: (a \rightarrow Word8 \rightarrow a) \rightarrow a \rightarrow Stream \rightarrow a$
 $foldlS' f z (Stream\ next\ s\ n) = loop\ z\ s$

where

$loop\ z\ s = \mathbf{case\ next\ s\ of}$

$Done \quad \rightarrow z$

$Yield\ x\ s' \rightarrow loop\ (f\ z\ x)\ s'$

$Skip\ s' \quad \rightarrow loop\ z\ s'$

FIND/SHORT-CIRCUITING

$findS :: (Word8 \rightarrow Bool) \rightarrow Stream \rightarrow Maybe Word8$

$findS\ p\ (Stream\ next\ s\ n) = loop\ s$

where

$loop\ s = \mathbf{case}\ next\ s\ \mathbf{of}$

$Done \quad \rightarrow Nothing$

$Yield\ x\ s' \mid p\ x \quad \rightarrow Just\ x$

$\quad \mid otherwise \rightarrow loop\ s'$

$Skip\ s' \quad \rightarrow loop\ s'$

FUSING WITH STREAMS

foldl' f z · map g · filter h

FUSING WITH STREAMS

foldl' f z · map g · filter h

= foldlS' f z · read · write · mapS g {inline *foldl'*, *map*
· read · write · filterS h · read and *filter*}

FUSING WITH STREAMS

$foldl' f z \cdot map g \cdot filter h$

$= foldlS' f z \cdot read \cdot write \cdot mapS g \quad \{\text{inline } foldl', \text{ map}\}$
 $\quad \cdot read \cdot write \cdot filterS h \cdot read \quad \{\text{and filter}\}$

$= foldlS' f z \cdot mapS g \cdot filterS h \cdot read \quad \{\text{read/write fusion}\}$

FUSING WITH STREAMS

$$\text{foldl}' f z \cdot \text{map } g \cdot \text{filter } h$$
$$= \text{foldlS}' f z \cdot \text{read} \cdot \text{write} \cdot \text{mapS } g \quad \{\text{inline } \text{foldl}', \text{map} \\ \cdot \text{read} \cdot \text{write} \cdot \text{filterS } h \cdot \text{read} \quad \text{and } \text{filter}\}$$
$$= \text{foldlS}' f z \cdot \text{mapS } g \cdot \text{filterS } h \cdot \text{read} \quad \{\text{read/write fusion}\}$$

- The original 3 loops, and 2 intermediate arrays, are *automatically* transformed into a single traversal.
- GHC then further inlines and combines transformers, eliminating *Step* values.
- 2.4 × faster

EXTENSIONS

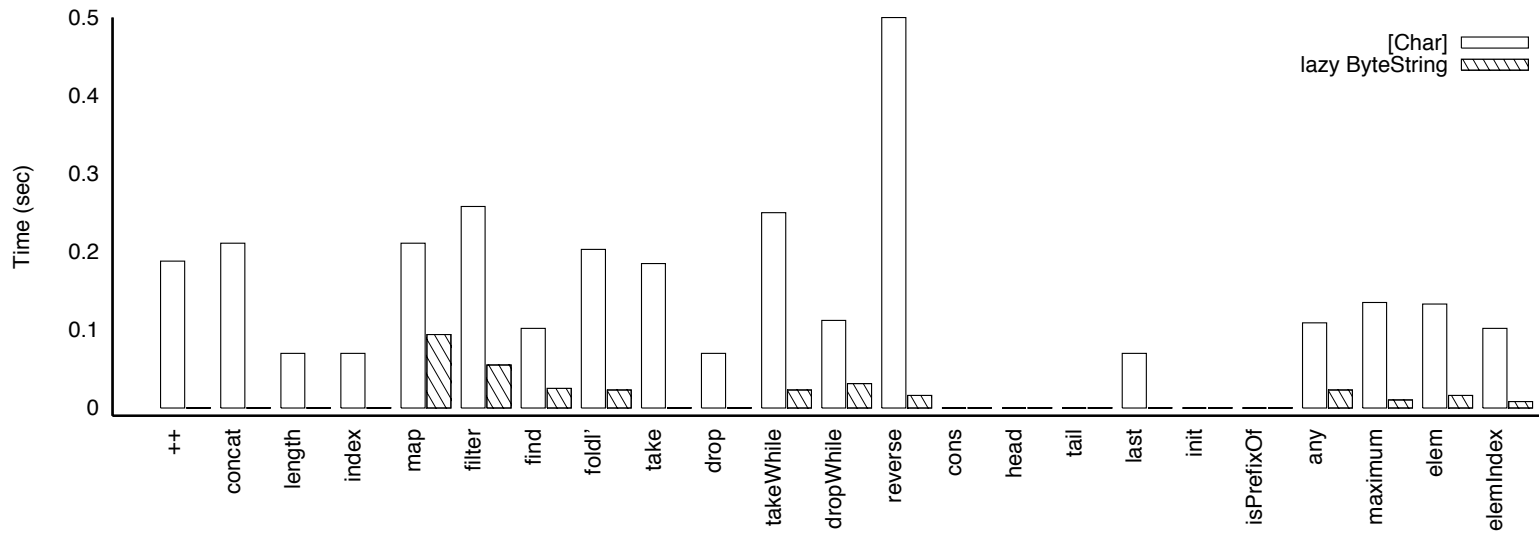
Over 40 functions in the ByteString library fuse, including:

- both up and down traversals (*foldl*, *foldr*, *scanl*, *scanr*)
- up and down loops will fuse with “bidirectional” loops like *map*
- *lazy* bytestrings, combining strict array chunks in the L2 cache, with a lazy spine

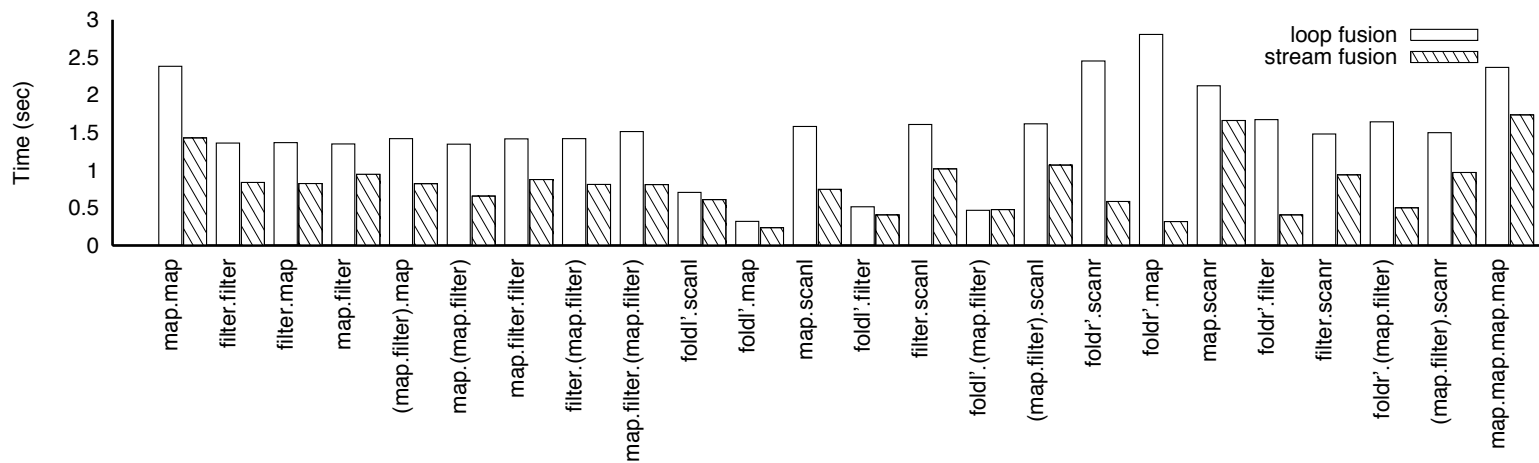
More details: see the paper.

RESULTS

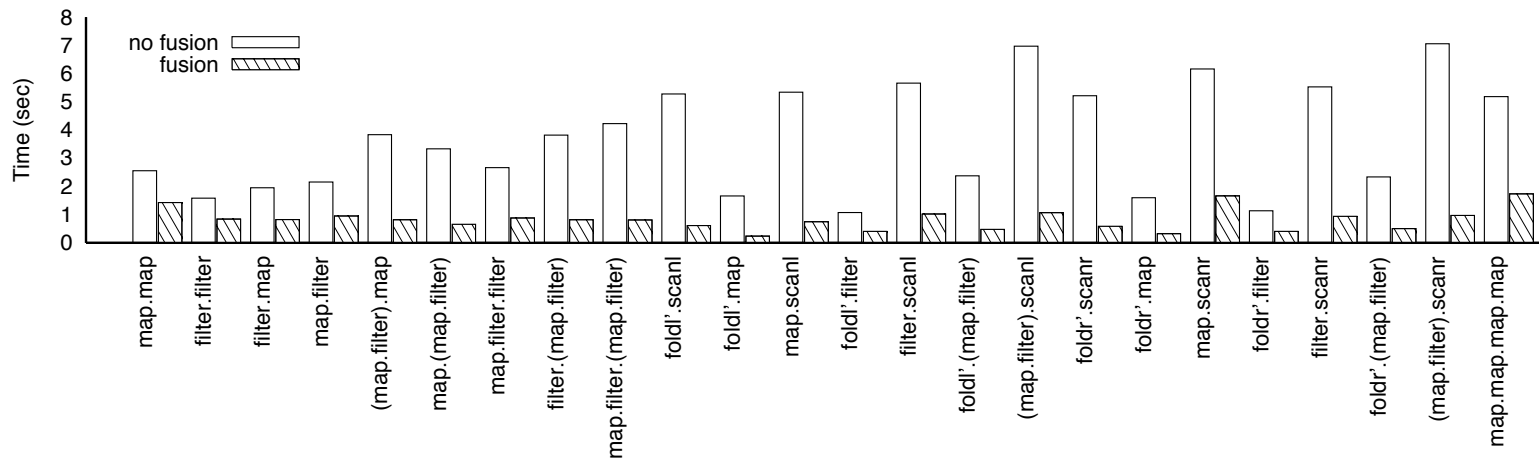
[Char] AND LAZY BYTESTRING : RUNNING TIME



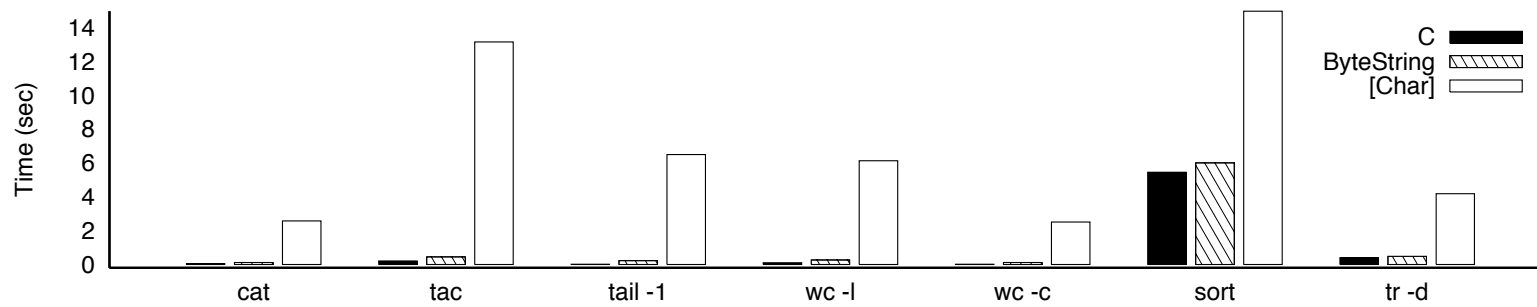
COMPARATIVE FUSION STRATEGIES



EFFECT OF FUSION: STREAMS WITH AND WITHOUT FUSION



C, BYTESTRING AND [CHAR] UNIX TOOLS



FUTURE

- Stream fusion for Haskell lists
- Improving the GHC backend for tight loops and branch prediction
- Efficiently fusing multiple traversals: `append`, `zip`, `concatMap`
- Fusible Binary serialisation
- Fusion across IO boundaries
- Establish the safety of streams in the presence of *seq*
- Rewrite Great Language Shootout entries, and take back 1st place! ;-)
- *Done*: Fusible *Storable a* \Rightarrow *Vector a* (Spencer Janssen for Summer of Code)

DATA.BYTESTRING



- Home page : <http://www.cse.unsw.edu.au/~dons/fps.html>
- Available with the current version of Hugs
- Available with GHC 6.6

