# The Art of
# Incremental Stream Processing

@markhibberd

I have a problem

```
q:Quit   d:Del   u:Undel   s:Save   m:Mail

INBOX [52]   10/6422    |   [   +]   04/18
archive 34982/52900     |   [   +]   04/18
sent [18]    38/ 840    |   [   +]   04/18
drafts        1/  54    |   [   F]   04/18
                        |   [   !]   04/18
                        |   [   +]   04/18
                        |   [    ]   04/18
                        |   [    ]   04/18
```

~150k Emails

~4GB Emails

**100s** of new email / day

~5 I want to Read

~2 I want to Reply To

Messages **delivered** where and when I need them

# Ability to **locate** important messages from days past

Lennart Augustsson wrote

```
> main = do
>     name:_ <- getArgs
>     file <- readFile name
>     print $ length $ lines file
```

Given the stance against top-level mutable variables, I have not
expected to see this Lazy IO code. After all, what could be more against
the spirit of Haskell than a `pure' function with observable side
effects. With Lazy IO, one indeed has to choose between correctness
and performance. The appearance of such code is especially strange
after the evidence of deadlocks with Lazy IO, presented on this list
less than a month ago. Let alone unpredictable resource usage and
reliance on finalizers to close files (forgetting that GHC does not
guarantee that finalizers will be run at all).

Is there an alternative?

Lennart Augustsson wrote

```
> main = do
>     name:_ <- getArgs
>     file <- readFile name
>     print $ length $ lines file
```

Given the stance against top-level mutable variables, I have not
expected to see this Lazy IO code. After all, what could be more against
the spirit of Haskell than a `pure' function with observable side
effects. With Lazy IO, one indeed has to choose between correctness
and performance. The appearance of such code is especially strange
after the evidence of deadlocks with Lazy IO, presented on this list
less than a month ago. Let alone unpredictable resource usage and
reliance on finalizers to close files (forgetting that GHC does not
guarantee that finalizers will be run at all).

Is there an alternative?

# Zeitgeist

```
13 find "$MAILDIR" -type f       | \
14   xargs -n 1 stat -f "%m|%N"   | \
15   sort -n                       | \
16   cut -d '|' -f 2               | \
17   xargs grep -l "$QUERY"        | \
18   head -5                       | \
19   xargs less
```

```haskell
 5 data Email =
 6   Email { date :: Int,  content :: String }
 7   deriving (Show, Eq)
 8
 9 search :: String -> [Email] -> [Email]
10 search term =
11   take 5
12   . filter (isInfixOf term . content)
13   . sortBy (compare `on` date)
```

# Reality Calling

"--mmap
    Use mmap(2) instead of read(2)
    to read input, which can result
    in better performance under
    some circumstances but can
    cause undefined behaviour."
                    — $(man grep)

```
275  struct file *
276  grep_open(const char *path)
277  {
278      struct file *f;
279
280      f = grep_malloc(sizeof *f);
281      memset(f, 0, sizeof *f);
282      if (path == NULL) {
283          /* Processing stdin implies --line-buffered.
284          */
285          lbflag = true;
286          f->fd = STDIN_FILENO;
287      } else if ((f->fd = open(path, O_RDONLY)) == -
288          1)
289          goto error1;
290
291      if (filebehave == FILE_MMAP) {
292          struct stat st;
293
294          if ((fstat(f->fd, &st) == -1) || (st.st_size >
295      OFF_MAX) ||
296              (!S_ISREG(st.st_mode)))
297              filebehave = FILE_STDIO;
298          else {
299              int flags = MAP_PRIVATE | MAP_NOCORE |
300              MAP_NOSYNC;
301  #ifdef MAP_PREFAULT_READ
302              flags |= MAP_PREFAULT_READ;
303  #endif
304              fsiz = st.st_size;
305              buffer = mmap(NULL, fsiz, PROT_READ, flags,
306                  f->fd, (off_t)0);
307              if (buffer == MAP_FAILED)
308                  filebehave = FILE_STDIO;
309              else {
310                  bufrem = st.st_size;
311                  bufpos = buffer;
312                  madvise(buffer, st.st_size, MADV_SEQUENTIAL)
313              ;
314              }
315          }
316      }
317
318      if ((buffer == NULL) || (buffer == MAP_FAILED))
319          buffer = grep_malloc(MAXBUFSIZ);
320
321      if (filebehave == FILE_GZIP &&
322          (gzbufdesc = gzdopen(f->fd, "r")) == NULL)
323          goto error2;
324
325  #ifndef WITHOUT_BZIP2
326      if (filebehave == FILE_BZIP &&
327          (bzbufdesc = BZ2_bzdopen(f->fd, "r")) ==
328          NULL)
329          goto error2;
330  #endif
331
332      /* Fill read buffer, also catches errors early
333      */
334      if (bufrem == 0 && grep_refill(f) != 0)
335          goto error2;
336
337      /* Check for binary stuff, if necessary */
338      if (binbehave != BINFILE_TEXT && memchr(bufpos,
339          '\0', bufrem) != NULL)
340      f->binary = true;
341
342      return (f);
343
344  error2:
345      close(f->fd);
346  error1:
347      free(f);
348      return (NULL);
349  }
350
351
```

"With Lazy IO, one indeed has
to choose between correctness
and performance."

— Oleg Kiselyov

```haskell
 5 type Maildir =
 6   FilePath
 7
 8 data Email =
 9   Email { date :: Int,  content :: String }
10   deriving (Show, Eq)
11
12 search :: String -> Maildir -> IO [Email]
13 search term =
14   {— oh noes! It's so horrible
15                  I can't even show it —}
```

Is there an alternative?

# Intuition 1: A Language

I Need To Produce Values

```
1 type In i
```

I Need To Consume Values

```
1 type In i
2
3 type Out o
```

I Need To **Transform** Values

```
1 type In i
2
3 type Out o
4
5 data Pipeline i o
```

I May Have **Effects**

```
1 type In i m
2
3 type Out o m
4
5 data Pipeline i o m
```

I May **Compute** A Value

```
1 type In i m a
2
3 type Out o m a
4
5 data Pipeline i o m a
```

# A (Simple) Interface

```
1 type In i m a = Pipeline i () m a
2
3 type Out o m a = Pipeline Void o m a
4
5 data Pipeline i o m a
```

```haskell
1 type In i m a = Pipeline i () m a
2
3 type Out o m a = Pipeline Void o m a
4
5 data Pipeline i o m a
6   = Done a
7   | Yield o (Pipeline i o a)
8   | Await (i -> Pipeline i o a)
```

# Intuition 2: Pipelines

Mail

Poll

Prepare

Classify

Deliver

Mail

Poll

Prepare

Classify

Deliver

# Getting Real

```
1 data Proxy i' i o' o m a
2
3 type Producer i m a = Proxy X () () o m a
4 type Consumer o m a = Proxy () i () X m a
5 type Pipe i o m a   = Proxy () i () o m a
```

```
1 type In i m a
2
3 type Out o m a
4
5 type Pipeline i o m a
```

Pipes

```
1 data Proxy i' i o' o m a
2
3 type Producer i m a = Proxy X () () o m a
4 type Consumer o m a = Proxy () i () X m a
5 type Pipe i o m a   = Proxy () i () o m a
```

**Explicit Input and Output**

**at Each Component**

Pipes

```
1 data Proxy i' i o' o m a
2
3 type Producer i m a = Proxy X () () o m a
4 type Consumer o m a = Proxy () i () X m a
5 type Pipe i o m a   = Proxy () i () o m a
```

Effects On Producers,

Consumers And Pipes

Pipes

```
1 data Proxy i' i o' o m a
2
3 type Producer i m a = Proxy X () () o m a
4 type Consumer o m a = Proxy () i () X m a
5 type Pipe i o m a   = Proxy () i () o m a
```

Can Terminate With A

Value Anywhere In Pipeline

Pipes

```
1 data Pipe l i o u m r
2
3 newtype ConduitM i o m r =
4   ConduitM { unConduitM :: Pipe i i o () m r }
5
6 type Source m o    = ConduitM () o    m ()
7 type Sink i m a    = ConduitM i  Void m a
8 type Conduit i m o = ConduitM i  o    m ()
```

```
1 type In i m a
2
3 type Out o m a
4
5 type Pipeline i o m a
```

Conduit

```
1 data Pipe l i o u m r
2
3 newtype ConduitM i o m r =
4    ConduitM { unConduitM :: Pipe i i o () m r }
5
6 type Source m o     = ConduitM () o    m ()
7 type Sink i m a     = ConduitM i  Void m a
8 type Conduit i m o = ConduitM i  o     m ()
```

Explicit Input and Output

at Each Component

Conduit

```
1 data Pipe l i o u m r
2
3 newtype ConduitM i o m r =
4   ConduitM { unConduitM :: Pipe i i o () m r }
5
6 type Source m o   = ConduitM () o    m ()
7 type Sink i m a   = ConduitM i  Void m a
8 type Conduit i m o = ConduitM i  o    m ()
```

Effects On Sources,

Sinks And Conduits

Conduit

```
1 data Pipe l i o u m r
2
3 newtype ConduitM i o m r =
4   ConduitM { unConduitM :: Pipe i i o () m r }
5
6 type Source m o    = ConduitM () o    m ()
7 type Sink i m a    = ConduitM i  Void m a
8 type Conduit i m o = ConduitM i  o    m ()
```

Can Only Terminate With A

Value On a Sink

Conduit

```
1  sealed abstract class Process[F[_],O]
2
3  type Process0[O]        = Process[Env[_,_]#Is, O]
4  type Process1[I, O]     = Process[Env[I,_]#Is, O]
5  type Sink[F[_], O]      = Process[F, O => F[Unit]]
```

```
1  type In i m a
2
3  type Out o m a
4
5  type Pipeline i o m a
```

Scalaz Streams

```
1 data Process m o
2
3 type Process0 o     = forall a. Process (Is a) o
4 type Process1 i o  = Process (Is i) o
5 type Sink m o       = Process m (o -> m ())
```

```
1 type In i m a
2
3 type Out o m a
4
5 type Pipeline i o m a
```

Scalaz Streams

```haskell
1 data Process m o
2
3 type Process0 o    = forall a. Process (Is a) o
4 type Process1 i o  = Process (Is i) o
5 type Sink m o      = Process m (o -> m ())
```

Model Request And Production

Rather Than Input and Output

Scalaz Streams

```
1 data Process m o
2
3 type Process0 o    = forall a. Process (Is a) o
4 type Process1 i o  = Process (Is i) o
5 type Sink m o      = Process m (o -> m ())
```

Effects Are Returned As

Values, Transducers are Pure

Scalaz Streams

```
1 data Process m o
2
3 type Process0 o    = forall a. Process (Is a) o
4 type Process1 i o  = Process (Is i) o
5 type Sink m o      = Process m (o -> m ())
6
7 runFoldMap :: (Monad m, Monoid b) =>
8                  Process m o -> (o -> m b) -> m b
```

Computation of Values

Modelled Externally

Scalaz Streams

:: Source Mail    :: Conduit Mail Features    :: Conduit Features Score    :: Sink Score

Poll    Prepare    Classify    Deliver

Conduit

`:: Process m Mail`

**Poll**

`:: Process1 Mail Features`

**Prepare**

`:: Process1 Features Score`

**Classify**

`:: Sink m Score`

**Deliver**

**Scalaz Stream**

# Horizontal Composition

# Mail Delivery

**Poll** **Prepare** **Classify** **Deliver**

Mail Delivery

Poll

Prepare

Classify

Deliver

:: In Event

:: Pipeline Mail Features

:: Pipeline Features Score

:: Out Event

Poll

Prepare

Classify

Deliver

(>|) :: Pipeline i o m a –> Pipeline o o' m a –> Pipeline i o' m a

:: Pipeline () Void

:: In Event

:: Pipeline Mail Features

:: Pipeline Features Score

:: Out Event

**Poll** >| **Prepare** >| **Classify** >| **Deliver**

(>|) :: Pipeline i o m a –> Pipeline o o' m a –> Pipeline i o' m a

eval :: Pipeline () Void m a -> m a

:: In Event

:: Pipeline Mail Features

:: Pipeline Features Score

:: Out Event

**Poll** >| **Prepare** >| **Classify** >| **Deliver**

(>|) :: Pipeline i o m a -> Pipeline o o' m a -> Pipeline i o' m a

:: Effect

:: Producer Event     :: Pipe Mail Features     :: Pipe Features Score     :: Consumer Score

Poll  >-->  Prepare  >-->  Classify  >-->  Deliver

Pipes

runEffect :: Effect m a -> m a

:: Producer Event

:: Pipe Mail Features

:: Pipe Features Score

:: Consumer Score

Poll >-> Prepare >-> Classify >-> Deliver

Pipes'

:: Source ()

:: Source Event

:: Conduit Mail Features

:: Conduit Features Score

:: Sink Score

Poll

=$=

Prepare

=$=

Classify

=$=

Deliver

Conduit

:: Source ()

:: Source Event

:: Conduit Mail Features

:: Conduit Features Score

:: Sink Score

Poll

$=

Prepare

=$=

Classify

=$

Deliver

Conduit'

`:: Process m ()`

`:: Process m Event`   `:: Process1 Mail Features`   `:: Process1 Features Score`   `:: Sink m Score`

**Poll** `|>` **Prepare** `|>` **Classify** `to` **Deliver**

Scalaz Stream

run :: Process m a -> m ()

:: Process m Event

:: Process1 Mail Features

:: Process1 Features Score

:: Sink m Score

**Poll** |> **Prepare** |> **Classify** to **Deliver**

Scalaz Stream'

# Is Composition About Combinators or Laws?

# Vertical Composition

Mail Delivery

Poll

Prepare

Classify

Deliver

Work

Throttle

Read

Home

Work

Throttle

Work

0

Throttle

2

Read

5

Home

Work

Throttle

**Work** 0

**Throttle** 1

**Read**

**Home** 5

**Work**

**Throttle**

Work

Home

Work

Home

Throttle

Read

Throttle

Work

Home **4**

Throttle **0**

Read

Work

Home

Throttle

Work

Home **4**

Throttle **0**

Read

Work

Home

Throttle

Work

Home

Throttle

4

Work

3

Throttle

Read

Home

Work

Throttle

`:: Producer Event`

`:: Pipe Event Event`

`:: Pipe Event Mail`

**Work** `>->` **Throttle** `>->` **Read**

`>>=`

**Home**

`forever`

**Work** **Throttle**

Pipes

**Work**

$=

**Throttle**

=$=

**Read**

>>=

**Home**

forever

**Work**

**Throttle**

**Conduit**

:: Process m Event

:: Process1 Event Event

:: Process1 Event Mail

**Work**

\>|

**Throttle**
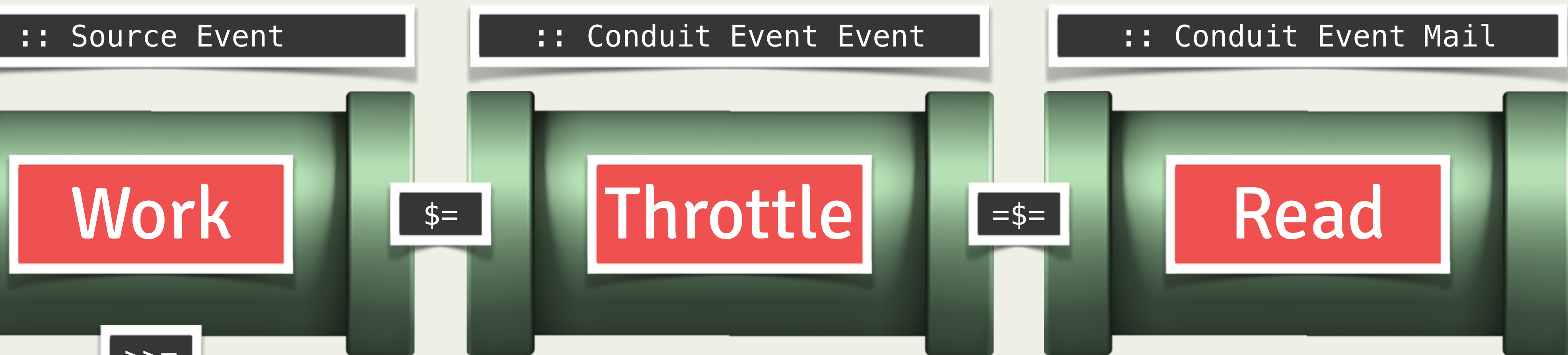
\>|

**Read**

fby

**Home**

repeat

**Work**

**Throttle**

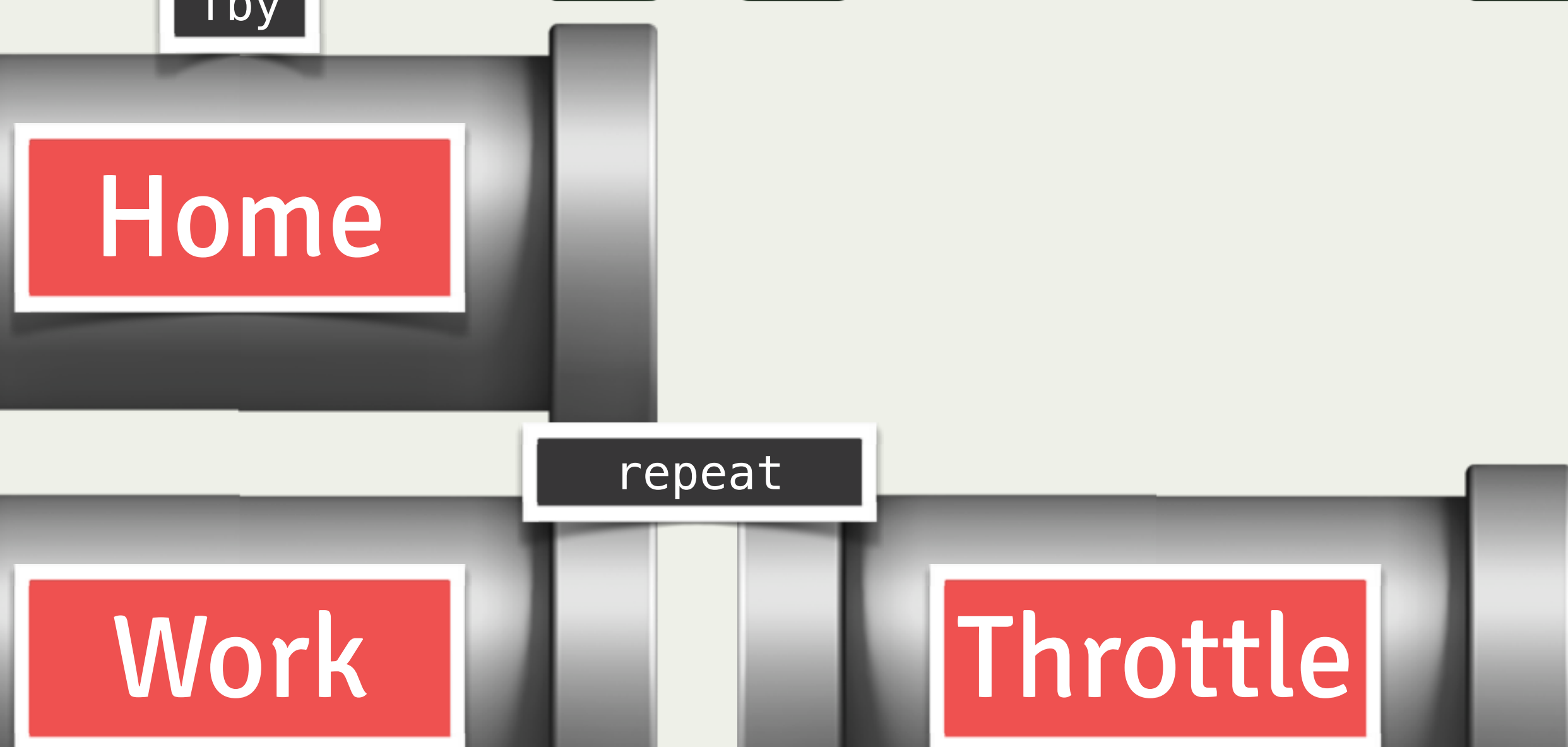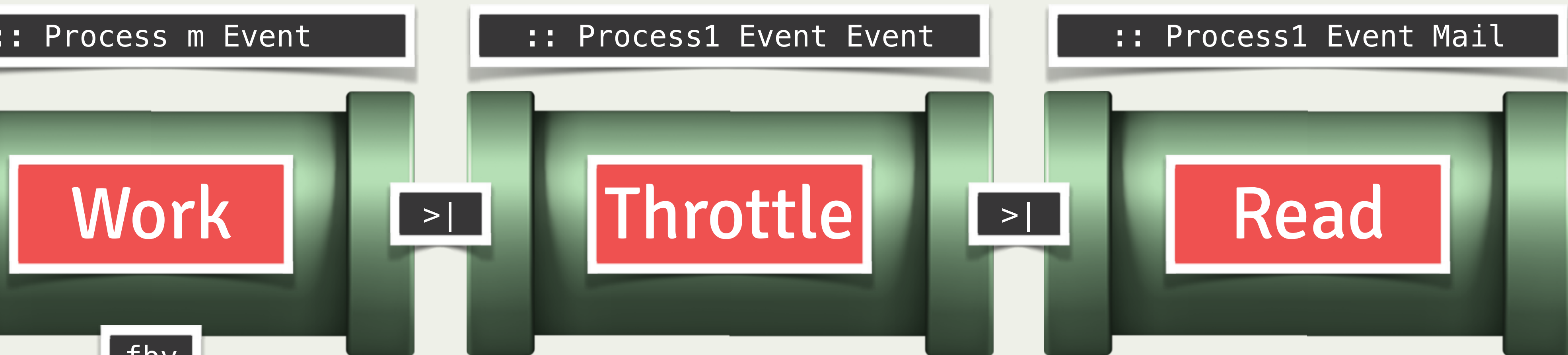**Scalaz Stream**

# Intuition 3: Parsers

```haskell
type In i m a = Pipeline i () m a

type Out o m a = Pipeline Void o m a

data Pipeline i o m a
  = Done a
  | Yield o (Pipeline i o a)
  | Await (i -> Pipeline i o a)
```

```haskell
type In i m a = Pipeline i () m a

type Out o m a = Pipeline Void o m a

data Pipeline i o m a
  = Done a
  | Yield o (Pipeline i o a)
  | Await (i -> Pipeline i o a)

yield :: o -> Pipeline i o m ()
yield = Yield o (Done ())

await :: Pipeline i o m i
await = Await Done
```

```
one :: Pipeline i i m ()
one = do
  i <- await
  yield i

cat :: Pipeline i i m ()
cat = forever one

pairs :: Pipeline i (i, i) m ()
pairs = forever $ do
  i1 <- await
  i2 <- await
  yield (i1, i2)
```

```haskell
counter :: Monad m => Pipeline i (Int, i) m ()
counter = flip evalStateT 0 . forever $ do
  i <- lift await
  n <- get
  lift . yield $ (n, i)

filter :: (i -> Bool) -> Pipeline i i m a
filter f = forever $ do
  i <- await
  when (f i) $ yield i
```

**Pipes**

```
1 yield :: o –> Pipe i o m ()
2
3 await :: Pipe i o m i
```

**Conduit**

```
1 yield :: o –> ConduitM i o m r
2
3 await :: ConduitM i o m (Maybe i)
4
5 awaitForever :: (\i –> ConduitM i o m a)
6                 –> ConduitM i o m ()
```

**Scalaz Stream**

```
1 emit :: o –> Process f o
2
3 await1 :: Process1 i i
```

**Subtlety** Fights Back

# Internal vs External Management of Resources

# Layered Streams

# Constant Memory Streaming

# How much does elegance cost?

to be continued...

@markhibberd