Polarized Data Parallel Data Flow <u>Ben Lippmeier</u>, Fil Mackay, Amos Robinson FHPC 2016/09/17

Image: A♥ / Aih. flickr. CC Generic.

do ls <- readLines "huge.txt"
 writeLines "out.txt" ls
 let xs = map read ls
 let total = sum xs
 print total</pre>

do ls <- readLines "huge.txt"
 writeLines "out.txt" ls
 print (sum (map read ls))</pre>





















readLines : Name -> m (S String)
writeLines : Name -> m (K String)

![](_page_13_Figure_0.jpeg)

readLines : Name -> m (S String)
writeLines : Name -> m (K String)
dup : S a -> K a -> S a

![](_page_14_Figure_0.jpeg)

readLines	:	Name -> m (S String)
writeLines	:	Name -> m (K String)
dup	:	<b>S</b> a -> <b>K</b> a -> <b>S</b> a
sum	•	m (K a, Ref Nat)

![](_page_15_Figure_0.jpeg)

data	Source	m	а	=	•••
data	Sin <mark>K</mark>	m	а	=	•••

readLines	: Name -> m (S String)
writeLines	: Name -> m (K String)
dup	: S a -> K a -> S a
sum	: m (K a, Ref Nat)
map : (a	-> b) -> K b -> K a

![](_page_16_Figure_0.jpeg)

readLines	: Name -> m (S String)
writeLines	: Name -> m (K String)
dup	: S a -> K a -> S a
sum	: m (K a, Ref Nat)
map : (a	-> b) -> K b -> K a
drain	: S a -> K a -> m ()

![](_page_17_Figure_0.jpeg)

```
data Source m a = ...
data SinK m a = ...
```

readLines	: Name -> m (S String)
writeLines	: Name -> m (K String)
dup	: S a -> K a -> S a
sum	: m (K a, Ref Nat)
map : (a	-> b) -> K b -> K a
drain	: S a -> K a -> m ()

```
do ls <- readLines "huge.txt"
  (ref, xs) <- sum
  k1 <- map read xs
  k2 <- writeLines "out.txt"
  s1 <- dup ls k1
  drain s1 k2
  total <- readRef ref
  print total</pre>
```

![](_page_18_Figure_0.jpeg)

![](_page_19_Figure_0.jpeg)

"Polarity Versions" :: a :: a +push to input pull from output induces induces map i map o push to output pull from input +"pushy" "pully" :: b :: b

dup\_iio + dup\_iio +

![](_page_22_Picture_1.jpeg)

![](_page_23_Figure_0.jpeg)

![](_page_23_Figure_1.jpeg)

![](_page_24_Figure_0.jpeg)

![](_page_25_Figure_0.jpeg)

![](_page_26_Figure_0.jpeg)

![](_page_27_Picture_0.jpeg)

![](_page_28_Figure_0.jpeg)

![](_page_29_Figure_0.jpeg)

![](_page_30_Figure_0.jpeg)

![](_page_31_Figure_0.jpeg)

![](_page_32_Picture_0.jpeg)

OK

![](_page_32_Picture_2.jpeg)

![](_page_32_Picture_3.jpeg)

![](_page_33_Figure_0.jpeg)

, kEject :: ix -> m () }

```
map_i :: Monad m
       => (ix -> a -> b)
       -> Sources ix m a -> m (Sources ix m b)
map i f (Sources n pullA)
 = return (Sources n pullB)
where
       pullB i eatB eject
        = pullA i eatA eject
        where
              eatA x = eatB (f i x)
```

![](_page_35_Figure_0.jpeg)

![](_page_36_Figure_0.jpeg)

![](_page_37_Figure_0.jpeg)

naturally concurrent input streams are contending for a shared output

uncontrolled order of consumption elements pushed in non-deterministic order naturally sequential read from the input streams one after the other

controlled order of consumption drain entire stream first, or round robin element-wise

![](_page_37_Picture_5.jpeg)

"Drain Fattening"

(funnel\_i s >>=  $\lambda$ s'. drainP s' k) => (funnel\_o (arity s) k >>=  $\lambda$ k'. drainP s k')

![](_page_38_Figure_2.jpeg)

## repa-flow package (on Hackage)

Used in production at Vertigo, financial data processing ~ few GBs.

Implementation uses chunked streams. Absolute performance primarily depends on the library used to process the chunks.

## Questions?

![](_page_41_Figure_0.jpeg)

![](_page_42_Figure_0.jpeg)

## "operator is in control"

![](_page_43_Figure_0.jpeg)

"operator is in control"

"context is in control"

![](_page_44_Figure_0.jpeg)

"operator is in control"

"context is in control"

# Comparison

Image: Leo Reynolds.flickr. CC-NC-SA.

## conduit - Michael Snoyman

![](_page_46_Figure_1.jpeg)

- Pipe is an instance of Monad.
- Data can flow both ways through the pipe, and yield a final result.
- Single stream, single element at a time.
- Individual Sources created by 'yield' action.
- Combine pipes/conduits with fusion operators.

#### pipes - Gabriel Gonzelez

![](_page_47_Figure_1.jpeg)

- = Request a' (a -> Proxy a' a b' b m r)
  | Respond b (b' -> Proxy a' a b' b m r)
  | M (m (Proxy a' a b' b m r))
  | Pure r
- Proxy / Pipe is an instance of Monad.
- Data can flow both ways through the pipe, and yield a final result.

#### machines - Edward Kmett

newtype MachineT m k o

- = MachineT
- { runMachine :: m (Step k o (MachineT m k o))

type Machine k o

= forall m. Monad m => MachineT m k o

type Process a b = Machine (Is a) b)

#### **type** Source b = forall k. Machine k b

- Like streams as used in Data.Vector stream fusion, except the step function returns a whole new Machine (stream)
- Clean and general API, but not sure if it supports array fusion. Machines library does not seem to attempt fusion.

#### repa-flow vs others

- Repa flow provides chunked, data parallel database-like operators with a straightforward API.
- Sources and Sinks are values rather than computations. The "Pipe" between them created implicitly in IO land.
- API focuses on simplicity and performance via stream and array fusion, rather than having the most general API.
- Suspect we could wrap single-stream Repa flow operators as either Pipes or Conduits, but neither of the former seem to naturally support data parallel flows.

![](_page_50_Picture_0.jpeg)

![](_page_51_Figure_0.jpeg)