

# Experimenting with Distributed Data Processing in Haskell

Utku Demir

YOW! Lambda Jam 2019

- ▶ 2014 - 2016: **Picus Security**
  - ▶ A distributed Haskell application to simulate cyber attacks
- ▶ 2016 - now: **Movio**
  - ▶ A SaaS to analyze moviegoer data
  - ▶ Billions of transactions, millions of people
  - ▶ Large-scale data analytics

# Distributed Data Processing

- ▶ “Big Data”
  - ▶ Terabytes?
- ▶ Unknown set of complex queries or transformations

## Resilient Distributed Datasets (Apache Spark)

- ▶ **Dataset** is a multiset of rows.
- ▶ A Dataset is stored as many **Partition**'s.
- ▶ Partitions can be processed in parallel on many **executor**'s
- ▶ **Driver** coordinates executors.
- ▶ High-level combinators for transforming a Dataset
  - ▶ map, filter, groupBy, reduce eg.

- ▶ Written in Haskell.
- ▶ Borrows ideas from Apache Spark.
- ▶ Composes nicely with the existing Haskell ecosystem.
- ▶ Can be used with “Function as a Service” / “Serverless” offerings.

## -XStaticPointers

- ▶ -XStaticPointers

```
static (\x -> x + 1) :: StaticPtr (Int -> Int)
```

- ▶ distributed-closure

```
data Closure a
```

```
closure :: StaticPtr a -> Closure a
```

```
cpure   :: Serializable a => a -> Closure a
```

```
cap     :: Closure (a -> b) -> Closure a -> Closure b
```

## In detail

```
mkPartition :: Closure (ConduitT () a (ResourceT IO) ())  
             -> Partition a
```

```
data Dataset a where
```

```
  DExternal  :: StaticSerialise a => [Partition a] -> Dataset a
```

```
  DPipe      :: (StaticSerialise a, StaticSerialise b)  
             => Closure (ConduitT a b (ResourceT IO) ())  
             -> Dataset a -> Dataset b
```

```
  DPartition :: (StaticHashable k, StaticSerialise a)  
             => Int  
             -> Closure (a -> k)  
             -> Dataset a -> Dataset a
```

# Aggregations

```
-- Aggregates many a's to a single b.
```

```
data Aggr a b = ...
```

```
instance StaticApply (Aggr a)
```

```
instance StaticProfunctor Aggr
```

```
aggrFromMonoid :: StaticSerialise a
```

```
    => Closure (Dict (Monoid a))
```

```
    -> Aggr a a
```

```
dAggr :: Aggr a b -> Dataset a -> IO b
```

```
dGroupedAggr :: StaticHashable k
```

```
    => (a -> k) -> Aggr a b
```

```
    -> Dataset (k, b)
```



# Backend

- ▶ LocalProcessBackend
- ▶ distributed-dataset-aws
  - ▶ Uses AWS Lambda to run executors and S3 to exchange information.
  - ▶ Scales well, cost-effective
  - ▶ No infrastructure necessary

## Example

```
ghArchive (fromGregorian 2018 1 1, fromGregorian 2018 12 31)
  & dConcatMap (static (\e ->
    let author = e ^. gheActor . ghaLogin
        commits = e ^.. gheType . _GHPushEvent
                . ghpepCommits . traverse . ghcMessage
    in map (author, ) commits
  ))
  & dFilter (static (\(_, commit) ->
    "cabal" `T.isInfixOf` T.toLower commit
  ))
  & dGroupedAggr 50 (static fst) dCount
  & dAggr (dTopK (static Dict) 20 (static snd))
  >>= mapM_ (liftIO . print)
```

- ▶ 126 GB compressed, 909 GB uncompressed JSON
- ▶ 2190 executors
- ▶ ~ 2 minutes, including the time to deploy infrastructure

## Alternatives in Haskell (as far as I know)

### Sparkle

*A library for writing resilient analytics applications in Haskell that scale to thousands of nodes, using Spark and the rest of the Apache ecosystem under the hood.*

### HSpark

*A port of Apache Spark to Haskell using distributed process*

# Thanks!

- ▶ <https://github.com/utdemir/distributed-dataset>
- ▶ [me@utdemir.com](mailto:me@utdemir.com)

Questions?



## An external data source

```
import Network.HTTP.Simple
import Data.Conduit.Zlib (ungzip)
import Data.Conduit.JSON.NewlineDelimited as NDJ

data GHEvent = ... deriving FromJSON

urlToPartition :: String -> Partition GHEvent
urlToPartition url' = mkPartition $ static (\url -> do
  req <- parseRequest url
  httpSource req getResponseBody
    .| ungzip
    .| NDJ.eitherParser @_ @GHEvent
    .| C.mapM (either fail return)
) `cap` cpure (static Dict) url'
```

## How to aggregate

```
input & groupedAggr 3 (static getColor) (dSum (static Dict))
```

- ▶ Input

- ▶ Partition 1: [3, 5, 2, 1]
- ▶ Partition 2: [3, 7, 2]
- ▶ Partition 3: [1, 2, 8]

- ▶ Aggregation Step 1:

- ▶ Partition 1: [5, 5, 1]
- ▶ Partition 2: [5, 7]
- ▶ Partition 3: [1, 2, 8]

- ▶ Shuffle!

- ▶ Partition 1: [5, 7, 2]
- ▶ Partition 2: [5, 1, 1]
- ▶ Partition 3: [5, 8]

- ▶ Aggregation Step 2:

- ▶ Partition 1: [14]
- ▶ Partition 2: [6, 1]
- ▶ Partition 3: [13]

## Composing Aggr's

```
dConstAggr :: (Typeable a, Typeable t)
            => Closure a -> Aggr t a
```

```
dSum :: StaticSerialise a
      => Closure (Dict (Num a)) -> Aggr a a
```

```
dCount :: Typeable a => Aggr a Integer
dCount = static (const 1) `staticLmap` dSum (static Dict)
```

```
dAvg :: Aggr Double Double
dAvg = dConstAggr (static (/))
      `staticApply` dSum (static Dict)
      `staticApply` staticMap (static realToFrac) dCount
```



# StaticFunctor

```
class Functor f where
  fmap :: (a -> b) -> f a -> f b

class Typeable f => StaticFunctor f where
  staticMap :: (Typeable a, Typeable b)
             => Closure (a -> b) -> f a -> f b
```

## Results

```
("peti",899)
("haskell-pushbot",535)
("bgamari",418)
("phadej",307)
("23Skidoo",208)
("alanz",174)
("edolstra",141)
("quasicomputational",136)
("jneira",135)
("hvr",133)
("Ericson2314",133)
("felixonmars-bot",130)
("DanielG",129)
("philderbeast",120)
("coreyoconnor",115)
("rcaballeromx",109)
("...",100)
```