# COMS 4995 W: Parallel Functional Programming Parallel PageRank with MapReduce

Xi Yang, Zefeng Liu - `xy2390, zl2715`

December 16, 2019

## 1 Introduction

PageRank (PR) is an algorithm used by Google Search to rank web pages in their search engine results.[1] PageRank works by counting the number and quality of links to a page to determine a rough estimate of how important the website is. The underlying assumption is that more important websites are likely to receive more links from other websites.[2]

PageRank algorithm can be generalized to measure the importance of any type of recursive documents. It can be viewed as a node weight metric for complex networks including social networks, transportation networks, electricity networks, species networks, etc. The computing of PageRank is, therefore, a fundamental yet nontrivial problem.

In this project, we propose a parallel PageRank calculation program based on the MapReduce framework.

## 2 Problem Formulation

The PageRank algorithm simulates a random surfer traveling within a directed graph. Given the initial weight configuration of nodes, the algorithm outputs the probability (weight) distribution which represents the likelihood of a person randomly traveling through the edges will arrive at any particular node.

Now we formulate the Map/Reduce version of the PageRank problem.

The mapper receives the pair of node and pagerank as key, and the list of adjacent nodes as value. It maps those key-value pairs to either the pairs of node and pagerank increment or the pairs of node and list of adjacent nodes. The intermediate pairs are aggregated by key and fed to the reducers.

The reducer receives the pairs emitted by the mappers and aggregates the pagerank increments and calculates the updated pagerank value.

---

[1] Wikipedia contributors. "PageRank." Wikipedia, The Free Encyclopedia. Wikipedia, The Free Encyclopedia, 15 Nov. 2019. Web. 21 Nov. 2019.

[2] "Facts about Google and Competition". Archived from the original on 4 November 2011. Retrieved 12 July 2014.

# 3   Implementation

## 3.1   Data Type Definitions

In this section, we would like to introduce some shared building blocks, the data types, upon which both our sequential and parallel solutions are implemented.

### 3.1.1   Nodes

Here we used *String* to represent a general node in the concerned graph. Intuitively, *Nodes* are a *set* of *Node*. An instance of *Nodes* might be {"a", "b", "c"}.

### 3.1.2   Edges

In our implementation, *Edges* are a map for which the key type is *Node*, while the value type is a *list* of nodes, representing the nodes connected to the corresponding key node. An instance of *Edges* might be {"a": ["b", "c"], "b": ["c"]}.

### 3.1.3   Graph

The *Graph* data type represents a directed graph for whose nodes we would like to calculate the page rank values. The fields of this data type are,

- *nodes.* A set of all the nodes in this graph.

- *inEdges.* A map from a node to a list of nodes from which it is linked.

- *outEdges.* A map from a node to a list of nodes to which it links.

And there are some utility functions for this data type,

- *parseLine :: Graph -> String -> Graph.* Formulate an *inEdge* and an *outEdge* from the given *String*, add them to the given *Graph*, then return the newly constructed *Graph*. Each line of the input file should conform to the format 'fromNode toNode'.

- *fromContent :: String -> Graph.* Given a file content, apply *parseLine* to every line of the file content to construct a Graph.

- *fromFile :: String -> IO Graph.* Given a file name, utilize *fromContent* to construct a *IO Graph*.

### 3.1.4   PageRank

A *PageRank* data type is a *map* from a *Node* to its current *PageRankValue*, which is a *double* in our case.

    It also has some utility functions, such as the mapper and the reducer functions to compute the page rank values for a given graph with MapReduce, and also a sequential method to compute page rank. We will explain more about these utility functions in the following sections.

## 3.2    Sequential Solution

The function type is defined as *PageRank -> Graph -> Int -> Double -> PageRank*. We can interpret it as, "given initial PageRank, the corresponding graph, a number of iterations to compute, and a damping factor, returns the resulting PageRank after those iterations of computation in a sequential way".

     Our sequential solution to compute the *PageRank* values for the next iteration works in this way,

---

**1**   **for** *each node n in all the Nodes of the Graph* **do**
**2**      $pr\_n \leftarrow 0$
**3**      **for** *each edge (m, n) of n's inEdges* **do**
**4**          *num_of_out_nodes_m ← the number of nodes to which m links*
**5**          *pr_previous_m ← the previous page rank value of m*
**6**          *pr_delta_m ← pr_previous_m / num_of_out_nodes*
**7**          *pr_n ← pr_n + pr_delta_m*
**8**      **end**
**9**      *update the new page rank value of node n in the new PageRank data*
**10**   **end**
**11**   *one iteration of computation is completed, return the updated PageRank data*

---

## 3.3    Parallel Solution with customized MapReduce

The function type is also defined as *PageRank -> Graph -> Int -> Double -> PageRank*. And the interpretation is also similar, despite that this time the page rank values for the next iteration will be calculated in a parallel way.

     The function type of the *mapper* is defined as *mapper :: (PageRankValue, [Node]) -> PageRank*. For each *Node* in the *Graph*, the *mapper* takes its current *PageRankValue* and the list of nodes in its *outEdges*, then produces a *map* for which the key is each of the node in its *outEdges*, and the value is its contribution to that node, defined as its current *PageRankValue* divided by the number of nodes in its *outEdges*.

     The function type of the *reducer* is defined as *reducer :: [PageRank] -> PageRank*. The *reducer* merges all the outputs that the *mapper* produces. The merging rule is a simple addition for each same node.

     With these definitions, our customized mapReduce function is implemented as,

```
1    mapReduce :: (a -> b) -> ([b] -> c) -> [a] -> c
2    mapReduce mapper reducer input = pseq mapResult reduceResult
3      where
4        mapResult = parMap rpar mapper input
5        reduceResult = runEval (rpar $ reducer mapResult)
```

Once the reducer completed its work in one iteration, we could simply update the page rank value for each node as $(base + d * pr)$, where $base = (1 - d)/num\_of\_nodes\_in\_graph$, $d$ is the damping factor, $pr$ is the corresponding value the reducer produced.

## 3.4   Benchmark based on External MapReduce Library

We wanted to have an external benchmark with which to compare and evaluate our MapReduce based parallel PageRank implementation.

A short web search yielded Haskell-MapReduce (`https://github.com/jdstmporter/Haskell-MapReduce`, `https://wiki.haskell.org/MapReduce_as_a_monad`) to be a promising general-purposed MapReduce library. Therefore we implemented a benchmark based on the mentioned library.

The library is implemented in a monadic fashion such that mappers and reducers can be viewed as generalized transformers of type signature `a -> ([(s,a)] -> [(s',b)])`. It provides a wrapper function `liftMR` that converts the map / reduce function into a monadic function.[3]

Given the aforementioned MapReduce library, we only need to implement conventional mapper and reducer.

According to the specification of the library, mapper should take the form of `[s] -> [(s', a)]`, where `s` is input data, `s'` is output data and `a` is output key. We implemented the mapper such that `s = (fromNode, (pageRankValue, toNodes))` and `s' = (toNode, pageRankIncrement)`. Each of the input data emits its pagerank increment contribution to all of its `toNodes`.

The reducer is implemented in a similar fashion, it takes input of the form `[(toNode, pageRankIncrement)]`. For a particular `toNode`, the pagerank increment contribution from all `fromNodes` are aggregated together, producing the pagerank value.

The evaluations to be given later in this report showed that this external benchmark has a vastly worse performance compared with our implementation.

# 4   Evaluation

## 4.1   Settings

We performed our experiments on a *MacBook Pro (15-inch, 2018)*, of which the processor is *2.2 GHz 6-core Intel Core i7*, and the memory is *16 GB 2400 MHz DDR4*.

## 4.2   Experiment Results

We performed our experiments by performing 10 iterations of page rank computation on two datasets with different sizes.

The first dataset is a larger fraction of the Wikipedia Note Network[4], which is 90Kb large with 11515 edges. Table 1 shows the experiment results of our MapReduce implementation. Table 2 shows the experiment results of our sequential implementation and the benchmark implementation.

---

[3]`https://github.com/jdstmporter/Haskell-MapReduce`
[4]`https://snap.stanford.edu/data/wiki-Vote.html`

Table 1: Experiment Result for a 90Kb Dataset (MapReduce)

| N | time(s) | converted | gc'd | fizzled | total |
|---|---------|-----------|------|---------|-------|
| 1 | 59.92 | 0 | 4412 | 33038 | 37450 |
| 2 | 37.68 | 23322 | 858 | 13270 | 37450 |
| 3 | 35.86 | 28519 | 535 | 8396 | 37450 |
| 4 | 35.8 | 30830 | 346 | 4757 | 37450 |
| 5 | 36.09 | 32348 | 346 | 4757 | 37450 |
| 6 | 34.13 | 33376 | 300 | 3774 | 37450 |
| 7 | 35.62 | 33595 | 295 | 3560 | 37450 |
| 8 | 38.03 | 34186 | 256 | 3008 | 37450 |
| 9 | 40.4 | 34642 | 243 | 2565 | 37450 |
| 10 | 44.66 | 34850 | 226 | 2374 | 37450 |
| 11 | 44.29 | 35192 | 212 | 2046 | 37450 |
| 12 | 48.84 | 35487 | 191 | 1772 | 37450 |

Table 2: Experiment Result for a 90Kb Dataset (Sequential & Benchmark)

| N | time(s) |
|---|---------|
| seq | 173.62 |
| benchmark-1 | 1923.26 |
| benchmark-6 | 906.79 |

The second dataset is a smaller fraction of the Wikipedia Note Network, which is 40Kb large with 5508 edges. Table 3 shows the experiment results of our MapReduce implementation. Table 4 shows the experiment results of our sequential implementation and the benchmark implementation.

Table 3: Experiment Result for a 40Kb Dataset (MapReduce)

| N | time(s) | converted | gc'd | fizzled | total |
|---|---------|-----------|------|---------|-------|
| 1 | 26.23 | 0 | 2371 | 24879 | 27250 |
| 2 | 21.96 | 15890 | 844 | 10516 | 27250 |
| 4 | 20.38 | 21898 | 407 | 4945 | 27250 |
| 6 | 19.71 | 23888 | 286 | 3076 | 27250 |
| 8 | 24.64 | 24641 | 238 | 2371 | 27250 |
| 10 | 26.32 | 25345 | 196 | 1709 | 27250 |
| 12 | 27.99 | 25883 | 141 | 1226 | 27250 |

## 4.3   Performance Analysis

From the results, we can conclude that our MapReduce implementation is much more efficient both than the sequential version and than the benchmark implementation.

Table 4: Experiment Result for a 40Kb Dataset (Sequential & Benchmark)

| N | time(s) |
|---|---|
| seq | 88.26 |
| benchmark-1 | 804.94 |
| benchmark-6 | 472.71 |

We can also observe that when $N = 6$, which is equal to the number of cores, the performance of our implementation is the best. If $N$ is set to be larger, even the convertion rate is increased, the overhead for parallelism is also increased, hence the consumed time becomes longer.

For furthur analysis, we scrutinized the event log for our MapReduce implementation running with the 40Kb dataset using *ThreadScope*. From the figure, we can observe that the bottleneck is the GC waiting time.
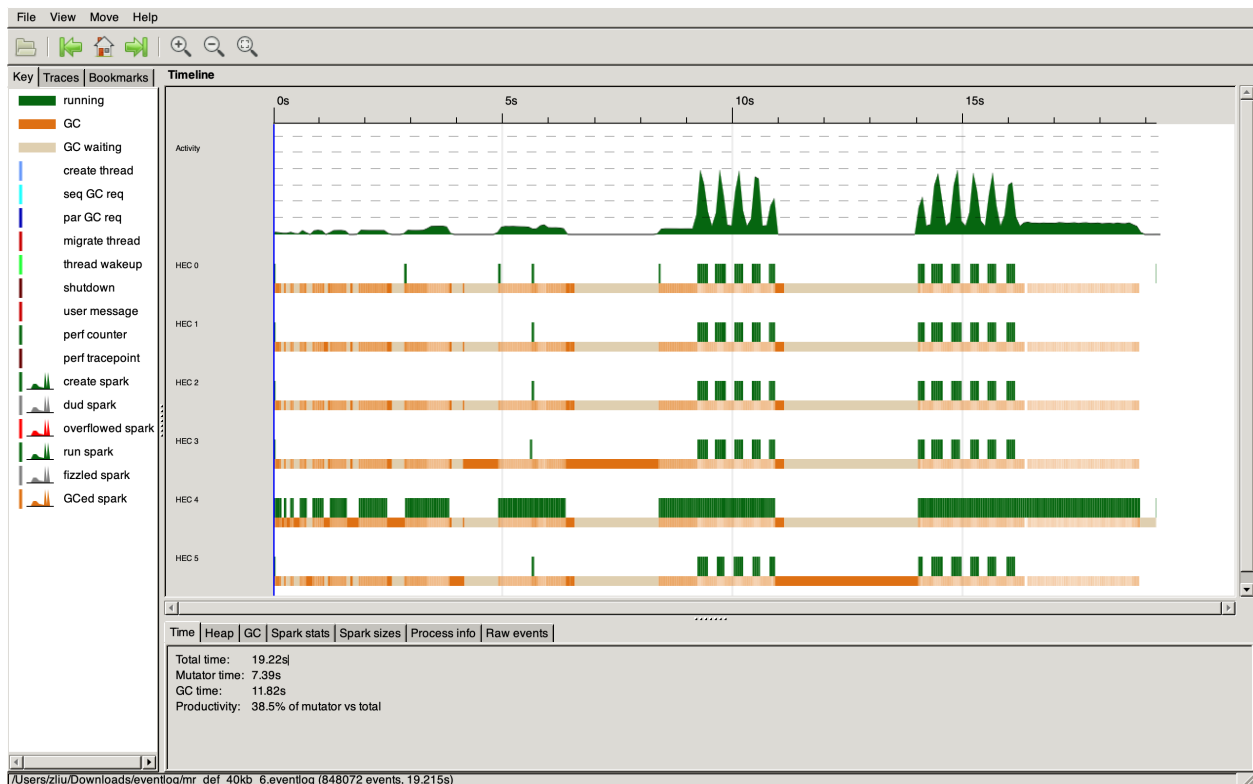


Figure 1: Eventlog for MapReduce experiment with 40Kb Dataset

# A    Code Listing

```haskell
1  {−
2
3  The main application program
4
5  Command line arguments: infilePath, outfilePath, itrs, [mode]
6
7  infilePath:   path of input file, which should be in the format of lines
8                    consisting of 'fromNode toNode'
9  outfilePath: path of output file
10 itrs:           number of iterations in the PageRank computation
11 mode:          optional, mode of PageRank computation, one of {seq, mr_def,
12                    mr_ext}, default to mr_def
13                    seq:    non−parallel sequential computation
14                    mr_def: parallel implementation based on default MapReduce
15                    mr_ext: benchmark parallel implementation based on external opensourced MapReduce library
16
17 −}
18
19 module Main (main) where
20
21 import Control.Monad (when)
22 import System.IO (openFile, IOMode(WriteMode), hPutStrLn, hClose)
23 import System.Environment (getArgs, getProgName)
24 import System.Exit
25 import Data.Map as M (toList)
26
27 import ProcessData (processData)
28 import PageRank (computePageRankSeq, computePageRankMR)
29 import PageRankExt (computePageRankMRext)
30
31 main :: IO()
32 main = do
33     progName <− getProgName
34     args  <− getArgs
35
36     when (length args /= 3 && length args /= 4) $
37         die $ "Usage: " ++ progName ++ " <infilePath> <outfilePath> <itrs> [mode], where\
38         \ mode is one of {seq, mr_def, mr_ext}, default to mr_def"
39
40     let infilePath : outfilePath : itrs : mode = args
41         computePageRank = case mode of
42             [] −> computePageRankMR
43             ["mr_def"] −> computePageRankMR
44             ["seq"] −> computePageRankSeq
45             ["mr_ext"] −> computePageRankMRext
46             _ −> error $ "Usage: " ++ progName ++ " <infilePath> <outfilePath> <itrs> [mode], where\
47             \ mode is one of {seq, mr_def, mr_ext}, default to mr_def"
```

```
48
49      (graph, pageRank) <- processData infilePath
50      let resPageRank = computePageRank pageRank graph (read itrs) 0.85
51      h <- openFile outfilePath WriteMode
52      mapM_ (hPutStrLn h) [ n ++ ": " ++ show pr | (n, pr) <- M.toList resPageRank ]
53      hClose h
```

Listing 1: app/Main.hs

```
1   {-
2
3       This module contains a utility  function,  which
4       1) reads  in  a graph from a input  file
5       2)  initializes   a PageRank data from the given graph
6       3) returns the  graph and the  initial  page rank
7
8   -}
9
10  module ProcessData
11  (processData) where
12
13  import Graph (Graph, fromFile)
14  import PageRank (PageRank, initFromGraph)
15
16  processData :: String -> IO (Graph, PageRank)
17  processData filename = do
18      graph <- fromFile filename
19      let pageRank = initFromGraph graph
20      return (graph, pageRank)
```

Listing 2: src/ProcessData.hs

```
1   {-
2
3   This module defines the Graph data type. The fields  are,
4   1) nodes:     set  of  all  the nodes in  this  graph
5   2) inEdges:  map from a node to a list  of nodes from which it  is  linked
6   3) outEdges: map from a node to a list  of nodes to  which it  links
7
8   This module also contains some utility  functions for  this  data type.
9
10  -}
11
12  module Graph
13  ( Graph(..)
14  , Node
```

```haskell
15    , fromFile
16  ) where
17
18  import qualified Data.Map as M (Map, insertWith, empty, keysSet)
19  import qualified Data.Set as S (Set, union, fromList, empty, toList, difference )
20  import System.IO (readFile)
21
22  type Node = String
23  type Nodes = S.Set Node
24  type InEdges = M.Map Node [Node]
25  type OutEdges = M.Map Node [Node]
26
27  data Graph = Graph { nodes    :: Nodes
28                     , inEdges  :: InEdges
29                     , outEdges :: OutEdges } deriving Show
30
31  -- Initial state  of an empty graph
32  empty :: Graph
33  empty = Graph S.empty M.empty M.empty
34
35  -- Read in a graph from a file
36  fromFile :: String -> IO Graph
37  fromFile filename = do
38      content <- readFile filename
39      return $ fromContent content
40
41  fromContent :: String -> Graph
42  fromContent content =
43      let ls = lines content
44      in postProcess $ foldl parseLine empty ls
45      where
46          postProcess :: Graph -> Graph
47          postProcess graph = foldl parseLine graph newLines
48              where
49                  ns = nodes graph
50                  sinkNodes = S. difference  ns $ M.keysSet $ outEdges graph
51                  newLines = [ n1 ++ " " ++ n2 |
52                  n1 <- S.toList sinkNodes, n2 <- S.toList ns, n1 /= n2 ]
53
54  {-
55    Parse each line of the input  file  as an edge in the  graph.
56    Each line of the input  file  should conform to the format 'fromNode toNode'.
57  -}
58  parseLine :: Graph -> String -> Graph
59  parseLine graph line =
60      let ws = words line
61      in case ws of
```

```
62        [fromNode, toNode] ->
63            Graph ns iEdges oEdges
64            where
65                ns = S.union (S.fromList ws) (nodes graph)
66                iEdges = M.insertWith (++) toNode [fromNode] (inEdges graph)
67                oEdges = M.insertWith (++) fromNode [toNode] (outEdges graph)
68        _ -> error "All lines of the input  file  \
69        \should be in the format of 'fromNode toNode'"
```

<div align="center">Listing 3: src/Graph.hs</div>

```
1  module MapReduce
2  ( mapReduce )
3
4  where
5
6  import Control.Parallel (pseq)
7  import Control.Parallel. Strategies (rpar, runEval, parMap)
8
9  mapReduce ::
10      (a -> b)          -- map function
11      -> ([b] -> c)     -- reduce function
12      -> [a]            -- list to map over
13      -> c
14  mapReduce mapper reducer input = pseq mapResult reduceResult
15    where mapResult = parMap rpar mapper input
16          reduceResult = runEval (rpar $ reducer mapResult)
```

<div align="center">Listing 4: src/MapReduce.hs</div>

```
1  {-
2
3  This module defines the PageRank data type, which is a map from a node to its
4  current page rank value.
5
6  This module also contains the empty definition and some utility functions for
7  this  data type, such as the mapper and the reducer functions to compute the page
8  rank values for a given graph with MapReduce, and also a sequential method to
9  compute page rank.
10
11 -}
12
13 module PageRank
14 ( PageRank
15 , initFromGraph
16 , computePageRankSeq
```

```
17  , computePageRankMR
18  ) where
19
20  import Graph (Graph(..), Node)
21  import MapReduce (mapReduce)
22  import qualified Data.Map as M (Map, empty, fromList, lookup, unionWith, toList)
23  import qualified Data.Set as S (toList, size )
24  import Data.Maybe (fromJust)
25
26  type PageRankValue = Double
27  type PageRank = M.Map Node PageRankValue
28
29  −− Initial state  of a PageRank data for an empty graph
30  empty :: PageRank
31  empty = M.empty
32
33  {−
34    Initial  state  of a PageRank data for a given graph, the page rank value
35    of each node is the  reciprocal  of  the  number of nodes in this  graph
36  −}
37  initFromGraph ::  Graph −> PageRank
38  initFromGraph graph =
39      let ns = nodes graph
40          pr = 1.0 / (fromIntegral $ S.size ns) in
41      M.fromList [ (n,  pr) | n <− S.toList ns ]
42
43  mapper :: (PageRankValue, [Node]) −> PageRank
44  mapper (pr, outNodes) =
45      let pr_ = pr / (fromIntegral $ length outNodes) in
46      M.fromList [ (n, pr_) | n <− outNodes ]
47
48  reducer  ::  [PageRank] −> PageRank
49  reducer  []  = empty
50  reducer  [x] = x
51  reducer  (x:xs) = M.unionWith (+) x (reducer xs)
52
53  {−
54    Given  initial  PageRank and the corresponding graph, a number of iterations
55    to compute, and a damping factor, returns the  resulting  PageRank after those
56    iterations  of computation in a  parallel  way with MapReduce
57  −}
58  computePageRankMR :: PageRank −> Graph −> Int −> Double −> PageRank
59  computePageRankMR pageRank _ 0 _ = pageRank
60  computePageRankMR pageRank graph itrs damping =
61      let nextPageRank = computeNextPageRankMR pageRank
62      in computePageRankMR nextPageRank graph (itrs−1) damping
63      where
```

11

```
64      computeNextPageRankMR :: PageRank -> PageRank
65      computeNextPageRankMR curPR =
66          let ns = S.toList $ nodes graph
67              input = map produceInput ns
68              produceInput n = (pr, outNodes)
69                  where
70                      pr = fromJust $ M.lookup n curPR
71                      outNodes = fromJust $ M.lookup n $ outEdges graph
72              mrResult = mapReduce mapper reducer input
73              base = (1 - damping) / (fromIntegral $ length ns)
74          in M.fromList [ (n, base + damping * pr) | (n, pr) <- M.toList mrResult ]

76  {-
77     Given initial PageRank and the corresponding graph, a number of iterations
78     to compute, and a damping factor, returns the resulting PageRank after those
79      iterations  of computation in a sequential way
80  -}
81  computePageRankSeq :: PageRank -> Graph -> Int -> Double -> PageRank
82  computePageRankSeq pageRank _ 0 _ = pageRank
83  computePageRankSeq pageRank graph itrs damping =
84      let nextPageRank = computeNextPageRank pageRank
85      in computePageRankSeq nextPageRank graph (itrs-1) damping
86      where
87          computeNextPageRank :: PageRank -> PageRank
88          computeNextPageRank curPR =
89              M.fromList [ (n, computePRValue n) | n <- ns ]
90              where
91                  ns = S.toList $ nodes graph
92                  iEdges = inEdges graph
93                  oEdges = outEdges graph
94                  computePRValue :: Node -> PageRankValue
95                  computePRValue n =
96                      let inNodes = fromJust $ M.lookup n iEdges
97                      in (1 - damping) / (fromIntegral $ length ns) + damping * (foldl sumUp 0 inNodes)
98                      where
99                          sumUp acc node =
100                             let numOutNodes = length $ fromJust $ M.lookup node oEdges
101                                 prValue = fromJust $ M.lookup node curPR
102                             in acc + prValue / (fromIntegral numOutNodes)
```

Listing 5: src/PageRank.hs

```
1  {-
2  External MapReduce Lib used to implement a benchmark
3  GitHub repository of the MapReduce library: https://github.com/jdstmporter/Haskell-MapReduce
4  -}
5
```

```
 6 {-# LANGUAGE MultiParamTypeClasses, FlexibleInstances #-}
 7
 8 -- | Module that defines the 'MapReduce' monad and exports the necessary functions.
 9 --
10 --  Mapper / reducers are generalised to functions of type
11 --  @a -> ([(s,a)] -> [(s',b)])@ which are combined using the monad's bind
12 --  operation.  The resulting monad is executed on initial data by invoking
13 --  'runMapReduce'.
14 --
15 --  For programmers only wishing to write conventional map / reduce algorithms,
16 --  which use functions of type @([s] -> [(s',b)])@ a wrapper function
17 --  'liftMR' is provided, which converts such a function into the
18 --  appropriate monadic function.
19 module MapReduceLibExt (
20 -- * Types
21         MapReduce,
22 -- * Functions
23 --
24 -- ** Monadic operations
25       return, (>>=),
26 -- ** Helper functions
27       run, distribute, lift ) where
28
29 import Data.List (nub)
30 import Control.Applicative ((<$>))
31 import Control.Monad (liftM)
32 import Control.DeepSeq (NFData)
33 import System.IO
34 import Prelude hiding (return,(>>=))
35 import Data.Digest.Pure.MD5
36 import Data.Binary
37 import qualified Data.ByteString.Lazy as B
38 import Control.Parallel.Strategies (parMap, rdeepseq)
39
40 -- | The parallel map function; it must be functionally identical to 'map',
41 --   distributing the computation across all available nodes in some way.
42 pMap :: (NFData b) => (a -> b) -- ^ The function to apply
43         -> [a]                 -- ^ Input
44         -> [b]                 -- ^ output
45 pMap = parMap rdeepseq
46
47 -- | Generalised version of 'Monad' which depends on a pair of 'Tuple's, both
48 --   of which change when '>>=' is applied.
49 class MonadG m where
50       return :: a                       -- ^ value.
51             -> m s x s a                -- ^ transformation that inserts the value
52                                         --   by replacing all
```

```
53                                         −−  the  key  values  with  the   specified
54                                         −−  value,  leaving  the  data  unchanged.
55
56
57         (>>=) :: (Eq b,NFData s'',NFData c) =>
58                 m s a s' b                −− ˆ Initial  processing  chain
59                 −> ( b −> m s' b s'' c )−− ˆ Transformation  to append  to  it
60                 −> m s a s'' c            −− ˆ Extended  processing  chain
61
62
63  −− | The  basic  type  that  provides  the  MapReduce  monad (strictly  a  generalised  monad).
64  −− In  the  definition
65  −− @(s,a)@  is the  type  of  the  entries  in  the   list   of  input  data  and  @(s',b)@
66  −− that  of  the  entries  in  the   list   of  output  data,  where  @s@  and  @s'@  are  data
67  −− and  @a@  and  @b@  are  keys.
68  −−
69  −− 'MapReduce'  represents  the  transformation  applied  to  data  by  one  or  more
70  −− MapReduce  staged.  Input  data  has  type  @[(s,a)]@  and  output  data  has  type
71  −− @[(s',b)]@  where  @s@  and  @s'@  are  data  types  and  @a@,  @b@  are  key  types.
72  −−
73  −− Its  structure  is  intentionally  opaque  to  application  programmers.
74  newtype MapReduce s a s' b = MR { runMR :: [(s,a)] −> [(s',b)] }
75
76  −− | Make  MapReduce  into  a  'MonadG'  instance
77  instance MonadG MapReduce where
78         return = ret
79         (>>=) = bind
80
81  −− | Insert  a  value  into  'MapReduce' by replacing all the  key  values  with the
82  −−   specified  value,  leaving  the  data  unchanged.
83  ret  ::  a                          −− ˆ value
84         −> MapReduce s x s a          −− ˆ transformation  that  inserts  the  value
85                                       −−  into  'MapReduce' by  replacing  all
86                                       −−  the  key  values  with  the   specified
87                                       −−  value,  leaving  the  data  unchanged.
88  ret  k = MR (\ss −> [(s,k) | s <− fst <$> ss])
89
90  −− ˆ Apply  a  generalised  mapper / reducer  to  the  end  of  a  chain  of  processing
91  −−   operations  to  extend  the  chain.
92  bind  :: (Eq b,NFData s'',NFData c) =>
93                 MapReduce s a s' b        −− ˆ Initial  state  of  the  monad
94         −> (b −> MapReduce s' b s'' c) −− ˆ Transformation  to  append  to  it
95         −> MapReduce s a s'' c          −− ˆ Extended  transformation  chain
96  bind  f  g = MR (\s −>
97         let
98                 fs = runMR f s
99                 gs = map g $ nub $ snd <$> fs
```

```haskell
100          in
101          concat $ pMap ('runMR' fs) gs
102
103  -- | Execute a MapReduce MonadG given specified initial data. Therefore, given
104  --   a 'MapReduce' @m@ and initial data @xs@ we apply the processing represented
105  --   by @m@ to @xs@ by executing
106  --
107  --   @run m xs@
108  run :: MapReduce s () s' b              -- ^ 'MapReduce' representing the required processing
109            -> [s]                        -- ^ Initial data
110            -> [(s',b)]                   -- ^ Result of applying the processing to the data
111  run m ss = runMR m [(s,()) | s <- ss]
112
113  -- | The hash_ function. Computes the MD5 hash_ of any 'Hashable' type
114  hash_ :: (Binary s) => s         -- ^ The value to hash_
115           -> Int                  -- ^ its hash_
116  hash_ s = sum $ map fromIntegral (B.unpack h)
117         where
118         h = encode (md5 $ encode s)
119
120  -- | Function used at the start of processing to determine how many threads of processing
121  --   to use.  Should be used as the starting point for building a 'MapReduce'.
122  --   Therefore a generic 'MapReduce' should look like
123  --
124  --   @'distribute' '>>=' f1 '>>=' ... '>>=' fn@
125  distribute :: (Binary s) => Int          -- ^ Number of threads across which to distribute initial data
126               -> MapReduce s () s Int -- ^ The 'MapReduce' required to do this
127  distribute n = MR (\ss -> [(s,hash_ s 'mod' n) | s <- fst <$> ss])
128
129  -- | The wrapper function that lifts mappers / reducers into the 'MapReduce'
130  --   monad. Application programmers can use this to apply MapReduce transparently
131  --   to their mappers / reducers without needing to know any details of the implementation
132  --   of MapReduce.
133  --
134  --   Therefore the generic 'MapReduce' using only traditional mappers and
135  --   reducers should look like
136  --
137  --   @'distribute' '>>=' 'lift' f1 '>>=' ... '>>=' 'lift' fn@
138  lift :: (Eq a) => ([s] -> [(s',b)])      -- traditional mapper / reducer of signature
139                                           -- @([s] -> [(s',b)]@
140          -> a                             -- the input key
141          -> MapReduce s a s' b            -- the mapper / reducer wrapped as an instance
142                                           -- of 'MapReduce'
143  lift f k = MR (\ss -> f $ fst <$> filter (\s -> k == snd s) ss)
```

Listing 6: src/MapReduceLibExt.hs

15

```haskell
1  {−
2  The benchmark PageRank computation implementation based on external opensourced MapReduce library
3  GitHub repository of the MapReduce library: https://github.com/jdstmporter/Haskell−MapReduce
4  −}
5
6  module PageRankExt
7  (
8  computePageRankMRext
9  ) where
10
11 import Graph (Graph(..), Node)
12 import qualified Data.Map as M (Map, empty, fromList, lookup, unionWith, toList)
13 import qualified Data.Set as S (toList, size )
14 import Data.Maybe (fromJust)
15 import MapReduceLibExt (run,distribute, lift ,(>>=))
16
17 type PageRankValue = Double
18 type PageRank = M.Map Node PageRankValue
19
20 empty :: PageRank
21 empty = M.empty
22
23 initFromGraph :: Graph −> PageRank
24 initFromGraph graph =
25      let ns = nodes graph
26          pr = 1.0 / (fromIntegral $ S.size ns) in
27      M.fromList [ (n, pr) | n <− S.toList ns ]
28
29 mr :: Double −> Double −> Int −> [(Node, (PageRankValue, [Node]))] −> [(Node, (PageRankValue, [Node]))]
30 mr damping numNodes n state = run f state
31      where
32          f = distribute  n MapReduceLibExt.>>= lift mapper MapReduceLibExt.>>= lift (reducer damping numNo
33
34 −− According to the specification  of Haskell−MapReduce lib
35 −− mapper should take the form of [s] −> [(s', a)]
36 −− where s is input data, s' is output data and a is output key
37 mapper :: [( Node, (PageRankValue, [Node]))] −> [((Node, (PageRankValue, [Node])), Node)]
38 mapper [] = []
39 mapper (x:xs) = parse x ++ mapper xs
40      where
41          parse (n, (pr, outNodes)) =
42              let pr_ = pr / (fromIntegral $ length outNodes)
43              in ((n, (0, outNodes)), n) : [ ((n_, (pr_,  [])),  n_) | n_ <− outNodes ]
44
45 −− According to the specification  of Haskell−MapReduce lib
46 −− reducer should take the form of [s '] −> [s'']
47 −− where s' is output data of mapper, s'' is output data of reducer
```

16

```
48 reducer :: Double -> Double -> [(Node, (PageRankValue, [Node]))] -> [(Node, (PageRankValue, [Node]))]
49 reducer _ _ [] = []
50 reducer damping numNodes xs@(x:_) =
51     [ foldl f (fst x, ((1 - damping) / numNodes, [])) xs ]
52     where f x y = (
53             fst x,
54             (
55                 (fst $ snd x) + damping * (fst $ snd y),
56                 (snd $ snd x) ++ (snd $ snd y)
57             )
58         )
59
60 computePageRankMRext :: PageRank -> Graph -> Int -> Double -> PageRank
61 computePageRankMRext pageRank _ 0 _ = pageRank
62 computePageRankMRext pageRank graph itrs damping =
63     let ns = S.toList $ nodes graph
64         numNodes = fromIntegral $ length ns
65         oEdges = outEdges graph
66         initMRinput = map toMRinput ns
67         toMRinput n =
68             let pr = fromJust $ M.lookup n pageRank
69                 outNodes = fromJust $ M.lookup n oEdges
70             in (n, (pr, outNodes))
71         mrOutput = mrItr initMRinput damping numNodes itrs
72     in M.fromList [ (n, pr) | (n, (pr, _)) <- mrOutput ]
73     where
74         mrItr :: [(Node, (PageRankValue, [Node]))] -> Double -> Double -> Int -> [(Node, (PageRankVa
75         mrItr input _ _ 0 = input
76         mrItr input damping numNodes itrs =
77             let output = mr damping numNodes 1 input
78             in mrItr output damping numNodes (itrs-1)
```

Listing 7: src/PageRankExt.hs

```
1 import Data.Map as M (toList)
2
3 import ProcessData (processData)
4 import PageRank (computePageRankMR)
5
6 main :: IO ()
7 main = do
8     (graph, pageRank) <- processData "../data/sample_input.txt"
9     let resPageRank = computePageRankMR pageRank graph 10 0.85
10    mapM_ putStrLn [ n ++ ": " ++ show pr | (n, pr) <- M.toList resPageRank ]
```

Listing 8: test/Spec.hs