

# ParVarys: Parallelizing Coflow Scheduling in Haskell

Project Report - COMS 4995 Parallel Functional Programming

Etesam Ansari, Yunlan Li

Columbia University

{ea2905, yl4387}@columbia.edu

## 1 Introduction

We parallelized the Varys [4] coflow scheduling algorithm in Haskell. We utilized the data parallelism in Varys using Haskell’s **Eval** monad and **Strategy**, and lazy data structures to parallelize Varys. Our parallel implementation achieved a max speed-up of 4x on a Macbook Pro M1 with 10 cores and 16GB memory.

In Section 2, we present some background on coflow scheduling and the Varys algorithm. In Section 3, we describe our sequential Haskell implementation of Varys, and then proceed to discuss how we parallelized it in Section 4, where we also present benchmark results that illustrate the effectiveness of our approaches and the pitfalls we encountered. Lastly, we summarize our takeaways from this project in Section 5.

## 2 Coflow Scheduling

Datacenter networks typically optimize for TCP *flow completion time*(FCT). However, nowadays, a datacenter application task could rarely be completed by sending a single request to another datacenter server. This has led to a mismatch between the network optimization objective(FCT) and application-level objective(Task Completion Time): minimizing FCT doesn’t necessarily contribute to better task completion time. To address this challenge, a new abstraction called *coflow* is proposed: a collection of flows that share a common performance goal.

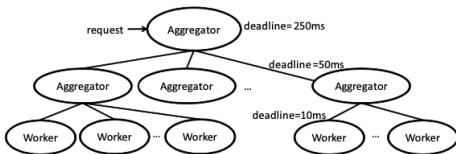


Figure 1: The partition/aggregate design pattern [3].

For example, web search workloads typically use a partition/aggregator model [3] where a user query will trigger multiple subtasks to search for results in each shard stored on different servers. The results from each shard are then aggregated to compute the final answer to be sent back to the user. In this example, all flows created for answering the user query would belong to the same coflow. Because latency observed by the end user is determined by the slowest of all flows of this coflow, it makes sense to optimize for this metric termed *coflow completion time*(CCT). By doing so, we could observe better application and end-user experience.

In the network community, there has been a large body of work to find a near-optimal coflow scheduling algorithm to minimize the average CCT.<sup>1</sup>

### 2.1 Offline Coflow Scheduling Problem (CSP)

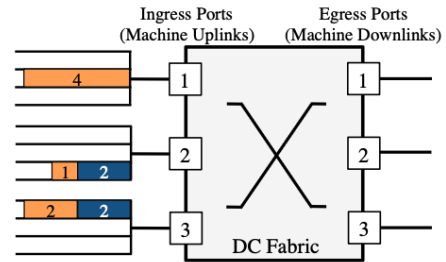


Figure 2: Coflow scheduling over a  $3 \times 3$  datacenter fabric with three ingress/egress ports. Flows in ingress ports are organized by destinations and color-coded by coflows -  $C_1$  in orange/light and  $C_2$  in blude/dark [4].

The offline coflow scheduling problem is defined as follows:

- the datacenter network fabric is abstracted into a single big switch consisting of  $m$  ingress ports (ToR switches)

<sup>1</sup>The average CCT of a collection of coflows.

and  $n$  egress ports (ToR switches)

- assume all coflows arrive simultaneously at time  $t = 0$ , and the information about each coflow (number of flows, size, source and destination port of each flow) is all known [4].

The goal is to find a schedule for the coflows (order, rate) to minimize the average CCT. This problem is NP-hard (via reduction from concurrent open-shop scheduling problem).

## 2.2 Varys

Varys uses a Smallest-Effective-Bottleneck-First (SEBF) heuristic to produce a near-optimal ordering of coflows, and then perform rate allocations to optimize average CCT.

$$\Gamma^C = \max \left( \max_i \frac{\sum_i d_{ij}}{Rem(P_i^{in})}, \max_j \frac{\sum_j d_{ij}}{Rem(P_j^{out})} \right) \quad (1)$$

---

**Pseudocode 1** Coflow Scheduling to Minimize CCT

```

1: procedure ALLOCBANDWIDTH(Coflows C, Rem(.), Bool cct)
2:   for all C' ∈ C do
3:     τ = ΓC (Calculated using Equation (1))
4:     if not cct then
5:       τ = DC
6:     end if
7:     for all dij ∈ C do           ▷ MADD
8:       rij = dij/τ
9:       Update Rem(Piin) and Rem(Pjout)
10:    end for
11:  end for
12: end procedure

13: procedure MINCCTOFFLINE(Coflows C, C, Rem(.))
14:  C' = SORT_ASC (C ∪ C') using SEBF
15:  allocBandwidth(C', Rem(.), true)
16:  Distribute unused bandwidth to C' ∈ C' ▷ Work conserv. (§5.3.4)
17:  return C'
18: end procedure

```

---

**Figure 3:** Varys Offline Coflow Scheduling Algorithm [4].

In equation 1

- $\Gamma^C$  represents the **shortest effective bottleneck** for coflow  $C$
- $d_{ij}$  represents the size of data that goes from ingress port  $i$  to egress port  $j$
- $Rem(\cdot)$  represents the remaining bandwidth of an ingress port or egress port
- $P_i^{in}$  represents ingress port  $i$ ,  $P_j^{out}$  represents egress port  $j$

The Shortest job first (SJF) scheduling discipline is optimal for flow scheduling. SEBF can be viewed as an approximation for shortest job first in the context of coflow scheduling where a coflow consist of one or more flows instead of just one.

## 3 Sequential Haskell Implementation

### 3.1 CSP Representation

In datacenter networks, a centralized controller running on commodity servers will be informed of the global view of network states, run Varys to perform coflow scheduling and finally disseminate the results to end switches. In an offline CSP problem, the network view consists of the state - switch ingress/egress bandwidth, and flow information - of all switches (see Figure 2).

In our implementation, a flow is abstracted into the **Flow** datatype, which holds the coflow id it belongs to, its flow size, and the egress port destination.

---

```

data Flow = Flow
{ coflowId    :: Int
, size        :: Int
, destinationId :: Int
}
deriving Show

```

---

We use **Switch** datatype to represent each ingress/egress port  $iId :: Int$ , which contains a number of `flows :: [Flow]` and has ingress/egress link rates specified by two **Int**.

---

```

data Switch = Switch
{ iId        :: Int
, flows      :: [Flow]
, iBandwidth :: Int
, eBandwidth :: Int
}
deriving Show

```

---

With these, we define the input to Varys - i.e. a CSP problem - as the datatype **CSP** which is a datacenter-wide view of network state consisting of ingress and egress switch states.

---

```

data CSP = CSP
{ ingressSwitches :: [Switch]
, egressSwitches  :: [Switch]
}
deriving Show

```

---

### 3.2 CSP Generator

We implement `generateProblem` to generate an offline CSP problem with `numIngress :: Int` ingress switches and `numEgress :: Int` egress switches. In addition, the function takes three specifications **RandomFlowSpec**, **RandomSwitchSpec**, and **RandomSwitchSpec** that defines the parameters for randomly generating the flows, ingress switches, and egress switches. In particular, for each ingress switch, we chose an **Int** in the range specified by `minFlows :: Int` and `maxFlows :: Int`

of `ingressSwitchSpec :: RandomSwitchSpec` as the number of flows. Each flow has a size and coflow id randomly generated with min and max threshold specified by `RandomFlowSpec`.

```
data RandomFlowSpec = RandomFlowSpec
  { minSwitchId :: Int
  , maxSwitchId :: Int
  , minCoflowId :: Int
  , maxCoflowId :: Int
  , minFlowSize :: Int
  , maxFlowSize :: Int
  }

data RandomSwitchSpec = RandomSwitchSpec
  { minFlows      :: Int
  , maxFlows      :: Int
  , ingressBandwidth :: Int
  , egressBandwidth :: Int
  }

generateProblem
  :: RandomFlowSpec
  -> RandomSwitchSpec
  -> RandomSwitchSpec
  -> Int
  -> Int
  -> IO CSP
```

The flow and switch specification can be defined by setting relevant command line options:

```
% stack exec ParVarys-exe -- -h
ParVarys: Parallel Varys Coflow Scheduling Using SEBF

Usage: ParVarys-exe [-t|--type STRING] [-n|--coflows NUMBER]
                  [-i|--ingress NUMBER] [-e|--egress NUMBER]
                  [-s|--min-flow-size NUMBER] [-S|--max-flow-size NUMBER]
                  [-f|--min-switch-flows NUMBER]
                  [-F|--max-switch-flows NUMBER]
                  [-b|--ingress-bandwidth NUMBER] [-B|--egress-bandwidth NUMBER]
                  [--seed NUMBER]

Generates an offline coflow scheduling problem, and uses the Varys Shortest
Effective Bottleneck First heuristic to order the coflows.

Available options:
  -h,--help                Show this help text
  -t,--type STRING          Varys mode: seq, parMap, chunk (default: "parMap")
  -n,--coflows NUMBER       Number of coflows (default: 4000)
  -i,--ingress NUMBER        Number of ingress switches (default: 1000)
  -e,--egress NUMBER         Number of egress switches (default: 1000)
  -s,--min-flow-size NUMBER  Smallest flow size in bytes (default: 0)
  -S,--max-flow-size NUMBER  Largest flow size in bytes (default: 1000)
  -f,--min-switch-flows NUMBER
                             Minimum number of flows arriving at an ingress switch
                             (default: 0)
  -F,--max-switch-flows NUMBER
                             Maximum number of flows arriving at an ingress switch
                             (default: 5000)
  -b,--ingress-bandwidth NUMBER
                             Ingress bandwidth (Gb/s) of a switch (default: 40)
  -B,--egress-bandwidth NUMBER
                             Egress bandwidth (Gb/s) of a switch (default: 40)
  --seed NUMBER              Seed for global pseudo-random number generator
                             (default: 4995)
```

We use Haskell’s standard pseudo-random number generator `StdGen` from the `System.Random` module. To reproduce a given CSP problem, one can specify an integer via the option `--seed` to seed the pseudo-random number generator.

### 3.3 Varys Controller

The Varys algorithm consists of two parts:

1. coflow ordering using the Shortest Effective Bottleneck First (SEBF) heuristic where the shortest effective bottleneck  $\Gamma$  of each coflow is computed with equation 1.

2. bandwidth allocation given the global coflow ordering in part 1.

For this project, we implemented part 1 and leave part 2 as future work.

The key insight is that the input to Varys `CSP` is a network-view centering around `Switch` entities, whereas  $\Gamma$  is a statistics derived from a coflow and switch bandwidth. This motivated us to first transform `CSP` into an intermediate representation that includes a coflow table `CoflowMap :: IntMap Coflow` and a switch bandwidth lookup table `BandwidthTable`.

```
type Switch2Flow = IntMap.IntMap [Flow]
data FlowDirection = Ingress | Egress deriving (Eq, Show)
data Coflow =
  Coflow Int -- coflow id
  [Flow] -- all flows belonging to this coflow
  Switch2Flow -- flows grouped by ingress switch
  Switch2Flow -- flows grouped by egress switch
  deriving (Show)

type CoflowMap = IntMap.IntMap Coflow
type BandwidthTable = Map.Map (Int, FlowDirection) Int
```

This transformation is enabled by two functions that traverse the ingress and egress switches of `CSP` to incrementally build the desired `CoflowMap` and `BandwidthTable`.

```
getSwitchBandwidth :: CSP -> BandwidthTable
toCoflows           :: CSP -> CoflowMap
```

With the abstraction of `Coflow`, computation of  $\Gamma$  becomes simple: the first fraction is a left-fold over the flows grouped by ingress switch with `Rem(.)` looked up using the `BandwidthTable`, and the second fraction is a left fold over the flows grouped by egress switch. The computation of  $\Gamma$  is implemented by `getGamma` which takes in a `BandwidthTable`, `Coflow` and produces a `Rational` that represents the shortest effective bottleneck for the given coflow.

```
getGamma :: BandwidthTable -> Coflow -> Rational
```

Lastly, we sort the coflows by their shortest effective bottleneck to produce a global coflow ordering schedule.

```
-- Given a Coflow Scheduling Problem,
-- use the SEBF heuristic to order the Coflows.
--
-- Returns [(coflow id, shortest effective bottleneck)]
sebf :: CSP -> [(Int, Rational)]
sebf csp = Key.sort snd $ map f coflows
  where
    switchLinkRates = getSwitchBandwidth csp
    coflows          = IntMap.toList $ toCoflows csp
    f (cid, coflow) = (cid, getGamma switchLinkRates coflow)
```

## 4 Parallelizing Varys

Unless otherwise specified, our experiments are run on an Apple Macbook Pro M1 with 10 cores and 16GB memory. We run both sequential and parallel version of our implementation of Varys on a CSP problem with 1000 ingress ports, 1000 egress ports, a maximum of 2000 coflows with a maximum flow size of 1000 bytes, and a maximum of 500 flows at each ingress port. In addition, both the ingress and egress bandwidth are set to 40Gbps. To ensure reproducibility of the generated CSP problem, we used the default seed of 4995.

We chose the parameters that generate the CSP problem so that it simulates a realistic workload in modern day datacenters where it's typical to have a fleet of thousands of servers, and thousands of coflows arriving per second.

### 4.1 Garbage Collector Optimization



**Figure 4:** Threadscope visualizing event traces of sequential Varys running on a single HEC, with and without tuning GHC runtime options of garbage collector.

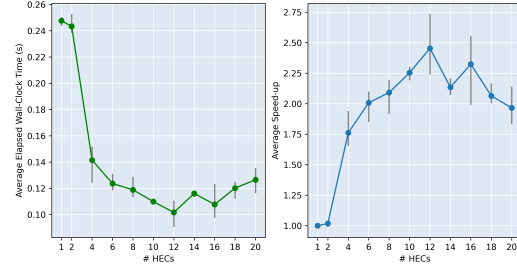
Figure 4(a) shows that when we run our sequential implementation of Varys with the default GHC settings, our application thread is frequently interleaved with the garbage collector thread (active throughout the duration of the program). This is explained by the scale of the problem we run our program against: thousands of switches and coflows whose data structures contains hundreds of flows, leaving a very large memory footprint.

Because of Amdahl's law, frequent interleaving garbage collection work would significantly limit the speedup we can achieve, we thus limit garbage collection by utilizing two GHC RTS options: `-H` for providing a suggested heap size for the garbage collector, and `-I` to increase the amount of time that must pass before an idle GC is performed [1]. We used 30s for the `-I` option, which is much larger than the time it takes for the program to run, and 8G (which is the maximum amount of memory typically available on the Macbook Pro M1) as the suggested heap size. Unless otherwise specified, these are the GHC RTS options we used for benchmarking.

This proved to be quite effective as can be seen from the

comparison in Figure 4. Garbage collection disappears after 1.55s, yielding a 1.8x speedup (from 471.91ms to 265.69ms).

### 4.2 Data Parallelism



**Figure 5:** Average wall-clock time and average speed-up of Parvarys vs # HECs.

There are two main opportunities for data parallelism in our sequential implementation of Varys, which we exploit using Haskell's **Eval** Monad and **Strategy**. We manage to obtain a max speedup of 2.75x on 12 cores (see Figure 6). This is much lower than what we had expected and could partly be attributed to long pause of application threads due to garbage collection (see Figure 6).

**Computation of  $\Gamma^C$**  Recall equation 1, this represents the shortest effective bottleneck for a coflow  $C$ . Importantly, it is independent of  $\Gamma^C$  of other coflows, meaning that  $\Gamma$  of all coflows can be computed in parallel. We fully evaluate `f (cid, coflow) :: (Int, Rational)` to normal form using `rdeepseq`. To prevent spark overflow, we set a limit of the spark pool size with `maxParSparks = 4000` using `parBuffer`.

---

```

parSebf :: CSP -> [(Int, Rational)]
parSebf csp = Key.sort
  snd
  (map f coflows `using` parBuffer maxParSparks rdeepseq)
  where
    switchLinkRates = parGetSwitchBandwidth csp
    coflows         = IntMap.toList $ parToCoflows csp
    f (cid, coflow) = (cid, getGamma switchLinkRates coflow)

```

---

**Building CoflowMap from Ingress Switches** This conversion is a left-fold over `ingressSwitches :: Switch` to produce a `coflowMap :: IntMap Coflow`. For each switch, the accumulator function iterates over each `flow :: Flow` of the switch, add updates the accumulated coflowMap by `cons` the flow to `flowsByIngress :: [Flow]` and `flowsByEgress :: [Flow]` of the coflow it belongs to. This accumulator function is associative, which allows us to transform it into a `map` and parallelize it. In particular, we

(1) split the ingress switches into chunks of 10, (2) build a partial coflowMap for each chunk, and (3) finally merge all the partial maps together using `unionsWith`. The construction of partial coflowMap for each chunk could then be parallelized using `parBuffer` `maxParSparks` `rdeepseq`.

```

data FlowDirection = Ingress | Egress deriving (Eq, Show)

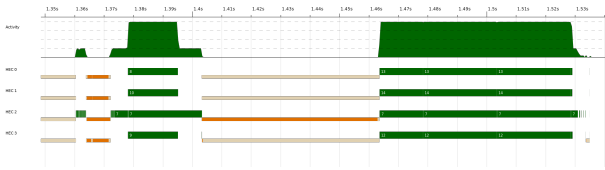
instance NFData FlowDirection where
  rnf dir = dir `seq` ()

-- coflow: id flows flowsByIngress flowsByEgress
data Coflow = Coflow Int [Flow] Switch2Flow Switch2Flow
  deriving Show

instance NFData Coflow where
  rnf (Coflow cid flows iFlows eFlows) =
    cid `seq` rnf flows `seq` rnf iFlows `seq` rnf eFlows

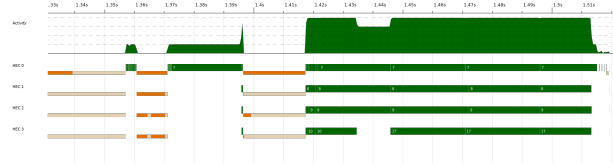
-- parallel version of toCoflows
parToCoflows :: CSP -> CoflowMap
parToCoflows csp = IntMap.unionsWith
  mergeCoflow
  (map f switches `using` parBuffer maxParSparks rdeepseq)
  where
    f = IntMap.unionsWith mergeCoflow . map (update
  ← IntMap.empty)
  switches = chunksOf 10 $ ingressSwitches csp

```



**Figure 6:** Threadscope visualizing event traces of parallel Varys running on 4 HECs, with data parallelism enabled.

The granularity of our data parallelism is fine enough to produce enough sparks to keep all HECs busy and is coarse enough to make the overhead of parallelization (spark creation, chunking, etc) worth it. For example, building a partial map from 10 ingress switches consist of a left-fold over thousands of flows. This is illustrated in Figure 6 where we can see work are spread evenly on all 4 HECs and constantly keeps the HECs busy. At time 1.39s to 1.41s, only HEC 2 busy, which is a result of the sequential `unionsWith` of partial coflowMaps in `parToCoflow :: CSP -> CoflowMap`. We mitigate this in Section 4.2 by using lazy maps. Lastly, near the end around 1.53s, again only HEC 3 is busy, this is due to sequentially sorting (`coflowId :: Int`, `gamma :: Rational`) pairs in `parSebf`. For thousands of coflows, this wasn't a bottleneck and thus we didn't attempt to parallelize it. However, when there are hundreds of thousands of coflows or more, this becomes the main performance bottleneck (see Figure 10 in Section 4.3).

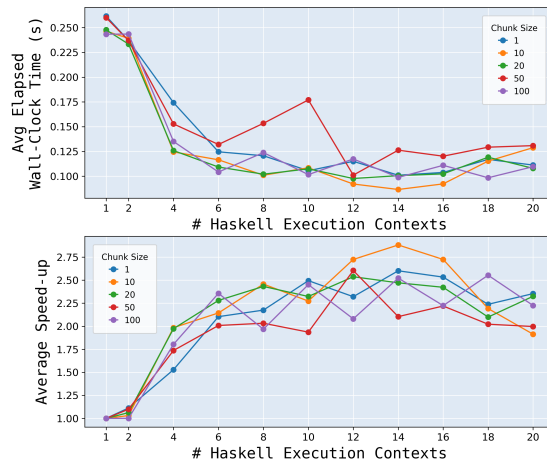


**Figure 7:** Threadscope visualizing event traces of parallel Varys running on 4 HECs, with data parallelism enabled and using lazy maps and chunking for building `BandwidthTable` and `CoflowMap`.

### 4.3 Lazy Coflow Transformation

To further improve speed-up, we reduce the amount of sequential work performed in `parToCoflows :: CSP -> CoflowMap` by using lazy `IntMap` from `Data.List.IntMap.Lazy` module.

A lazy map is strict in its keys but lazy in its values [2]. This allows us to delay the computation of map values when they are needed, i.e. during the calculation of shortest-effective-bottleneck  $\Gamma^C$  for each coflow  $C$ . Since we parallelized the calculation of  $\Gamma^C$ , the evaluation of the map value of type `Coflow` to normal form could thus be parallelized. With this change, we no longer spot the sequential merge of partial coflowMaps in threadscope (Figure 7), where only a single HEC is busy.

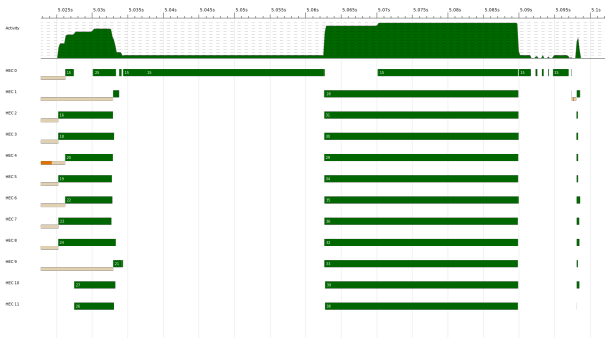


**Figure 8:** Comparison of average elapsed wall-clock time and speed-ups for building partial lazy `CoflowMaps` using chunks of  $N$  ingress switches and then perform a sequential union of these lazy `IntMaps`.

Initially, we thought chunk size we use to split ingress switches plays a big role in speedup, but experiment results

say otherwise as Figure 8 shows. Initially, we thought that the amount of time the sequential `IntMap.unionsWith` is related to the number of maps it merges, so using a larger chunk size would take less time to complete. We later realized that this was incorrect since the time complexity of unions is determined by the number of keys, which would be the same regardless of the chunk size we choose. This explains the negligible difference in performance for different chunk sizes used.

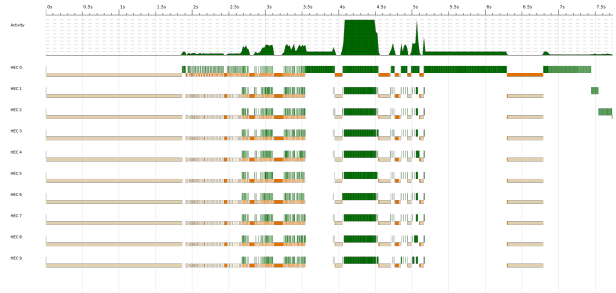
#### 4.4 Amdahl’s Law Strikes



**Figure 9:** Threadscope visualizing event traces of parallel Varys running on 12 HECs, with all parallelization optimizations enabled.

Before parallelizing Varys, transforming the problem into `coflowMap` and calculating the shortest-effective-bottleneck for each coflow are the performance bottleneck. After parallelizing Varys, the more HECs we use, the less time it takes to perform the aforementioned calculations. What previously account for a small percentage of the total computation time now start to dominate and become the new bottleneck that prevent further speed-up (see Figure 9). This gives us an alternative view of Amdahl’s Law, which is performance bottleneck shifts as we parallelize the algorithm, and eventually the sequential portion becomes the bottleneck that prevents further possibility for speed-up.

Initially, we thought Varys is straightforward to parallelize and easy to achieve linear speed-ups. However, as can be seen from our discussion so far, this is far from the case. Issues such as garbage collection due to the memory-intensive nature of the algorithm, map merging that’s inherently sequential make Varys hard to parallelize. To make things worse, when there are hundreds of thousands coflows or more, sorting coflows by their shortest-effective-bottleneck now dominates the computation time as Figure 10 shows.



**Figure 10:** Threadscope visualizing event traces of parallel Varys running on 10 HECs solving a CSP problem of up to 500K coflows, 1K ingress ports and a maximum of 5000 flows per each ingress port.

## 5 Takeaways

Parallelizing Varys was harder than we expected and the amount of speed-ups achieved was also surprisingly far from our expectation going into the project. However, we found this project to be an interesting one and were able to take away something meaningful:

1. garbage collection can be really expensive for memory-intensive applications
2. alternative view of Amdahl’s law is that performance bottleneck shifts
3. lazy data structures can be helpful for parallelization, and Haskell’s paradigm of defining a structure for holding computation and a strategy to evaluate the computation is very elegant

## References

- [1] ghc runtime options. [https://downloads.haskell.org/ghc/latest/docs/users\\_guide/runtime\\_control.html#rts-options-to-control-the-garbage-collector](https://downloads.haskell.org/ghc/latest/docs/users_guide/runtime_control.html#rts-options-to-control-the-garbage-collector).
- [2] Haskell lazy map. <https://hackage.haskell.org/package/containers-0.6.6/docs/Data-Map-Lazy.html>.
- [3] Mohammad Alizadeh, Albert Greenberg, David A. Maltz, Jitendra Padhye, Parveen Patel, Balaji Prabhakar, Sudepta Sengupta, and Murari Sridharan. Data center tcp (dctcp). In *Proceedings of the ACM SIGCOMM 2010 Conference, SIGCOMM ’10*, page 63–74, New York, NY, USA, 2010. Association for Computing Machinery.
- [4] Mosharaf Chowdhury, Yuan Zhong, and Ion Stoica. Efficient coflow scheduling with varys. In *Proceedings of the*

2014 ACM Conference on SIGCOMM, SIGCOMM '14,  
page 443–454, New York, NY, USA, 2014. Association  
for Computing Machinery.

## Appendix

src/Generator.hs

```
1 {-# LANGUAGE FlexibleContexts #-}
2 {-# LANGUAGE GADTs #-}
3 {-# LANGUAGE NamedFieldPuns #-}
4
5 module Generator
6   ( RandomFlowSpec(..)
7   , RandomSwitchSpec(..)
8   , Flow(..)
9   , Switch(..)
10  , CSP(..)
11  , generateProblem
12  ) where
13
14 import           Control.DeepSeq                ( NFData
15                                                , rnf
16                                                )
17 import           System.Random                 ( randomRIO )
18
19 data Flow = Flow
20   { coflowId      :: Int
21   , size          :: Int
22   , destinationId :: Int
23   }
24   deriving Show
25
26 instance NFData Flow where
27   rnf (Flow cid size dId) = cid `seq` size `seq` dId `seq` ()
28
29 data Switch = Switch
30   { iId          :: Int
31   , flows        :: [Flow]
32   , iBandwidth   :: Int
33   , eBandwidth   :: Int
34   }
35   deriving Show
36
37 -- Coflow Scheduling Problem
38 data CSP = CSP
39   { ingressSwitches :: [Switch]
40   , egressSwitches  :: [Switch]
41   }
42   deriving Show
43
44 data RandomFlowSpec = RandomFlowSpec
45   { minSwitchId :: Int
46   , maxSwitchId :: Int
47   , minCoflowId :: Int
48   , maxCoflowId :: Int
49   , minFlowSize :: Int
50   , maxFlowSize :: Int
```



```

51 }
52
53 data RandomSwitchSpec = RandomSwitchSpec
54 { minFlows      :: Int
55 , maxFlows      :: Int
56 , ingressBandwidth :: Int
57 , egressBandwidth :: Int
58 }
59
60 -- FUNCTIONS:
61 -- Generates random Integer from lb to ub (inclusive? Yes)
62 generateRandomNum :: Int -> Int -> IO Int
63 generateRandomNum lb ub = do
64   randomRIO (lb, ub)
65
66 generateFlows :: RandomFlowSpec -> Int -> IO [Flow]
67 generateFlows spec n = if n <= 0
68   then do
69     return []
70   else do
71     flows      <- generateFlows spec $ n - 1
72     coflowId    <- generateRandomNum (minCoflowId spec) (maxCoflowId spec)
73     egressSwitchId <- generateRandomNum (minSwitchId spec) (maxSwitchId spec)
74     flowSize    <- generateRandomNum (minFlowSize spec) (maxFlowSize spec)
75
76     return $ Flow coflowId flowSize egressSwitchId : flows
77
78 generateSwitches
79   :: RandomFlowSpec -> RandomSwitchSpec -> Int -> Int -> IO [Switch]
80 generateSwitches flowSpec switchSpec minId maxId = if maxId - minId < 0
81   then do
82     return []
83   else do
84     numOfFlows <- generateRandomNum (minFlows switchSpec) (maxFlows switchSpec)
85     flows      <- generateFlows flowSpec numOfFlows
86     let switch = Switch minId
87         flows
88             (ingressBandwidth switchSpec)
89             (egressBandwidth switchSpec)
90
91     switches <- generateSwitches flowSpec switchSpec (minId + 1) maxId
92     return $ switch : switches
93
94 generateProblem
95   :: RandomFlowSpec
96   -> RandomSwitchSpec
97   -> RandomSwitchSpec
98   -> Int
99   -> Int
100  -> IO CSP
101 generateProblem flowSpec ingressSwitchSpec egressSwitchSpec numIngress numEgress
102 = do
103   let (minIngressId, maxIngressId) = (1, numIngress)
104       (minEgressId , maxEgressId ) = (numIngress + 1, numIngress + numEgress)

```

```
105
106 iSwitches <- generateSwitches flowSpec
107                               ingressSwitchSpec
108                               minIngressId
109                               maxIngressId
110 eSwitches <- generateSwitches flowSpec
111                               egressSwitchSpec
112                               minEgressId
113                               maxEgressId
114
115 return $ CSP iSwitches eSwitches
```

```

1  {-# LANGUAGE NamedFieldPuns #-}
2
3  module Controller
4    ( Coflow(..)
5    , toCoflows
6    , parToCoflows
7    , getSwitchBandwidth
8    , getGamma
9    , sebf
10   , parSebf
11   ) where
12
13  import           Control.DeepSeq          ( NFData
14                                           , rnf
15                                           )
16  import           Control.Parallel.Strategies ( parBuffer
17                                           , rdeepseq
18                                           , using
19                                           )
20  import qualified Data.IntMap.Lazy         as IntMap
21  import qualified Data.List.Key           as Key
22  import           Data.List.Split         ( chunksOf )
23  import qualified Data.Map.Lazy           as Map
24  import           Data.Maybe              ( fromMaybe )
25  import           Data.Ratio              ( (%) )
26
27  import           Generator                ( CSP(..)
28                                           , Flow(..)
29                                           , Switch(..)
30                                           )
31
32
33  data FlowDirection = Ingress | Egress deriving (Eq, Show)
34
35  instance NFData FlowDirection where
36    rnf dir = dir `seq` ()
37
38  instance Ord FlowDirection where
39    a <= b = case (a, b) of
40      (Egress, Ingress) -> False
41      _                  -> True
42
43  -- coflow: id flows flowsByIngress flowsByEgress
44  data Coflow = Coflow Int [Flow] Switch2Flow Switch2Flow
45    deriving Show
46
47  instance NFData Coflow where
48    rnf (Coflow cid flows iFlows eFlows) =
49      cid `seq` rnf flows `seq` rnf iFlows `seq` rnf eFlows
50
51
52  type Switch2Flow = IntMap.IntMap [Flow]

```

```

53 type CoflowMap = IntMap.IntMap Coflow
54 type BandwidthTable = Map.Map (Int, FlowDirection) Int
55
56
57 maxParSparks :: Int
58 maxParSparks = 4000
59
60 updateMap :: Int -> Flow -> Switch2Flow -> Switch2Flow
61 updateMap k v = IntMap.alter f k
62   where
63     f pv = case pv of
64       Nothing -> Just [v]
65       Just vs -> Just $ v : vs
66
67 addFlow :: CoflowMap -> (Int, Flow) -> CoflowMap
68 addFlow currMap (ingressPort, flow) = IntMap.alter f (coflowId flow) currMap
69   where
70     egressPort = destinationId flow
71     f val = case val of
72       Nothing -> Just $ Coflow (coflowId flow)
73         [flow]
74         (IntMap.singleton ingressPort [flow])
75         (IntMap.singleton egressPort [flow])
76       Just (Coflow cid coflow flowsByISwitch flowsByESwitch) -> Just $ Coflow
77         cid
78         (flow : coflow)
79         (updateMap ingressPort flow flowsByISwitch)
80         (updateMap egressPort flow flowsByESwitch)
81
82 update :: CoflowMap -> Switch -> CoflowMap
83 update currMap switch =
84   foldl addFlow currMap $ zip (repeat $ iId switch) (flows switch)
85
86 -- Assumes that the coflowId of the two coflows passed in are the same
87 mergeCoflow :: Coflow -> Coflow -> Coflow
88 mergeCoflow (Coflow cid flows ingress egress) (Coflow _ flows' ingress' egress')
89   = Coflow cid
90     (flows ++ flows')
91     (IntMap.unionWith (++) ingress ingress')
92     (IntMap.unionWith (++) egress egress')
93
94 toCoflows :: CSP -> CoflowMap
95 toCoflows csp = foldl update IntMap.empty $ ingressSwitches csp
96
97 parToCoflows :: CSP -> CoflowMap
98 parToCoflows csp = IntMap.unionsWith
99   mergeCoflow
100   (map f switchess `using` parBuffer maxParSparks rdeepseq)
101   where
102     f      = IntMap.unionsWith mergeCoflow . map (update IntMap.empty)
103     switchess = chunksOf 10 $ ingressSwitches csp
104
105 getSwitchBandwidth :: CSP -> BandwidthTable
106 getSwitchBandwidth csp = Map.fromList $ concatMap f switches where

```

```

107 f (Switch sid _ iBw eBw) = [(sid, Ingress), iBw], ((sid, Egress), eBw)]
108 switches = ingressSwitches csp ++ egressSwitches csp
109
110 parGetSwitchBandwidth :: CSP -> BandwidthTable
111 parGetSwitchBandwidth csp = Map.fromList $ concat
112   (map (concatMap f) switchess `using` parBuffer maxParSparks rdeepseq)
113   where
114     f (Switch sid _ iBw eBw) = [(sid, Ingress), iBw], ((sid, Egress), eBw)]
115     switchess = chunksOf 20 $ ingressSwitches csp ++ egressSwitches csp
116
117 getGamma :: BandwidthTable -> Coflow -> Rational
118 getGamma bwTbl (Coflow _ _ ingressFlows egressFlows) = max
119   (maximum ingressTimes)
120   (maximum egressTimes)
121   where
122     sumFlows :: (Int, [Flow]) -> (Int, Int)
123     sumFlows (switchId, flows) = (switchId, foldl (\a el -> a + size el) 0 flows)
124
125     calcTime :: FlowDirection -> (Int, Int) -> Rational
126     calcTime flowDir (switchId, flowSize) = fromIntegral flowSize
127       % fromIntegral bandwidth
128       where bandwidth = fromMaybe 0 $ Map.lookup (switchId, flowDir) bwTbl
129
130     ingressTimes = map (calcTime Ingress . sumFlows) $ IntMap.toList ingressFlows
131     egressTimes  = map (calcTime Egress . sumFlows) $ IntMap.toList egressFlows
132
133
134 -- Given a Coflow Scheduling Problem, use Shortest Effective Bottleneck First
135 -- heuristic to order the Coflows.
136 --
137 -- Returns [(coflow id, effective bottleneck)]
138 sebf :: CSP -> [(Int, Rational)]
139 sebf csp = Key.sort snd $ map f coflows
140   where
141     switchLinkRates = getSwitchBandwidth csp
142     coflows         = IntMap.toList $ toCoflows csp
143     f (cid, coflow) = (cid, getGamma switchLinkRates coflow)
144
145 -- Parallel version of sebf
146 parSebf :: CSP -> [(Int, Rational)]
147 parSebf csp = Key.sort
148   snd
149   (map f coflows `using` parBuffer maxParSparks rdeepseq)
150   where
151     switchLinkRates = parGetSwitchBandwidth csp
152     coflows         = IntMap.toList $ parToCoflows csp
153     f (cid, coflow) = (cid, getGamma switchLinkRates coflow)

```

```

1  {-# LANGUAGE NamedFieldPuns #-}
2
3  import      Control.DeepSeq          ( force )
4  import      Control.Exception        ( evaluate )
5  import      Control.Monad            ( join )
6  import      Formatting                ( fprintfLn )
7  import      Formatting.Clock          ( timeSpecs )
8  import      Generator                 ( RandomFlowSpec(..)
9                                          , RandomSwitchSpec(..)
10                                         , generateProblem
11                                         )
12 import      Options.Applicative
13 import      System.Clock
14 import      System.Exit                ( die )
15 import      System.Random              ( mkStdGen
16                                         , setStdGen
17                                         )
18
19 import      Controller                 ( parSebf
20                                         , sebf
21                                         )
22
23 -- Arg Parser template adapted from:
24 -- https://ro-che.info/articles/2016-12-30-optparse-applicative-quick-start
25 main :: IO ()
26 main = join . customExecParser (prefs showHelpOnError) $ info
27   (helper <*> parser)
28   ( fullDesc
29   <> header "ParVarys: Parallel Varys Coflow Scheduling Using SEBF "
30   <> progDesc
31     ( "Generates an offline coflow scheduling problem, and uses the "
32     ++ "Varys Shortest Effective Bottleneck First heuristic to order "
33     ++ "the coflows."
34     )
35   )
36 where
37   parser :: Parser (IO ())
38   parser =
39     work
40     <$> strOption
41       ( long "type"
42       <> short 't'
43       <> metavar "STRING"
44       <> help "Varys mode: seq, parMap"
45       <> value "parMap"
46       <> showDefault
47       )
48     <*> option
49       auto
50       ( long "coflows"
51       <> short 'n'
52       <> metavar "NUMBER"

```

```

53         <> help "Number of coflows"
54         <> value 4000
55         <> showDefault
56     )
57 <*> option
58     auto
59     ( long "ingress"
60     <> short 'i'
61     <> metavar "NUMBER"
62     <> help "Number of ingress switches"
63     <> value 1000
64     <> showDefault
65     )
66 <*> option
67     auto
68     ( long "egress"
69     <> short 'e'
70     <> metavar "NUMBER"
71     <> help "Number of egress switches"
72     <> value 1000
73     <> showDefault
74     )
75 <*> option
76     auto
77     ( long "min-flow-size"
78     <> short 's'
79     <> metavar "NUMBER"
80     <> help "Smallest flow size in bytes"
81     <> value 0
82     <> showDefault
83     )
84 <*> option
85     auto
86     ( long "max-flow-size"
87     <> short 'S'
88     <> metavar "NUMBER"
89     <> help "Largest flow size in bytes"
90     <> value 1000
91     <> showDefault
92     )
93 <*> option
94     auto
95     ( long "min-switch-flows"
96     <> short 'f'
97     <> metavar "NUMBER"
98     <> help "Minimum number of flows arriving at an ingress switch"
99     <> value 0
100    <> showDefault
101    )
102 <*> option
103     auto
104     ( long "max-switch-flows"
105     <> short 'F'
106     <> metavar "NUMBER"

```

```

107         <> help "Maximum number of flows arriving at an ingress switch"
108         <> value 5000
109         <> showDefault
110     )
111     <*> option
112         auto
113         ( long "ingress-bandwidth"
114         <> short 'b'
115         <> metavar "NUMBER"
116         <> help "Ingress bandwidth (Gb/s) of a switch"
117         <> value 40
118         <> showDefault
119         )
120     <*> option
121         auto
122         ( long "egress-bandwidth"
123         <> short 'B'
124         <> metavar "NUMBER"
125         <> help "Egress bandwidth (Gb/s) of a switch"
126         <> value 40
127         <> showDefault
128         )
129     <*> option
130         auto
131         ( long "seed"
132         <> metavar "NUMBER"
133         <> help "Seed for global pseudo-random number generator"
134         <> value 4995
135         <> showDefault
136         )
137
138
139 work
140     :: String
141     -> Int
142     -> Int
143     -> Int
144     -> Int
145     -> Int
146     -> Int
147     -> Int
148     -> Int
149     -> Int
150     -> Int
151     -> IO ()
152 work mode numCoflows numIngress numEgress minFlowSize maxFlowSize minFlows maxFlows
153     ↪ ingressBandwidth egressBandwidth seed
154     = do
155     sebfImpl <- case mode of
156     "seq"     -> return sebf
157     "parMap" -> return parSebf
158     _        ->
159     die
160         $ "Unrecognized varys mode: "

```



```

160     ++ "expect one of seq, parMap, got "
161     ++ show mode
162
163 let flowSpec = RandomFlowSpec { minSwitchId = numIngress + 1
164                               , maxSwitchId = numIngress + numEgress
165                               , minCoflowId = 1
166                               , maxCoflowId = numCoflows
167                               , minFlowSize
168                               , maxFlowSize
169                               }
170
171 ingressSwitchSpec = RandomSwitchSpec { minFlows
172                                       , maxFlows
173                                       , ingressBandwidth
174                                       , egressBandwidth
175                                       }
176
177 egressSwitchSpec = RandomSwitchSpec { minFlows           = 0
178                                       , maxFlows           = 0
179                                       , ingressBandwidth
180                                       , egressBandwidth
181                                       }
182
183 -- seed the global pseudo-random number generator
184 -- for reproducibility of CSP problems
185 setStdGen $ mkStdGen seed
186 problem <- generateProblem flowSpec
187                               ingressSwitchSpec
188                               egressSwitchSpec
189                               numIngress
190                               numEgress
191
192 start      <- getTime Monotonic
193 coflowOrder <- evaluate $ force $ sebfImpl problem
194 end        <- getTime Monotonic
195
196 print coflowOrder
197 putStr "Calculation Time: "
198 fprintfLn timeSpecs start end

```

```

1  import      Data.Ratio          ( (%) )
2  import      System.Random      ( mkStdGen
3                                  , setStdGen
4                                  )
5  import      Test.HUnit
6
7  import      Controller          ( parSebf
8                                  , sebf
9                                  )
10 import      Generator           ( CSP(..)
11                                  , Flow(..)
12                                  , RandomFlowSpec(..)
13                                  , RandomSwitchSpec(..)
14                                  , Switch(..)
15                                  , generateProblem
16                                  )
17
18 testSebf1 :: Test
19 testSebf1 = TestCase
20   (do
21     let csp = CSP
22         { ingressSwitches =
23           [ Switch
24             { iId      = 1
25               , iBandwidth = 1
26               , eBandwidth = 1
27               , flows = [Flow { coflowId = 1, size = 4, destinationId = 5 } ]
28             }
29           , Switch
30             { iId      = 2
31               , iBandwidth = 1
32               , eBandwidth = 1
33               , flows = [ Flow { coflowId = 1, size = 1, destinationId = 6 }
34                           , Flow { coflowId = 2, size = 2, destinationId = 6 }
35                         ]
36             }
37           , Switch
38             { iId      = 3
39               , iBandwidth = 1
40               , eBandwidth = 1
41               , flows = [ Flow { coflowId = 1, size = 2, destinationId = 4 }
42                           , Flow { coflowId = 2, size = 2, destinationId = 4 }
43                         ]
44             }
45         ]
46
47         , egressSwitches = [ Switch { iId      = n
48                                   , iBandwidth = 1
49                                   , eBandwidth = 1
50                                   , flows      = []
51                                   }
52           | n <- [4 .. 6]

```

```

53         ]
54     }
55
56     assertEquals "" [(2, 2 % 1), (1, 4 % 1)] $ sebf csp
57 )
58
59 testSebf2 :: Test
60 testSebf2 = TestCase
61     (do
62         let csp = CSP
63             { ingressSwitches =
64                 [ Switch
65                     { iId      = 1
66                       , iBandwidth = 2
67                       , eBandwidth = 1
68                       , flows = [Flow { coflowId = 1, size = 4, destinationId = 5 } ]
69                     }
70                 , Switch
71                     { iId      = 2
72                       , iBandwidth = 1
73                       , eBandwidth = 1
74                       , flows = [ Flow { coflowId = 1, size = 1, destinationId = 6 }
75                                   , Flow { coflowId = 2, size = 2, destinationId = 4 }
76                                 ]
77                     }
78                 , Switch
79                     { iId      = 3
80                       , iBandwidth = 1
81                       , eBandwidth = 1
82                       , flows = [ Flow { coflowId = 1, size = 2, destinationId = 4 }
83                                   , Flow { coflowId = 2, size = 2, destinationId = 4 }
84                                 ]
85                     }
86             ]
87
88             , egressSwitches =
89                 [ Switch { iId = 4, iBandwidth = 1, eBandwidth = 1, flows = [] }
90                 , Switch { iId = 5, iBandwidth = 1, eBandwidth = 4, flows = [] }
91                 , Switch { iId = 6, iBandwidth = 1, eBandwidth = 1, flows = [] }
92                 ]
93         }
94
95     assertEquals "" [(1, 2 % 1), (2, 4 % 1)] $ sebf csp
96 )
97
98 -- validate the correctness of parallel implementation using
99 -- the sequential implementation of SEBF
100 testParSebf1 :: Test
101 testParSebf1 = TestCase
102     (do
103         let seed      = 91845734
104             flowSpec = RandomFlowSpec { minSwitchId = 201
105                                         , maxSwitchId = 400
106                                         , minCoflowId = 1

```

```

107         , maxCoflowId = 1000
108         , minFlowSize = 0
109         , maxFlowSize = 100
110     }
111     ingressSpec = RandomSwitchSpec 0 100 40 40
112     egressSpec  = RandomSwitchSpec 0 0 40 40
113
114     setStdGen $ mkStdGen seed
115     problem <- generateProblem flowSpec ingressSpec egressSpec 200 200
116     assertEquals "" (sebf problem) (parSebf problem)
117 )
118
119
120 main :: IO Counts
121 main = runTestTT $ TestList
122   [ "SEBF Sequential (varys_paper_fig1)" ~: testSebf1
123   , "SEBF Sequential (variable_link_rates)" ~: testSebf2
124   , "SEBF Parallel (parMap)" ~: testParSebf1
125   ]

```