

# GRNPar: Parallel Inference of Gene Regulatory Networks Using Boolean Network Models and Mutual Information

Anushka Gupta (ag4351), William Das (whd2108)

December 22, 2022

## 1 Introduction

Gene regulatory networks (GRNs) are graph-like representations of genes prone to activating or inhibiting the expression of other genes in a localized network. Boolean networks are often used to model GRNs—these networks consist of a series of nodes with connections between genes represented as boolean functions. Each node is a boolean variable that encodes its state (0 or 1). Each node’s state is determined by a boolean expression that takes as input the states of a subset of other nodes in the network. Researchers have developed algorithms for inferring these boolean expressions to model gene regulation activity.

For our project, we implemented a parallelized algorithm for inferring gene regulatory networks from gene expression time-series data using the boolean network model. Specifically, we modified a state-of-the-art approach put forth by Barman & Kwon (2017), which used a mutual information based approach coupled with a "swapping" subroutine to select dependent genes in the network and infer logical relations between input and output nodes.

In particular, we implemented and parallelized: (1) a mutual information approach for inferring dependencies in a boolean network, (2) boolean expressions mapping input nodes to output nodes, optimizing for a "gene-wise dynamics consistency" criterion, and a (3) method for visualizing and graphing the inferred boolean network. We parallelized each step, and specifically optimized key operations in the computation of possible boolean expressions that satisfy the various gene expression interactions in the network, where we saw the most speed-up.

### 1.1 Boolean Networks

The boolean network model  $G(V, A)$  proposed by Barman & Kwon consists of nodes:  $V = \{v_1, \dots, v_n\}$  and a set of interactions, or directed edges between input and output nodes:  $A = \{(v_i, v_j) \mid v_i, v_j \in V\}$ , where  $v_j$  is a target node whose state depends on that of  $v_i$ .

We further extend this model to account for boolean logic between the inferred input nodes and output nodes.

## 2 GRNPar Algorithm

### 2.1 Overview

We infer gene regulatory networks from time-series data by:

1. First, finding the  $k$  input nodes with the highest mutual information in relation to every target node.
2. Second, finding an optimal boolean expression that takes the  $k$  inferred input nodes with the highest dynamics criterion with respect to the target node. We choose the expression among a combination of  $2^{k-1}$  possible boolean expressions.

We also plot the inferred network to visualize our results as a final feature.

To determine the  $k$  most "closely related" input nodes for each node in the network, we used a mutual inference feature selection (MIFS) method based on entropy values (2.2.1). The testing of boolean expressions consists of finding possible boolean expressions which can be evaluated using the top  $k$  input nodes mapping to each target node, and choosing the one with highest gene-wise dynamics consistency (2.2.2).

Let  $V = \{\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_n\}$  where each vector  $\mathbf{v}_i$  corresponds to a gene and  $\mathbf{v}_i(t)$  corresponds to its value, 0 or 1, at time  $t$ .

---

#### Algorithm 1 GRNPar algorithm

---

```

procedure GRNPAR( $V$ )
    network  $\leftarrow$  generateNetwork( $k, V$ )
    optimalExpressions  $\leftarrow$  getOptimalBoolExpressions(network,  $|\mathbf{v}_0|$ )
    networkImage  $\leftarrow$  plotBoolNetwork(network)
    return network, optimalExpressions, networkImage
end procedure

```

---

$|\mathbf{v}_0|$  in *getOptimalBoolExpressions* corresponds to a target node—given the network and connections generated, we iterate through each node and test for boolean expressions of its top  $k$  input nodes.

Finally, we plot the network as a directed graph, and return the inferred network, optimal boolean expressions for each node, and the image.

## 2.2 Criterions

### 2.2.1 Mutual Information

The mutual information metric used in Algorithm 2  $I(X; Y)$  is based on the entropy of two variables. First, the entropy  $H(X)$  of a discrete random variable (gene) X is defined to measure the uncertainty of X over all time steps. A joint entropy  $H(X, Y)$  between X and Y measures the joint probability distribution between X and Y. The mutual information metric is then calculated based on these two entropy values.

$$\begin{aligned} H(X) &= - \sum_{x \in X} (p(x) \log p(x)) \\ H(X, Y) &= - \sum_{x \in X} \sum_{y \in Y} (p(x) \log p(x)) \\ I(X; Y) &= H(X) + H(Y) - H(X, Y) \end{aligned}$$

### 2.2.2 Gene-wise Dynamics Consistency

Given the length of the time-series  $T$ , and the observed and inferred boolean states of a node  $v$  and  $v'$  across a time-series, respectively, the **gene-wise dynamics consistency** is:

$$E(v, v') = \frac{\sum_{t=2}^T I(v(t) = v'(t))}{T - 1} \quad (1)$$

where  $I(v(t) = v'(t))$  here is equal to 1 if the condition  $v(t) = v'(t)$  is true, and 0 if false.  $v'(t)$ , is obtained from a possible boolean expression that maps the inferred input nodes at time  $t$  with the target node at  $v'(t)$ .

In essence, this criterion exposes the proportion of nodes whose states are predicted correctly given a boolean expression that outputs  $v'(t)$  at every timestep in the series. We used this criterion to choose an optimal boolean expression that maps the states of inferred input nodes to a target node.

### 2.2.3 searchUpdateRule

The `searchUpdateRule` function in **src/BDDUtils.hs** we implemented is where we search for possible boolean expressions that can be evaluated using the top k input nodes for each target node.

The authors from the original study included a second subroutine that iteratively swaps the selected  $k$  nodes with unselected nodes to see if different input nodes could yield higher dynamics consistency metrics—in their algorithm, the resulting number of input nodes could be less than or equal to  $k$ —whereas our implementation, due to time constraints, does not implement this swapping routine and selects a fixed  $k$  input nodes for each target node.

The essence of their swapping routine was a "search\_update\_rule" that searched through

---

**Algorithm 2** Generation of Boolean Network

---

```
procedure GENERATENETWORK( $k, V$ )
    for each node  $v_i \in V$  do
        connections  $\leftarrow \{\}$ 
        create empty BoolNetwork boolNetwork
        for each node  $w \in W (W = V \setminus v_i)$  do
            mutualInfo  $\leftarrow I(v_i, w)$ 
            if mutualInfo in highest  $k$  values then
                connections +=  $w$ 
            end if
        end for
    end for
    return boolNetwork
end procedure
```

---

---

**Algorithm 3** Optimal Boolean Expression Inference

---

```
procedure GETOPTIMALBOOLEXPRESSIONS( $boolNetwork, timeLength$ )
    for each node  $v_i \in$  the nodes in  $boolNetwork$  do
        regulatoryNodes  $\leftarrow boolNetwork(v_i)$ 
        allValidExpressions  $\leftarrow$  createExpressions(regulatoryNodes)
        for each expression  $exp \in$  allValidExpressions do
            find expression with highest geneWiseDynamicCosnsitency value and return
        end for
    end for
end procedure
```

---

possible boolean expressions between the selected input nodes and updated the rule for each target node with the boolean expression that yielded the highest dynamics consistency. Our implementation of this lies in *getOptimalBoolExpressions* (Algorithm 3).

## 3 Haskell Implementation

### 3.1 Data Types

#### 3.1.1 NodeState

*NodeState* represents the states of each node across the time-series. It consists of a *String* attribute as the name of the node/gene, and an *[Int]* to represent the boolean states/expression of that particular node across the time-series from time  $1..T$ , where  $T$  is the length of the time-series.

#### 3.1.2 BoolEdge

*BoolEdge* consists of two *NodeState* attributes— $v\_i$  and  $v\_j$ , which represents a directed edge dependency in the network from node  $v\_i$  to  $v\_j$ .

#### 3.1.3 BoolNetwork

*BoolNetwork* represents a boolean network as a directed graph. It contains a an attribute of *[NodeState]* consisting of the different nodes and their states across the time series, and an attribute of *[BoolEdge]* to represented the directed edges.

#### 3.1.4 BDD

*BDD* is a recursive data type similar to a binary decision diagram. We created a variation to represent and evaluate boolean expressions. *BDD* consists of several different states:

- *Name String*: Represents he name of a paricular boolean variable.
- *State Int*: Represents the state, 1 or 0, of a particular boolean variable.
- *AND BDD BDD*: Represents an AND logical operation between two *BDD* types.
- *OR BDD BDD*: Represents an OR logical operation between two *BDD* types.
- *XOR BDD BDD*: Represents an XOR logical operation between two *BDD* types.
- *NOT BDD*: Represents a NOT operation to be applied to a *BDD*.

The idea with *BDD* is to map out a boolean expression specifying variable names, and supply a *[(String, Int)]* which consists of Int boolean values for each variable name to a *BDD* data type that can evaluate the expression—*evaluateFunc* in `src/BDDUtils.hs` performs this task, like so:

```
-- Define a boolean expression with variable names x, y, a, z:
f = XOR (AND (Name "x") (OR (Name "y") (Name "a")))) (NOT (Name "z"))

-- Evaluate function by supplying values for each variable:
fromEnum $ evaluateFunc f [("x", 1), ("y", 0), ("a", 1), ("z", 1)]
```

This will output 1 when evaluated. For simplicity, based on the original algorithm, we only look at a combination of conjunctive and disjunctive operations (AND/OR) evaluated using the top  $k$  input nodes, sorted in descending order based on each node's mutual information in relation to the target node.

## 3.2 Generating Random Time-Series Data

We generate random time series gene expression data in `src/generate_data.py`, which takes as input the number of nodes and timesteps to be randomly generated at each timestep.

## 3.3 Running GRNPar

This program takes in five arguments:

```
1 stack exec GRNPar-exe <csvFilename> <k> <genExpressions> <genImage> <mode>
```

Specific implementation details can be found in the README. The output of the program is a .png file that has the boolean network mapped out in `src/output_files`. This section will demonstrate parts of the code that contain essential steps in GRNPar.

First, we find the top  $k$  input nodes for each node in `genBoolEdgesSeq`, combining it into one list of `BoolEdge` and creating a `BoolNetwork` with the observed `nodeStates` and all of the generated `BoolEdge` connections (set of interactions).

```
genNetworkSeq :: [NodeState] -> Int -> BoolNetwork
genNetworkSeq nodeStates k = BoolNetwork nodeStates (combineEdgesSeq
  ↳ nodeStates k)

combineEdgesSeq :: [NodeState] -> Int -> [BoolEdge]
combineEdgesSeq nodeStates k = concatMap (genBoolEdgesSeq nodeStates k)
  ↳ nodeStates

genBoolEdgesSeq :: [NodeState] -> Int -> NodeState -> [BoolEdge]
genBoolEdgesSeq nodeStates k targetNode =
  let topKMutual = getMutualInfoSeq targetNode (filter (/= targetNode)
    ↳ nodeStates) k
  in map (`BoolEdge` targetNode) topKMutual
```

Next, we determined the optimal boolean expressions (`getOptimalBoolExpressions` - line 223 `BDDUtils.hs`) for each input node. Using the `boolNetwork` and top  $k$  connections generated for each node, we generate a set of  $2^{k-1}$  possible boolean expressions and use

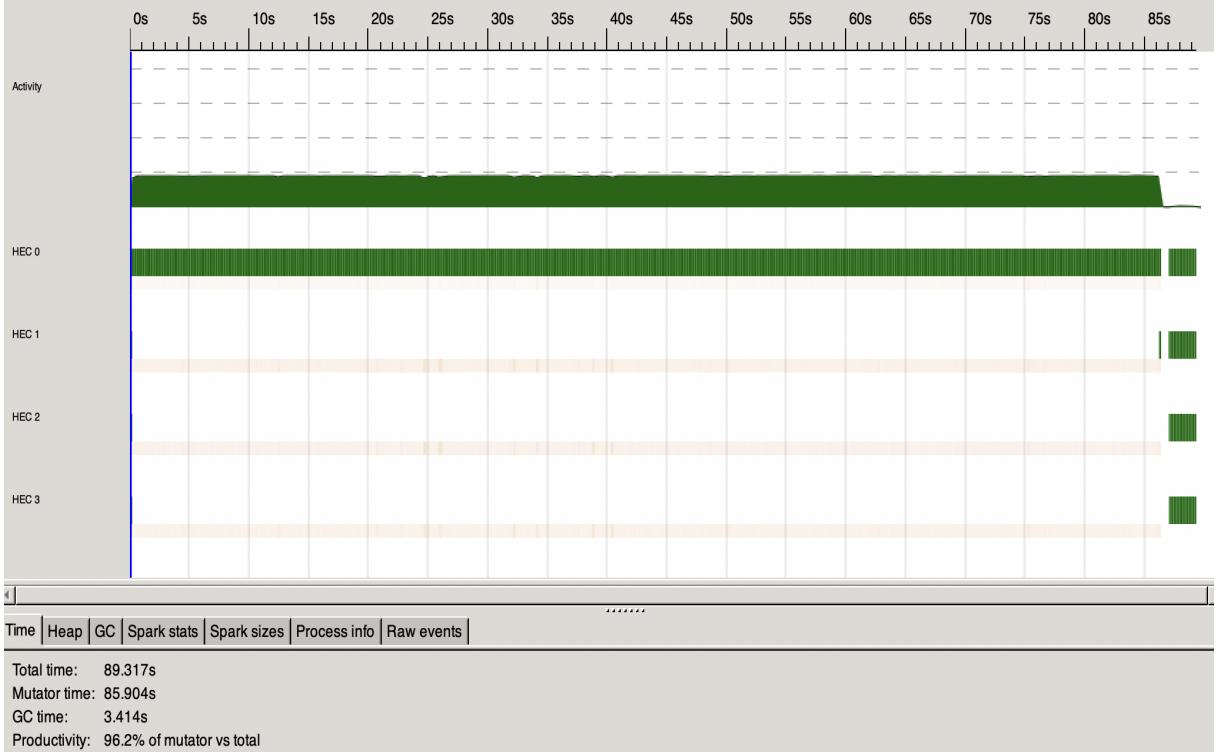


Figure 1: Sequential program for a randomized sample where 300 nodes are generated for 300 time steps and k=5

each to infer a time-series to compare with the observed series (`searchUpdateRule` - line 169 `BDDUtil.hs`). We choose the boolean expression that maximizes the genewise dynamics consistency. (`searchUpdateRule` - line 169 `BDDUtils.hs`).

Lastly, we plotted the BoolNetwork using the `Data.Graph.DGraph` library and `graphviz` library.

```
plotBoolNetworkPng :: BoolNetwork -> FilePath -> Bool -> IO FilePath
plotBoolNetworkPng network fname labelEdges = plotDGPng networkDG fname
  <-> labelEdges
  where
    networkDG = boolNetworkToDG network labelEdges
```

## 4 Parallel Solution

We used `genNetworkPar`, `getOptimalBoolExpressionsPar`, and `plotBoolNetworkPngPar` in `Main.hs` to run parallel strategies. We made `NFData` instances for `NodeState`, `BoolEdge`, and `BoolNetwork` to perform these computations.

There are three layers of parallelism applied to the implementation to improve the performance

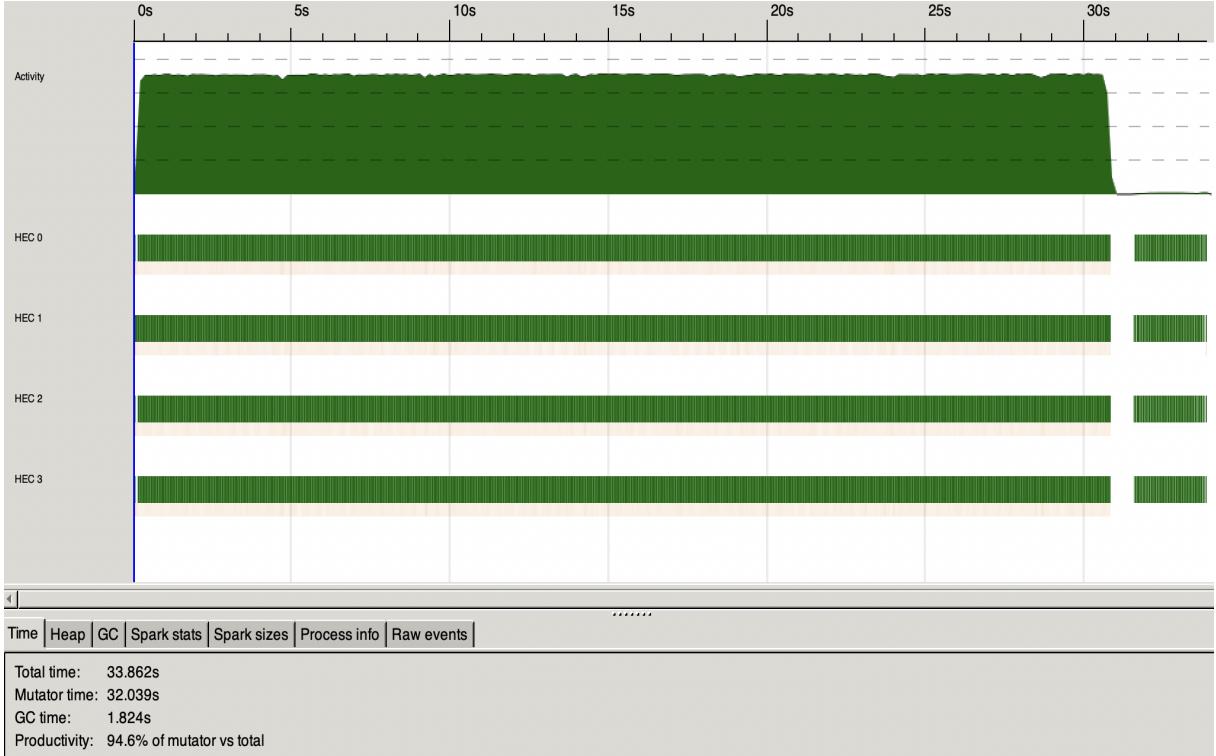


Figure 2: Parallel program for a randomized sample where 300 nodes are generated for 300 time steps and  $k=5$

of the program. The first layer uses the `parMap` `deepSeq` functions from the package `Control.Parallel.Strategies`. The second and third layer use `parBuffer` and `parChunkList`, respectively.

We used `parMap rdeepseq` to parallelize map computations. We used `parMap` to combine the results of generating the top  $k$  connections into one list, shown in the new `combineEdgesPar` function in `GRPPar.hs`. Similarly, when parsing through all the possible nodes to determine the top  $k$  nodes based on mutual information, we applied `parMap rdeepseq` and ran the calculations in parallel.

```
allMutualInfo =
  map (\inp ->
    (inp, mutualInformation (timeStates inp) (timeStates targetNode)
    - sum (parMap rdeepseq (mutualInformation (timeStates inp) . timeStates
      $ map fst regNodes))) inpNodes `using` parBuffer 3 rdeepseq
```

Here, we use `parBuffer` to iterate through the input nodes and calculate mutual information.

Figure 2 shows the eventlog when running parallel on a randomly generated time-series containing 300 nodes and 300 timesteps. The test was run on 4 cores with  $k = 5$ . The average runtime was 33.9 seconds.

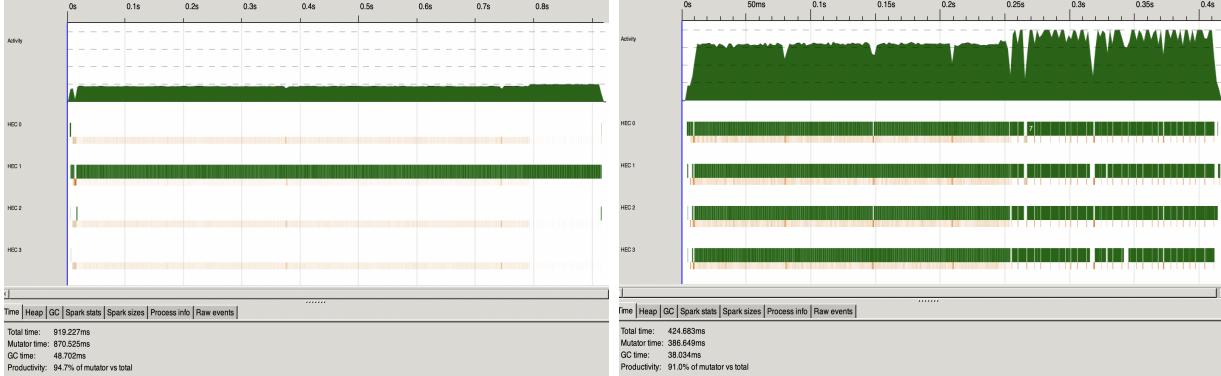


Figure 3: Sequential (left) and Parallel (right) program for GRNPar on Insilico E. coli Dataset using 4 cores.  $k = 8$ .

## 5 Performance Analysis

### 5.1 Results

#### 5.1.1 Real-world Datasets

We tested GRNPar with three real gene expression datasets: an E. coli network with 10 genes and 20 timesteps, an Insilico E.coli network with 100 genes and 20 timesteps, and a fission yeast cell cycle network with 10 genes and 10 timesteps. We discretized the InSilico data using a rudimentary KMeans clustering algorithm, and the other datasets had been discretized to boolean states using KMeans clustering prior by Barman & Kwon.

We ran the Insilico E. coli dataset through GRNPar, only generating the network without computing expressions or outputting an image, both in parallel and sequentially—the results are shown in Figure 3. There was a 2.16x speedup in total time and 2.25x reduction in mutator time from the sequential to parallel run.

When running the Insilico data through GRNPar, there was an extra period for generating the image, in both the sequential and parallel versions of the algorithm. This is seen in the extra time taken to complete the algorithm in Figure 4. In both the sequential and parallel algorithms, it took 1.5 seconds extra to run as compared with the runs when the graph image was not generated.

The graph generated for the Insilico set is shown in figure 5. There is a tendency for the central nodes to be in the center of the graph, as they have more dependencies. Figure 6 shows the graphs generated by testing GRNPar on the other E. coli dataset and fission yeast cell cycle network—similar to the Insilico set, we observe 3.33x speedup between the sequential and parallel usage of GRNPar.

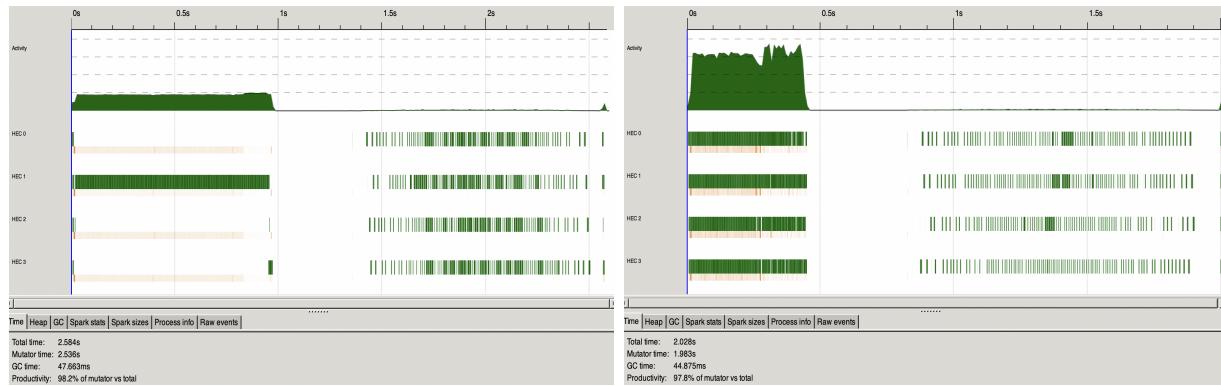


Figure 4: Generation of network + plotting image  $-i$ , both the sequential (left) and parallel (right) implementation took 1.5 extra seconds, compared to figure 3

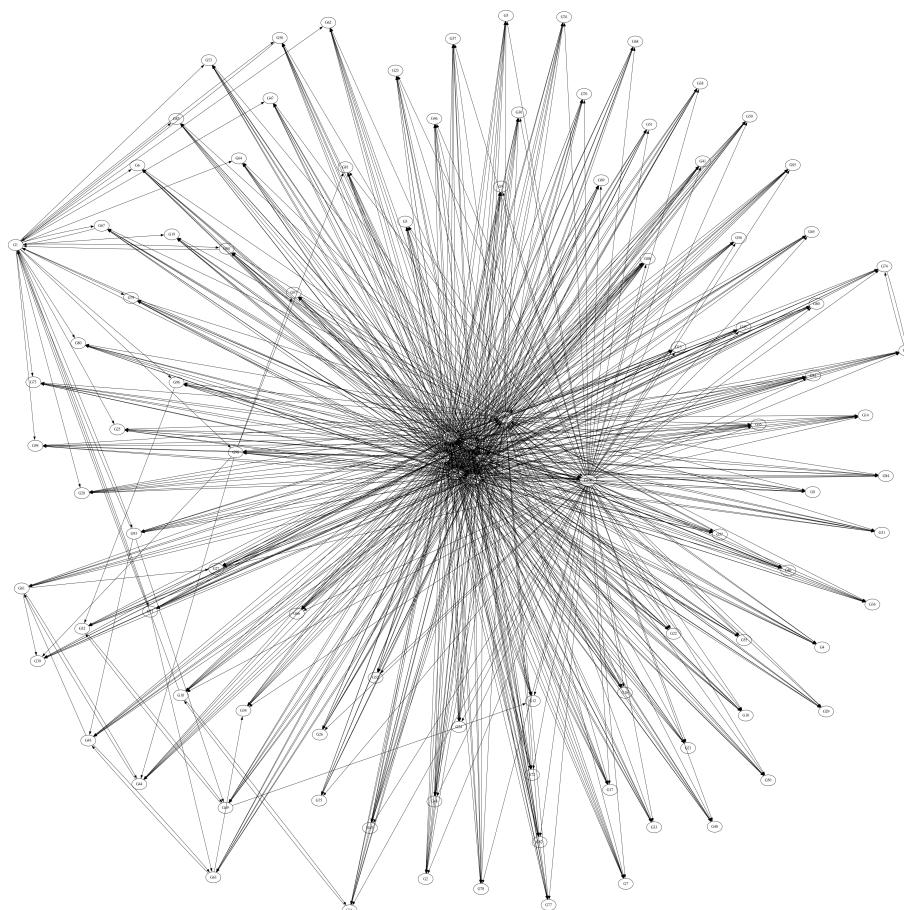


Figure 5: Boolean Network Inferred for InSilico Ecoli Protein Trajectories

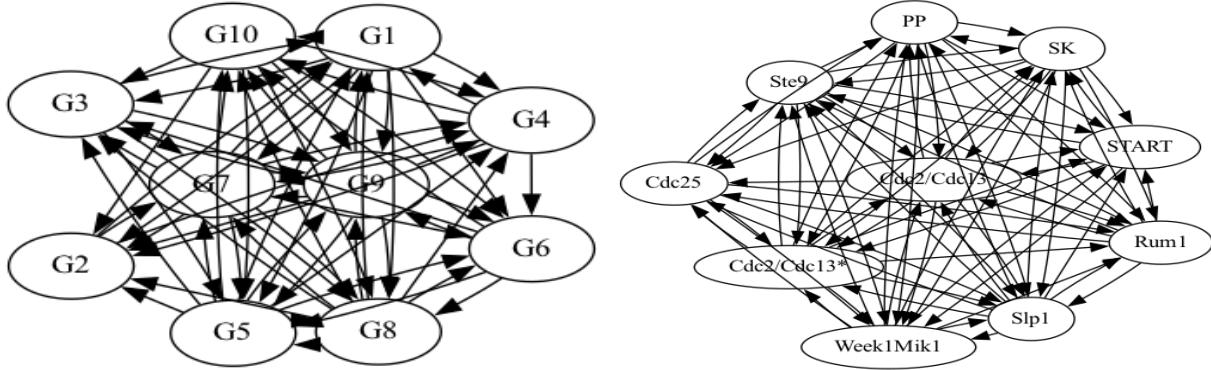


Figure 6: E. coli gene regulatory network(left) & a fission yeast cell cycle network (right)

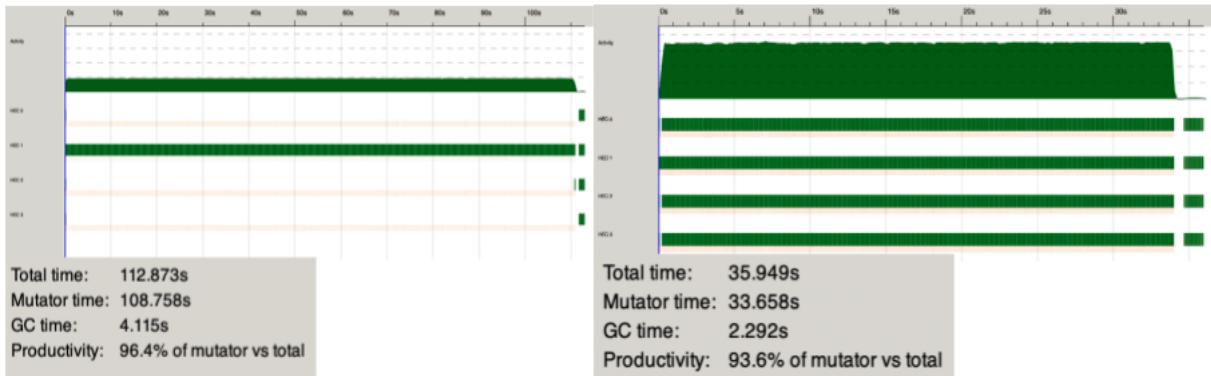


Figure 7: Randomly Generated Data: 500 Nodes over 300 time steps (4 cores) in a sequential and parallel method

## 5.2 Analysis over Randomly Generated Data

In this section, we analyzed how changing the number of nodes (genes), cores, and k-values affected performance of GRNPar on randomly generated time-series expression data. Figure 8 shows three tables worth of runtime data. In Table 1, as the number of nodes (genes) increases, the runtime for sequential GRNPar increases at a much larger rate than parallel GRNPar. However, when compared with the parallel runtime in Table 2 (Figure 8), the sequential runtime increases as the number of cores increases, since unnecessary garbage collection occurs during the sequential run as more cores are added on.

## 6 References

Barman S, Kwon YK (2017) A novel mutual information-based Boolean network inference method from time-series gene expression data. PLOS ONE 12(2): e0171097.  
<https://doi.org/10.1371/journal.pone.0171097>

Prill RJ, Marbach D, Saez-Rodriguez J, Sorger PK, Alexopoulos LG, Xue X, et al. Towards a Rigorous Assessment of Systems Biology Models: The DREAM3 Challenges. PLOS One.

Nodes = 300, k = 5			
Cores	runtime		
	par	seq	
4	12.714	40.147	
5	11.586	40.825	
6	9.948	42.407	
7	9.096	44.005	
8	9.095	44.932	

Nodes = 300, Cores = 4		
k	runtime	
	par	seq
3	12.493	41.234
4	18.758	62.14
5	27.429	87.863
6	39.134	122.933

Figure 8: From left to right, Table 1, 2, 3. Table 1: Change in runtime for a parallel and sequential implementation as the number of nodes increases. Table 2: Change in runtime for a parallel and sequential implementation as the number of cores increases. Table 3: Change in runtime for a parallel and sequential implementation as the k value increases

## 7 Code Listing

### 7.1 app/Main.hs

```

1 {-
2 Main application program:
3 -}
4
5 module Main (main) where
6
7 import System.IO ()
8 import System.Environment (getArgs)
9 import Control.Monad
10
11 import GRNPar (genNetworkSeq, genNetworkPar)
12 import GraphUtils (plotBoolNetworkPng, plotBoolNetworkPngPar, NodeState(..))
13 import BDDUtils (getOptimalBoolExpressions, getOptimalBoolExpressionsPar)
14 import ProcessData (csvToNodeStates)
15
16 import Control.DeepSeq
17
18 {-
19 Main script:
20   1. Infer + generate boolean network from input file
21   2. Determine optimal boolean expression
22   3. Output network to png file
23
24 Usage: GRNPar-exe <csvfilename> <k> <outputFile> <genExpressions> <genImage>
25   ↳ <mode>
26 -}
27 main :: IO ()
28 main = do
29   args <- getArgs

```

```

29     case args of
30         [fname, k, genExpressions, genImage, mode] -> do
31             putStrLn fname
32             putStrLn "Parsing data..."
33             nodeStates@(_x) <- csvToNodeStates fname 1
34
35             -- Generate network with dependencies
36             let k'          = read k :: Int
37                 timeLength = length $ timeStates x
38                 network     = if mode == "par"
39                             then force $ genNetworkPar nodeStates k'
40                         else force $ genNetworkSeq nodeStates k'
41             putStrLn "Generated network."
42
43             when (genExpressions == "1") $ do
44                 -- Get optimal boolean expressions for each node
45                 let optimalExpressions = if mode == "par"
46                     then force $
47                         getOptimalBoolExpressionsPar network timeLength
48                         else force $
49                         getOptimalBoolExpressions network timeLength
50                         putStrLn "Optimal boolean expressions for each variable:"
51                         mapM_ (\(n, b, c) -> putStrLn $ name n ++ " = " ++ show b ++ "\n" ++ ". Gene-wise consistency: " ++ show c) optimalExpressions
52                         --print optimalExpressions
53
54             when (genImage == "1") $ do
55                 -- Print image
56                 let outputFile = "src/output_files/" ++ substring (length
57                     "src/data/") (length fname - 4) fname
58                     print outputFile
59                     imgFilepath <- if mode == "par"
60                         then plotBoolNetworkPngPar network
61                         else plotBoolNetworkPng network outputFile
62                     outputFile False
63                     putStrLn $ "Printed to " ++ imgFilepath ++ "."
64
65             putStrLn "Finished."
66             _ -> putStrLn "Usage: GRNPar-exe <csvFilename> <k> <outputFile>
67             <genExpressions> <genImage> <mode>"
68
69     substring :: Int -> Int -> String -> String
70     substring x y s = take (y - x) (drop x s)

```

## 7.2 src/GRNPar.hs

```
1 {-# LANGUAGE ParallelListComp #-}  
2  
3 module GRNPar  
4     ( genNetworkSeq  
5     , genNetworkPar  
6     ) where  
7  
8  
9 import Control.Monad ()  
10  
11 import qualified Data.Map as Map  
12 import Data.Ord (comparing)  
13 import Data.List (maximumBy, sortBy)  
14  
15 import GraphUtils (NodeState(..), BoolEdge(..), BoolNetwork(..))  
16  
17 import Control.Parallel.Strategies (parMap, rdeepseq, parBuffer, using)  
18 import Control.DeepSeq ()  
19  
20 {-  
21 -- Testing with sample nodes:  
22 x_1 :: NodeState  
23 x_1 = NodeState "x_1" [1, 1, 0, 0, 0, 1, 0, 1]  
24  
25 x_2 :: NodeState  
26 x_2 = NodeState "x_2" [1, 0, 1, 0, 1, 0, 1, 0]  
27  
28 x_3 :: NodeState  
29 x_3 = NodeState "x_3" [1, 0, 1, 0, 1, 0, 1, 0]  
30  
31 x_4 :: NodeState  
32 x_4 = NodeState "x_4" [1, 0, 0, 0, 1, 0, 1, 1]  
33  
34 x_5 :: NodeState  
35 x_5 = NodeState "x_5" [1, 0, 1, 0, 0, 0, 0, 0]  
36  
37 x_6 :: NodeState  
38 x_6 = NodeState "x_6" [1, 1, 0, 0, 1, 1, 1, 0]  
39  
40 sampleNodes :: [NodeState]  
41 sampleNodes = [x_1,x_2,x_3,x_4,x_5,x_6]  
42 -}  
43  
44 {-  
45 Parallel pipeline for inferring boolean network given NodeState data.  
46 -}
```

```

47 genNetworkPar :: [NodeState] -> Int -> BoolNetwork
48 genNetworkPar nodeStates k = BoolNetwork nodeStates (combineEdgesPar nodeStates
   ↪ k)
49
50 {-
51 Sequential pipeline for inferring boolean network given NodeState data.
52 -}
53 genNetworkSeq :: [NodeState] -> Int -> BoolNetwork
54 genNetworkSeq nodeStates k = BoolNetwork nodeStates (combineEdgesSeq nodeStates
   ↪ k)
55
56 {-
57 Combine all BoolEdge connections for each target node into one array.
58 -}
59 combineEdgesSeq :: [NodeState] -> Int -> [BoolEdge]
60 combineEdgesSeq nodeStates k = concatMap (genBoolEdgesSeq nodeStates k)
   ↪ nodeStates
61
62 {-
63 Generate BoolEdge connections in network given NodeState and a targetNode
64 -}
65 genBoolEdgesSeq :: [NodeState] -> Int -> NodeState -> [BoolEdge]
66 genBoolEdgesSeq nodeStates k targetNode =
67   let topKMutual = getMutualInfoSeq targetNode (filter (/= targetNode)
   ↪ nodeStates) k
68   in map (`BoolEdge` targetNode) topKMutual
69
70 combineEdgesPar :: [NodeState] -> Int -> [BoolEdge]
71 combineEdgesPar nodeStates k = concat $ parMap rdeepseq (genBoolEdgesPar
   ↪ nodeStates k) nodeStates
72
73 genBoolEdgesPar :: [NodeState] -> Int -> NodeState -> [BoolEdge]
74 genBoolEdgesPar nodeStates k targetNode = map (`BoolEdge` targetNode) topKMutual
75   where
76     topKMutual = getMutualInfoPar targetNode (filter (/= targetNode) nodeStates)
   ↪ k
77
78 {-
79 Get top k input nodes with the highest mutual information relative to a target
   ↪ node.
80
81 Return input nodes sorted descending based on mutual information with target
   ↪ node.
82
83 * inputNodes should not contain targetNode when calling getMutualInfo
84 -}
85 getMutualInfoSeq :: NodeState -> [NodeState] -> Int -> [NodeState]

```

```

86  getMutualInfoSeq targetNode inputNodes k = map fst
87          $ sortBy (flip (comparing snd))
88          $ getMutualInfo' [maxMutualNodeInfo]
89          (filter (/= maxMutual) inputNodes) k
90
91  where
92      maxMutualNodeInfo@(maxMutual, _) = maximumBy (comparing snd)
93          $ map (\inp -> (inp, mutualInformation
94          ← (timeStates inp) (timeStates targetNode))) inputNodes
95      getMutualInfo' :: [(NodeState, Double)] -> [NodeState] -> Int -> [(NodeState,
96          Double)]
97      getMutualInfo' regNodes inpNodes k'
98          | length regNodes == k' = regNodes
99          | otherwise             =
100          let newMaxInfo@(newMax, _) = maximumBy (comparing snd)
101              $ map (\inp -> (inp, mutualInformation
102              ← (timeStates inp) (timeStates targetNode)
103                  - sum (map (mutualInformation
104                  ← (timeStates inp) . timeStates) $ map fst regNodes))) inpNodes
105              newInpNodes           = filter (/= newMax) inpNodes
106              newRegNodes          = regNodes ++ [newMaxInfo]
107          in getMutualInfo' newRegNodes newInpNodes k'
108
109
110  getMutualInfoPar :: NodeState -> [NodeState] -> Int -> [NodeState]
111  getMutualInfoPar targetNode inputNodes k = map fst
112          $ sortBy (flip (comparing snd))
113          $ getMutualInfo' [maxMutualNodeInfo]
114          (filter (/= maxMutual) inputNodes) k
115
116  where
117      mutualInfo'           = map (\inp -> (inp, mutualInformation
118          ← (timeStates inp) (timeStates targetNode)))
119          inputNodes `using` parBuffer 3
120
121      rdeepseq
122      maxMutualNodeInfo@(maxMutual, _) = maximumBy (comparing snd) mutualInfo'
123      getMutualInfo' :: [(NodeState, Double)] -> [NodeState] -> Int -> [(NodeState,
124          Double)]
125      getMutualInfo' regNodes inpNodes k'
126          | length regNodes == k' = regNodes
127          | otherwise             =
128          let allMutualInfo        = map (\inp -> (inp, mutualInformation
129          ← (timeStates inp) (timeStates targetNode)
130              - sum (parMap rdeepseq
131              ← (mutualInformation (timeStates inp) . timeStates) $ map fst regNodes)))
132                  inputNodes `using` parBuffer 3
133          in rdeepseq
134          newMaxInfo@(newMax, _) = maximumBy (comparing snd) allMutualInfo
135          newInpNodes           = filter (/= newMax) inpNodes
136          newRegNodes          = regNodes ++ [newMaxInfo]

```

```

123     in getMutualInfo' newRegNodes newInpNodes k'
124
125 -- Compute the entropy of a list of numerical values
126 entropy :: (Ord a) => [a] -> Double
127 entropy xs = sum [ -p * logBase 2 p | p <- ps ]
128 where
129     counts = Map.fromListWith (+) [(x, 1 :: Int) | x <- xs]
130     values = Map.elems counts
131     total = sum values
132     ps = [fromIntegral v / fromIntegral total | v <- values]
133
134 -- Compute the mutual information of discrete variables x and y
135 mutualInformation :: (Ord a) => [a] -> [a] -> Double
136 mutualInformation xs ys = entropy xs + entropy ys - entropy (zip xs ys)

```

### 7.3 src/GraphUtils.hs

```

1 module GraphUtils
2     ( NodeState(..)
3     , BoolEdge(..)
4     , BoolNetwork(..)
5     , plotBoolNetworkPng
6     , plotBoolNetworkPngPar
7     , plotDGPng
8     ) where
9
10 import Data.GraphViz
11 import Data.GraphViz.Attributes.Complete
12 import Data.Hashable
13
14 import Control.Parallel.Strategies (parListChunk, rdeepseq, using)
15 import Control.DeepSeq ( NFData(..) )
16
17 import qualified Data.Text.Lazy      as TL
18 import qualified Data.Graph.DGraph as DG
19 import qualified Data.Graph.Types   as GT
20
21 {-
22 Data types representing boolean states of genes and boolean network consisting of
23 → NodeStates and directed edges.
24
25 - NodeState:
26     - name :: String -> name of gene
27     - timeStates :: [Int] -> expression of gene across time-series in order
28 -}
29 data NodeState = NodeState { name :: String

```

```

30 , timeStates :: [Int] } deriving (Eq, Ord)
31
32 instance Show NodeState where
33   show (NodeState n _) = n
34
35 instance NFData NodeState where
36   rnf (NodeState n ts) = rnf n `seq` rnf ts
37
38 {-#
39 - BoolEdge
40   - v_i :: NodeState -> source node
41   - v_j :: NodeState -> target node
42 #-}
43 data BoolEdge = BoolEdge { v_i :: NodeState
44                           , v_j :: NodeState } deriving (Eq)
45
46 instance Show BoolEdge where
47   show (BoolEdge i j) = "(" ++ name i ++ " -> " ++ name j ++ ")"
48
49 instance NFData BoolEdge where
50   rnf (BoolEdge i j) = rnf i `seq` rnf j
51
52 {-#
53 - BoolNetwork
54   - nodes :: [NodeState] -> all nodes/genes in network
55   - connections :: [BoolEdge] -> directed edges between genes
56 #-}
57 data BoolNetwork = BoolNetwork { nodes :: [NodeState]
58                                   , connections :: [BoolEdge] } deriving (Eq, Show)
59 instance NFData BoolNetwork where
60   rnf (BoolNetwork n c) = rnf n `seq` rnf c
61
62 {-#
63 Modified from graphite lib source code:
64 - Functions for plotting directed graph and BoolNetwork to png file
65 - graphviz needs to be installed on local machine to see output:
66   - brew install graphviz
67     OR
68   - sudo apt install graphviz
69
70 Example: Plot DGraph foundationUniverse to "foundation.png"
71
72 foundationUniverse :: DG.DGraph String Double
73 foundationUniverse = DG.fromArcsList
74   [ GT.Arc "Helicon" "Nishaya" 200.00
75   , GT.Arc "Helicon" "Wencory" 382.20
76   , GT.Arc "Nishaya" "Wencory" 820.32

```



```

119 labeledArcs :: (Hashable v, Eq v, Show e) => DG.DGraph v e -> [(v, v, String)]
120 labeledArcs g = (\(GT.Arc v1 v2 attr) -> (v1, v2, show attr)) <$> DG.arcs g
121
122 toDirectedDot :: (Hashable v, Ord v, Show v, Show e)
123   => Bool -- ^ Label edges
124   -> DG.DGraph v e
125   -> DotGraph v
126 toDirectedDot labelEdges g = graphElemsToDot params (labeledNodes g) (labeledArcs
127   -> g)
128   where params = sensibleDotParams True labelEdges
129
130 sensibleDotParams
131   :: Bool -- ^ Directed
132   -> Bool -- ^ Label edges
133   -> GraphvizParams t l String () l
134 sensibleDotParams directed edgeLabeled = nonClusteredParams
135   { isDirected = directed
136   , globalAttributes =
137     [ GraphAttrs [Overlap ScaleOverlaps]
138       , EdgeAttrs [FontColor (X11Color DarkGreen)]
139     ]
140   , fmtEdge = edgeFmt
141   }
142   where
143     edgeFmt (_, _, l) = if edgeLabeled
144       then [Label $ StrLabel $ TL.pack l]
145       else []

```

## 7.4 src/BDDUtils.hs

```

1 {-# LANGUAGE TupleSections #-}
2
3 module BDDUtils
4   ( BDD(..)
5   , evaluateFunc
6   , getOptimalBoolExpressions
7   , getOptimalBoolExpressionsPar
8   , getRegulatoryNodes
9   , searchUpdateRule
10  , getBDDFromFunc
11  ) where
12
13 import GraphUtils (NodeState(..), BoolEdge(..), BoolNetwork(..))
14
15 import Data.Ord (comparing)
16 import Data.List (maximumBy)
17

```

```

18 import qualified Data.Matrix as M
19 import qualified Data.Vector as Vec
20
21 import Control.Parallel.Strategies (parMap, rdeepseq, using, parListChunk,
22                                     ↳ parBuffer)
22 import Control.DeepSeq(NFData(..))
23
24 {-
25   Functions for evaluating and representing binary decision diagrams.
26 -}
27
28 {-
29   Binary decision diagram data type used for evaluating boolean expressions.
30
31 Example:
32 - f = XOR (AND (Name "x") (OR (Name "y") (Name "a")))) (NOT (Name "z"))
33 - fromEnum $ evaluateFunc f [("x", 1), ("y", 0), ("a", 1), ("z", 1)]
34   - Output: 1
35 -}
36
37 -- Define binary decision diagram
38 data BDD = Name String -- name of node (gene name)
39           | State Int    -- state of node (0 or 1)
40           | AND BDD BDD -- an AND node
41           | OR BDD BDD  -- an OR node
42           | XOR BDD BDD -- an XOR node
43           | NOT BDD      -- a NOT node
44           deriving (Eq)
45
46 instance Show BDD where
47   show (Name x) = x
48   show (State x) = show x
49   show (AND x y) = "(" ++ show x ++ " AND " ++ show y ++ ")"
50   show (OR x y) = "(" ++ show x ++ " OR " ++ show y ++ ")"
51   show (XOR x y) = "(" ++ show x ++ " XOR " ++ show y ++ ")"
52   show (NOT x) = "(NOT " ++ show x ++ ")"
53
54 instance NFData BDD where
55   rnf (Name x) = rnf x
56   rnf (State x) = rnf x
57   rnf (AND x y) = rnf x `seq` rnf y
58   rnf (OR x y) = rnf x `seq` rnf y
59   rnf (XOR x y) = rnf x `seq` rnf y
60   rnf (NOT x) = rnf x
61
62 {-
63   Evaluate a boolean expression given a BDD data type and a list of tuples
64   ↳ specifying values for each variable.

```

```

64
65 Params:
66 - BDD data type
67 - func: list of tuples [(varName, value)], where varName is a boolean variable
   ↳ name in the BDD, and value is the value assigned to it
68
69 Example:
70 - f = XOR (AND (Name "x") (OR (Name "y") (Name "a")))) (NOT (Name "z"))
71 - fromEnum $ evaluateFunc f [("x", 1), ("y", 0), ("a", 1), ("z", 1)]
72   - Output: 1
73
74 -}
75 -- evaluate a BDD for a given assignment of the variables
76 evaluateFunc :: BDD -> [(String, Int)] -> Bool
77 evaluateFunc (Name x) func = case lookup x func of
78   Just a -> a == 1
79   Nothing -> error "Variable not in assignment"
80 evaluateFunc (AND x y) func = evaluateFunc x func && evaluateFunc y func
81 evaluateFunc (OR x y) func = evaluateFunc x func || evaluateFunc y func
82 evaluateFunc (XOR x y) func
83   | andTrue = False
84   | otherwise = evaluateFunc (OR x y) func
85 where
86   andTrue = evaluateFunc x func && evaluateFunc y func
87 evaluateFunc (NOT b) func = not (evaluateFunc b func)
88 evaluateFunc (State x) _
89   | x == 1 = True
90   | otherwise = False
91
92
93 {-
94 Get gene wise dynamics consistency metric given a predicted time series and the
   ↳ actual time series.
95
96 Params:
97 - predictedStates :: [Int]
98 - actualStates :: [Int]
99
100 -}
101 geneWiseDynamicsConsistency :: [Int] -> [Int] -> Double
102 geneWiseDynamicsConsistency predictedStates actualStates = sumPredictions /
   ↳ fromIntegral timeLength
103 where
104   sumPredictions = sum $ zipWith (\x y -> if x == y then 1 :: Double else 0 ::
   ↳ Double) predictedStates actualStates
105   timeLength = length predictedStates
106

```

```

107 {-  

108 Get all combinations of conjunctive and disjunctive boolean expressions.  

109  

110 Params:  

111 - n :: Int -> number of operators to evaluate  

112  

113 Returns:  

114 - [[Int]] -> list of Int lists, where each list corresponds to a sequence of  

    ↳ AND/OR operations.  

115  

116 Each list, e.g. [1, 0, 1], signifies a sequence of AND/OR operators, where 1 =  

    ↳ AND and 0 = OR  

117  

118 Example: given a list of truth values, such as [1, 0, 0, 1], and a boolean  

    ↳ expression combination represented as [1, 0, 1]:  

119 - Using the BDD data type, we would evaluate this as: ((1 AND 0) OR 0) AND 1  

120 -}  

121 getConjDisjCombos :: Int -> [[Int]]  

122 getConjDisjCombos 1 = [[0], [1]]  

123 getConjDisjCombos n = [[mod x 2^i | i <- [0..n-1]] | x <- [0..(2^n)-1]]  

124  

125 {-  

126 Construct a BDD given boolean variable names and a combination of  

    ↳ conjunctive/disjunctive operations.  

127  

128 Params:  

129 - (x:xs) :: [String] -> boolean variable names  

130 - ops :: [Int] -> conjunctive/disjunctive combinations as one element outputted  

    ↳ by getConjDisjCombos  

131  

132 Returns: BDD  

133  

134 The number of variables should always be 1 more than the number of operations.  

135  

136 Example:  

137 - getBDDFromFunc ["v1", "v2", "v3"] [1, 0] returns a BDD of: ((v1 AND v2) OR v3)  

138 - bdd = getBDDFromFunc ["v1", "v2", "v3"] [1, 0]  

139 - fromEnum $ evaluateFunc bdd [("v1", 1), ("v2", 0), ("v3", 1)]  

140 - Output: 1  

141 -}  

142  

143 getBDDFromFunc :: [String] -> [Int] -> BDD  

144 getBDDFromFunc [] _      = error "Invalid input."  

145 getBDDFromFunc [_] _      = error "Invalid input."  

146 getBDDFromFunc (x:xs) ops = foldl (\acc (y, ys) -> if ys == 1 then AND acc (Name  

    ↳ y) else OR acc (Name y)) (Name x) tailZipped  

    where

```

```

148     tailZipped = zip xs ops
149
150 {-
151 Get regulatory genes of node given target node and boolean network.
152
153 Params:
154 - targetNode: target gene
155 - network : BoolNetwork
156 -}
157 getRegulatoryNodes :: NodeState -> BoolNetwork -> [NodeState]
158 getRegulatoryNodes targetNode network = map v_i $ filter (\(BoolEdge _ out) ->
159   → out == targetNode) $ connections network
160
161 {-
162 Compute genewise dynamics consistency metric for possible boolean expressions in
163   → input nodes sorted by mutual information.
164
165 Params:
166 - inpNodes :: [NodeState] -> top k input nodes with highest mutual information
167   → with respect to targetNode
168 - targetNode :: NodeState
169 - timeLength :: Int -> length of time series
170 -}
171 searchUpdateRule :: [NodeState] -> NodeState -> Int -> [(Double, BDD)]
172 searchUpdateRule inpNodes targetNode timeLength =
173   case allTargetStates of
174     []          -> error "Invalid target states."
175     [_]         -> error "Invalid target states."
176     (_:targetStates) -> map (\(p, r) -> (geneWiseDynamicsConsistency p
177       → targetStates, r)) predStates
178   where
179     -- Get target states of target node
180     allTargetStates = timeStates targetNode
181
182     inpNodeNames = map name inpNodes
183     inpMatrix    = M.fromLists $ map (\xs -> map (name xs,) (timeStates xs))
184     → inpNodes
185
186     -- Get boolean expression combos
187     ruleCombos   = map (getBDDFromFunc inpNodeNames) $ getConjDisjCombos (length
188       → inpNodes - 1)
189
190     -- Predict states using boolean expression
191     predStates   = map (\ruleBDD ->
192       let p = map (fromEnum . \t -> evaluateFunc ruleBDD
193       → (filter (\(nn, _) -> nn `elem` inpNodeNames)

```

```

188                                         $ Vec.toList (M.getCol t
189                                         ↵  inpMatrix)))
190                                         [1..(timeLength - 1)]
191                                         in (p, ruleBDD))
192                                         ruleCombos
193
194 searchUpdateRulePar :: [NodeState] -> NodeState -> Int -> [(Double, BDD)]
195 searchUpdateRulePar inpNodes targetNode timeLength =
196   case allTargetStates of
197     []          -> error "Invalid target states."
198     [_]         -> error "Invalid target states."
199     (_:targetStates) -> map (\(p, r) -> (geneWiseDynamicsConsistency p
200     ↵  targetStates, r)) predStates
201   where
202     -- Get target states of target node
203     allTargetStates = timeStates targetNode
204
205     inpNodeNames      = map name inpNodes
206     inpMatrix        = M.fromLists $ map (\xs -> map (name xs,) (timeStates xs)
207     ↵  `using` parBuffer 50 rdeepseq) inpNodes
208
209     -- Get boolean expression combos
210     ruleCombos       = map (getBDDFromFunc inpNodeNames) (getConjDisjCombos
211     ↵  (length inpNodes)) `using` parBuffer 100 rdeepseq
212
213     -- Predict states using boolean expression
214     predStates        = parMap rdeepseq (\ruleBDD ->
215       let p = map (fromEnum . \t -> evaluateFunc ruleBDD
216       ↵  (filter (\(nn, _) -> nn `elem` inpNodeNames)
217           $ Vec.toList (M.getCol t inpMatrix)))
218           [1..(timeLength - 1)]
219           `using` parBuffer 50 rdeepseq
220           in (p, ruleBDD))
221           ruleCombos
222
223 {-
224 Get optimal boolean expressions for each node in network that optimizes genewise
225 dynamics consistency.
226 -}
227 getOptimalBoolExpressions :: BoolNetwork -> Int -> [(NodeState, BDD, Double)]
228 getOptimalBoolExpressions inferredNetwork timeLength = map (\targetNode ->
229   let inpNodes
230     = getRegulatoryNodes targetNode inferredNetwork
231     -- Get
232     consistencyMetrics for each boolean expression
233
234     consistencyMetrics      = searchUpdateRule inpNodes targetNode
235     timeLength

```

```

226                                         -- Get boolean
227     ← expression with max consistency metric
228     ← optimalBDD) = maximumBy (comparing fst) consistencyMetrics
229                                         (maxConsistency,
230                                         in (targetNode,
231                                         ← optimalBDD, maxConsistency))
232                                         $ nodes inferredNetwork
233
234                                         getOptimalBoolExpressionsPar :: BoolNetwork -> Int -> [(NodeState, BDD, Double)]
235                                         getOptimalBoolExpressionsPar inferredNetwork k = map (\targetNode ->
236                                         let inpNodes
237                                         ← = getRegulatoryNodes targetNode inferredNetwork
238                                         consistencyMetrics
239                                         ← = searchUpdateRulePar inpNodes targetNode k
240                                         (maxConsistency,
241                                         ← optimalBDD) = maximumBy (comparing fst) consistencyMetrics
242                                         in (targetNode,
243                                         ← optimalBDD, maxConsistency))
244                                         (nodes inferredNetwork)
245                                         `using` parListChunk 50 rdeepseq

```

## 7.5 src/ProcessData.hs

```

1 module ProcessData
2   ( csvToNodeStates
3   ) where
4
5 import qualified Data.ByteString.Char8 as C
6
7 import GraphUtils (NodeState(..))
8 import qualified Data.Matrix as M
9 import qualified Data.Vector as Vec
10
11
12 -- Get CSV Data
13 -- Return (header, values as multi-D list)
14 getCSVData :: String -> IO ([String], [[Int]])
15 getCSVData fname = do
16   inp <- C.readFile fname
17   let inpLines = C.lines inp
18   case inpLines of
19     []          -> error "Invalid csv file."
20     [_]         -> error "Invalid csv file."
21     (header:csvLines) -> do
22       let headerFormatted = map C.unpack (C.split ',' header)
23       csvLinesFormatted = map (map (\l2 -> read (C.unpack l2) :: Int) .
24     ← C.split ',') csvLines

```

```

24         return (headerFormatted, csvLinesFormatted)
25
26 {-  

27 Params:
28   - x: csvData returned by getCSVData
29   - axis: 0 = row, 1 = column
30 -}  

31 parseCSVData :: IO ([String], [[Int]]) -> Int -> IO [NodeState]
32 parseCSVData x axis = do
33   csvData <- x
34   let (header, csvList) = csvData
35   csvMatrix           = M.fromLists csvList
36   nodeStates          = foldl (\acc (i, hName) -> acc ++ [NodeState hName
37   ↵  (Vec.toList $  

38   ↵  if axis == 0  

39   ↵    then M.getRow i  

40   ↵  ↵  else M.getCol i  

41   ↵  ↵  )])  

42   ↵  []
43   ↵  $ filter (\(_ , n) -> n /= "time") (zip [0..]
44   ↵  ↵  header)
45
46 {-  

47 Params:
48   - fname: csv filename
49   - axis: 0 = row, 1 = column
50
51 Example:
52 t <- csvToNodeStates "../test2.csv" 1
53 -}  

54 csvToNodeStates :: String -> Int -> IO [NodeState]
55 csvToNodeStates fname = parseCSVData (getCSVData fname)

```

## 7.6 src/generate\_data.py

```

1 # Generate large test samples
2 """
3 Main script for generating random large datasets of input/output gene expression
4   ↵ samples
5 """
6 import random
7 import pandas as pd

```

```

8 import numpy as np
9 import argparse
10
11 def generate_io_pairs(num_nodes, time):
12     nodes = ["v" + str(i) for i in range(1, num_nodes + 1)] + ["v" + str(i) + "!""
13     ↪ for i in range(1, num_nodes + 1)]
14     state_series = []
15
16     inp_seen = []
17     out_seen = []
18
19     for t in range(1, time + 1):
20         rand_states_inp = [random.randint(0, 1) for _ in range(num_nodes)]
21         rand_states_out = [random.randint(0, 1) for _ in range(num_nodes)]
22
23         inp_done = False
24         out_done = False
25
26         while not (inp_done and out_done):
27             if rand_states_inp not in inp_seen:
28                 inp_seen.append(rand_states_inp)
29                 inp_done = True
30             else:
31                 rand_states_inp = [random.randint(0, 1) for _ in
32                     ↪ range(num_nodes)]
33
34             if rand_states_out not in out_seen:
35                 out_seen.append(rand_states_out)
36                 out_done = True
37             else:
38                 rand_states_out = [random.randint(0, 1) for _ in
39                     ↪ range(num_nodes)]
40
41         state_series.append(rand_states_inp + rand_states_out)
42
43     state_series_df = pd.DataFrame(state_series, columns = nodes)
44     return state_series_df
45
46
47 def generate_data_time_series(num_nodes, time):
48     nodes = ["v" + str(i) for i in range(1, num_nodes + 1)]
49     state_series = []
50     states_seen = []
51     for t in range(1, time + 1):
52         rand_states = [random.randint(0, 1) for _ in range(len(nodes))]
53         while True:
54             if rand_states not in states_seen:

```

```

52         states_seen.append(rand_states)
53         break
54     else:
55         rand_states = [random.randint(0, 1) for _ in range(len(nodes))]
56
57     state_series.append([t] + rand_states)
58
59 state_series_df = pd.DataFrame(state_series, columns = ["time"] + nodes)
60 return state_series_df
61
62 def output_data_to_file(state_series_df, fname):
63     state_series_df.to_csv(fname, index=False)
64     print("Printed data to {}".format(fname))
65
66 def parse_args():
67     parser = argparse.ArgumentParser()
68
69     parser.add_argument("--numNodes", type=int, help="Number of nodes")
70     parser.add_argument("--time", type=int, help="Length of time-series")
71     parser.add_argument("--outputFile", help="Output file")
72
73     args = parser.parse_args()
74
75     return args
76
77 # Example: python generate_data.py --numNodes 100 --time 300 --outputFile
78 #           "test.csv"
79 if __name__ == "__main__":
80     args = parse_args()
81     data = generate_data_time_series(args.numNodes, args.time)
82     print(data)
83     output_data_to_file(data, args.outputFile)

```