

ECS 158 Final Project: An Attempt At Parallelizing R's Phylobase::ShortestPath()

Raymond S. Chan, Alicia Luu, Bryan Ng 997544611, 999999999, 999999999
raschan@ucdavis.edu, ajuu@ucdavis.edu, bng@ucdavis.edu

Abstract

This report attempts parallelize CRAN's phylobase package's shortestPath function in RSnow, OpenMP and CUDA. The function takes in a phylogenetic tree, two nodes in the tree and produces the shortest path of nodes inbetween them. The RSnow implementation built on top of the existant code and parallelized the descendants function. The OpenMP implementation did a similar approach, except with **ALICIA WHAT DID YOU DO**. The CUDA implementation took the brute force approach to the shortestPath problem. Overall, **SENTENCE ABOUT TEST RESULTS**. **CLOSE WITH A MORE GENERAL STATEMENT BUILDING OFF RESULTS OF THE TESTS**.

1 Introduction and Motivation

Alongside the rise of "big data" in the recent years, bioinformatics has gained considerable momentum. But the a consistent issues remain: what do we do with all the data and how do we make sense of it at a reasonable rate? The R community has taken a stab at those issues. For this report, we are examining the R's phylobase package, which provides the base class and functions for phylogenetic or evolutionary structures and comparative data. (CITE) We will be focusing our efforts in making the "treewalk" utility functions, such as finding descendants/ancestors and shortest pathes, fast through three different parallel programming models (RSnow, OpenMP and CUDA).

2 Approaches

Only pseudocode or key chunks of code of each implementation are shown or described, see Appendix A for code details.

2.1 Original

The original R implementation calculated the shortest path between the two nodes of interest by first calculating their Most Recent Common Ancestor (MRCA). Then that MRCA's descendants are calculated and compared to the the two nodes of interest's ancestors. Any overlap is stored and that is the shortest path.

The C version of the descendants function, called Cdescendant(), works by first marking given node in a preordered list of edges. Then the direct descendants (or children) of the given node is marked. Cdescendant() then iterates through the other edges and marks each marked node's direct descendants (children).

The C version of the ancestors function works the same as Cdescendants, expect direct *ancestors* (or parents) are marked instead of direct *descendants* (children).

Pseudocode (source code in original phylobase package):

```

1 descendants(tree, given_node){
2     #let x be the edges of tree listed in PREORDER,
3     #with the older node occupying the first column
4     # C function call
5     isDescendant <- Cdescendants(x[,1], x[,2], given_node)
6     retval <- getNode(tree, isDescendant)
7 }
8
9 ancestors(tree, node1){
10     #let x be the edges of tree listed in POSTORDER,
11     #with the older node occupying the first column
12     # C function call
13     isAncestor <- Cancestors(x[,1], x[,2], given_node)
14     retval <- getNode(tree, isAncestor)
15 }
16
17 MRCA(tree, node1, node2 ... noden){
18     nodes <- unique(node1, node2, ..., noden)
19     ances <- lapply(nodes, ancestors, phy=phy, type="ALL")
20     retval <- getNode(phy, max(Reduce(intersect, ances)))
21 }
22 shortestPath(tree, node1, node2){
23     t1 <- getNode(tree, node1)
24     t2 <- getNode(tree, node2)
25
26     # most recent common ancestor
27     comAnc <- MRCA(tree, t1, t2)
28     desComAnc <- descendants(tree, comAnc)
29
30     # path: common ancestor to t1
31     ancT1 <- ancestors(x, t1)
32     path1 <- intersect(desComAnc, ancT1)
33
34     # path: common ancestor to t2
35     ancT2 <- ancestors(x, t2)
36     path2 <- intersect(desComAnc, ancT2)
37
38     # union of the path above paths
39     retval <- union(path1, path2)
40 }

```

2.2 RSnow

The RSnow implemenation builds on top of original R version by paralleizing the descendants function. In order to make independent subproblems, for a given node, every other node keeps marching upwards to its ancestors until they either reached the root or encountered the given node. We were

unable to parallelize the ancestor function. The original R version seems to have taken the most efficient serial approach for calculating a given node’s ancestors.

Pseudocode: (see code in Appendix A.1)

```

1 descendants(tree , given_node){
2     #let x be the list of all nodes except the given_node
3     for i from 1 to height of tree
4         if x == given_node -> append node to retval
5         update x with its corresponding ancestor
6     return retval
7 }

```

2.3 OpenMP

ALICIA WRITE STUFF HERE

See code in Appendix A.2.

2.4 CUDA

Our CUDA implementation of the shortestPath function utilizes the GPU to find all ancestors of a given pair of nodes and then construct the shortest path between them. Our implementation assumes CSIF’s pc43’s resources, which are 1024 threads per block and 1 GB of global memory. We assume the given data can fit in our GPUs global memory. This assumption may limit the test we will be able to perform. Our solution utilizes the fact that the shortest path between two nodes in a tree must converge at the lowest common ancestor of both nodes. In cases, where one node is an ancestor of another, then the shortest path is then found by traversing the parents of the child node. We parallelized our code by finding both sets of ancestors of the given nodes at the same time. Since neither node needs to know about the other to find its own ancestors, this problem can be done independently of each other. Both sets of ancestors are then traversed to find the shortest path. We were unable to parallelize this part of the solution since each list of ancestors must be checked to find overlapping elements.

See code in Appendix A.3.

3 Experiment Results

These are the results of running the above scripts with a simplified internet (i.e. $n = 6$).

```

>> i = [ 2 6 3 4 4 5 6 1 1];
>> j = [ 1 1 2 2 3 3 3 4 6];
>> n = 6;
>> G = sparse(i,j,1,n,n);
>> Finaltimetest
pagerank1          pagerank2          pagerank3A          pagerank3B          pagerankpow

```

0.0002	0.0001	0.0002	0.0003	0.0004
--------	--------	--------	--------	--------

These are the results of running the above scripts with the Harvard500 dataset (i.e. $n = 500$).

```
>> load Harvard500
>> Finaltimetest
```

pagerank1	pagerank2	pagerank3A	pagerank3B	pagerankpow
0.0024	0.0329	0.0226	0.0011	0.0255

These results are consistent with the dicussion above.

4 Discussion

4.1 RSnow

TALK ABOUT BIG O's.

4.2 OpenMP

TALK ABOUT BIG O's.

4.3 CUDA

THIS IS JUST A GUESS.

Compared to the serial version, the cuda implementation performed slower in most test cases. While the cuda version can compute both given nodes ancestors at the same time, it must also load the entire tree into the GPUs memory.

5 Conclusion

The PageRank algorithm was the starting point of Google's rise to fame. It was able to numerically quantify the "quality" on links/web pages of the internet. Pagerank is a Markov chain for which we solve for the dominant eigenvector of its transition probability matrix. There are two main methods of solving such a system of linear equations. The Power method is shown here to be the best method because of its efficiency in run time and memory usage. The run times are decent, as Hopcraft stated, it varies logarithmically with the size of input (n web paes). For the current day, the sheer amount of data that needs to processed is daunting. Further studies on PageRank could be done in further optimizing its space usage.

6 Acknowledgements

We would like to thank Professor Norman Matloff for his guidance and knowledge presented during lectures. This work is the result of a final project for ECS 158 Winter Quarter 2015. His open-source textbook, blog and various tutorial were an essential part of our learning. We would like to

thank the teaching assistance Shengren Li for offering invaluable advice and feedback on our codes (especially our CUDA) throughout the quarter.

7 Appendix

A Codes

INSERT ALL CODES HERE ALONG WITH A PARAGRAPH EXPLAINING IT

A.1 RSnow Code

```
1 SNOW <- function(x, size, root, type=c("descendants")){
2   ans <- rep(0, size)
3   mystart <- (myid-1)*length(x)+1
4   myend <- myid*length(x)
5
6   type <- match.arg(type)
7   if (type == "descendants"){
8     v1 <- descendant
9     v2 <- ancestor
10    #initialization
11    temp <- v1[mystart:myend]
12
13    #second and beyond iteration
14    for (j in 1:size){
15      if (node %in% temp){
16        setthese <- which(temp == node) + mystart-1
17        ans[setthese] <- 1
18      }
19      blah <- rep(-1, length(temp))
20      for (i in (1:length(temp))){
21        matched_pos <- which(v1 == temp[i])
22        if (length(matched_pos) != 0){
23          blah[which(temp == temp[i])] <- matched_pos
24        }
25        else{#matched_pos == 0
26          ## R is 1 INDEXED!
27          if (type == "descendants"){
28            blah[i] <- 1
29          }
30        }
31      }#for i
32      #"go to your parents set"
33      difference <- length(temp) - length(v2[blah])
34      temp <- v2[blah]
35      if (difference > 0){
36        temp <- c(rep(0, difference), temp)
```

```

37         }
38         if (node %in% temp){
39             setthese <- which(temp == node) + mystart-1
40             ans[setthese] <- 1
41         }
42     }#j loop
43 }#new endif for type==descendants
44 return(ans)
45 }# end SNOW
46
47 setmyid <- function(i){
48     myid <- i
49 }
50
51 ## get descendants with RSnow
52 RSnowdescendants <- function (phy, node, type=c("tips","children","all"),cls) {
53     type <- match.arg(type)
54
55     ## look up nodes, warning about and excluding invalid nodes
56     oNode <- node
57     node <- getNode(phy, node, missing="warn")
58     isValid <- !is.na(node)
59     node <- as.integer(node[isValid])
60
61     if (type == "children") {
62         res <- lapply(node, function(x) children(phy, x))
63         ## if just a single node, return as a single vector
64         if (length(res)==1) res <- res[[1]]
65     } else {
66         ## edge matrix must be in preorder for the C function!
67         #if (phy@order=="preorder") {
68             edge <- phy@edge
69         #} else {
70             # edge <- reorder(phy, order="postorder")@edge
71         #}
72         ## extract edge columns
73         ancestor <- as.integer(edge[, 1])
74         descendant <- as.integer(edge[, 2])
75
76         ## return indicator matrix of ALL descendants (including self)
77         #isDes <- .Call("descendants", node, ancestor, descendant)
78         clusterExport(cls, c("node", "ancestor", "descendant", "setmyid", "SNOW"), env=
79             dextrgs <- splitIndices(length(ancestor), length(cls))
80             rootdex <- which(phy@edge[,1] == 0)
81             clusterApply(cls, 1:length(cls), setmyid)
82             newisDes <- clusterApply(cls, dextrgs, SNOW, length(ancestor), rootdex,
"descendants")
83             isDes <- (matrix(Reduce( '+', newisDes ), nrow=length(ancestor), ncol=1))

```

```

84     storage.mode(isDes) <- "logical"
85     ## for internal nodes only, drop self (not sure why this rule?)
86     int.node <- intersect(node, nodeId(phy, "internal"))
87     isDes[cbind(match(int.node, descendant),
88                 match(int.node, node))] <- FALSE
89
90     ## if only tips desired, drop internal nodes
91     if (type=="tips") {
92         isDes[descendant %in% nodeId(phy, "internal"),] <- FALSE
93     }
94     ## res <- lapply(seq_along(node), function(n) getNode(phy,
95     ##             descendant[isDes[,n]]))
96     res <- getNode(phy, descendant[isDes[, seq_along(node)]]])
97 }
98 ## names(res) <- as.character(oNode[isValid])
99
100     res
101 }
102
103 #####
104 # shortestPath
105 #####
106
107 RSnowshortestPath <- function(phy, node1, node2, cls){
108
109     ## conversion from phylo, phylo4 and phylo4d
110     if (class(phy) == "phylo4d") {
111         x <- extractTree(phy)
112     }
113     else if (class(phy) != "phylo4"){
114         x <- as(phy, "phylo4")
115     }
116     ## some checks
117     t1 <- getNode(x, node1)
118     t2 <- getNode(x, node2)
119     if(any(is.na(c(t1,t2)))) stop("wrong_node_specified")
120     if(t1==t2) return(NULL)
121
122     ## main computations
123     comAnc <- MRCA(x, t1, t2) # common ancestor
124     desComAnc <- RSnowdescendants(x, comAnc, type="all", cls)
125     ancT1 <- ancestors(x, t1, type="all")
126     path1 <- intersect(desComAnc, ancT1) # path: common anc -> t1
127
128     ancT2 <- ancestors(x, t2, type="all")
129     path2 <- intersect(desComAnc, ancT2) # path: common anc -> t2
130
131     res <- union(path1, path2) # union of the path

```

```

132     ## add the common ancestor if it differs from t1 or t2
133     if(!comAnc %in% c(t1,t2)){
134         res <- c(comAnc,res)
135     }
136
137     res <- getNode(x, res)
138
139     return(res)
140 } # end shortestPath

```

A.2 OpenMP Code

1 ALICIA CODE GOES HERE

A.3 CUDA Code

1 BRYAN CODE GOES HERE

B Who Did What

Alicia wrote the OpenMP implementation. Bryan wrote the CUDA implemenation. Raymond wrote the RSnow implementation. We worked on running tests and writing the report in L^AT_EX.

References

- [1] C. Moler, *Numerical Computing with MATLAB Revised Reprint* 2004.
- [2] L. Page, S. Brin, R. Motwani, and T. Wingograd, *The PageRank Citation Ranking: Bringing Order to the Web*, available at <http://ilpubs.stanford.edu:8090/422/1/1999-66.pdf>, 1998
- [3] J. Hopcraft and R. Kannan, *Foundations of Data Science*, available at <http://www.cs.cornell.edu/jeh/NOSOLUTIONS90413.pdf>, 2011