

Project Report: Hascade

Yiming Fang
Email: yf2484@columbia.edu

December 22, 2022

1 Introduction

Influence Maximization (IM) is an exciting and well-researched topic, and it has practical applications to commercial marketing and social network management. In any given social network, it often is the case that some nodes are more influential than others. Sometimes, identifying a set of influential nodes can greatly help decision-makers make marketing decisions. However, since the problem is proven to be NP-Hard, most existing algorithms use greedy heuristics that run sequentially. Although the approximation algorithms have theoretical error bounds, the computational cost of IM solvers is usually very high due to their sequential nature.

This project presents a parallel implementation of the IM solver. In particular, we focus on the IM problem in the context of the Independent Cascade model. The following sections are going to be divided as follows. Section 2 formulates the problem. Section 3 describes the greedy algorithm we used. Section 4 discusses the choices we made for parallelization. Section 5 presents the empirical results of the performance.

2 Problem Formulation

Given a graph $G = (V, E)$, the task of the IM problem identifying a "most influential seed set" $S \subseteq V$ of size k , such that

$$S = \arg \max_{|S|=k} E[f(S)]$$

where S is called the seeds, and $f(S)$ is the total expected cascade size resulting from S .

In the Independent Cascade model, each edge $e_{(u,v)}$ has an influence probability $p_{u,v}$, such that u has a one-shot opportunity of influencing v with $p_{u,v}$ probability.

The influence process runs as follows. In the beginning, only the seeds are activated. At each timestep, the activated vertices have the opportunity to influence their neighbors, and if they succeed in influencing a vertex, the new vertex will join the set of activated vertices in the next iteration. The process terminates when there all vertices are either activated or have been tried to be activated by some neighbor.

3 Implementation

The baseline implementation of IM is a sequential, greedy algorithm [1] [2]. The algorithm adds vertices to the "most influential set" one-by-one, by always choosing the vertex that brings the highest expected increase of influence to the current set. The algorithm terminates when all k slots are populated. We present this process in Algorithm 1.

Algorithm 1 Greedy

```
1: Input: Graph  $G = (V, E)$ .
2: Output: Most Influential Set  $S$ .
3:
4:  $S_0 \leftarrow \{\}$ 
5: for  $i = 1, \dots, k$ , do
6:    $u \leftarrow \arg \max_u f(S_{i-1} \cup \{u\})$ 
7:   Activate  $u$ 
8:    $S_i \leftarrow S_{i-1} \cup \{u\}$ 
9: end
10:
11: return  $S_k$ 
12: =0
```

The most important and time-consuming part of this algorithm is line 6, which computes an approximation of the expected influence of a set of activated vertices. This approximation can be done through a simple Monte Carlo-style search that runs sequentially.

Algorithm 2 Monte-Carlo

```
1: Input: Graph  $G = (V, E)$  Vertex set  $S_i = S_{i-1} \cup \{u\}$ , number of trails  $N$ .
2: Output: The expected influence of  $S_i$ 
3:
4: count  $\leftarrow 0$ 
5: for  $j = 1, \dots, N$ , do
6:   Simulate independent cascade on  $G, S_i$ 
7:   count += number of influenced vertex
8: end
9:
10: return count /  $N$ .
11: =0
```

One obvious thing to notice is that we are running the independent function calls over the exact same input vertex set for a large number of iterations, and this can be easily made parallel, as discussed in the next section.

The simulation of the independent cascade model is done by finding the all neighbors of every vertex in the input set, and trying to activate them by generating a random number and comparing it with the influence probability $p_{u,v}$. Activated vertices are added to the input set of the recursive call to the next simulation, and vertices that have been attempted but not successfully activated are removed from the graph given to the recursive call.

Algorithm 3 Independent-Cascade

```
1: Input: Graph  $G = (V, E)$  Vertex set  $S_i$ .
2: Output: The simulated influence of  $S_i$ 
3:
4: neighbors  $\leftarrow \cup_{v \in S} G.lookup(v)$ 
5: neighbors = neighbors  $\setminus S$ 
6: activated  $\leftarrow \text{tryActivate}(\text{neighbors})$ 
7: activated = activated  $\setminus S$ 
8: failed = neighbors  $\setminus$  activated
9:  $G.keys \leftarrow G.keys \setminus$  failed
10:
11: return length(activated) + Independent-Cascade( $G, S \cup$  activated)
12: =0
```

4 Parallelization

As we observed in the last section, the most appropriate place to introduce parallelism is in Algorithm 3, because the structure of the sequential algorithm can be modified very slightly to bring parallelism to the overall algorithm. Since we used lazy data structures in our code, we want to ensure that our expressions are all fully evaluated to the normal form. We can accomplish this by using the "rdeepseq" strategy. The modified parallel Monte Carlo algorithm is presented as follows.

Algorithm 4 Monte-Carlo Par

```
1: Input: Graph  $G = (V, E)$  Vertex set  $S_i = S_{i-1} \cup \{u\}$ , number of trails  $N$ .
2: Output: The expected influence of  $S_i$ 
3:
4: count  $\leftarrow 0$ 
5: chunks  $\leftarrow \text{split numCores } N$ 
6: for  $chunk \in \text{chunks}$ , do
7:   for  $j = 1, \dots, N/\text{numCores}$ , do
8:     Simulate independent cascade on  $G, S_i$ 
9:     count += number of influenced vertex
10:  end
11: end 'using' parList rdeepseq
12: return count /  $N$ .
13: =0
```

A further parallelization trick that we used to make the algorithm more efficient is static chunking. While it is in general more beneficial to use dynamic partitioning to initialize sparks, we observe that in the situation of Monte Carlo, static partitioning suffices, and even outperforms dynamic partitioning.

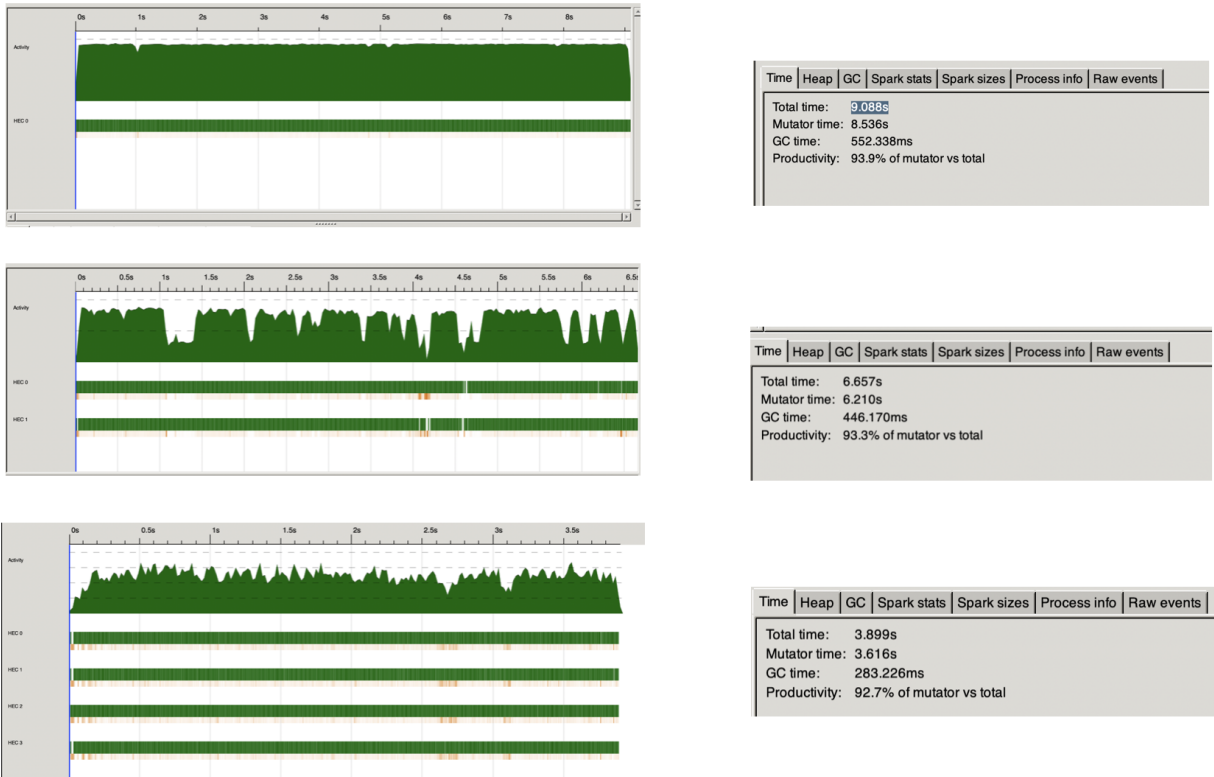
For each function call to Monte-Carlo, we are going to split the number of trials into numCores chunks, and start a spark for each chunk.

The rationale for doing so is the following. In most cases, whenever we do Monte Carlo simulation, we would like to ensure that the approximation we get is both stable and accurate. To this end,

the number of trials for Monte Carlo is going to be always a large number. Therefore, when we divide the workload into chunks, although each Independent-Cascade simulation can differ a lot in terms of workload and compute time, the chunk of simulation should all have similar workloads. This property ensures that the sparks started for the same Monte Carlo function call should finish roughly at the same time, minimizing the idle time of cores waiting for other cores to finish.

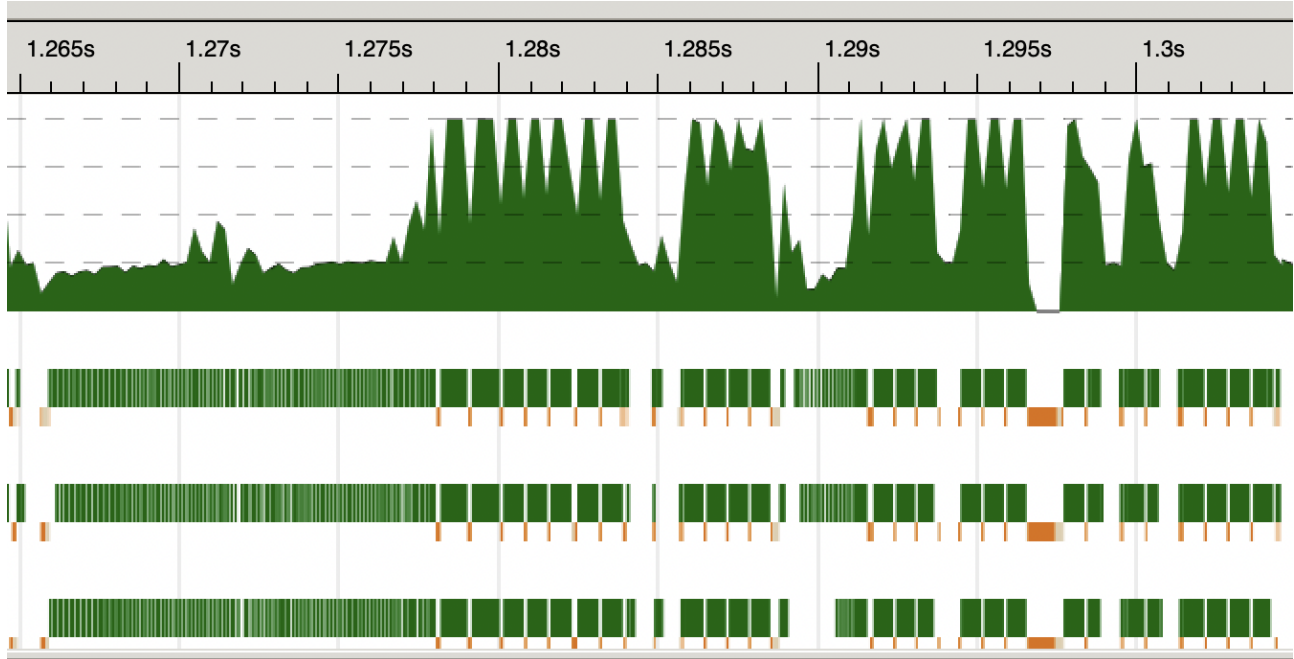
5 Performance

For testing, we used a dataset consisting of 1000 vertices, and set the hyper-parameters of Monte Carlo trials to be 1000. I used a 2018 MacBook Pro with 4 Intel cores and 8GB RAM, and I tested for sequential, 2 cores, and 4 cores respectively. The results is summarized in the following screenshot from ThreadScope.



As one can see from the picture, the multi-core tests achieved reasonable speedup ratios over the sequential version of the algorithm. For 2 cores, the speedup is **1.47x**, and for 4 cores the speedup is **2.33x**.

From the trace of the sparks, one can also see that the tasks are well-balanced, although there are periods of time where the balance is not so great and the performance becomes near sequential. To see exactly what happened during these sections, we can zoom into look at a more microscopic image:



In this test run with 4 cores, one can observe that the beginning section is most likely doing some sequential operations, such as chunking, combining, and folding. During this time, the cores are not fully occupied, and the stack trace tells us that the threads are blocked by some other threads, waiting for some relevant execution to finish. At around 1.275 seconds, the cores begin to do more meaningful work, interrupted periodically by garbage collections.

Given the nature of this algorithm, the degree of parallelization can be affected by a lot of factors, and can also differ across input. For example, a dense graph would cause the algorithm to spend more time doing Independent-Cascade simulation, resulting in a better ratio of parallelization. Moreover, the influence probability is also positively related to the ratio of speedup. To make a fair comparison, we have used a real-world graph from the SNAP datasets, representing the wikipedia community [3]

6 Conde Listing

```

1 {-# OPTIONS_GHC -Wno-missing-export-lists #-}
2 module Main where
3
4 import           BasicTypes                ( UnweightedGraph )
5 import qualified Data.Map.Strict          as Map
6 import qualified Data.Set                 as Set
7 import           Solver                    ( greedySolver )
8 import           System.Environment        ( getArgs
9                                           , getProgName
10                                          )
11 import           System.Exit              ( die )
12
13
14 main :: IO ()
15 main = do

```

```

16  args <- getArgs
17  case args of
18    [k, filename] -> do
19      contents <- readFile filename
20      let inputGraph = constructGraph contents
21          print $ greedySolver inputGraph Set.empty (read k) 0.1 100
22    _ -> do
23      pn <- getProgName
24      die $ "Usage: " ++ pn ++ "<num_cores> <filename>"
25
26
27  constructGraph :: String -> UnweightedGraph
28  constructGraph = Map.fromList . map extractLine . lines
29
30
31  extractLine :: String -> (Int, [Int])
32  extractLine str = (node, neighbors)
33  where
34    (node, neighbors) = case words str of
35      (this : others) -> (read this, map read others)
36      []                -> (-1, [])

```

Listing 1: Main.hs

```

1  {-# OPTIONS_GHC -Wno-unrecognised-pragmas #-}
2  {-# HLINT ignore "Use newtype instead of data" #-}
3  module BasicTypes
4    ( Vertex
5    , Weight
6    , UnweightedGraph
7    , WeightedGraph
8    ) where
9
10
11  import qualified Data.Map.Strict          as Map
12
13  type Vertex = Int
14  type Weight = Float
15
16  type UnweightedGraph = Map.Map Vertex [Vertex]
17  type WeightedGraph = Map.Map Vertex [(Vertex, Weight)]

```

Listing 2: BasicTypes.hs

```

1  {-# LANGUAGE BlockArguments #-}
2  module Solver
3    ( greedySolver
4    ) where
5
6  import           BasicTypes          ( UnweightedGraph
7                                       , Vertex
8                                       )
9  import           Control.Monad        ( replicateM )
10 import           Control.Parallel.Strategies ( parList
11                                              , rdeepseq
12                                              -- , rpar
13                                              -- , rseq
14                                              , using
15                                              )

```

```

16 import qualified Data.Map.Strict           as Map
17 import           Data.Maybe               ( fromMaybe )
18 import qualified Data.Set                 as Set
19 import           Data.Set                 ( Set
20                                           , (\\)
21                                           )
22 import           System.IO.Unsafe        ( unsafePerformIO )
23 import           System.Random            ( randomIO )
24
25
26
27 greedySolver
28   :: UnweightedGraph -> Set Vertex -> Int -> Float -> Int -> Set Vertex
29 greedySolver graph vSet k thresh mcTrials
30   | k == 0
31   = vSet
32   | otherwise
33   = let runMC :: Vertex -> (Float, Vertex)
34       runMC = monteCarlo graph vSet buffer
35
36       runChunk :: [Vertex] -> [(Float, Vertex)]
37       runChunk vs = map runMC vs
38
39       findMaxV :: [(Float, Vertex)] -> (Float, Vertex) -> Vertex
40       findMaxV [] acc = snd acc
41       findMaxV (x : xs) acc | fst x > fst acc = findMaxV xs x
42                               | otherwise      = findMaxV xs acc
43
44       buffer      = replicate mcTrials thresh
45       candidateVs = Map.keys graph
46
47       -- candidateChunks = split 10 candidateVs
48       -- scores          = map runChunk candidateChunks 'using' parList rdeepseq
49       -- vMax             = findMaxV (concat scores) (0, -1)
50
51       scores      = map runMC candidateVs -- 'using' parList rdeepseq
52       vMax        = findMaxV scores (0, -1)
53
54       vSet'       = Set.insert vMax vSet
55   in greedySolver graph vSet' (k - 1) thresh mcTrials
56
57
58 split :: Int -> [a] -> [[a]]
59 split numChunks xs = chunk (length xs 'quot' numChunks) xs
60
61
62 chunk :: Int -> [a] -> [[a]]
63 chunk _ [] = []
64 chunk n xs = let (as, bs) = splitAt n xs in as : chunk n bs
65
66
67 monteCarloV1
68   :: UnweightedGraph -> Set Vertex -> [Float] -> Int -> (Float, Vertex)
69 monteCarloV1 graph vSet ps vNew = (mean, vNew)
70   where
71     vs = Set.insert vNew vSet
72     lens = map (independentCascade graph vs 0) ps 'using' parList rdeepseq
73     mean = sum lens / realToFrac (length lens)
74

```

```

75
76 monteCarlo
77   :: UnweightedGraph -> Set Vertex -> [Float] -> Vertex -> (Float, Vertex)
78 monteCarlo graph vSet ps vNew = (mean, vNew)
79   where
80     meansWithSizes = map mc pss 'using' parList rdeepseq
81
82     vs              = Set.insert vNew vSet
83     pss             = split 4 ps
84
85     mc              = monteCarloChunk graph vs
86
87     totalSum        = foldr ((+) . multSize) 0 meansWithSizes
88     totalSize       = foldr ((+) . snd) 0 meansWithSizes
89
90     multSize (a, b) = a * realToFrac b
91     mean = totalSum / realToFrac totalSize
92
93
94 monteCarloChunk :: UnweightedGraph -> Set Vertex -> [Float] -> (Float, Vertex)
95 monteCarloChunk graph vSet ps = (mean, length ps)
96   where
97     lens = map (independentCascade graph vSet 0) ps -- 'using' parList rseq
98     mean = sum lens / realToFrac (length lens)
99
100
101 independentCascade :: UnweightedGraph -> Set Vertex -> Int -> Float -> Float
102 independentCascade graph vSet depth thresh = if null activatedSet'
103   then 0
104   else
105     let nextCascade =
106         independentCascade graph' activeSet (depth + 1) thresh
107         thisCascade = realToFrac $ length activeSet
108     in thisCascade + nextCascade
109
110   where
111     graph'          = graph Map.\ \ setMap
112     setMap          = Map.fromSet ('Map.lookup' graph) neighborSet'
113
114     neighborSet     = getNeighborSet graph vSet
115     neighborSet'    = neighborSet \ \ vSet
116
117     activatedSet    = tryActivate neighborSet' thresh
118     activatedSet'   = activatedSet \ \ vSet
119
120     activeSet       = Set.union vSet activatedSet'
121
122
123 randSeq :: Int -> [Float]
124 randSeq k = unsafePerformIO (replicateM k (randomIO :: IO Float))
125
126
127 tryActivate :: Set Vertex -> Float -> Set Vertex
128 tryActivate vs thresh = Set.fromList newActiveVs
129   where
130     strengths      = randSeq 1
131     threshs        = replicate 1 thresh
132     l               = length vs
133

```



```

134     diff      = zipWith (-) threshs strengths
135     threshVs  = zip diff (Set.toList vs)
136     newActive = filter ((< 0) . fst) threshVs
137
138     newActiveVs = map snd newActive
139
140
141 getNeighborSet :: UnweightedGraph -> Set Vertex -> Set Vertex
142 getNeighborSet graph vSet = Set.fromList newSets
143     where
144     findChildren :: Int -> [Int]
145     findChildren v = Data.Maybe.fromMaybe [] (Map.lookup v graph)
146
147     newSets = concatMap findChildren $ Set.toList vSet

```

Listing 3: Solver.hs

References

- [1] <https://snap-stanford.github.io/cs224w-notes/network-methods/influence-maximization>
- [2] https://hautahi.com/im_greedyself
- [3] <https://snap.stanford.edu/data/wiki-Talk.html>