

Parallelizing NFA to DFA Conversions

Alexis Gadonniex, Nikhil Mehta

December 22, 2022

1 Problem Statement

We seek to parallelize the subset construction algorithm that is applied to a nondeterministic finite automaton (NFA) to convert it into a deterministic finite automaton (DFA). This is particularly useful in the context of string processing and compilers where one wants to check whether a given string matches a regular expression or not. Since regexes can be represented as NFAs, this is the same as asking the question of whether a given string will be accepted by an NFA. At the surface this might seem like a tricky problem since NFA's are, as their name suggests, non-deterministic. However, NFA's can always be converted into DFAs through the subset construction algorithm. This algorithm, described below, is one that has to consider a variety of possible intermediary states in the NFA at the same time. Thus, the algorithm is ripe for parallelization. To generate NFAs in the first place we will use three techniques. First, we will implement non-parallelized functionality to convert a regex to an NFA using Thompson's algorithm. Second, we will create another type of NFA directly from a dictionary of all accepted strings in a language. Third, we will directly generate NFA's with random edges and transitions.

Thus our project requires several parts. First, we need to create a data structure to represent automata whether they are deterministic or not. Then, we need write single-threaded algorithms to generate NFAs. At the core of our project, we need to write a single-threaded version of the subset construction algorithm and then work to parallelize it. Last, we opted to write an algorithm to check if strings are accepted by a given automaton.

2 Algorithm

2.1 Subset construction

The subset construction algorithm is used to convert an NFA into a DFA. The subset construction algorithm works by considering sets of states. That is to say, for a given NFA consider the set of all states reached from the starting state. That constitutes the starting state for your DFA. Then consider all the states (if any) accessible from the initial set of states that are reached when you pass in a given character. That constitutes the next state of the DFA. Then, consider the set of all states reached when you pass in a different character. The algorithm continues to reapply this step to find all possible sets of states in the NFA for all possible sequences of characters. Each set of NFA states then corresponds to a unique state in the DFA.

Below we present pseudocode of the algorithm:

```
Init: Add epsilon-closure of the start states to DStates
```

```
While there is a non-visited state T in DStates do
  mark T as visited
  for each symbol c do
    compute S the set of states that can be reached from T through c
    compute SE the epsilon-closure of S
    if SE is not in DStates do
      add SE to DStates as a non-visited state
    mark SE as final if it contains a final state of the NFA
```

```

end
Transition [T, c] = SE
end
end

```

The algorithm has a worst-case complexity of $O(2^n)$ since the number of subsets of an n -state NFA is 2^n . This worst-case complexity can be achieved with fairly simple languages and this will be useful to test our code against difficult instances.

2.2 Thompson's rules

While not a parallelized feature of our code, we use Thompson's algorithm to generate NFA's from regexes. Thompson's algorithm is a set of conversions that can be applied to each regex rule to generate an NFA. Thompson's algorithm can be applied recursively to a syntax tree of a regular expression. Below we describe the four rules for Thompson's algorithm.

To represent two sub-regexes $N(s)$ and $N(t)$ that are concatenated with one another one can use the rule shown in Fig. 1. To represent two sub-regexes $N(s)$ and $N(t)$ that are joined via a union, sometimes referred to as "or" and typically denoted by $|$, one can use the rule shown in Fig. 2. To represent a sub regex $N(s)$ that is encased in a Kleene Enclosure, normally denoted by a $*$, one can use the rule shown in Fig. 3. Finally, to represent the transition of an accepted symbol a , one can apply Thompson's symbol rule shown in Fig. 4

We further discuss Thompson's algorithm and how we implemented it in section 4.

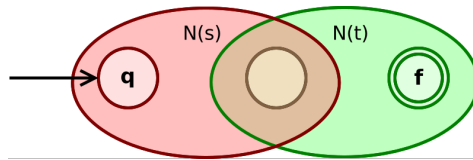


Figure 1: Thompson's concatenation rule

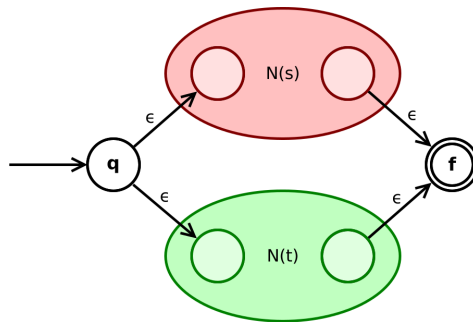


Figure 2: Thompson's union rule

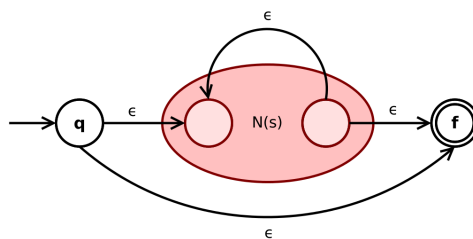


Figure 3: Thompson's symbol Enclosure rule

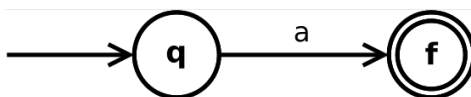


Figure 4: Thompson’s Kleene Enclosure rule

3 Implementation

3.1 Automaton Data Structure

To represent NFAs and DFAs in our project, we wrote an Automaton data type in Haskell in `Automaton.hs`. States in an Automaton are represented as unique integers. A transition is denoted as an Edge which consists of a tuple of a Label and the next State. A Label is either an Epsilon or a Label Int where the integer denotes the character in the alphabet that is labeling a given transition. An Automaton consists of an Adjacency List which is a Map with each State serving as the key and their subsequent list of Edges serving as the value. An Automaton also takes a list of Labels to represent the alphabet. Last, an Automaton takes a Set of final states and a Set of start states.

3.2 Subset Construction

From the very beginning, we tried to implement the subset construction algorithm in a way that could be easily parallelized in the future. Thus, we tried to isolate some map patterns.

We used a layer-by-layer approach, similar to a BFS, that we divided into two phases:

- An exploration phase, where we explore all the transitions leaving from the previous layer. We create the cartesian product between the alphabet and the DFA states of the previous layer, and for each pair, we compute the target DFA State (set of NFA States) by performing look-ups to the NFA’s transition table and ε -closure. We parallelized this phase by running these computations using `parMap` and the `rdeepseq` strategy.
- A Sequential phase where we will try to merge all the newly discovered states from the parallel phase with our partially created DFA. In order to keep track of the mapping between a set of NFA’s states (set of integers) and the corresponding DFA state (an integer), we use a Haskell Map with Sets of Ints as keys and Ints as values. A newly discovered DFA state can lead to several results: it is a new state and should be used in the next exploration phase; it is empty so we can ignore it; it already exists elsewhere in the DFA (a previous layer or the same one) and we just have to add the edge to our adjacency list.

We also tried to explore several layers at once during the exploration phase to create more sparks but it was not very efficient because the number of states in each layer is growing very quickly and most of them are actually redundant. It is more effective to sort the new states from the already existing one at each layer. We illustrate our algorithm working on a sample NFA in Fig. 5

A clear bottleneck in our implementation is the sequential phase which limits how much we can get from parallelism (Amdahl’s law). The most expensive operation during this phase is the look-up in the “NFA state sets to DFA state”, especially if the DFA States are made of a large number of NFA states. A possible improvement would be to use some kind of hash function for a set of integers to produce smaller keys for the Map. This hash could be computed during the parallel phase. Unfortunately, we didn’t have time to explore this possibility.

Reducing the time spent in this phase was one of our main concerns when considering the various options we had to generate NFAs.

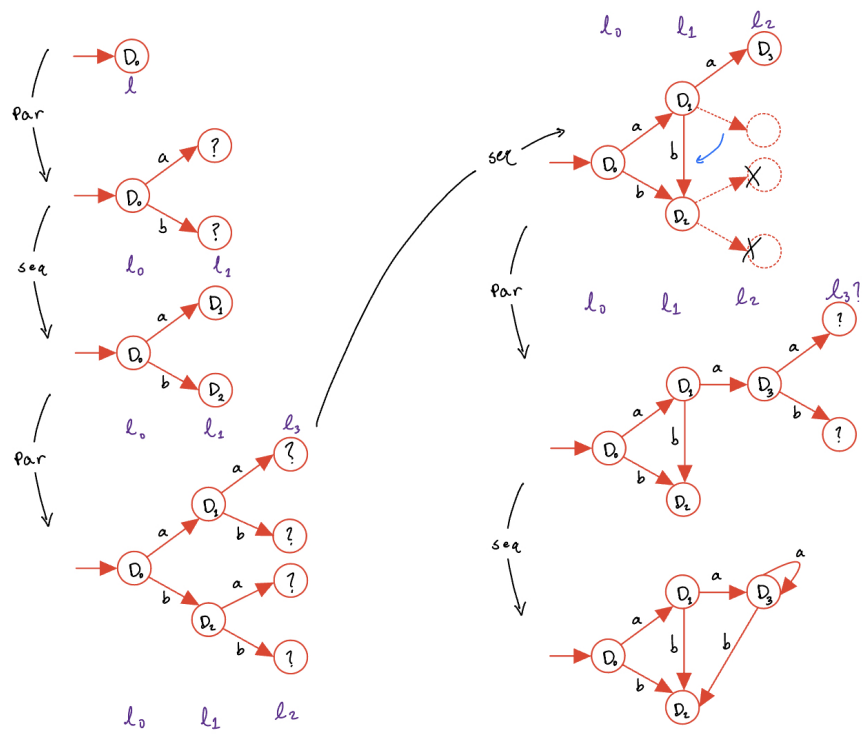


Figure 5: Subset construction implementation diagram

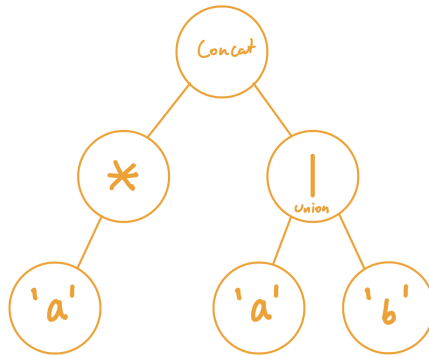


Figure 6: Syntax tree generated from "a*(a|b)"

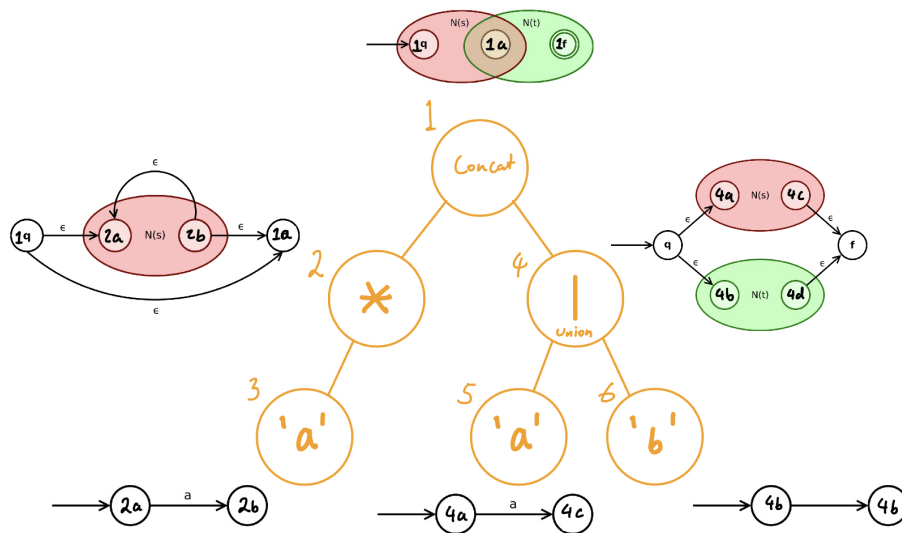


Figure 7: Traversal of syntax tree from Fig. 6 with corresponding Thompson's construction at each node.

4 NFA Generation

4.1 From regular expressions

We used three techniques for generating NFAs. Our first method involved using Thompson's algorithm to generate an NFA from a regex. To achieve this we wrote a parser `regexParser` in `Thompson.hs` using the Haskell `Parsec` library to take in a regex as a string and output a syntax tree. We used the standard notation of `*` for kleene enclosure and `|` for union. Our parser supports the use of parentheses to enforce a certain precedence. Characters of the accepted alphabet for the regex could be denoted as themselves. Concatenation is implicit when elements of the regex succeed each other. An example regex could be: "a*(a|b)". We show the resultant syntax tree for this regex in Fig. 6.

The resultant syntax tree is then fed into our implementation of Thompson's Algorithm in `thompsons` which is called by `makeThompsonNFA` in `Thompson.hs` which traverses the syntax tree using a depth first search and builds up the adjacency table for NFA recursively. To generate the alphabet for the NFA we used a separate DFS to search for all characters, a.k.a. leaf nodes on the syntax tree. We show an example DFS with the each node in the tree from Fig. 6 labeled according to the order it was visited in Fig. 7 We then show the resultant Thompson's construction for that node.

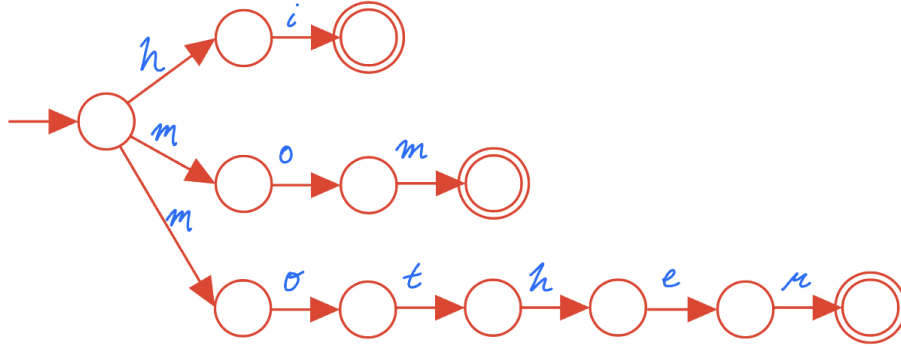


Figure 8: Sample NFA created from a dictionary of all accepted words in a language.

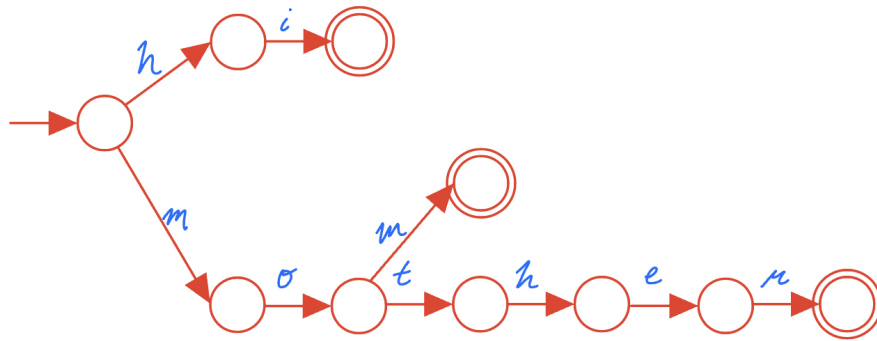


Figure 9: DFA converted from NFA above from dictionary of all accepted words in a language.

We quickly realized that, in order to create sufficiently large problems for our algorithm, we would need to create random regular expressions. One way of doing that could be to chose a number of terminal states and, in a recursive fashion, randomly split the remaining nodes between the two branches of an operator (concatenation or $|$). However, we were afraid that the structure of this type of NFA would be too linear for our algorithm to produce interesting results. We chose to spend time on other methods.

4.2 From a list of words

Our second method of generating NFAs was to build it directly from a dictionary of all the accepted strings in the language. We implemented the function `makeDictNFA` in `Dictionary.hs` which takes in a `[Char]` of all the strings from a dictionary. Then, each string is added onto the NFA from the starting state via an epsilon transition. For example, the following input to `["hello", "world", "mom"]` would generate the NFA shown in Fig. 8.

This is a pretty straightforward way of generating NFAs with a very regular structure. The results are discussed below.

4.3 From a given size and density

The last method we tried was to generate a fully random NFA, without any particular meaning, from a set of parameters. The parameters we chose to expose are:

- The number of states

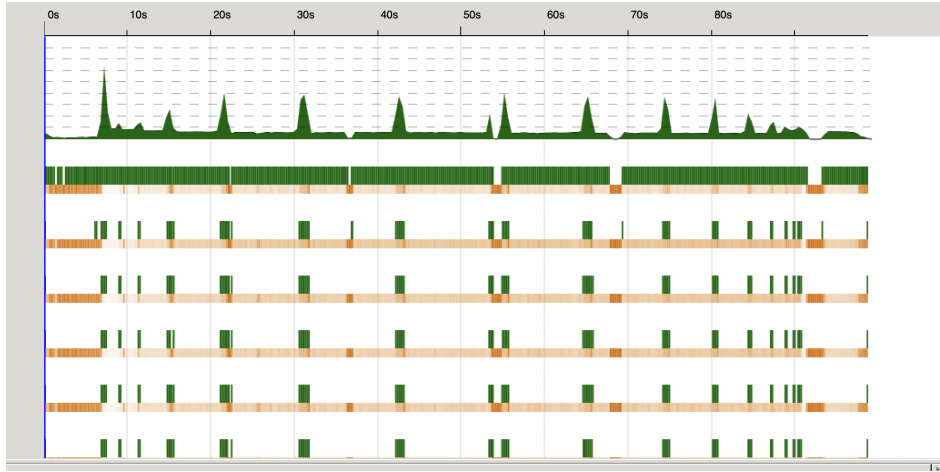


Figure 10: Most time spent in sequential phase when using an NFA created from a list of words

- The size of the alphabet
- The probability that, given two states s_1 and s_2 and an element of the alphabet c , there is an edge from s_1 to s_2 with a label c .
- (The number of initial states)
- (The number of final states)

The last two are not very important since they don't have any noticeable influence on the complexity of the algorithm.

Since our implementation is a bit naive, we have to use high probability values (it depends on the size of the alphabet) in order to keep our graph connected and avoid isolated states in the NFA.

5 Results

5.1 List of words

The NFAs built from lists of words looked like a perfect use-case for our algorithm. But it turned out that the result were a lot worse than we initially expected. You can see in the threadscope graph in Fig. 10 that the algorithm was running on a single core most of the time. The main reason why it spends all this time in the sequential phase is that each DFA state is made of a lot of NFA states. To actually see the effect of parallelism, we had to consider very long lists of words (200000 words). But the problem is that it leads to huge NFA States (for example: about 10000 words start with the letter a). In this case, the look-ups are extremely expensive and the memory usage is extremely high (maximum heap size 3Gb; total allocated 200Gb !!). We can also note that there are 20 million sparks created, half were converted and the other half overflowed. We tried using `IntSet` as keys but it didn't significantly improve the results.

5.2 Random NFA

When using the random NFAs we described above, we were able to tune more precisely the parameters to find NFAs that fit our algorithm.

The best results were obtained for NFAs with 600 states, an alphabet of size 20 and a probability of 0.4. We obtained a speed-up of 3 using 8 cores (on an M1 Macbook Pro with 8Gb of RAM). You can see on Fig. 11 that most of the time is spent in the parallel phase. The speed-up for different number of threads is plotted on Fig. 12. Note that we could get a slightly better speed-up by generating the random NFAs in advance and parsing them from a text file since it takes about 10 percent

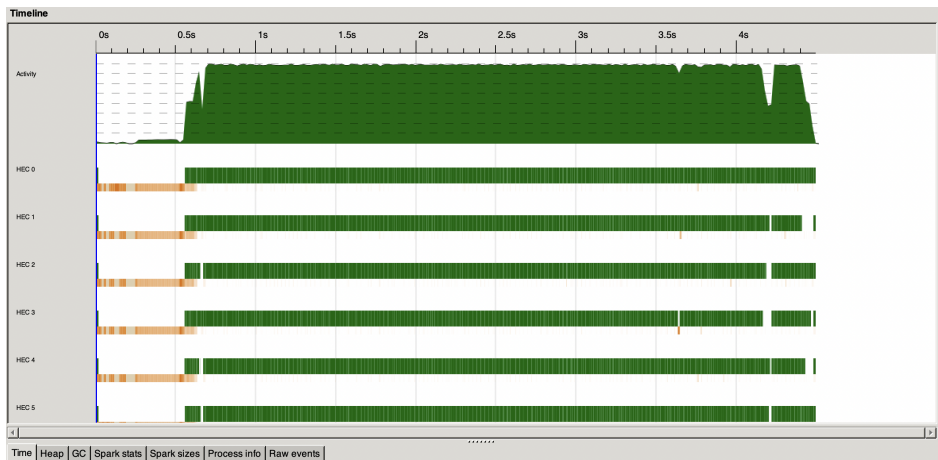


Figure 11: Threadscope on 8 threads for a random NFA with 600 states, alphabet of size 20 and probability of 0.4

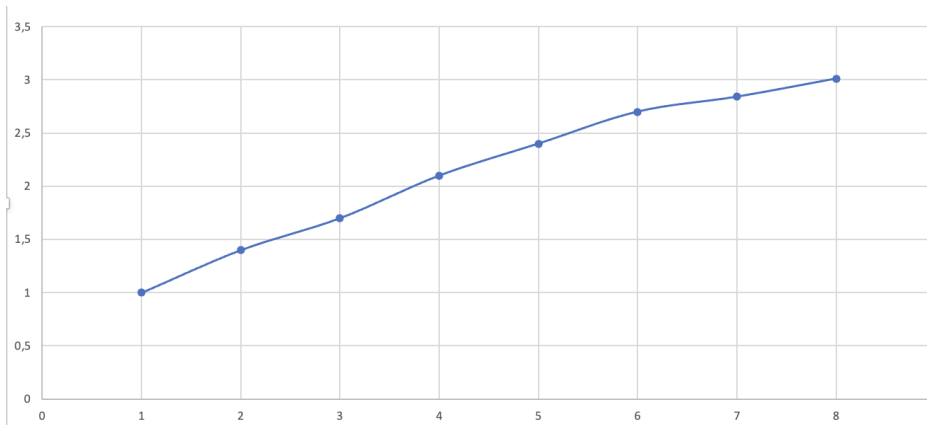


Figure 12: Speed-ups from 1 to 8 threads. Measured on a sample of 10 random NFAs with 600 states, alphabet of size 20 and probability of 0.4

of the total time and is fully sequential.

These are relatively small NFAs but highly interconnected. This leads to pretty small DFA states (in term of number of NFA states) but a lot of different combinations of those. The average size of a state is about 60 NFA states, a lot smaller than the NFAs from a list of words (several hundreds in average and up to several thousands for about the same computation time).

Changing the number of states of the NFA has a relatively small impact on the results. Below 500 the algorithm finishes too quickly for the parallelism to make a real difference. Above (we tested up to 2000 states), the computation time increases very quickly but the speed-up stays the same. It takes about 10 times longer when going from 500 to 1000 states.

We tried to modify the probability but it is a very sensitive parameter and shows the limitations of our modelization. If we decrease the probability too much the graph becomes poorly connected and we end up with a small number NFA states that are truly relevant. In the contrary, if we increase the probability too much, the graph becomes too connected and there are only a small number of DFA states to explore.

The alphabet size is also tricky since increasing it leads to an exponential number of additional edges.

We chose to stick with numbers between 20 and 30 since it is the range of the latin alphabet.

6 Usage

Our code is intended to run with GHC 8.10.7, HLS 1.8.0.0, Stack 2.9.1, and cabal 3.6.2.0.

To build our code run `$ stack build` in the base directory. To run our code using the following command: `$ stack exec subset-construction-exe -- +RTS -ls -s -lf -N8`

7 References

We cite the following sources and thank the respective owners:

1. https://en.wikipedia.org/wiki/Thompson%27s_construction for the Thompson's construction graphics
2. Professor Stephen Edwards PLT slides for sample Regexes used in figures and testing

8 Appendix

/src/Automaton.hs

```
1 module Automaton (Automaton(..), AdjacencyList, DfaStatesMap, State, Edge, Label(..),
2   exampleAutomaton, exampleAutomaton2, ioDumbAutomaton) where
3 import Control.DeepSeq (NFData (rnf))
4 import qualified Data.List as List
5 import qualified Data.Map as Map
6 import qualified Data.Set as Set
7
8 -- A node's "id"
9 type State = Int
10 -- A labeled edge pointing to a node
11 data Label = Epsilon | Label !Int deriving (Eq, Show, Ord)
12
13 instance NFData Label where
14   rnf (Label x) = rnf x
15   rnf Epsilon = ()
16
17 type Edge = (Label, State)
18 -- Successors map, initial states and final states
19 type AdjacencyList = Map.Map State [Edge]
20 data Automaton = Automaton !AdjacencyList ![Label] !(Set.Set State) !(Set.Set State)
21   deriving (Show)
22
23 type DfaStatesMap = (Map.Map (Set.Set State) State, Int)
24
25 -- Example taken from https://en.wikipedia.org/wiki/Powerset_construction (5 states
26   NFA generating a 16 states DFA through the algorithm)
27 alphabet :: [Label]
28 alphabet = [Label 0, Label 1]
29 initStates :: Set.Set State
30 initStates = Set.fromList [0]
31 finalStates :: Set.Set State
32 finalStates = Set.fromList [4]
33 successors :: AdjacencyList
34 successors = Map.fromList [(0, [(Label 0, 0), (Label 1, 0), (Label 1, 1)]), (1, [(
35   Label 0, 2), (Label 1, 2)]), (2, [(Label 0, 3), (Label 1, 3)]), (3, [(Label 0, 4),
36   (Label 1, 4)]), (4, [])]
```

```

37 finalStates2 = Set.fromList [2]
38 successors2 :: AdjacencyList
39 successors2 = Map.fromList [(0, [(Label 0, 1)]), (1, [(Label 1, 2)]), (2, [(Epsilon,
    0)])]
40 exampleAutomaton2 :: Automaton
41 exampleAutomaton2 = Automaton successors2 alphabet initState2 finalStates2
42
43 intAlphabet :: [Label]
44 intAlphabet = List.map Label [0..50]
45
46 dumbAutomaton :: Int -> Automaton
47 dumbAutomaton nStates = Automaton adj intAlphabet (Set.singleton 0) (Set.singleton $
    nStates - 1)
48     where adj = Map.insert 0 [(Label 0, 0), (Label 1, 0), (Label 1, 1)]
49           allButFirstSucc
50           allButFirstSucc = Map.insert (nStates-1) [] allButFirstAndLast
51           allButFirstAndLast = Map.fromList $ List.map (\n -> (n, [(l, r) | l
    <- intAlphabet, r <- [n+1,n, n-1]])) [1..nStates-2]
52
53 ioDumbAutomaton :: Int -> IO Automaton
54 ioDumbAutomaton n = return $ dumbAutomaton n

```

/src/Checks.hs

```

1 module Checks (checkAccept, checkAlphabet) where
2 import Automaton (Automaton (..), Label (..), State)
3 import Data.Char (ord)
4 import Data.List (find)
5 import qualified Data.Map as Map
6
7 -- checkAlphabet: Check if word agrees with alphabet for a given Automaton
8 checkAlphabet :: [Char] -> Automaton -> Bool
9 checkAlphabet (x:xs) dfa@(Automaton _ alph _) = n && checkAlphabet xs dfa
10     where n = Label (ord x) `elem` alph
11 checkAlphabet [] _ = True
12
13 -- checkAccept: check if word is in a language. Automaton MUST be a DFA
14 checkAccept :: [Char] -> Automaton -> State -> Bool
15 checkAccept (x:xs) dfa c = case findTransition x dfa c of
16     Just e -> checkAccept xs dfa n where (_,n) = e
17     Nothing -> False
18 checkAccept [] (Automaton _ _ _ end) c = c `elem` end
19
20 findTransition :: Char -> Automaton -> Int -> Maybe (Label, State)
21 findTransition x (Automaton lst _ _ _) c = case Map.lookup c lst of
22     Just es -> find (\y -> fst y == Label (ord x)) es
23     Nothing -> Nothing

```

/src/Dictionary.hs

```

1 module Dictionary (dictNfa) where
2
3 import Automaton (AdjacencyList, Automaton (..), Label (..), State)
4 import Data.Char (ord)
5 import qualified Data.Map as Map
6 import qualified Data.Set as Set
7 import WordList (buildList)
8
9 -- addWord: (non thompsons) helper method for dictNFA
10 addWord :: [Char] -> AdjacencyList -> Set.Set Label -> State -> Set.Set State -> State
11     -> (AdjacencyList, Set.Set Label, Set.Set State, State)
12 addWord (x:xs) adjList alph fromState finalStates toState = addWord xs newAdjList
13     newAlph toState finalStates newToState where
14     t = Label (ord x)
15     newToState = toState + 1
16     newAlph = Set.insert t alph
17     newAdjList = Map.insertWith (++) fromState [(t, toState)] adjList
18 addWord [] adjList alph lastState finalStates toState = (adjList, alph, newFinals,
19     toState) where
20     newFinals = Set.insert lastState finalStates
21
22 -- dictNFA: (non thompsons) helper method for buildDictNFA

```

```

20 dictNFA :: [[Char]] -> AdjacencyList -> Set.Set Label -> State -> Set.Set State ->
    State -> (AdjacencyList, Set.Set Label, Set.Set State)
21 dictNFA (x:xs) lst alph st fi l = dictNFA xs nlst nalph st nfi nl where
22   (nlst, nalph, nfi, nl) = addWord x lst alph 0 fi il
23   il = l+1
24 dictNFA [] lst alph _ fi _ = (lst, alph, fi)
25
26 -- makeDictNFA: Build NFA from dict of accepted strings in a language
27 makeDictNFA :: [[Char]] -> Automaton
28 makeDictNFA l@(_:_ ) = Automaton lst alph start fi where
29   alph = Epsilon : Set.toList salph
30   (lst, salph, fi) = dictNFA l (Map.insert 0 [] Map.empty) Set.empty 0 Set.empty 0
31   start = Set.singleton 0
32 makeDictNFA [] = error "empty dictionary"
33
34 dictNfa :: FilePath -> IO Automaton
35 dictNfa fp = do
36   wordList <- buildList fp
37   return $ makeDictNFA wordList

```

/src/RandomNfa.hs

```

1 module RandomNfa (ioRandomNfa) where
2 import Automaton (Label (...))
3 import qualified Automaton as A
4 import qualified Data.List as List
5 import qualified Data.Map as Map
6 import qualified Data.Set as Set
7 import qualified System.Random as Random
8
9
10
11 -- Inspired from https://hackage.haskell.org/package/random-1.2.1.1/docs/System-Random.html
12 rolls :: Int -> Int -> Int -> [Int]
13 rolls n maxInt seed = take n . List.unfoldr (Just . Random.uniformR (0, maxInt)) $
    Random.mkStdGen seed
14
15 -- Generate a random NFA with a given number of states, an alphabet size, a number of
    final states and a probability
16 -- (probability that there is an edge with a given label between 2 given states)
17 randomNFA :: Int -> Int -> Int -> Int -> A.Automaton
18 randomNFA numStates alphabetSize nbFinals proba =
19   A.Automaton transitions alphabet initStates finalStates
20   where
21
22     intForGen = numStates + alphabetSize + nbFinals + proba
23
24     -- list of states
25     states = [0..(numStates-1)]
26     -- final states
27     finalStates = Set.fromList [(numStates - nbFinals)..(numStates-1)]
28     -- alphabet
29     alphabet = [Label i | i <- [0..(alphabetSize - 1)]]
30
31     allTransitions = [
32       (state1, symbol, state2) |
33         state1 <- states,
34         state2 <- states,
35         symbol <- alphabet ]
36
37     keepTransitions = [ rdInt <= proba | rdInt <- rolls (length allTransitions)
    100 intForGen]
38
39     transitionsList = [ v | (v,keep) <- List.zip allTransitions keepTransitions,
    keep]
40
41     -- turn the list into a map
42     transitions :: A.AdjacencyList
43     transitions =
44       List.foldl'
45         (\m (fromS, label, toS) -> Map.insertWith (++) fromS [(label, toS)] m

```

```

)
46         Map.empty
47         transitionsList
48
49         -- Generate a random start state
50         initState = Set.singleton 0
51
52 ioRandomNfa :: Int -> Int -> Int -> Int -> IO A.Automaton
53 ioRandomNfa nbStates alphabetSize nbFinals probability =
54     return $ randomNFA nbStates alphabetSize nbFinals probability

```

/src/SubsetConstruction.hs

```

1 module SubsetConstruction (nfaToDfa) where
2 import qualified Automaton          as A
3 import           Control.Parallel.Strategies (parMap, rdeepseq)
4 import qualified Data.List          as List
5 import qualified Data.Map           as Map
6 import           Data.Maybe         (fromMaybe)
7 import qualified Data.Set           as Set
8
9 exploreLabelFromNFAState :: A.AdjacencyList -> A.Label -> A.State -> [A.State]
10 exploreLabelFromNFAState nfaAdjacency label state = [s | (l, s) <- edges, l == label]
11     where edges = fromMaybe [] (Map.lookup state nfaAdjacency)
12
13 exploreLabelFromDFAState :: A.AdjacencyList -> A.Label -> Set.Set A.State -> Set.Set A
14     .State
15 exploreLabelFromDFAState nfaAdjacency label =
16     Set.fromList
17     . List.concatMap (exploreLabelFromNFAState nfaAdjacency label)
18     . Set.toList
19
20 epsilonClosure :: A.AdjacencyList -> Set.Set A.State -> Set.Set A.State
21 epsilonClosure nfaAdjacency nfaStates | Set.size nfaStates == Set.size explored =
22     nfaStates
23     | otherwise = epsilonClosure nfaAdjacency
24     explored
25     where explored = Set.union nfaStates (exploreLabelFromDFAState
26     nfaAdjacency A.Epsilon nfaStates)
27
28 nextStates :: A.AdjacencyList -> [A.Label] -> [Set.Set A.State] -> [(A.Label, Set.Set
29     A.State, Set.Set A.State)]
30 nextStates nfaAdjacency alphabet dfaStates =
31     parMap rdeepseq
32     (\(l, s) -> (l, s, (epsilonClosure nfaAdjacency . exploreLabelFromDFAState
33     nfaAdjacency l) s))
34     [(l,s) |
35     l <- alphabet,
36     s <- dfaStates
37     ]
38
39 addDfaEdge :: Set.Set A.State -> (A.DfaStatesMap, A.AdjacencyList, Set.Set A.State, [
40     Set.Set A.State])
41     -> (A.Label, Set.Set A.State, Set.Set A.State)
42     -> (A.DfaStatesMap, A.AdjacencyList, Set.Set A.State, [Set.Set A.State
43     ])
44 addDfaEdge nfaFinals ((dfaSM, maxIdx), dfaA, dfaF, tV) (l, originS, destS) = case Map.
45     lookup destS dfaSM of
46     Nothing | Set.null destS -> ((dfaSM, maxIdx), dfaA, dfaF, tV)
47     Nothing -> ((Map.insert destS newState dfaSM, newState), Map.
48     insertWith (++) originState [(l, newState)] dfaA, newDfaF, destS : tV)
49     where newState = maxIdx + 1
50           isFinal = List.any ('Set.member' nfaFinals) (Set.toList
51     destS)
52           newDfaF = if isFinal then Set.insert newState dfaF else dfaF
53     Just s -> ((dfaSM, maxIdx), Map.insertWith (++) originState [(l,s)]
54     dfaA, dfaF, tV)
55     where originState = case Map.lookup originS dfaSM of
56     Just st -> st
57     Nothing -> error "Couldn't find origin state"

```

```

48 explore :: A.AdjacenclyList -> Set.Set A.State -> [A.Label] -> A.DfaStatesMap -> A.
    AdjacenclyList -> Set.Set A.State -> [Set.Set A.State] -> (A.AdjacenclyList, Set.Set
    A.State)
49 explore _ _ _ dfaAdjacency dfaFinals [] = (
    dfaAdjacency, dfaFinals)
50 explore nfaAdjacency nfaFinals alphabet dfaStatesMap dfaAdjacency dfaFinals toVisit =
51 explore nfaAdjacency nfaFinals alphabet newDfaSM newDfaA newDfaFinals newToVisit
52     where (newDfaSM, newDfaA, newDfaFinals, newToVisit) =
53         List.foldl'
54             (addDfaEdge nfaFinals)
55             (dfaStatesMap, dfaAdjacency, dfaFinals, [])
56             nStates
57     nStates = nextStates nfaAdjacency alphabet toVisit
58
59 nfaToDfa :: A.Automaton -> A.Automaton
60 nfaToDfa (A.Automaton nfaAdjacency alphabet inits nfaFinals) = A.Automaton
    newAdjacency alphabet dfaInits dfaFinals
61     where (newAdjacency, dfaFinals) = explore nfaAdjacency nfaFinals alphabet (
    initDfaStatesMap, 0) initDfaAdjacency Set.empty [inits]
62     initDfaStatesMap = Map.fromList [(inits, 0)]
63     initDfaAdjacency = Map.fromList [(0, [])]
64     dfaInits = Set.fromList [0]

```

/src/Thompson.hs

```

1 module Thompson (regexParser, makeThompsonNFA) where
2 import Automaton (AdjacencyList,
3                  Automaton (..), Label (..),
4                  State)
5 import Control.Monad (msum)
6 import Data.Char (ord)
7 import qualified Data.Map as Map
8 import qualified Data.Set as Set
9 import qualified Text.ParserCombinators.Parsec as P
10 import qualified Text.ParserCombinators.Parsec.Expr as PE
11
12 data Node = Concat !Node !Node | Star !Node | Or !Node !Node | Character !Int deriving
    (Show)
13
14 regexParser :: P.Parser Node
15 regexParser = PE.buildExpressionParser opTable base
16     where
17         opTable = [ [ PE.Postfix (P.char '*' >> return Star)]
18                   , [ PE.Infix (return Concat) PE.AssocLeft]
19                   , [ PE.Infix (P.char '|' >> return Or) PE.AssocLeft]
20                   ]
21         base = msum [Character . ord <$> P.noneOf "()*|", parens regexParser]
22         parens = P.between (P.char '(') (P.char ')')
23
24 -- Build NFA from regex AST following Thompson's Algorithm
25 makeThompsonNFA :: Either a Node -> Automaton
26 makeThompsonNFA ((Right ast)) = Automaton table alphabet start end where
27     (table, _) = thompsons ast 1 0 1
28     alphabet = Set.toList (buildAlph ast Set.empty)
29     start = Set.singleton 1
30     end = Set.singleton 0
31 makeThompsonNFA (Left _) = error "Bad AST"
32
33 -- buildAlph: Build alphabet from regex AST. Helper method for makeThompsonNFA
34 buildAlph :: Node -> Set.Set Label -> Set.Set Label
35 buildAlph (Concat r l) alph = Set.union (buildAlph r alph) (buildAlph l alph)
36 buildAlph (Star l) alph = buildAlph l alph
37 buildAlph (Or r l) alph = Set.union (buildAlph r alph) (buildAlph l alph)
38 buildAlph (Character x) alph = Set.insert (Label x) alph
39
40 -- thompsons: Build adjacency list from regex AST. Helper method for makeThompsonNFA
41 thompsons :: Node -> State -> State -> State -> (AdjacencyList, State)
42 thompsons (Character x) q f l = (Map.fromList [(q, [(Label x, f)])], l)
43 thompsons (Concat s t) q f l = (Map.union smap tmap, lt) where
44     (smap, ls) = thompsons s q i ln
45     (tmap, lt) = thompsons t i f ls
46     i = l + 1

```

```

47   ln = l + 1
48 thompsons (Or s t) q f l = (Map.union (Map.union smap tmap) omap, lt) where
49   omap = Map.fromList [(q, [(Epsilon, si), (Epsilon, ti)]), (sf, [(Epsilon, f)]), (tf,
      [(Epsilon, f)])]
50   (smap, ls) = thompsons s si sf ln
51   (tmap, lt) = thompsons t ti tf ls
52   si = l+1; sf = l+2
53   ti = l+3; tf = l+4
54   ln = l+4
55 thompsons (Star s) q f l = (Map.union smap stmap, ls) where
56   stmap = Map.fromList [(q, [(Epsilon, si), (Epsilon, f)]), (sf, [(Epsilon, si), (
      Epsilon, f)])]
57   (smap, ls) = thompsons s si sf ln
58   si = l+1; sf = l+2
59   ln = l+2

```

/src/WordList.hs

```

1 module WordList (buildList) where
2 import Data.Char (isAlpha)
3 import qualified Data.Set as Set
4 import qualified Data.Text as T
5 import Data.Text.IO as TIO (readFile)
6
7 buildList :: FilePath -> IO [[Char]]
8 buildList filename = do
9   h <- TIO.readFile filename
10  let l = T.words h
11      lnorm = map normalize l
12          lnormset = Set.fromList lnorm
13          lnormunique = Set.toList lnormset
14      return lnormunique
15
16 normalize :: T.Text -> [Char]
17 normalize string = [ x | x <- a, isAlpha x ] where a = T.unpack (T.toLower string)

```

/Setup.hs

```

1 import Distribution.Simple
2 main = defaultMain

```

/stack.yaml

```

1 # This file was automatically generated by 'stack init'
2 #
3 # Some commonly used options have been documented as comments in this file.
4 # For advanced use and comprehensive documentation of the format, please see:
5 # https://docs.haskellstack.org/en/stable/yaml_configuration/
6
7 # Resolver to choose a 'specific' stackage snapshot or a compiler version.
8 # A snapshot resolver dictates the compiler version and the set of packages
9 # to be used for project dependencies. For example:
10 #
11 # resolver: lts-3.5
12 # resolver: nightly-2015-09-21
13 # resolver: ghc-7.10.2
14 #
15 # The location of a snapshot can be provided as a file or url. Stack assumes
16 # a snapshot provided as a file might change, whereas a url resource does not.
17 #
18 # resolver: ./custom-snapshot.yaml
19 # resolver: https://example.com/snapshots/2018-01-01.yaml
20 system-ghc: true
21 resolver: ghc-8.10.7
22
23 # User packages to be built.
24 # Various formats can be used as shown in the example below.
25 #
26 # packages:
27 # - some-directory
28 # - https://example.com/foo/bar/baz-0.0.2.tar.gz
29 #   subdirs:
30 #     - auto-update

```

```

31 # - wai
32 packages:
33 - .
34 # Dependency packages to be pulled from upstream that are not in the resolver.
35 # These entries can reference officially published versions as well as
36 # forks / in-progress versions pinned to a git hash. For example:
37 #
38 extra-deps:
39 - random-1.2.1.1
40 - splitmix-0.1.0.4
41 - parallel-3.2.2.0
42 # - acme-missiles-0.3
43 # - git: https://github.com/commercialhaskell/stack.git
44 #   commit: e7b331f14bcffb8367cd58fbfc8b40ec7642100a
45 #
46 # extra-deps: []
47
48 # Override default flag values for local packages and extra-deps
49 # flags: {}
50
51 # Extra package databases containing global packages
52 # extra-package-dbs: []
53
54 # Control whether we use the GHC we find on the path
55 # system-ghc: true
56 #
57 # Require a specific version of stack, using version ranges
58 # require-stack-version: -any # Default
59 # require-stack-version: ">=2.9"
60 #
61 # Override the architecture used by stack, especially useful on Windows
62 # arch: i386
63 # arch: x86_64
64 #
65 # Extra directories used by stack for building
66 # extra-include-dirs: [/path/to/dir]
67 # extra-lib-dirs: [/path/to/dir]
68 #
69 # Allow a newer minor version of GHC than the snapshot specifies
70 # compiler-check: newer-minor

```

/stack.yaml.lock

```

1 # This file was autogenerated by Stack.
2 # You should not edit this file by hand.
3 # For more information, please see the documentation at:
4 #   https://docs.haskellstack.org/en/stable/lock_files
5
6 packages:
7 - completed:
8   hackage: random-1.2.1.1@sha256:
9     dea1f11e5569332dc6c8efaad1cb301016a5587b6754943a49f9de08ae0e56d9,6541
10    pantry-tree:
11      sha256: 646ee77fe01178837ee928b61ae8653dcf190e9b5353ebabd094079c77a18b76
12      size: 1528
13    original:
14      hackage: random-1.2.1.1
15 - completed:
16   hackage: splitmix-0.1.0.4@sha256:804
17     e2574bc7e32d08cbab91e47ee6287b4df7d50851d73f9e778f94a9a7814c7,6521
18     pantry-tree:
19       sha256: b56f706c092dc0ac4875ef45bd18719386358c56667f6b604a733b66f9e4657f
20       size: 1519
21     original:
22       hackage: splitmix-0.1.0.4
23 - completed:
24   hackage: parallel-3.2.2.0@sha256:6
25     edd5a06938cea3d28b406d5231683f89737e854af144a8800aa69e1eee785e0,1821
26     pantry-tree:
27       sha256: 6ec36425356925d6d9042769a29ab4ec2aa69c2a7161c49ff18a9a77c1d957b1
28       size: 392
29     original:

```

```
27     package: parallel-3.2.2.0
28 snapshots: []
```

/package.yaml

```
1 name: subset-construction
2 version: 0.1.0.0
3 github: "AlexisGado/subset-construction"
4 license: BSD3
5 author: "Alexis Gadonneix"
6 maintainer: "ag4625@columbia.edu"
7 copyright: "2022 Alexis Gadonneix"
8
9 extra-source-files:
10 - README.md
11 - CHANGELOG.md
12 - assets/**
13
14 # Metadata used when publishing your package
15 # synopsis:          Short description of your package
16 # category:         Web
17
18 # To avoid duplicated efforts in documentation and dealing with the
19 # complications of embedding Haddock markup inside cabal files, it is
20 # common to point users to the README.md file.
21 description: Please see the README on GitHub at <https://github.com/githubuser/subset-
    construction#readme>
22
23 dependencies:
24 - base >= 4.7 && < 5
25 - containers
26 - random
27 - parsec
28 - parallel
29 - deepseq
30 - text
31 - mtl
32
33 ghc-options:
34 - -Wall
35 - -Wcompat
36 - -Widentities
37 - -Wincomplete-record-updates
38 - -Wincomplete-uni-patterns
39 - -Wmissing-export-lists
40 - -Wmissing-home-modules
41 - -Wpartial-fields
42 - -Wredundant-constraints
43
44 library:
45   source-dirs: src
46
47 executables:
48   subset-construction-exe:
49     main: Main.hs
50     source-dirs: app
51     ghc-options:
52       - -threaded
53       - -rtsopts
54       - -eventlog
55       - -with-rtsopts=-N
56       - -O2
57     dependencies:
58       - subset-construction
59
60 tests:
61   subset-construction-test:
62     main: Spec.hs
63     source-dirs: test
64     ghc-options:
65       - -threaded
66       - -rtsopts
```



```
67 - -with-rtsopts=-N
68 dependencies:
69 - subset-construction
```

/app/Main.hs

```
1 module Main (main) where
2 import Automaton (Automaton (..))
3 import qualified Data.Map as Map
4 import RandomNfa (ioRandomNfa)
5 import qualified SubsetConstruction as SC
6
7 usedAutomaton :: IO Automaton
8 -- random NFA with control over density and size
9 usedAutomaton = ioRandomNfa 600 20 10 50
10
11 -- NFA from a list of words
12 -- usedAutomaton = dictNfa "assets/words.txt"
13
14 -- NFA from a regular expression
15 -- usedAutomaton = return $ makeThompsonNFA $ P.parse regexParser "" "a*(a|b)b*"
16
17 -- Some simple examples
18 -- usedAutomaton = return exampleAutomaton
19 -- usedAutomaton = return exampleAutomaton2
20
21 -- a simple linear automaton to test scale
22 -- usedAutomaton = ioDumbAutomaton 300
23
24
25 main :: IO ()
26 main = do
27     nfa <- usedAutomaton
28     let Automaton nfaAdj _ _ = nfa
29         Automaton dfaAdj _ _ = SC.nfaToDfa nfa
30     print $ Map.size nfaAdj
31     print $ Map.size dfaAdj
```