

## Handout 3 (Finite Automata)

Every formal language and compiler course I know of bombards you first with automata and then to a much, much smaller extent with regular expressions. As you can see, this course is turned upside down: regular expressions come first. The reason is that regular expressions are easier to reason about and the notion of derivatives, although already quite old, only became more widely known rather recently. Still, let us in this lecture have a closer look at automata and their relation to regular expressions. This will help us with understanding why the regular expression matchers in Python, Ruby, Java and so on are so slow with certain regular expressions. On the way we will also see what are the limitations of regular expressions. Unfortunately, they cannot be used for *everything*.

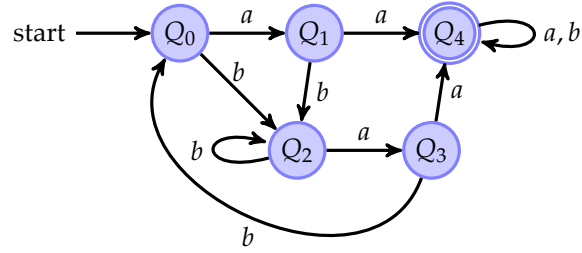
### Deterministic Finite Automata

Lets start...the central definition is:

A *deterministic finite automaton* (DFA), say  $A$ , is given by a five-tuple written  $\mathcal{A}(\Sigma, Qs, Q_0, F, \delta)$  where

- $\Sigma$  is an alphabet,
- $Qs$  is a finite set of states,
- $Q_0 \in Qs$  is the start state,
- $F \subseteq Qs$  are the accepting states, and
- $\delta$  is the transition function.

I am sure you have seen this definition already before. The transition function determines how to “transition” from one state to the next state with respect to a character. We have the assumption that these transition functions do not need to be defined everywhere: so it can be the case that given a character there is no next state, in which case we need to raise a kind of “failure exception”. That means actually we have *partial* functions as transitions—see the Scala implementation for DFAs later on. A typical example of a DFA is



In this graphical notation, the accepting state  $Q_4$  is indicated with double circles. Note that there can be more than one accepting state. It is also possible that a DFA has no accepting state at all, or that the starting state is also an accepting state. In the case above the transition function is defined everywhere and can also be given as a table as follows:

$(Q_0, a)$	$\rightarrow$	$Q_1$
$(Q_0, b)$	$\rightarrow$	$Q_2$
$(Q_1, a)$	$\rightarrow$	$Q_4$
$(Q_1, b)$	$\rightarrow$	$Q_2$
$(Q_2, a)$	$\rightarrow$	$Q_3$
$(Q_2, b)$	$\rightarrow$	$Q_0$
$(Q_3, a)$	$\rightarrow$	$Q_4$
$(Q_3, b)$	$\rightarrow$	$Q_0$
$(Q_4, a)$	$\rightarrow$	$Q_4$
$(Q_4, b)$	$\rightarrow$	$Q_4$

Please check that this table represents the same transition function as the graph above.

We need to define the notion of what language is accepted by an automaton. For this we lift the transition function  $\delta$  from characters to strings as follows:

$$\begin{aligned}\widehat{\delta}(q, []) &\stackrel{\text{def}}{=} q \\ \widehat{\delta}(q, c::s) &\stackrel{\text{def}}{=} \widehat{\delta}(\delta(q, c), s)\end{aligned}$$

This lifted transition function is often called *delta-hat*. Given a string, we start in the starting state and take the first character of the string, follow to the next state, then take the second character and so on. Once the string is exhausted and we end up in an accepting state, then this string is accepted by the automaton. Otherwise it is not accepted. This also means that if along the way we hit the case where the transition function  $\delta$  is not defined, we need to raise an error. In our implementation we will deal with this case elegantly by using Scala's Try. Summing up: a string  $s$  is in the *language accepted by the automaton*  $\mathcal{A}(\Sigma, Q, Q_0, F, \delta)$  iff

$$\widehat{\delta}(Q_0, s) \in F$$

I let you think about a definition that describes the set of all strings accepted by a deterministic finite automaton.

My take on an implementation of DFAs in Scala is given in Figure ?? . As you can see, there are many features of the mathematical definition that are quite closely reflected in the code. In the DFA-class, there is a starting state, called *start*, with the polymorphic type *A*. There is a partial function *delta* for specifying the transitions—these partial functions take a state (of polymorphic type *A*) and an input (of polymorphic type *C*) and produce a new state (of type *A*). For the moment it is OK to assume that *A* is some arbitrary type for states and the input is just characters. (The reason for not having concrete types, but polymorphic types for the states and the input of DFAs will become clearer later on.)

The DFA-class has also an argument for specifying final states. In the implementation it is not a set of states, as in the mathematical definition, but a function from states to booleans (this function is supposed to return true whenever a state is final; false otherwise). While this boolean function is different from the sets of states, Scala allows us to use sets for such functions (see Line 41 where the DFA is initialised). Again it will become clear later on why I use functions for final states, rather than sets.

The most important point in the implementation is that I use Scala’s partial functions for representing the transitions; alternatives would have been Maps or even Lists. One of the main advantages of using partial functions is that transitions can be quite nicely defined by a series of case statements (see Lines 29 – 39 for an example). If you need to represent an automaton with a sink state (catch-all-state), you can use Scala’s pattern matching and write something like

```
abstract class State
...
case object Sink extends State

val delta : (State, Char) => State =
{ case (S0, 'a') => S1
  case (S1, 'a') => S2
  case _ => Sink
}
```

I let you think what the corresponding DFA looks like in the graph notation. Also, I suggest you to tinker with the Scala code in order to define the DFA that does not accept any string at all.

Finally, I let you ponder whether this is a good implementation of DFAs or not. In doing so I hope you notice that the  $\Sigma$  and  $Qs$  components (the alphabet and the set of *finite* states, respectively) are missing from the class definition. This means that the implementation allows you to do some “fishy” things you are not meant to do with DFAs. Which fishy things could that be?

```

1 // DFAs in Scala using partial functions
2 import scala.util.Try
3
4 // a type abbreviation for partial functions
5 type =>[A, B] = PartialFunction[A, B]
6
7 // main DFA class
8 case class DFA[A, C](start: A,           // the starting state
9                      delta: (A, C) => A, // transitions (partial fun)
10                     fins: A => Boolean) { // the final states
11
12     def deltas(q: A, s: List[C]) : A = s match {
13         case Nil => q
14         case c::cs => deltas(delta(q, c), cs)
15     }
16
17     def accepts(s: List[C]) : Boolean =
18         Try(fins(deltas(start, s))).getOrElse(false)
19 }
20
21 // the example shown in the handout
22 abstract class State
23 case object Q0 extends State
24 case object Q1 extends State
25 case object Q2 extends State
26 case object Q3 extends State
27 case object Q4 extends State
28
29 val delta : (State, Char) => State =
30     { case (Q0, 'a') => Q1
31       case (Q0, 'b') => Q2
32       case (Q1, 'a') => Q4
33       case (Q1, 'b') => Q2
34       case (Q2, 'a') => Q3
35       case (Q2, 'b') => Q2
36       case (Q3, 'a') => Q4
37       case (Q3, 'b') => Q0
38       case (Q4, 'a') => Q4
39       case (Q4, 'b') => Q4 }
40
41 val dfa = DFA(Q0, delta, Set[State](Q4))
42 dfa.accepts("aaabbb".toList) // true

```

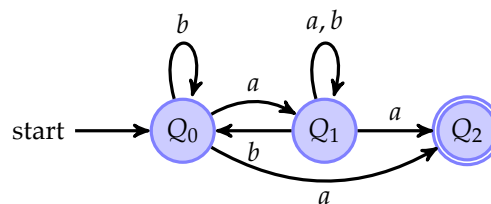
**Figure 1:** An implementation of DFAs in Scala using partial functions. Note some subtleties: `deltas` implements the  $\delta$ -hat construction by lifting the (partial) transition function to lists of characters. Since `delta` is given as a partial function, it can obviously go “wrong” in which case the `Try` in `accepts` catches the error and returns `false`—that means the string is not accepted. The example `delta` in Line 22–43 implements the DFA example shown earlier in the handout.

## Non-Deterministic Finite Automata

Remember we want to find out what the relation is between regular expressions and automata. To do this with DFAs is a bit unwieldy. While with DFAs it is always clear that given a state and a character what the next state is (potentially none), it will be convenient to relax this restriction. That means we allow states to have several potential successor states. We even allow more than one starting state. The resulting construction is called a *Non-Deterministic Finite Automaton* (NFA) given also as a five-tuple  $\mathcal{A}(\Sigma, Q_s, Q_{0s}, F, \rho)$  where

- $\Sigma$  is an alphabet,
- $Q_s$  is a finite set of states
- $Q_{0s}$  is a set of start states ( $Q_{0s} \subseteq Q_s$ )
- $F$  are some accepting states with  $F \subseteq Q_s$ , and
- $\rho$  is a transition relation.

A typical example of a NFA is



This NFA happens to have only one starting state, but in general there could be more than one. Notice that in state  $Q_0$  we might go to state  $Q_1$  or to state  $Q_2$  when receiving an  $a$ . Similarly in state  $Q_1$  and receiving an  $a$ , we can stay in state  $Q_1$  or go to  $Q_2$ . This kind of choice is not allowed with DFAs. The downside of this choice in NFAs is that when it comes to deciding whether a string is accepted by a NFA we potentially have to explore all possibilities. I let you think which strings the above NFA accepts.

There are a number of additional points you should note about NFAs. Every DFA is a NFA, but not vice versa. The  $\rho$  in NFAs is a transition *relation* (DFAs have a transition *function*). The difference between a function and a relation is that a function has always a single output, while a relation gives, roughly speaking, several outputs. Look again at the NFA above: if you are currently in the state  $Q_1$  and you read a character  $b$ , then you can transition to either  $Q_0$  or  $Q_2$ . Which route, or output, you take is not determined. This non-determinism can be represented by a relation.

My implementation of NFAs in Scala is shown in Figure ?? . Perhaps interestingly, I do not actually use relations for my NFAs, but use transition functions that return sets of states. DFAs have partial transition functions that return

a single state; my NFAs return a set of states. I let you think about this representation for NFA-transitions and how it corresponds to the relations used in the mathematical definition of NFAs. An example of a transition function in Scala for the NFA shown above is

```
val nfa_delta : (State, Char) => Set[State] =
  { case (Q0, 'a') => Set(Q1, Q2)
    case (Q0, 'b') => Set(Q0)
    case (Q1, 'a') => Set(Q1, Q2)
    case (Q1, 'b') => Set(Q0, Q1) }
```

Like in the mathematical definition, `starts` is in NFAs a set of states; `fins` is again a function from states to booleans. The `next` function calculates the set of next states reachable from a single state `q` by a character `c`. In case there is no such state—the partial transition function is undefined—the empty set is returned (see function `applyOrElse` in Lines 11 and 12). The function `nexts` just lifts this function to sets of states.

Look very careful at the `accepts` and `deltas` functions in NFAs and remember that when accepting a string by a NFA we might have to explore all possible transitions (recall which state to go to is not unique any more with NFAs...we need to explore potentially all next states). The implementation achieves this exploration through a *breadth-first search*. This is fine for small NFAs, but can lead to real memory problems when the NFAs are bigger and larger strings need to be processed. As result, some regular expression matching engines resort to a *depth-first search* with *backtracking* in unsuccessful cases. In our implementation we can implement a depth-first version of `accepts` using Scala's `exists`-function as follows:

```
def search(q: A, s: List[C]) : Boolean = s match {
  case Nil => fins(q)
  case c::cs => next(q, c).exists(search(_, cs))
}

def accepts2(s: List[C]) : Boolean =
  starts.exists(search(_, s))
```

This depth-first way of exploration seems to work quite efficiently in many examples and is much less of a strain on memory. The problem is that the backtracking can get “catastrophic” in some examples—remember the catastrophic backtracking from earlier lectures. This depth-first search with backtracking is the reason for the abysmal performance of some regular expression matchings in Java, Ruby and Python. I like to show you this in the next two sections.

## Epsilon NFAs

In order to get an idea what calculations are performed by Java & friends, we need a method for transforming a regular expression into a corresponding automaton. This automaton should accept exactly those strings that are accepted by the regular expression. The simplest and most well-known method for this is

```

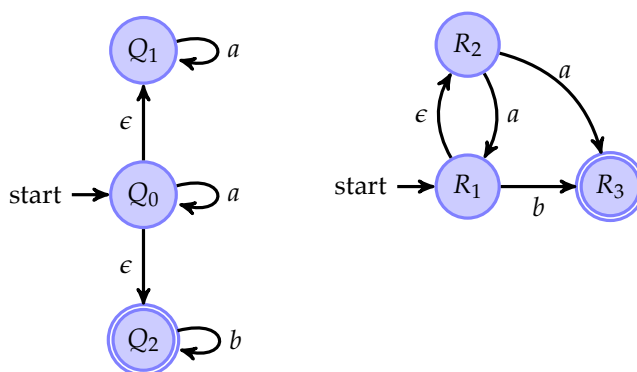
1 // NFAs in Scala using partial functions (returning
2 // sets of states)
3 //
4
5 // load dfas
6 import $file.dfa, dfa._
7
8
9 // return an empty set, when a partial function is not defined
10 import scala.util.Try
11 def applyOrElse[A, B](f: A => Set[B], x: A) : Set[B] =
12   Try(f(x)).getOrElse(Set[B]())
13
14
15 // NFAs
16 case class NFA[A, C](starts: Set[A],           // the starting states
17                      delta: (A, C) => Set[A], // the transition function
18                      fins: A => Boolean) {      // the final states
19
20   // given a state and a character, what is the set of
21   // next states? if there is none => empty set
22   def next(q: A, c: C) : Set[A] =
23     applyOrElse(delta, (q, c))
24
25   def nexts(qs: Set[A], c: C) : Set[A] =
26     qs.flatMap(next(_, c))
27
28   // given some states and a string, what is the set
29   // of next states?
30   def deltas(qs: Set[A], s: List[C]) : Set[A] = s match {
31     case Nil => qs
32     case c::cs => deltas(nexts(qs, c), cs)
33   }
34
35   // is a string accepted by an NFA?
36   def accepts(s: List[C]) : Boolean =
37     deltas(starts, s).exists(fins)
38
39   // depth-first version of accepts
40   def search(q: A, s: List[C]) : Boolean = s match {
41     case Nil => fins(q)
42     case c::cs => next(q, c).exists(search(_, cs))
43   }

```

**Figure 2:** A Scala implementation of NFAs using partial functions. Notice that the function `accepts` implements the acceptance of a string in a breadth-first search fashion. This can be a costly way of deciding whether a string is accepted or not in applications that need to handle large NFAs and large inputs.

called the *Thompson Construction*, after the Turing Award winner Ken Thompson. This method is by recursion over regular expressions and depends on the non-determinism in NFAs described in the previous section. You will see shortly why this construction works well with NFAs, but is not so straightforward with DFAs.

Unfortunately we are still one step away from our intended target though—because this construction uses a version of NFAs that allows “silent transitions”. The idea behind silent transitions is that they allow us to go from one state to the next without having to consume a character. We label such silent transition with the letter  $\epsilon$  and call the automata  $\epsilon$ NFAs. Two typical examples of  $\epsilon$ NFAs are:



Consider the  $\epsilon$ NFA on the left-hand side: the  $\epsilon$ -transitions mean you do not have to “consume” any part of the input string, but “silently” change to a different state. In this example, if you are in the starting state  $Q_0$ , you can silently move either to  $Q_1$  or  $Q_2$ . You can see that once you are in  $Q_1$ , respectively  $Q_2$ , you cannot “go back” to the other states. So it seems allowing  $\epsilon$ -transitions is a rather substantial extension to NFAs. On first appearances,  $\epsilon$ -transitions might even look rather strange, or even silly. To start with, silent transitions make the decision whether a string is accepted by an automaton even harder: with  $\epsilon$ NFAs we have to look whether we can do first some  $\epsilon$ -transitions and then do a “proper”-transition; and after any “proper”-transition we again have to check whether we can do again some silent transitions. Even worse, if there is a silent transition pointing back to the same state, then we have to be careful our decision procedure for strings does not loop (remember the depth-first search for exploring all states).

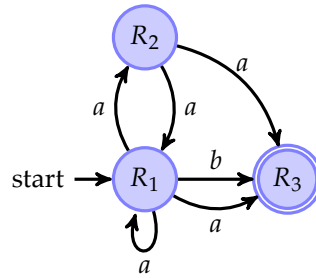
The obvious question is: Do we get anything in return for this hassle with silent transitions? Well, we still have to work for it...unfortunately. If we were to follow the many textbooks on the subject, we would now start with defining what  $\epsilon$ NFAs are—that would require extending the transition relation of NFAs. Next, we would show that the  $\epsilon$ NFAs are equivalent to NFAs and so on. Once we have done all this on paper, we would need to implement  $\epsilon$ NFAs. Let’s try to take a shortcut instead. We are not really interested in  $\epsilon$ NFAs; they are only a convenient tool for translating regular expressions into automata. So we are



not going to implementing them explicitly, but translate them immediately into NFAs (in a sense  $\epsilon$ NFAs are just a convenient API for lazy people ;o). How does this translation work? Well we have to find all transitions of the form

$$q \xrightarrow{\epsilon} \dots \xrightarrow{\epsilon} \xrightarrow{a} \xrightarrow{\epsilon} \dots \xrightarrow{\epsilon} q'$$

where somewhere in the “middle” is an  $a$ -transition (for a character, say,  $a$ ). We replace them with  $q \xrightarrow{a} q'$ . Doing this to the  $\epsilon$ NFA on the right-hand side above gives the NFA



where the single  $\epsilon$ -transition is replaced by three additional  $a$ -transitions. Please do the calculations yourself and verify that I did not forget any transition.

So in what follows, whenever we are given an  $\epsilon$ NFA we will replace it by an equivalent NFA. The Scala code for this translation is given in Figure ???. The main workhorse in this code is a function that calculates a fixpoint of function (Lines 6–12). This function is used for “discovering” which states are reachable by  $\epsilon$ -transitions. Once no new state is discovered, a fixpoint is reached. This is used for example when calculating the starting states of an equivalent NFA (see Line 28): we start with all starting states of the  $\epsilon$ NFA and then look for all additional states that can be reached by  $\epsilon$ -transitions. We keep on doing this until no new state can be reached. This is what the  $\epsilon$ -closure, named in the code `ec1`, calculates. Similarly, an accepting state of the NFA is when we can reach an accepting state of the  $\epsilon$ NFA by  $\epsilon$ -transitions.

Also look carefully how the transitions of  $\epsilon$ NFAs are implemented. The additional possibility of performing silent transitions is encoded by using `Option[C]` as the type for the “input”. The `Somes` stand for “proper” transitions where a character is consumed; `None` stands for  $\epsilon$ -transitions. The transition functions for the two  $\epsilon$ NFAs from the beginning of this section can be defined as

```

1 // epsilon NFAs...immediately translated into NFAs
2
3 import $file.dfa, dfa._
4 import $file.nfa, nfa._
5
6 // a fixpoint construction
7 import scala.annotation.tailrec
8 @tailrec
9 def fixpT[S](f: S => S, x: S): S = {
10   val fx = f(x)
11   if (fx == x) x else fixpT(f, fx)
12 }
13
14 // translates eNFAs directly into NFAs
15 def eNFA[A, C](starts: Set[A], // the starting states
16               delta: (A, Option[C]) => Set[A], // the epsilon-transitions
17               fins: A => Boolean) : NFA[A, C] = { // the final states
18
19   // epsilon transitions
20   def enext(q: A) : Set[A] =
21     applyOrElse(delta, (q, None))
22
23   def enexts(qs: Set[A]) : Set[A] =
24     qs | qs.flatMap(enext(_)) // | is the set-union in Scala
25
26   // epsilon closure
27   def ecl(qs: Set[A]) : Set[A] =
28     fixpT(enexts, qs)
29
30   // "normal" transitions
31   def next(q: A, c: C) : Set[A] =
32     applyOrElse(delta, (q, Some(c)))
33
34   def nexts(qs: Set[A], c: C) : Set[A] =
35     ecl(ecl(qs).flatMap(next(_, c)))
36
37   // result NFA
38   NFA(ecl(starts),
39       { case (q, c) => nexts(Set(q), c) },
40       q => ecl(Set(q)) exists fins)
41 }

```

**Figure 3:** A Scala function that translates  $\epsilon$ NFA into NFAs. The transition function of  $\epsilon$ NFA takes as input an `Option[C]`. `None` stands for an  $\epsilon$ -transition; `Some(c)` for a “proper” transition consuming a character. The functions in Lines 19–24 calculate all states reachable by one or more  $\epsilon$ -transition for a given set of states. The NFA is constructed in Lines 30–34. Note the interesting commands in Lines 7 and 8: their purpose is to ensure that `fixpT` is the tail-recursive version of the fixpoint construction; otherwise we would quickly get a stack-overflow exception, even on small examples, due to limitations of the JVM.

```

val enfa_trans1 : (State, Option[Char]) => Set[State] =
  { case (Q0, Some('a')) => Set(Q0)
    case (Q0, None) => Set(Q1, Q2)
    case (Q1, Some('a')) => Set(Q1)
    case (Q2, Some('b')) => Set(Q2) }

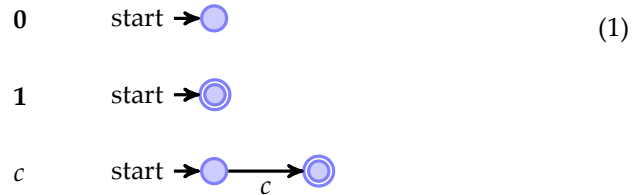
val enfa_trans2 : (State, Option[Char]) => Set[State] =
  { case (R1, Some('b')) => Set(R3)
    case (R1, None) => Set(R2)
    case (R2, Some('a')) => Set(R1, R3) }

```

I hope you agree now with my earlier statement that the  $\epsilon$ NFAs are just an API for NFAs.

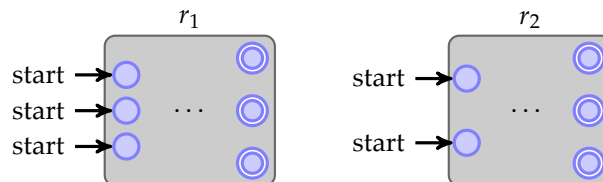
## Thompson Construction

Having the translation of  $\epsilon$ NFAs to NFAs in place, we can finally return to the problem of translating regular expressions into equivalent NFAs. Recall that by equivalent we mean that the NFAs recognise the same language. Consider the simple regular expressions  $0$ ,  $1$  and  $c$ . They can be translated into equivalent NFAs as follows:



I let you think whether the NFAs can match exactly those strings the regular expressions can match. To do this translation in code we need a way to construct states “programmatically”...and as an additional constraint Scala needs to recognise that these states are being distinct. For this I implemented in Figure ?? a class `TState` that includes a counter and a companion object that increases this counter whenever a new state is created.<sup>1</sup>

The case for the sequence regular expression  $r_1 \cdot r_2$  is a bit more complicated: Say, we are given by recursion two NFAs representing the regular expressions  $r_1$  and  $r_2$  respectively.



<sup>1</sup>You might have to read up what *companion objects* do in Scala.

```

1  // Thompson Construction
2  //=====
3
4  import $file.dfa, dfa._
5  import $file.nfa, nfa._
6  import $file.enfa, enfa._
7
8
9  // states for Thompson construction
10 case class TState(i: Int) extends State
11
12 object TState {
13     var counter = 0
14
15     def apply() : TState = {
16         counter += 1;
17         new TState(counter)
18     }
19 }
20
21 ...
22
23 // NFA that does not accept any string
24 def NFA_ZERO(): NFAt = {
25     val Q = TState()
26     NFA(Set(Q), { case _ => Set() }, Set())
27 }
28
29 // NFA that accepts the empty string
30 def NFA_ONE() : NFAt = {
31     val Q = TState()
32     NFA(Set(Q), { case _ => Set() }, Set(Q))
33 }
34
35 // NFA that accepts the string "c"
36 def NFA_CHAR(c: Char) : NFAt = {
37     val Q1 = TState()
38     val Q2 = TState()
39     NFA(Set(Q1), { case (Q1, d) if (c == d) => Set(Q2) }, Set(Q2))
40 }

```

**Figure 4:** The first part of the Thompson Construction. Lines 10–19 implement a way of how to create new states that are all distinct by virtue of a counter. This counter is increased in the companion object of TState whenever a new state is created. The code in Lines 38–45 constructs NFAs for the simple regular expressions 0, 1 and c. Compare this code with the pictures given in (??) on Page ??.

```

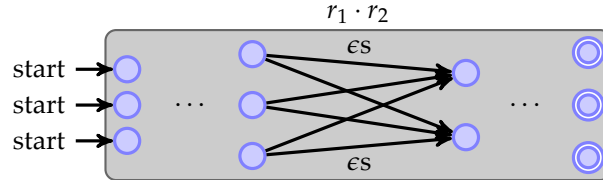
48 // for composing an eNFA transition with an NFA transition
49 // | is for set union
50 extension (f: eNFATrans) {
51   def +++(g: NFATrans) : eNFATrans =
52     { case (q, None) => applyOrElse(f, (q, None))
53       case (q, Some(c)) => applyOrElse(f, (q, Some(c))) | applyOrElse(g, (q, c)) }
54 }
55
56
57 // sequence of two NFAs
58 def NFA_SEQ(nfa1: NFAT, nfa2: NFAT) : NFAT = {
59   val new_delta : eNFATrans =
60     { case (q, None) if nfa1.fins(q) => nfa2.starts }
61
62   eNFA(nfa1.starts,
63        new_delta +++ nfa1.delta +++ nfa2.delta,
64        nfa2.fins)
65 }
66
67 // alternative of two NFAs
68 def NFA_ALT(nfa1: NFAT, nfa2: NFAT) : NFAT = {
69   val new_delta : NFATrans = {
70     case (q, c) => applyOrElse(nfa1.delta, (q, c)) |
71                   applyOrElse(nfa2.delta, (q, c)) }
72   val new_fins = (q: TState) => nfa1.fins(q) || nfa2.fins(q)
73
74   NFA(nfa1.starts | nfa2.starts, new_delta, new_fins)
75 }
76
77 // star of a NFA
78 def NFA_STAR(nfa: NFAT) : NFAT = {
79   val Q = TState()
80   val new_delta : eNFATrans =
81     { case (Q, None) => nfa.starts
82       case (q, None) if nfa.fins(q) => Set(Q) }
83
84   eNFA(Set(Q), new_delta +++ nfa.delta, Set(Q))
85 }

```

**Figure 5:** The second part of the Thompson Construction implementing the composition of NFAs according to  $\cdot$ ,  $+$  and  $*$ . The extension (Lines 48–54) about rich partial functions implements the infix operation  $+++$  which combines an eNFA transition with an NFA transition (both are given as partial functions—but with different type!).

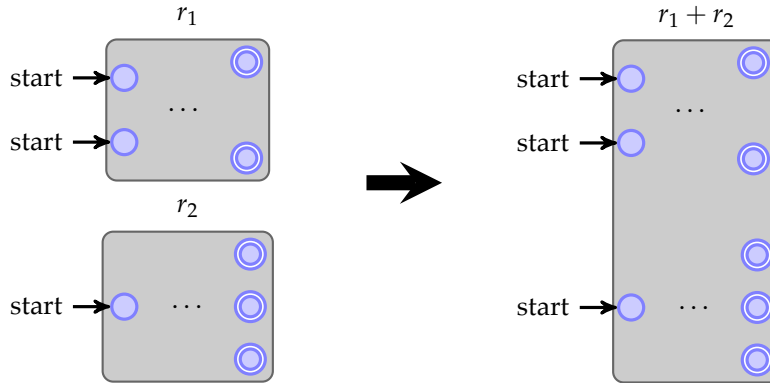
---

The first NFA has some accepting states and the second some starting states. We obtain an  $\epsilon$ NFA for  $r_1 \cdot r_2$  by connecting the accepting states of the first NFA with  $\epsilon$ -transitions to the starting states of the second automaton. By doing so we make the accepting states of the first NFA to be non-accepting like so:



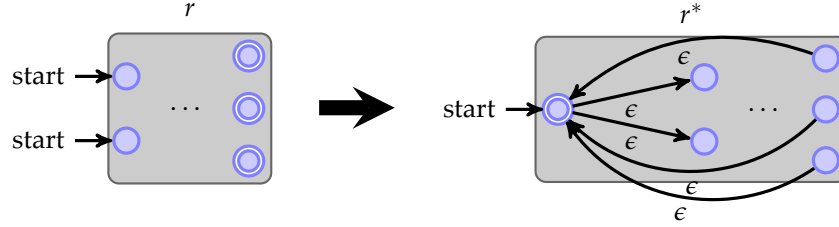
The idea behind this construction is that the start of any string is first recognised by the first NFA, then we silently change to the second NFA; the ending of the string is recognised by the second NFA...just like matching of a string by the regular expression  $r_1 \cdot r_2$ . The Scala code for this construction is given in Figure ?? in Lines 57–65. The starting states of the  $\epsilon$ NFA are the starting states of the first NFA (corresponding to  $r_1$ ); the accepting function is the accepting function of the second NFA (corresponding to  $r_2$ ). The new transition function is all the “old” transitions plus the  $\epsilon$ -transitions connecting the accepting states of the first NFA to the starting states of the second NFA (Lines 59 and 60). The  $\epsilon$ NFA is then immediately translated in a NFA.

The case for the alternative regular expression  $r_1 + r_2$  is slightly different: We are given by recursion two NFAs representing  $r_1$  and  $r_2$  respectively. Each NFA has some starting states and some accepting states. We obtain a NFA for the regular expression  $r_1 + r_2$  by composing the transition functions (this crucially depends on knowing that the states of each component NFA are distinct—recall we implemented for this to hold by some bespoke code for TStates). We also need to combine the starting states and accepting functions appropriately.



The code for this construction is in Figure ?? in Lines 67–75.

Finally for the  $*$ -case we have a NFA for  $r$  and connect its accepting states to a new starting state via  $\epsilon$ -transitions. This new starting state is also an accepting state, because  $r^*$  can recognise the empty string.



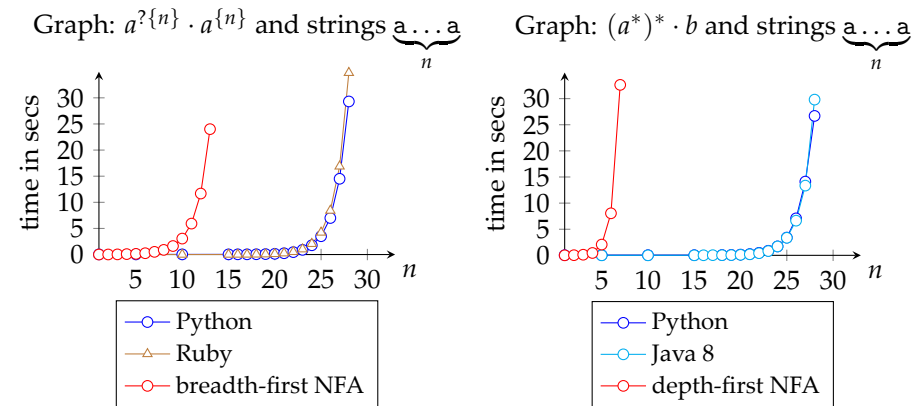
The corresponding code is in Figure ?? in Lines 77–85)

To sum up, you can see in the sequence and star cases the need for silent  $\epsilon$ -transitions. Otherwise this construction just becomes awkward. Similarly the alternative case shows the need of the NFA-nondeterminism. It looks non-obvious to form the ‘alternative’ composition of two DFAs, because DFA do not allow several starting and successor states. All these constructions can now be put together in the following recursive function:

```
def thompson(r: Rexp) : NFAt = r match {
  case ZERO => NFA_ZERO()
  case ONE  => NFA_ONE()
  case CHAR(c) => NFA_CHAR(c)
  case ALT(r1, r2) => NFA_ALT(thompson(r1), thompson(r2))
  case SEQ(r1, r2) => NFA_SEQ(thompson(r1), thompson(r2))
  case STAR(r1) => NFA_STAR(thompson(r1))
}
```

It calculates a NFA from a regular expression. At last we can run NFAs for our evil regular expression examples. The graph on the left shows that when translating a regular expression such as  $a^{?{n}} \cdot a^{{n}}$  into a NFA, the size can blow up and then even the relative fast (on small examples) breadth-first search can be slow...the red line maxes out at about 15 ns.

The graph on the right shows that with ‘evil’ regular expressions also the depth-first search can be abysmally slow. Even if the graphs not completely overlap with the curves of Python, Ruby and Java, they are similar enough.



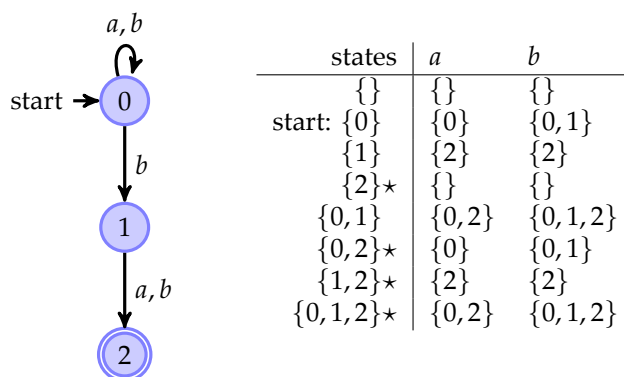
OK...now you know why regular expression matchers in those languages are sometimes so slow. A bit surprising, don't you agree? Also it is still a mystery why Rust, which because of the reasons above avoids NFAs and uses DFAs instead, cannot compete in all cases with our simple derivative-based regular expression matcher in Scala. There is an explanation for this as well...remember there the offending examples are of the form  $r^{\{n\}}$ . Why could they be a problem in Rust?

## Subset Construction

So of course, some clever developers of regular expression matchers are aware of these problems with sluggish NFAs and try to address them. One common technique for alleviating the problem I like to show you in this section. This will also explain why we insisted on polymorphic types in our DFA code (remember I used `A` and `C` for the types of states and the input, see Figure ?? on Page ??).

To start, remember that we did not bother with defining and implementing  $\epsilon$ NFAs: we immediately translated them into equivalent NFAs. Equivalent in the sense of accepting the same language (though we only claimed this and did not prove it rigorously). Remember also that NFAs have non-deterministic transitions defined as a relation, or alternatively my Scala implementation used transition functions returning sets of states. This non-determinism is crucial for the Thompson Construction to work (recall the cases for  $\cdot$ ,  $+$  and  $*$ ). But this non-determinism makes it harder with NFAs to decide when a string is accepted or not; whereas such a decision is rather straightforward with DFAs: recall their transition function is a “real” function that returns a single state. So with DFAs we do not have to search at all. What is perhaps interesting is the fact that for every NFA we can find a DFA that also recognises the same language. This might sound a bit paradoxical: NFA  $\rightarrow$  decision of acceptance hard; DFA  $\rightarrow$  decision easy. But this *is* true...but of course there is always a caveat—nothing ever is for free in life. Let's see what this caveat is.

There are actually a number of methods for transforming a NFA into an equivalent DFA, but the most famous one is the *subset construction*. Consider the following NFA where the states are labelled with 0, 1 and 2.





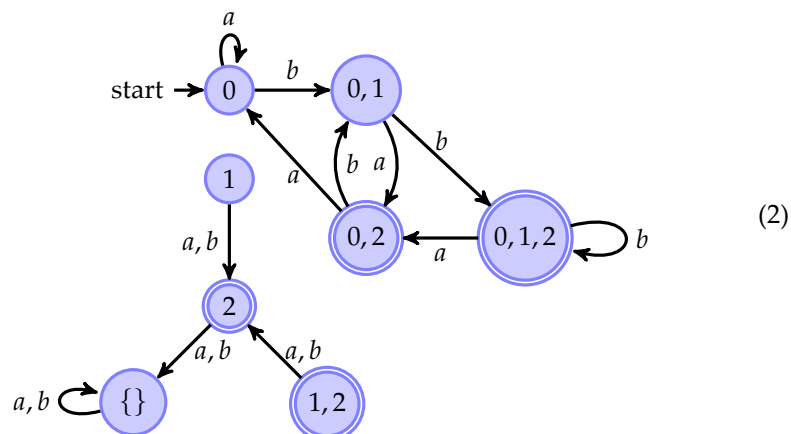
The states of the corresponding DFA are given by generating all subsets of the set  $\{0,1,2\}$  (seen in the states column in the table on the right). The other columns define the transition function for the DFA for inputs  $a$  and  $b$ . The first row states that  $\{\}$  is the sink state which has transitions for  $a$  and  $b$  to itself. The next three lines are calculated as follows:

- Suppose you calculate the entry for the  $a$ -transition for state  $\{0\}$ . Look for all states in the NFA that can be reached by such a transition from this state; this is only state 0; therefore from state  $\{0\}$  we can go to state  $\{0\}$  via an  $a$ -transition.
- Do the same for the  $b$ -transition; you can reach states 0 and 1 in the NFA; therefore in the DFA we can go from state  $\{0\}$  to state  $\{0,1\}$  via an  $b$ -transition.
- Continue with the states  $\{1\}$  and  $\{2\}$ .

Once you filled in the transitions for 'simple' states  $\{0\} \dots \{2\}$ , you only have to build the union for the compound states  $\{0,1\}$ ,  $\{0,2\}$  and so on. For example for  $\{0,1\}$  you take the union of Line  $\{0\}$  and Line  $\{1\}$ , which gives  $\{0,2\}$  for  $a$ , and  $\{0,1,2\}$  for  $b$ . And so on.

The starting state of the DFA can be calculated from the starting states of the NFA, that is in this case  $\{0\}$ . But in general there can of course be many starting states in the NFA and you would take the corresponding subset as *the* starting state of the DFA.

The accepting states in the DFA are given by all sets that contain a 2, which is the only accepting state in this NFA. But again in general if the subset contains any accepting state from the NFA, then the corresponding state in the DFA is accepting as well. This completes the subset construction. The corresponding DFA for the NFA shown above is:



Please check that this is indeed a DFA. The big question is whether this DFA can recognise the same language as the NFA we started with? I let you ponder about this question.

There are also two points to note: One is that very often in the subset construction the resulting DFA contains a number of “dead” states that are never reachable from the starting state. This is obvious in the example, where state  $\{1\}$ ,  $\{2\}$ ,  $\{1,2\}$  and  $\{\}$  can never be reached from the starting state. But this might not always be as obvious as that. In effect the DFA in this example is not a *minimal* DFA (more about this in a minute). Such dead states can be safely removed without changing the language that is recognised by the DFA. Another point is that in some cases, however, the subset construction produces a DFA that does *not* contain any dead states...this means it calculates a minimal DFA. Which in turn means that in some cases the number of states can by going from NFAs to DFAs exponentially increase, namely by  $2^n$  (which is the number of subsets you can form for sets of  $n$  states). This blow up in the number of states in the DFA is again bad news for how quickly you can decide whether a string is accepted by a DFA or not. So the caveat with DFAs is that they might make the task of finding the next state trivial, but might require  $2^n$  times as many states than a NFA.

To conclude this section, how conveniently we can implement the subset construction with our versions of NFAs and DFAs? Very conveniently. The code is just:

```
def subset[A, C](nfa: NFA[A, C]) : DFA[Set[A], C] = {
  DFA(nfa.starts,
    { case (qs, c) => nfa.nexts(qs, c) },
    _.exists(nfa.fins))
}
```

The interesting point in this code is that the state type of the calculated DFA is `Set[A]`. Think carefully that this works out correctly.

The DFA is then given by three components: the starting states, the transition function and the accepting-states function. The starting states are a set in the given NFA, but a single state in the DFA. The transition function, given the state `qs` and input `c`, needs to produce the next state: this is the set of all NFA states that are reachable from each state in `qs`. The function `nexts` from the NFA class already calculates this for us. The accepting-states function for the DFA is true whenever at least one state in the subset is accepting (that is true) in the NFA.

You might be able to spend some quality time tinkering with this code and time to ponder about it. Then you will probably notice that it is actually a bit silly. The whole point of translating the NFA into a DFA via the subset construction is to make the decision of whether a string is accepted or not faster. Given the code above, the generated DFA will be exactly as fast, or as slow, as the NFA we started with (actually it will even be a tiny bit slower). The reason is that we just re-use the `nexts` function from the NFA. This function implements the non-deterministic breadth-first search. You might be thinking: This is cheating! ... Well, not quite as you will see later, but in terms of speed we still need to work a bit in order to get sometimes(!) a faster DFA. Let's do this next.

## DFA Minimisation

As seen in (??), the subset construction from NFA to a DFA can result in a rather “inefficient” DFA. Meaning there are states that are not needed. There are two kinds of such unneeded states: *unreachable* states and *non-distinguishable* states. The first kind of states can just be removed without affecting the language that can be recognised (after all they are unreachable). The second kind can also be recognised and thus a DFA can be *minimised* by the following algorithm:

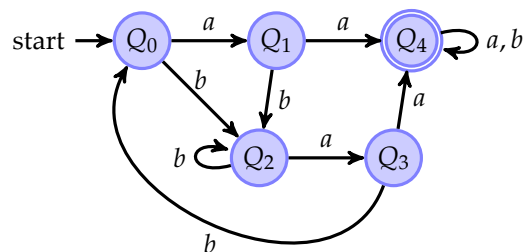
1. Take all pairs  $(q, p)$  with  $q \neq p$
2. Mark all pairs that accepting and non-accepting states
3. For all unmarked pairs  $(q, p)$  and all characters  $c$  test whether

$$(\delta(q, c), \delta(p, c))$$

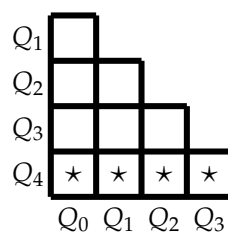
are marked. If there is one, then also mark  $(q, p)$ .

4. Repeat last step until no change.
5. All unmarked pairs can be merged.

Unfortunately, once we throw away all unreachable states in (??), all remaining states are needed. In order to illustrate the minimisation algorithm, consider the following DFA.



In Step 1 and 2 we consider essentially a triangle of the form

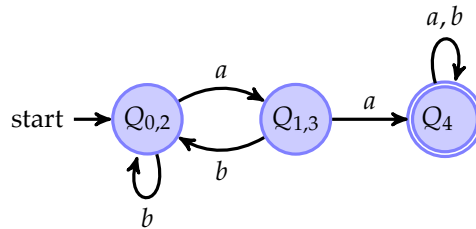


where the lower row is filled with stars, because in the corresponding pairs there is always one state that is accepting ( $Q_4$ ) and a state that is non-accepting (the other states).

In Step 3 we need to fill in more stars according to whether one of the next-state pairs are marked. We have to do this for every unmarked field until there is no change any more. This gives the triangle

$Q_1$	*			
$Q_2$		*		
$Q_3$	*		*	
$Q_4$	*	*	*	*
	$Q_0$	$Q_1$	$Q_2$	$Q_3$

which means states  $Q_0$  and  $Q_2$ , as well as  $Q_1$  and  $Q_3$  can be merged. This gives the following minimal DFA



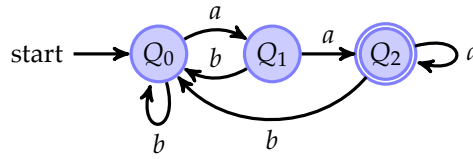
This minimised DFA is certainly fast when it comes deciding whether a string is accepted or not. But this is not universally the case. Suppose you count the nodes in a regular expression (when represented as tree). If you look carefully at the Thompson Construction you can see that the constructed NFA has states that grow linearly in terms of the size of the regular expression. This is good, but as we have seen earlier deciding whether a string is matched by an NFA is hard. Translating an NFA into a DFA means deciding whether a string is matched by a DFA is easy, but the number of states can grow exponentially, even after minimisation. Say a NFA has  $n$  states, then in the worst case the corresponding minimal DFA that can match the same language as the NFA might contain  $2^n$  of states. Unfortunately in many interesting cases this worst case bound is the dominant factor.

By the way, we are not bothering with implementing the above minimisation algorithm: while up to now all the transformations used some clever composition of functions, the minimisation algorithm cannot be implemented by just composing some functions. For this we would require a more concrete representation of the transition function (like maps). If we did this, however, then many advantages of the functions would be thrown away. So the compromise is to not being able to minimise (easily) our DFAs. We want to use regular expressions directly anyway.

## Brzowski's Method

I know this handout is already a long, long rant: but after all it is a topic that has been researched for more than 60 years. If you reflect on what you have read so far, the story is that you can take a regular expression, translate it via the Thompson Construction into an  $\epsilon$ NFA, then translate it into a NFA by removing all  $\epsilon$ -transitions, and then via the subset construction obtain a DFA. In all steps we made sure the language, or which strings can be recognised, stays the same. Of course we should have proved this in each step, but let us cut corners here. After the last section, we can even minimise the DFA (maybe not in code). But again we made sure the same language is recognised. You might be wondering: Can we go into the other direction? Can we go from a DFA and obtain a regular expression that can recognise the same language as the DFA?

The answer is yes. Again there are several methods for calculating a regular expression for a DFA. I will show you Brzowski's method because it calculates a regular expression using quite familiar transformations for solving equational systems. Consider the DFA:



for which we can set up the following equational system

$$Q_0 = 1 + Q_0 b + Q_1 b + Q_2 b \quad (3)$$

$$Q_1 = Q_0 a \quad (4)$$

$$Q_2 = Q_1 a + Q_2 a \quad (5)$$

There is an equation for each node in the DFA. Let us have a look how the right-hand sides of the equations are constructed. First have a look at the second equation: the left-hand side is  $Q_1$  and the right-hand side  $Q_0 a$ . The right-hand side is essentially all possible ways how to end up in node  $Q_1$ . There is only one incoming edge from  $Q_0$  consuming an  $a$ . Therefore the right hand side is this state followed by character—in this case  $Q_0 a$ . Now let's have a look at the third equation: there are two incoming edges for  $Q_2$ . Therefore we have two terms, namely  $Q_1 a$  and  $Q_2 a$ . These terms are separated by  $+$ . The first states that if in state  $Q_1$  consuming an  $a$  will bring you to  $Q_2$ , and the second that being in  $Q_2$  and consuming an  $a$  will make you stay in  $Q_2$ . The right-hand side of the first equation is constructed similarly: there are three incoming edges, therefore there are three terms. There is one exception in that we also "add" 1 to the first equation, because it corresponds to the starting state in the DFA.

Having constructed the equational system, the question is how to solve it? Remarkably the rules are very similar to solving usual linear equational systems. For example the second equation does not contain the variable  $Q_1$  on the

right-hand side of the equation. We can therefore eliminate  $Q_1$  from the system by just substituting this equation into the other two. This gives

$$Q_0 = \mathbf{1} + Q_0 b + Q_0 a b + Q_2 b \quad (6)$$

$$Q_2 = Q_0 a a + Q_2 a \quad (7)$$

where in Equation (6) we have two occurrences of  $Q_0$ . Like the laws about  $+$  and  $\cdot$ , we can simplify Equation (6) to obtain the following two equations:

$$Q_0 = \mathbf{1} + Q_0 (b + a b) + Q_2 b \quad (8)$$

$$Q_2 = Q_0 a a + Q_2 a \quad (9)$$

Unfortunately we cannot make any more progress with substituting equations, because both (8) and (9) contain the variable on the left-hand side also on the right-hand side. Here we need to now use a law that is different from the usual laws about linear equations. It is called *Arden's rule*. It states that if an equation is of the form  $q = q r + s$  then it can be transformed to  $q = s r^*$ . Since we can assume  $+$  is symmetric, Equation (9) is of that form:  $s$  is  $Q_0 a a$  and  $r$  is  $a$ . That means we can transform (9) to obtain the two new equations

$$Q_0 = \mathbf{1} + Q_0 (b + a b) + Q_2 b \quad (10)$$

$$Q_2 = Q_0 a a (a^*) \quad (11)$$

Now again we can substitute the second equation into the first in order to eliminate the variable  $Q_2$ .

$$Q_0 = \mathbf{1} + Q_0 (b + a b) + Q_0 a a (a^*) b \quad (12)$$

Pulling  $Q_0$  out as a single factor gives:

$$Q_0 = \mathbf{1} + Q_0 (b + a b + a a (a^*) b) \quad (13)$$

This equation is again of the form so that we can apply Arden's rule ( $r$  is  $b + a b + a a (a^*) b$  and  $s$  is  $\mathbf{1}$ ). This gives as solution for  $Q_0$  the following regular expression:

$$Q_0 = \mathbf{1} (b + a b + a a (a^*) b)^* \quad (14)$$

Since this is a regular expression, we can simplify away the  $\mathbf{1}$  to obtain the slightly simpler regular expression

$$Q_0 = (b + a b + a a (a^*) b)^* \quad (15)$$

Now we can unwind this process and obtain the solutions for the other equations. This gives:

$$Q_0 = (b + a b + a a (a^*) b)^* \quad (16)$$

$$Q_1 = (b + a b + a a (a^*) b)^* a \quad (17)$$

$$Q_2 = (b + a b + a a (a^*) b)^* a a (a)^* \quad (18)$$

Finally, we only need to “add” up the equations which correspond to a terminal state. In our running example, this is just  $Q_2$ . Consequently, a regular expression that recognises the same language as the DFA is

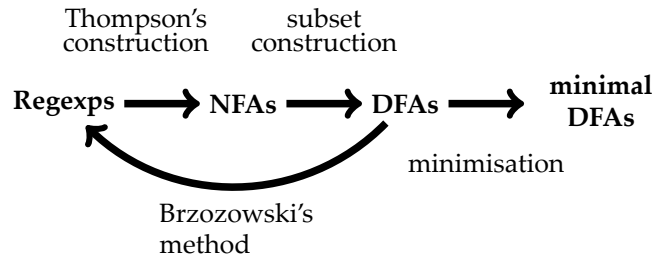
$$(b + a b + a a (a^*) b)^* a a (a)^*$$

You can somewhat crosscheck your solution by taking a string the regular expression can match and see whether it can be matched by the DFA. One string for example is *aaa* and *voila* this string is also matched by the automaton.

We should prove that Brzozowski’s method really produces an equivalent regular expression. But for the purposes of this module, we omit this. I guess you are relieved.

## Regular Languages

Given the constructions in the previous sections we obtain the following overall picture:



By going from regular expressions over NFAs to DFAs, we can always ensure that for every regular expression there exists a NFA and a DFA that can recognise the same language. Although we did not prove this fact. Similarly by going from DFAs to regular expressions, we can make sure for every DFA there exists a regular expression that can recognise the same language. Again we did not prove this fact.

The fundamental conclusion we can draw is that automata and regular expressions can recognise the same set of languages:

A language is *regular* iff there exists a regular expression that recognises all its strings.

or equivalently

A language is *regular* iff there exists an automaton that recognises all its strings.

Note that this is not a statement for a particular language (that is a particular set of strings), but about a large class of languages, namely the regular ones.

As a consequence for deciding whether a string is recognised by a regular expression, we could use our algorithm based on derivatives or NFAs or DFAs. But let us quickly look at what the differences mean in computational terms. Translating a regular expression into a NFA gives us an automaton that has  $O(n)$  states—that means the size of the NFA grows linearly with the size of the regular expression. The problem with NFAs is that the problem of deciding whether a string is accepted or not is computationally not cheap. Remember with NFAs we have potentially many next states even for the same input and also have the silent  $\epsilon$ -transitions. If we want to find a path from the starting state of a NFA to an accepting state, we need to consider all possibilities. In Ruby, Python and Java this is done by a depth-first search, which in turn means that if a “wrong” choice is made, the algorithm has to backtrack and thus explore all potential candidates. This is exactly the reason why Ruby, Python and Java are so slow for evil regular expressions. An alternative to the potentially slow depth-first search is to explore the search space in a breadth-first fashion, but this might incur a big memory penalty.

To avoid the problems with NFAs, we can translate them into DFAs. With DFAs the problem of deciding whether a string is recognised or not is much simpler, because in each state it is completely determined what the next state will be for a given input. So no search is needed. The problem with this is that the translation to DFAs can explode exponentially the number of states. Therefore when this route is taken, we definitely need to minimise the resulting DFAs in order to have an acceptable memory and runtime behaviour. But remember the subset construction in the worst case explodes the number of states by  $2^n$ . Effectively also the translation to DFAs can incur a big runtime penalty.<sup>2</sup>

But this does not mean that everything is bad with automata. Recall the problem of finding a regular expressions for the language that is *not* recognised by a regular expression. In our implementation we added explicitly such a regular expressions because they are useful for recognising comments. But in principle we did not need to. The argument for this is as follows: take a regular expression, translate it into a NFA and then a DFA that both recognise the same language. Once you have the DFA it is very easy to construct the automaton for the language not recognised by a DFA. If the DFA is completed (this is important!), then you just need to exchange the accepting and non-accepting states. You can then translate this DFA back into a regular expression and that will be

---

<sup>2</sup>Therefore the clever people in Rust try to *not* do such calculations upfront, but rather delay them and in this way can avoid much of the penalties...in many practical relevant places. As a result, they make the extraordinary claim that their time complexity is in the worst case  $O(m \times n)$  where  $m$  is proportional to the size of the regex and  $n$  is proportional to the size of strings. Does this claim hold water?



the regular expression that can match all strings the original regular expression could *not* match.

It is also interesting that not all languages are regular. The most well-known example of a language that is not regular consists of all the strings of the form

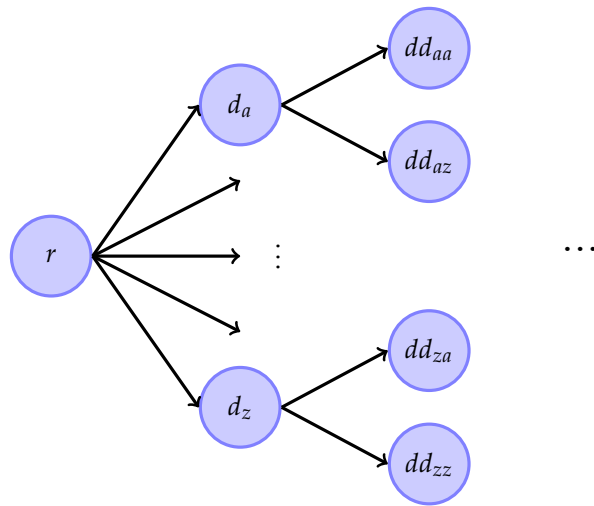
$$a^n b^n$$

meaning strings that have the same number of *as* and *bs*. You can try, but you cannot find a regular expression for this language and also not an automaton. One can actually prove that there is no regular expression nor automaton for this language, but again that would lead us too far afield for what we want to do in this module.

## Where Have Derivatives Gone?

By now you are probably fed up with this text. It is now way too long for one lecture, but there is still one aspect of the automata-regular-expression-connection I like to describe:

Where have the derivatives gone? Did we just forget them? Well, they actually do play a role of generating a DFA from a regular expression. And we can also see this in our implementation...because there is one flaw in our representation of automata and transitions as partial functions....remember I said something about fishy things. Namely, we can represent automata with infinite states, which is actually forbidden by the definition of what an automaton is. We can exploit this flaw as follows: Suppose our alphabet consists of the characters *a* to *z*. Then we can generate an “automaton” (it is not really one because it has infinitely many states) by taking as starting state the regular expression *r* for which we want to generate an automaton. There are *n* next-states which corresponds to the derivatives of *r* according to *a* to *z*. Implementing this in our slightly “flawed” representation is not too difficult. This will give a picture for the “automaton” looking something like this, except that it extends infinitely far to the right:



You might want to implement this “automaton”. What do you get?

While this makes all sense (modulo the flaw with the infinite states), does this automaton teach us anything new? The answer is no, because it boils down to just another implementation of the Brzozowski algorithm from Lecture 2. There *is* however something interesting in this construction which Brzozowski already cleverly found out, because there is a way to restrict the number of states to something finite. Meaning it would give us a real automaton. However, this would lead us far, far away from what we want discuss here. The end.