



COMP 181

Lecture 4
More scanning

September 14, 2006



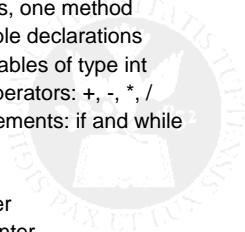
Programming assignment

- Big picture
 - Develop a compiler for *Decaf*, a subset of Java
 - 4 parts – we'll skip the lexer
 - PA1: an interpreter
 - PA2: parser
 - PA3: light version of semantic checker/code generator
 - PA4: full version of code generator
- Today: start PA1



Tufts University Computer Science

2



PA1: Interpreter

- *SkimDecaf*: subset of subset of Java
 - One class, one method
 - No variable declarations
 - Only variables of type int
 - Only 4 operators: +, -, *, /
 - Two statements: if and while
- Your Job:
 - Interpreter
 - Pretty printer




Tufts University Computer Science

3

Starter kit

- What I'll give you
 - Code skeleton
 - Parser: produces AST for *SkimDecaf* programs
- What you do:
 - Learn about the AST structure
 - Fill in code skeleton
 - Develop some test cases



Tufts University Computer Science

4



Environment

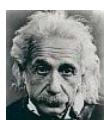
- JDK
 - Need access to java and javac
 - I have installed jdk 1.5 for x86/Linux:
`/g/181/2006f/jdk1.5.0_08`
- Eclipse
 - This assignment relies on the Eclipse parser
 - I have installed Eclipse for x86/Linux:
`/g/181/2006f/eclipse`
- Demo...




Tufts University Computer Science

5

Prelude



"God does not play dice with the universe."

What does this quote refer to?

- Newtonian physics: nature is deterministic
- Quantum mechanics: nature is non-deterministic
- Einstein didn't like this
- We will take Einstein's view...



Tufts University Computer Science

6

Scanner construction

[0] Define tokens as regular expressions

[1] Construct NFA for all REs

- Connect REs with ϵ transitions
- *Thompson's construction*

[2] Convert NFA into a DFA

- DFA is a simulation of NFA
- Possible much larger than NFA

[3] Minimize the DFA

- *Hopcroft's algorithm*

[4] Generate implementation



Tufts University Computer Science



7

[1] Thompson's construction

• **Goal:**

Systematically convert regular expressions for our language into a finite state automaton

• **Key idea:**

- FA "pattern" for each RE operator
- Start with atomic REs, build up a big NFA

• Idea due to Ken Thompson in 1968



Tufts University Computer Science



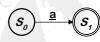
8

Thompson's construction

By induction on RE structure

• **Base case:**

Construct FA that recognizes atomic regular expressions:



• **Induction:**

Given FAs for two regular expressions, x and y , build a new FA that recognizes:

- xy
- $x|y$
- x^*



Tufts University Computer Science

9

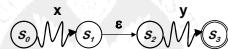


Thompson's construction

• Given:



• Build xy



• Why can't we do this?



Tufts University Computer Science



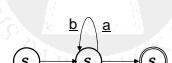
10

Need for ϵ transitions

• What if x and y look like this:

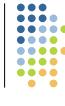


• Then xy ends up like this:



Tufts University Computer Science

11



Thompson's construction

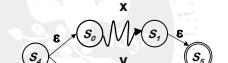
• Given:



• xy



• $x|y$



• x^*



Tufts University Computer Science

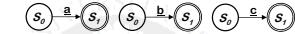


12

Example

Regular expression: $\underline{a} (\underline{b} \mid \underline{c})^*$

- \underline{a} , \underline{b} , & \underline{c}



- $\underline{b} \mid \underline{c}$



- $(\underline{b} \mid \underline{c})^*$

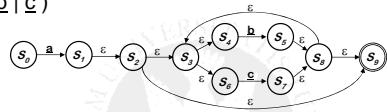


Tufts University Computer Science

13

Example

- $\underline{a} (\underline{b} \mid \underline{c})^*$



• Note: a human could design something simpler...



Tufts University Computer Science

14

Problem

- How to implement NFA scanner code?

- Non-determinism is a problem
- Explore all possible paths?

- Observation:

- We can build a DFA that simulates the NFA
- Accepts the same language
 - Explores all paths simultaneously



Tufts University Computer Science

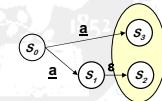
15

[2] NFA to DFA

- Subset construction

- Algorithm

- Intuition: each DFA state represents the possible states reachable after one input in the NFA



- Two key functions

- $\text{next}(s_i, a)$ – the set of states reachable from s_i on a
- $\epsilon\text{-closure}(s_i)$ – the set of states reachable from s_i on ϵ



Tufts University Computer Science

16

Subset construction

- Algorithm

Build a set of subsets of NFA states

```

 $s_0 \leftarrow \epsilon\text{-closure}(q_0)$ 
 $S \leftarrow \{s_0\}$ 
 $worklist \leftarrow \{s_0\}$ 
 $while (worklist is not empty)$ 
     $remove s from worklist$ 
     $for each \underline{a} \in \Sigma$ 
         $t \leftarrow \epsilon\text{-closure}(\text{next}(s, a))$ 
         $if (t \notin S)$ 
             $add t to S$ 
             $add t to worklist$ 
             $add transition (s, a, t)$ 

```

- Start with the ϵ -closure over the initial state
- Initialize the worklist to this one subset
- While there are more subsets on the worklist, remove the next subset
- Apply each input symbol to the set of NFA states to produce a new set.
- If we haven't seen that subset before, add it to S and the worklist, and record the set-to-set transition

17



Tufts University Computer Science

NFA to DFA example

$\underline{a} (\underline{b} \mid \underline{c})^*$:



$\epsilon\text{-closure}(\text{next}(s, a))$

Subsets S	NFA states	\underline{a}	\underline{b}	\underline{c}
s_0	q_0	q_1, q_2, q_3 q_4, q_5, q_6	none	none
s_1	q_1, q_2, q_3 q_4, q_5, q_6	none	q_5, q_6, q_7 q_3, q_4, q_6	q_7, q_8, q_9 q_3, q_4, q_6
s_2	q_5, q_6, q_7 q_3, q_4, q_6	none	(also s_2)	(also s_3)
s_3	q_7, q_8, q_9 q_3, q_4, q_6	none	(also s_2)	(also s_3)

Tufts University Computer Science

Accepting states

18

NFA to DFA example

- Convert each subset in S into a state:



- All transitions are deterministic
- Use same code engine to execute
- Smaller than NFA, but still bigger than necessary



Tufts University Computer Science

19

Does it work?

- The algorithm halts
 - S contains no duplicate subsets
 - $2^{|NFA|}$ is finite
 - Main loop adds to S , but does not remove
It is a monotone function
- S contains all the reachable NFA states
Tries all input symbols, builds all NFA configurations
- Note:** important class of compiler algorithms
 - Fixpoint** computation
 - Monotonic update function
 - Convergence is guaranteed



Tufts University Computer Science

20

[3] DFA minimization

- Hopcroft's algorithm
 - Discover sets of *equivalent* states in DFA
 - Represent each set with a single state
- Two states are equivalent iff:
 - The set of paths leading to them are the same
 - For all input symbols, transitions lead to equivalent states
- This is the key to the algorithm



Tufts University Computer Science

21

DFA minimization

- A partition P of the states S
 - Each $s \in S$ is in exactly one set $p_i \in P$
 - Idea:
If two states s and t transition to different partitions, then they must be in different partitions
- Algorithm: iteratively partition the DFA's states
 - Group states into maximal size sets, *optimistically*
 - Iteratively subdivide those sets, as needed
 - States that remain grouped together are equivalent



Tufts University Computer Science

22

DFA minimization

- Details:
 - Given DFA $(S, \Sigma, \delta, s_0, F)$
 - Initial partition: $P_0 = \{F, S-F\}$
 - Intuition: final states are always different
- Splitting a set around symbol a
 - Assume $s_a \in p_i$ and $\delta(s_a, a) = s_x$ & $\delta(s_b, a) = s_y$
 - Split p_i if:
 - If s_x & s_y are not in the same set
 - If s_a has a transition on a , but s_b does not
 - Intuition: one state in DFA cannot have two transitions on a



Tufts University Computer Science

23

Examples



Tufts University Computer Science

24

DFA minimization algorithm

```

 $P \leftarrow \{F, \{Q-F\}\}$ 
while (P is still changing)
   $T \leftarrow \{\}$ 
  for each set  $S \in P$ 
    for each  $a \in \Sigma$ 
      partition S by  $a$ 
      into  $S_1$  and  $S_2$ 
     $T \leftarrow T \cup S_1 \cup S_2$ 
  if  $T \neq P$  then
     $P \leftarrow T$ 
  
```

This is a fixed-point algorithm!



Tufts University Computer Science



25

Does it work?

- Algorithm halts

- Partition $P \in 2^S$
- Start off with 2 subsets of $S - \{F\}$ and $\{S-F\}$
- While loop takes $P_i \rightarrow P_{i+1}$ by splitting 1 or more sets
- P_{i+1} is at least one step closer to partition with $|S|$ sets
- Maximum of $|S|$ splits

- Note that

- Partitions are never combined
- Initial partition ensures that final states are intact



Tufts University Computer Science

26

DFA minimization

Refining the algorithm

- As written, it examines every $S \in P$ on each iteration
 - This does a lot of unnecessary work
 - Only need to examine S if some T , reachable from S , has been split
- Reformulate the algorithm using a “worklist”
 - Start worklist with initial partition, F and $\{Q-F\}$
 - When it splits S into S_1 and S_2 , place S_2 on worklist

→ This version looks at each $S \in P$ many fewer times
Well-known, widely used algorithm due to John Hopcroft



Tufts University Computer Science

27

Hopcroft's algorithm

```

 $W \leftarrow \{F, Q-F\}; P \leftarrow \{F, \{Q-F\}\}; // W is the worklist, P the current partition$ 
while (  $W$  is not empty ) do begin
  select and remove  $S$  from  $W$ ; // S is a set of states
  for each  $\alpha \in \Sigma$  do begin
    let  $I_\alpha \leftarrow \delta_\alpha^{-1}(S)$ ; //  $I_\alpha$  is set of all states that can reach  $S$  on  $\alpha$ 
    for each  $R$  in  $P$  such that  $R \cap I_\alpha$  is not empty
      and  $R$  is not contained in  $I_\alpha$  do begin
        partition  $R$  into  $R_1$  and  $R_2$  such that  $R_1 \leftarrow R \cap I_\alpha; R_2 \leftarrow R - R_1$ ;
        replace  $R$  in  $P$  with  $R_1$  and  $R_2$ ;
        if  $R \in W$  then replace  $R$  with  $R_1$  in  $W$  and add  $R_2$  to  $W$ ;
        else if  $|R_1| \leq |R_2|$ 
          then add  $R_1$  to  $W$ ;
        else add  $R_2$  to  $W$ ;
      end
    end
  end
end

```

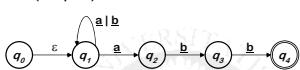


Tufts University Computer Science

28

Full example

- Consider $(a \mid b)^* abb$



- Applying the subset construction:

Iter.	State	Contains	ϵ -closure(move(s, a))	ϵ -closure(move(s, b))
0	S_0	q_0, q_1	q_1, q_2	q_1
1	S_1	q_1, q_2	q_1, q_2	q_1, q_3
2	S_2	q_1	q_1, q_2	q_1
3	S_3	q_1, q_3	q_1, q_2	q_1, q_3
	S_4	$q_1, q_3 \leftarrow$	q_1, q_2	q_1

contains q_4 (final state)

- Iteration 3 adds nothing to S , so the algorithm halts

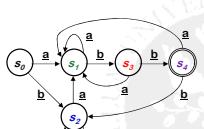


Tufts University Computer Science

29

Full example

- DFA for $(a \mid b)^* abb$



δ	a	b
S_0	S_1	S_2
S_1	S_1	S_3
S_2	S_1	S_2
S_3	S_1	S_4
S_4	S_1	S_2

- About the same size as NFA

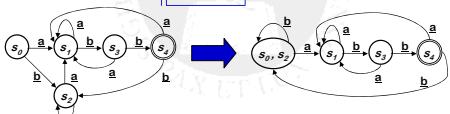


Tufts University Computer Science

30

Full example – minimization

	Current Partition	Worklist	s	Split on a	Split on b
P_0	{ s_4 } { s_0, s_1, s_2, s_3 }	{ s_4 }		none	{ s_0, s_1, s_2 } { s_3 }
P_1	{ s_3 } { s_2 } { s_0, s_1, s_2 }	{ s_0, s_1, s_2 } { s_3 }		none	{ s_0, s_2 } { s_1 }
P_2	{ s_4 , s_3 } { s_1 } { s_0, s_2 }	{ s_0, s_2 } { s_1 }		none	none

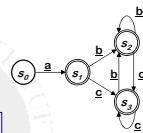


31

Another example

- What about $a \ (b \mid c)^*$

	Current Partition	a	b	c
P_0	{ s_0, s_2 , s_3 } { s_1 }	none	none	none



Previously, we observed that a human would design a simpler automaton than Thompson's construction & the subset construction did.
Minimizing that DFA produces the one that a human would design!

Tufts University Computer Science

32

We're done with the theory

- Take a nice, deep breath...



Tufts University Computer Science

33

Application

- What do our DFAs recognize?

Given a string, is it a token

- Problem:

Input consists of many tokens, not just one

`elsex = 0;`

- Option 1: "else", "x", "=", "0"
- Option 2: "elsex", "=", "0"

Tufts University Computer Science

34

Practical scanners

- Solution:

- List of regular expressions
- Longest matching RE wins
- Ties are broken by order of Res

- How do we implement that?

Lookahead

- Keep reading characters as long as the DFA is satisfied
- If no transition and we're in a final state for one of the REs:
 - Return that token
 - Go back to start state
- If no transition and we're **not** in a final state → error



Tufts University Computer Science

35

C scanner

Declarations

```
%{
#include "parser.tab.h"
%}
```

Short-hand

```
identifier ((a-zA-Z_)[0-9a-zA-Z_]*)
octal_escape ((0-7[^\\\"\\n])*)
any_white ([\\t\\v\\n\\r\\f\\b\\s])
%%
```

REs and actions

```
(any_white)+ {
    for { lval.tok = get_pos(); return ctokFOR; }
    if { lval.tok = get_pos(); return ctokIF; }
    { identifier } { lval.tok = get_pos();
        lval.idN = new idNode(cbtext, cbval.tok);
        if (is_typename(cbtext)) return TYPEDEFname;
        else return IDENTIFIER; }
    { decimal_constant } { lval.exprN = atoi(cbtext);
        return INTEGERconstant; }
}
%% ...any special code...
```

Tufts University Computer Science

36

Implementation

- Table driven
 - Read and classify character
 - Select action
 - Find the next state, assign to state variable
 - Repeat
- Alternative: direct coding
 - Each state is a chunk of code
 - Transitions test and branch directly
 - Very ugly code – but who cares?
 - Very efficient

This is how lex/flex work: states are encoded as cases in a giant switch statement

37



Tufts University Computer Science



Building scanners

- The point
 - Theory lets us automate construction
 - Implementer writes down regular expressions
 - Generator does: RE \rightarrow NFA \rightarrow DFA \rightarrow code
 - Reliably produces fast, robust scanners
- Works for most modern languages

Think twice about language features that defeat the DFA-based scanners



Tufts University Computer Science

38



Next time...

- Grammars and parsing
- Project 1:
 - SkimDecaf Interpreter
 - I'll post everything after class



Tufts University Computer Science

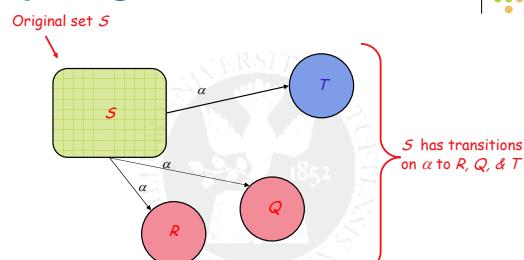
39



40



Splitting S around α



The algorithm partitions S around α

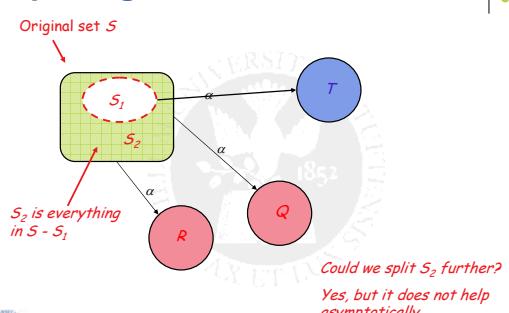
41



Tufts University Computer Science



Splitting S around α



42



Tufts University Computer Science

