

present:



Wednesday, October 1, 6:00-7:30pm Gore Hall Room 103

Bring your resume for a chance to win great prizes!! FREE FOOD will be provided!



Lexical Analysis: DFA Minimization & Wrap Up

Automating Scanner Construction

PREVIOUSLY

RE→NFA (Thompson's construction)

- Build an NFA for each term
- Combine them with ϵ -moves

NFA \rightarrow DFA (subset construction)

Build the simulation

TODAY

Some Project Related Material

 $DFA \rightarrow Minimal DFA$

Hopcroft's algorithm







- Report errors for lexicographically malformed inputs
 reject illegal characters, or meaningless character sequences
- Return an abstract representation of the code
 - \rightarrow character sequences (e.g., "if" or "loop") turned into tokens.
- Resulting sequence of tokens will be used by the parser
- Makes the design of the parser a lot easier.



- A scanner specification (e.g., for JLex), is list of (typically short) regular expressions.
- Each regular expressions has an action associated with it.
- Typically, an action is to return a token.
- On a given input string, the scanner will:
 - → find the longest prefix of the input string, that matches one of the regular expressions.
 - → will execute the action associated with the matching regular expression highest in the list.
- Scanner repeats this procedure for the remaining input.
- If no match can be found at some point, an error is reported.



- Consider the following scanner specification.
 - 1. aaa { return T1 }
 - 2. a*b { return T2 }
 - 3. b { return S }
- Given the following input string into the scanner aaabbaaa

the scanner as specified above would output

T2 T2 T1

- Note that the scanner will report an error for example on the string 'aa'.
- Rule 3 is redundant!

Special Return Tokens



- Sometimes one wants to extract information out of what prefix of the input was matched.
- Example:

"[a-zA-Z0-9]*" { return STRING(yytext()) }

- Above RE matches every string that
 - \rightarrow starts and ends with quotes, and
 - \rightarrow has any number of alpha-numerical chars between them.
- Associated action returns a string token, which is the exact string that the RE matched.
- Note that yytext() will also include the quotes.
- Note this RE doesn't handle escaped characters, special characters (e.g., punctuation characters), etc.

Quiz



- Consider the following scanner specification.
 - 1. bab { return T1 }
 - 2. ba⁺ { return T2 }
 - 3. b $\{ return S1 \}$
 - 4. a { return S2 }
- Given the following input string into the scanner baabbabaa

the scanner as specified above would output

What?



- Consider the following scanner specification.
 - 1. bab { return T1 }
 - 2. ba⁺ { return T2 }
 - 3. b $\{ return S1 \}$
 - 4. a { return S2 }
- Given the following input string into the scanner baabbabaa

the scanner as specified above would output

T2 51 T1 52 52

ELAWARE 17 4 3 *

Details of the algorithm

- Group states into maximal size sets, optimistically
- Iteratively subdivide those sets, as needed
- States that remain grouped together are equivalent

Initial partition, P_0 , has two sets: $\{D_F\}$ & $\{D-D_F\}$ (DFA = $(Q, \Sigma, \delta, q_0, F)$)

Splitting a set ("partitioning a set by \underline{a} ")

- Assume q_i , & $q_j \in p$, and $\delta(q_i,\underline{a}) = q_x$, & $\delta(q_j,\underline{a}) = q_y$
- If $q_x \& q_y$ are not in the same set, then p must be split $\rightarrow q_i$ has transition on a, q_j does not $\Rightarrow \underline{a}$ splits p
- One state in the final DFA cannot have two transitions on <u>a</u>

DFA Minimization

The algorithm

 $P \leftarrow \{ D_{F'} \{ D - D_F \} \}$ while (P is still changing) $T \leftarrow \emptyset$ for each set $p \in P$ $T \leftarrow T \cup Split(p)$ $P \leftarrow T$ Split(S) for each $\alpha \in \Sigma$ if α splits S into s_1 and s_2 then return $\{s_1, s_2\}$ return S

Why does this work?

- Partition $P \in 2^{D}$
- Starts with 2 subsets of D
 {D_F} and {D-D_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 the partition with |D| sets
- Maximum of |D| splits

Note that

- Partitions are <u>never</u> combined
- Initial partition ensures that final states are intact

This is a fixed-point algorithm!





The algorithm partitions S around α







First, the subset construction:

		ϵ -closure(Delta(s,*))		
	NFA states	<u>a</u>	<u>b</u>	<u>c</u>
S ₀	<i>q</i> ₀	<mark>q</mark> 1, q2, q3, q4, q6, q9	none	none
S ₁	$q_1, q_2, q_3, q_4, q_6, q_9$	none	q 5, q 8, q 9, q 3, q 4, q 6	q ₇ , q ₈ , q ₉ , q ₃ , q ₄ , q ₆
S ₂	$q_{5}, q_{8}, q_{9}, q_{3}, q_{4}, q_{6}$	none	S ₂	S ₃
S 3	<i>q</i> ₇ , <i>q</i> ₈ , <i>q</i> ₉ , <i>q</i> ₃ , <i>q</i> ₄ , <i>q</i> ₆	none	S ₂	S ₃
Final states				





To produce the minimal DFA



In lecture 4, 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!

Abbreviated Register Specification

Start with a regular expression

r0 | r1 | r2 | r3 | r4 | r5 | r6 | r7 | r8 | r9



RE→NFA → DFA → DFA



Abbreviated Register Specification

The subset construction builds



This is a DFA, but it has a lot of states ...

minimal →NFA DFA DFA



Abbreviated Register Specification

The DFA minimization algorithm builds



This looks like what a skilled compiler writer would do!

minimal DFA ►RE →NFA →DFA -



Limits of Regular Languages

Advantages of Regular Expressions

- Simple & powerful notation for specifying patterns
- Automatic construction of fast recognizers
- Many kinds of syntax can be specified with REs

Example — an expression grammar

 $Term \rightarrow [a-zA-Z] ([a-zA-z] | [0-9])^*$ $Op \rightarrow \pm | - | \pm | /$ $Expr \rightarrow (Term Op)^* Term$

Of course, this would generate a DFA ...

If REs are so useful ...

Why not use them for everything?



Limits of Regular Languages

Not all languages are regular $RL's \subset CFL's \subset CSL's$

You cannot construct DFA's to recognize these languages

• $L = \{ p^k q^k \}$

This is not a regular language

(nor an RE)

(parenthesis languages)

But, this is a little subtle. You <u>can</u> construct DFA's for

- Strings with alternating 0's and 1's $(\epsilon \mid 1)(01)^*(\epsilon \mid 0)$
- Strings with an even number of 0's and 1's

RE's can count bounded sets and bounded differences



Poor language design can complicate scanning

- Reserved words are important
 if then then then = else; else else = then
- Insignificant blanks
 do 10 i = 1,25
 do 10 i = 1.25

ELAWARE 1743

(Fortran & Algol68)

(PL/I)

- String constants with special characters (C, C++, Java, ...) newline, tab, quote, comment delimiters, ...
- Finite closures
 - → Limited identifier length
 - → Adds states to count length

(Fortran 66 & Basic)

Building Faster Scanners from the DFA

Table-driven recognizers waste effort

- Read (& classify) the next character
- Find the next state
- Assign to the state variable
- Trip through case logic in *action()*
- Branch back to the top

We can do better

- Encode state & actions in the code
- Do transition tests locally
- Generate ugly, spaghetti-like code
- Takes (many) fewer operations per input character



char ← next character; state ← s_{0;} call action(state,char); while (char ≠ <u>eof</u>) state ← δ(state,char); call action(state,char); char ← next character;

if T(state) = <u>final</u> then report acceptance; else report failure;

Building Faster Scanners from the DFA

A direct-coded recognizer for \underline{r} Digit Digit*

- Many fewer operations per character
- Almost no memory operations
- Even faster with careful use of fall-through cases

```
goto s_{0};
s_{0}: word \leftarrow \emptyset;
char \leftarrow next character;
if (char = 'r')
then goto s_{1};
else goto s_{e};
s_{1}: word \leftarrow word + char;
char \leftarrow next character;
if ('0' \leq char \leq '9')
then goto s_{2};
else goto s_{e};
```

 $s2: word \leftarrow word + char;$ $char \leftarrow next character;$ $if ('0' \leq char \leq '9')$ $then goto s_{2};$ else if (char = eof) then report success; $else goto s_{e};$ $s_{e}: print error message;$ return failure;



Building Faster Scanners

Hashing keywords versus encoding them directly

- Some (well-known) compilers recognize keywords as identifiers and check them in a hash table
- Encoding keywords in the DFA is a better idea
 - \rightarrow O(1) cost per transition
 - → Avoids hash lookup on each identifier

It is hard to beat a well-implemented DFA scanner



The point



- Implementer writes down the regular expressions
- Scanner generator builds NFA, DFA, minimal DFA, and then writes out the (table-driven or direct-coded) code
- This reliably produces fast, robust scanners

For most modern language features, this works

- You should think twice before introducing a feature that defeats a DFA-based scanner
- The ones we've seen (e.g., insignificant blanks, non-reserved keywords) have not proven particularly useful or long lasting

