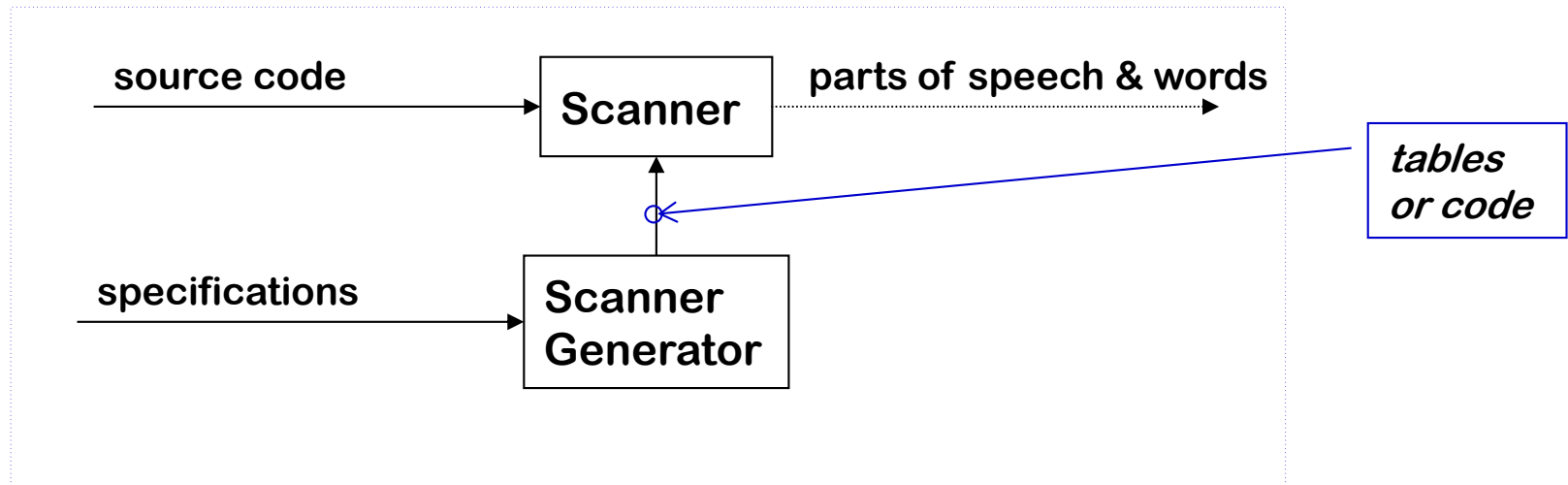# CS415 Compilers

# Lexical Analysis
## Part 3

These slides are based on slides copyrighted by
Keith Cooper, Ken Kennedy & Linda Torczon at Rice
University

- Homework solutions for homeworks 1 and 2 have been posted on canvas under "Files" tab

- Third homework will be posted after exam.

- First project (local instruction scheduler) has been posted
    Deadline for code: March 2
    Deadline for report: March 4

- First midterm: This Wednesday, February 23
    In class exam, 60 minutes,
    Topics: ILOC, instruction scheduling, register allocation

source code → **Scanner** ⋯ parts of speech & words ⋯→

tables or code

specifications → **Scanner Generator**

→ The scanner is the first stage in the front end

→ Specifications can be expressed using regular expressions

→ Build tables and code from a DFA

- We will show how to construct a finite state automaton to recognize any RE

- Overview:

  → Direct construction of a nondeterministic finite automaton (NFA) to recognize a given RE

    ▪ Requires $\varepsilon$-transitions to combine regular subexpressions

  → Construct a deterministic finite automaton (DFA) to simulate the NFA

    ▪ Use a set-of-states construction

  → Minimize the number of states

    ▪ Hopcroft state minimization algorithm

  → Generate the scanner code

    ▪ Additional specifications needed for details

- All strings of 1s and 0s ending in a <u>1</u>

  ( <u>0</u> | <u>1</u> )<sup>*</sup> <u>1</u>

- All strings over lowercase letters where the vowels (a,e,i,o, & u) occur exactly once, in ascending order

  *Cons* → (<u>b</u>|<u>c</u>|<u>d</u>|<u>f</u>|<u>g</u>|<u>h</u>|<u>j</u>|<u>k</u>|<u>l</u>|<u>m</u>|<u>n</u>|<u>p</u>|<u>q</u>|<u>r</u>|<u>s</u>|<u>t</u>|<u>v</u>|<u>w</u>|<u>x</u>|<u>y</u>|<u>z</u>)

- All strings of <u>1</u>s and <u>0</u>s that do not contain three <u>0</u>s in a row:

# More Regular Expressions

- All strings of 1s and 0s ending in a <u>1</u>

  ( <u>0</u> | <u>1</u> )$^*$ <u>1</u>

- All strings over lowercase letters where the vowels (a,e,i,o, & u) occur exactly once, in ascending order

  *Cons* → (<u>b|c|d|f|g|h|j|k|l|m|n|p|q|r|s|t|v|w|x|y|z</u>)
  *Cons*$^*$ <u>a</u> *Cons*$^*$ <u>e</u> *Cons*$^*$ <u>i</u> *Cons*$^*$ <u>o</u> *Cons*$^*$ <u>u</u> *Cons*$^*$

- All strings of <u>1</u>s and <u>0</u>s that do not contain three <u>0</u>s in a row:

- All strings of 1s and 0s ending in a <u>1</u>

  $( \underline{0} \mid \underline{1} )^* \underline{1}$

- All strings over lowercase letters where the vowels (a,e,i,o, & u) occur exactly once, in ascending order

  *Cons* $\rightarrow$ (b|c|d|f|g|h|j|k|l|m|n|p|q|r|s|t|v|w|x|y|z)
  *Cons*$^*$ <u>a</u> *Cons*$^*$ <u>e</u> *Cons*$^*$ <u>i</u> *Cons*$^*$ <u>o</u> *Cons*$^*$ <u>u</u> *Cons*$^*$
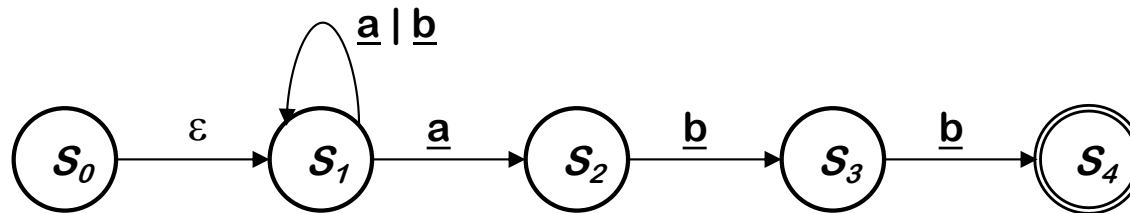
- All strings of <u>1</u>s and <u>0</u>s that do not contain three <u>0</u>s in a row:

  $( \underline{1}^* ( \varepsilon \mid \underline{01} \mid \underline{001} ) \underline{1}^* )^* ( \varepsilon \mid \underline{0} \mid \underline{00} )$

Each RE corresponds to a *deterministic finite automaton* (DFA)
- May be hard to directly construct the right DFA

What about an RE such as ( <u>a</u> | <u>b</u> )* <u>abb</u> ?



This is a little different

- $S_0$ has a transition on $\varepsilon$
- $S_1$ has two transitions on <u>a</u>

This is a *non-deterministic finite automaton* (NFA)

- An NFA accepts a string $x$ iff $\exists$ a path though the transition graph from $s_0$ to a final state such that the edge labels spell $x$

- Transitions on $\varepsilon$ consume no input

- To "run" the NFA, start in $s_0$ and *guess* the right transition at each step
  - → Always guess correctly
  - → If some sequence of correct guesses accepts x then accept

Why study NFAs?

- They are the key to automating the RE→DFA construction

- We can paste together NFAs with $\varepsilon$-transitions

DFA is a special case of an NFA

- DFA has no $\varepsilon$ transitions
- DFA's transition function is single-valued
- Same rules will work

DFA can be simulated with an NFA
  → *Obviously*

NFA can be simulated with a DFA *(less obvious)*
- Simulate sets of possible states
- Possible exponential blowup in the state space
- Still, one state transition per character in the input stream

To convert a specification into code:

1 Write down the RE for the input language
2 Build a big NFA
3 Build the DFA that simulates the NFA
4 Systematically shrink the DFA
5 Turn it into code

Scanner generators

- Lex and Flex work along these lines
- Algorithms are well-known and well-understood
- Key issue is interface to parser
- You could build one in a weekend!

RE→ NFA  *(Thompson's construction)*
- Build an NFA for each term

- Combine them with ε-moves

NFA → DFA *(subset construction)*
- Build the simulation

DFA → Minimal DFA
- Hopcroft's algorithm
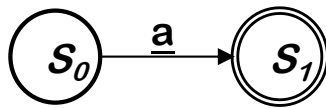
DFA →RE *(Not part of the scanner construction)*
- All pairs, all paths problem
- Take the union of all paths from $s_0$ to an accepting state

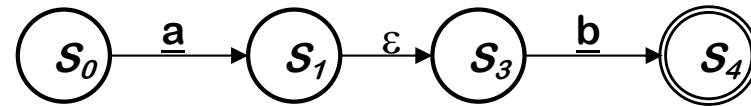*The Cycle of Constructions*

RE → NFA → DFA → *minimal* DFA

Key idea
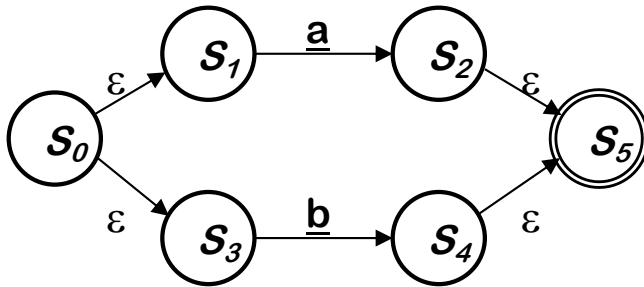
- NFA pattern for each symbol and each operator
- Each NFA has a single start and accept state
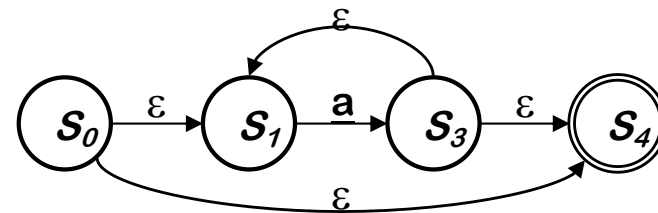- Join them with ε moves in precedence order



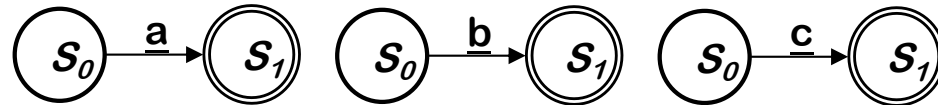NFA for <u>a</u>



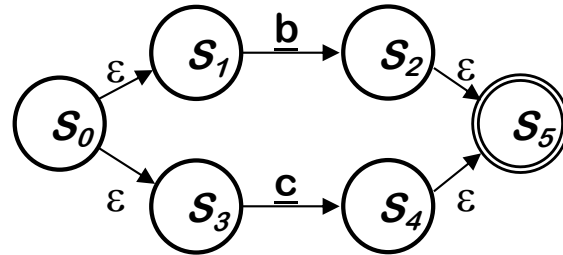NFA for <u>ab</u>



NFA for <u>a</u> | <u>b</u>



NFA for <u>a</u>$^*$
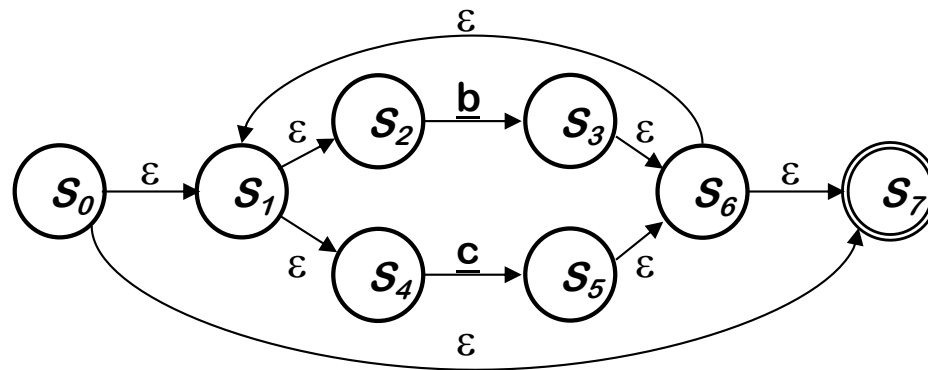
Ken Thompson, CACM, 1968

Let's try a ( b | c )*
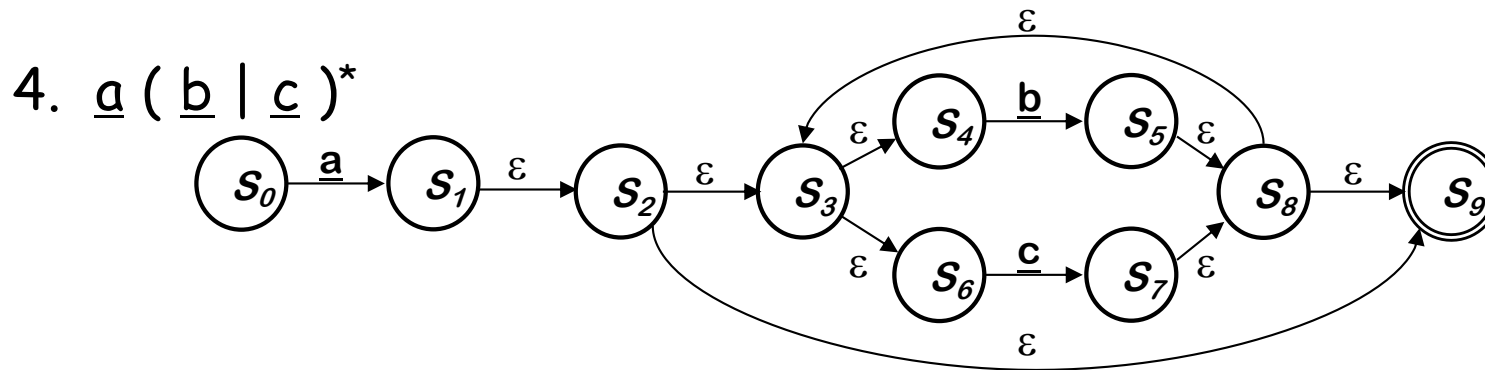
1. a, b, & c

$$S_0 \xrightarrow{a} S_1 \qquad S_0 \xrightarrow{b} S_1 \qquad S_0 \xrightarrow{c} S_1$$

2. b | c



3. ( b | c )*

4. $\underline{a} \, ( \, \underline{b} \mid \underline{c} \, )^{*}$



Of course, a human would design something simpler ...



> But, we can automate production of the more complex one ...

Need to build a simulation of the NFA

*Note: $s_i$ are sets of states of the NFA, which together constitute a single state in the simulating DFA*

Two key functions

- *move($s_i$, $\underline{a}$)*    is set of states reachable from $s_i$ by $\underline{a}$

- *ε-closure($s_i$)*   is set of states reachable from $s_i$ by *ε*

The algorithm (sketch):
- Start state derived from $s_0$ of the NFA
- Take its ε-closure $S_0$ = ε-closure($s_0$)
- For each state S, compute move(S, a) for each a ∈ Σ, and take its ε-closure
- Iterate until no more states are added

*Sounds more complex than it is...*

**The algorithm:**

$s_0 \leftarrow \varepsilon\text{-closure}(\{q_0\})$

*add $s_0$ to $S$*

*while ( $S$ is still changing )*

  *for each $s_i \in S$*

    *for each $a \in \Sigma$*

      $s_? \leftarrow \varepsilon\text{-closure}(move(s_i, a))$

      *if ( $s_? \notin S$ ) then*

        *add $s_?$ to $S$ as $s_j$*

        $T[s_i, a] \leftarrow s_j$

      *else*

        $T[s_i, a] \leftarrow s_?$

*Let's think about why this works*

**The algorithm halts:**

1. *$S$ contains no duplicates (test before adding)*

2. *$2^Q$ is finite*

3. *while loop adds to $S$, but does not remove from $S$ (monotone)*

$\Rightarrow$ the loop halts

**$S$ contains all the reachable NFA states**

*It tries each symbol in each $s_i$.*

*It builds every possible NFA configuration.*

$\Rightarrow$ *$S$ and $T$ form the DFA*

Example of a *fixed-point* computation

- Monotone construction of some finite set
- Halts when it stops adding to the set
- Proofs of halting & correctness are similar
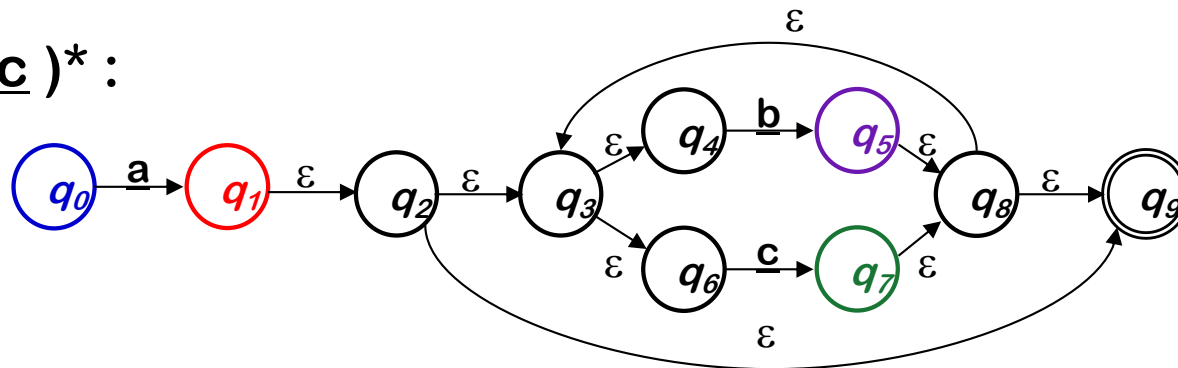- These computations arise in many contexts

Other fixed-point computations

- Canonical construction of sets of LR(1) items
  → Quite similar to the subset construction
- Classic data-flow analysis
  → Solving sets of simultaneous set equations
- DFA minimization algorithm (coming up!)

*We will see many more fixed-point computations*

**a ( b | c )\* :**



Applying the subset construction:

**a ( b | c )\* :**



Applying the subset construction:

|  | NFA states | ε-closure(move(s,*)) | | |
|---|---|---|---|---|
|  |  | **a** | **b** | **c** |
| $s_0$ | $q_0$ | $q_1, q_2, q_3, q_4, q_6, q_9$ | 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 | $s_2$ | $s_3$ |
| $s_3$ | $q_7, q_8, q_9, q_3, q_4, q_6$ | none | $s_2$ | $s_3$ |

**Final states**

The DFA for <u>a</u> ( <u>b</u> | <u>c</u> )*



| δ | <u>a</u> | <u>b</u> | <u>c</u> |
|-------|-------|-------|-------|
| $s_0$ | $s_1$ | - | - |
| $s_1$ | - | $s_2$ | $s_3$ |
| $s_2$ | - | $s_2$ | $s_3$ |
| $s_3$ | - | $s_2$ | $s_3$ |

- Ends up smaller than the NFA
- All transitions are deterministic

**More Lexical Analysis**

**Syntax Analysis (top-down parsing)**

Read EaC: Chapter 3.1 - 3.3