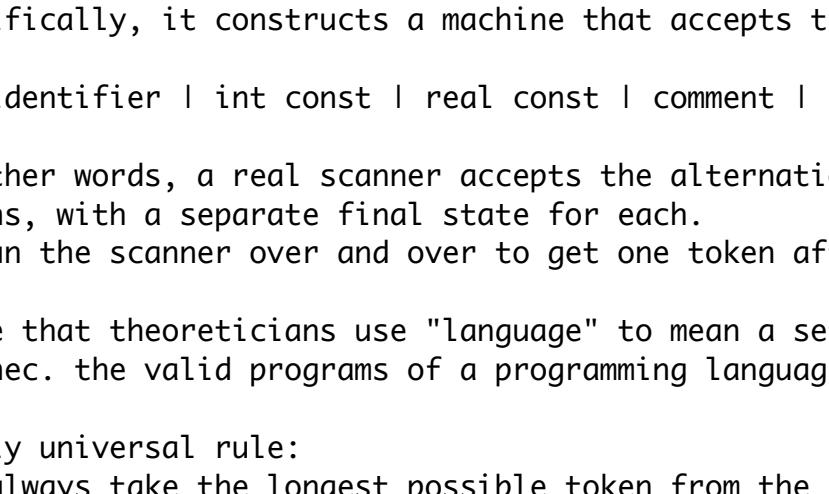


=====

SCANNING

Scanner is responsible for
 tokenizing source
 removing comments
 saving text of identifiers, numbers, strings
 saving source locations (file, line, column) for error messages

a DFA for identifiers:



Can be built by hand (ad hoc) or automatically from regular expressions (REs)
 ad-hoc generally yields the fastest, most compact code
 by doing lots of special-purpose things
 automatically-generated scanners can come close, though
 and are easy to develop and change

A scanner generator builds a DFA automatically from a set of REs
 Specifically, it constructs a machine that accepts the "language"

identifier | int const | real const | comment | symbol | ...

In other words, a real scanner accepts the alternation of a language's tokens, with a separate final state for each.

We run the scanner over and over to get one token after another.

(Note that theoreticians use "language" to mean a set of strings -- not nec. the valid programs of a programming language.)

Nearly universal rule:

always take the longest possible token from the input
 thus foobar is foobar and never f or foo or foob
 more to the point, 3.14159 is a real const and never 3, ., and 14159

An RE generates a regular language; a DFA _recognizes_ it.

The standard Unix lex (flex) outputs C code.

Some other tools produce numeric tables that are read by a separate driver.

The table is the transition function

two-dimensional array indexed by current state and input character
 entries specify
 next state
 whether to keep scanning, return a token, or announce an error

Longest-possible token rule means we return only when the next character can't be used to continue the current token.

That next character must generally be saved for the next token.

In some cases you may need to peek at more than one character of lookahead in order to know whether to proceed.

In Pascal, when you had a 3 and you saw a dot, did you proceed (in hopes of getting 3.14) or did you stop (in fear of getting 3..5)?

In messier cases, you may not be able to get by with any fixed amount of lookahead. In Fortran IV (c. 1962), for example, one had

DO 5 I = 1,25 loop
 DO 5 I = 1.25 assignment
 DO 5,I = 1,25 alternate syntax for loop, f77

For most languages it suffices to remember we were in a potentially final state, and save enough information that we can back up to it if we get stuck later.

For some languages (famously, Fortran), that isn't enough.

Sometimes need semantic information in order to scan (yuck).

Building a scanner from regular expressions

multi-step process

(1) write REs by hand, including for whitespace and comments, but with identifiers and reserve words (keywords) combined

(2) build NFA from REs

(3) build DFA from NFA

(4) minimize DFA

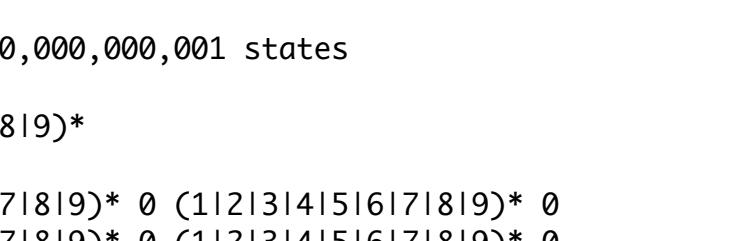
(5) add extra logic to
 implement the longest-possible-token rule, with backup
 discard white space and comments (i.e., start over when you
 realize that's what you found)
 distinguish reserve words from identifiers
 save text of "interesting" tokens
 tag returned tokens with location and text
 return an extra \$\$ token at end-of-file

Step (2) is inductive. It starts with a trivial DFA to accept a single character:

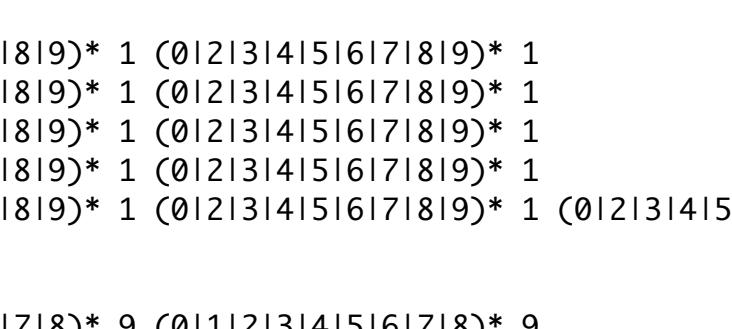


Note that this machine has a single start state, a single final state, no transitions into the start state, and no transitions out of the final state. We'll maintain these invariants in three inductive steps:

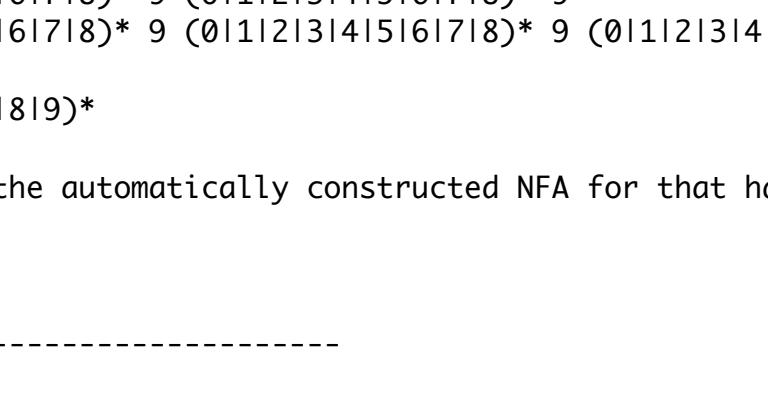
Concatenation (A B):



Alternation (A | B):



Kleene closure (A*):



Step (3) uses what's called a "set of subsets" construction.

Step (4) divides the states of an initial DFA into progressively finer equivalence classes, until it can prove that additional refinement makes no difference.

Example 1 (in the book): real numbers (no exponential notation)

RN = d*(.d|d.)d*

14-state NFA results from construction

5-state subset DFA

4-state minimal DFA

Example 2: character strings with optional backslash-escaped quotes

S = " ([^"] | \a) * " "a" for anything

11-state NFA results from construction

6-state subset DFA

4-state minimal DFA

Step (3) uses what's called a "set of subsets" construction.

Step (4) divides the states of an initial DFA into progressively finer equivalence classes, until it can prove that additional refinement makes no difference.

Example 1 (in the book): real numbers (no exponential notation)

RN = d*(.d|d.)d*

14-state NFA results from construction

5-state subset DFA

4-state minimal DFA

Example 2: character strings with optional backslash-escaped quotes

S = " ([^"] | \a) * " "a" for anything

11-state NFA results from construction

6-state subset DFA

4-state minimal DFA

Step (3) uses what's called a "set of subsets" construction.

Step (4) divides the states of an initial DFA into progressively finer equivalence classes, until it can prove that additional refinement makes no difference.

Example 1 (in the book): real numbers (no exponential notation)

RN = d*(.d|d.)d*

14-state NFA results from construction

5-state subset DFA

4-state minimal DFA

Example 2: character strings with optional backslash-escaped quotes

S = " ([^"] | \a) * " "a" for anything

11-state NFA results from construction

6-state subset DFA

4-state minimal DFA

Step (3) uses what's called a "set of subsets" construction.

Step (4) divides the states of an initial DFA into progressively finer equivalence classes, until it can prove that additional refinement makes no difference.

Example 1 (in the book): real numbers (no exponential notation)

RN = d*(.d|d.)d*

14-state NFA results from construction

5-state subset DFA

4-state minimal DFA

Example 2: character strings with optional backslash-escaped quotes

S = " ([^"] | \a) * " "a" for anything

11-state NFA results from construction

6-state subset DFA

4-state minimal DFA

Step (3) uses what's called a "set of subsets" construction.

Step (4) divides the states of an initial DFA into progressively finer equivalence classes, until it can prove that additional refinement makes no difference.

Example 1 (in the book): real numbers (no exponential notation)

RN = d*(.d|d.)d*

14-state NFA results from construction

5-state subset DFA

4-state minimal DFA

Example 2: character strings with optional backslash-escaped quotes

S = " ([^"] | \a) * " "a" for anything

11-state NFA results from construction

6-state subset DFA

4-state minimal DFA

Step (3) uses what's called a "set of subsets" construction.

Step (4) divides the states of an initial DFA into progressively finer equivalence classes, until it can prove that additional refinement makes no difference.

Example 1 (in the book): real numbers (no exponential notation)

RN = d*(.d|d.)d*

14-state NFA results from construction

5-state subset DFA

4-state minimal DFA

Example 2: character strings with optional backslash-escaped quotes

S = " ([^"] | \a) * " "a" for anything

11-state NFA results from construction

6-state subset DFA

4-state minimal DFA

Step (3) uses what's called a "set of subsets" construction.

Step (4) divides the states of an initial DFA into progressively finer equivalence classes, until it can prove that additional refinement makes no difference.

Example 1 (in the book): real numbers (no exponential notation)

RN = d*(.d|d.)d*

14-state NFA results from construction

5-state subset DFA

4-state minimal DFA

Example 2: character strings with optional backslash-escaped quotes

S = " ([^"] | \a) * " "a" for anything

11-state NFA results from construction

6-state subset DFA

4-state minimal DFA

Step (3) uses what's called a "set of subsets" construction.

Step (4) divides the states of an initial DFA into progressively finer equivalence classes, until it can prove that additional refinement makes no difference.

Example 1 (in the book): real numbers (no exponential notation)

RN = d*(.d|d.)d*

14-state NFA results from construction

5-state subset DFA

4-state minimal DFA

Example 2: character strings with optional backslash-escaped quotes

S = " ([^"] | \a) * " "a" for anything

11-state NFA results from construction

6-state subset DFA

4-state minimal DFA

Step (3) uses what's called a "set of subsets" construction.

Step (4) divides the states of an initial DFA into progressively finer equivalence classes, until it can prove that additional refinement makes no difference.

Example 1 (in the book): real numbers (no exponential notation)

RN = d*(.d|d.)d*

14-state NFA results from construction

5-state subset DFA

4-state minimal DFA

Example 2: character strings with optional backslash-escaped quotes

S = " ([^"] | \a) * " "a" for anything

11-state NFA results from construction

6-state subset DFA

4-state minimal DFA

Step (3) uses what's called a "set of subsets" construction.

Step (4) divides the states of an initial DFA into progressively finer equivalence classes, until it can prove that additional refinement makes no difference.

Example 1 (in the book): real numbers (no exponential notation)

RN = d*(.d|d.)d*

14-state