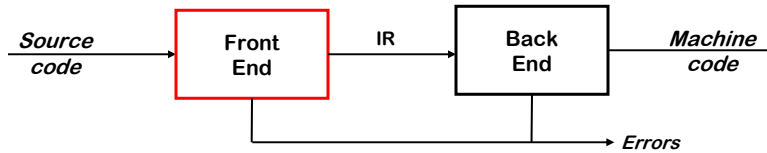


The Front End



The purpose of the front end is to deal with the input language

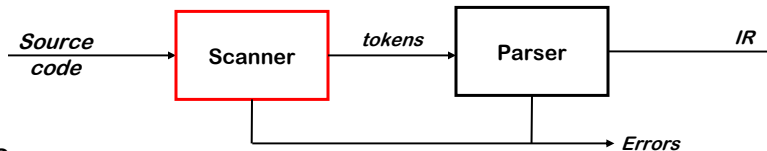
- Perform a membership test: $\text{code} \in \text{source language?}$
- Is the program well-formed (semantically) ?
- Build an IR version of the code for the rest of the compiler

The front end is not monolithic

CMSC430 Spring 2007

1

The Front End



Scanner

- Maps stream of characters into words
 - > Basic unit of syntax
 - > $x = x + y ;$ becomes
 $\langle \text{id}, x \rangle \langle \text{eq}, = \rangle \langle \text{id}, x \rangle \langle \text{pl}, + \rangle \langle \text{id}, y \rangle \langle \text{sc}, ; \rangle$
- Characters that form a word are its *lexeme*
- Its *part of speech* (or *syntactic category*) is called its *token type*
- Scanner discards white space & (often) comments

Speed is an issue in scanning
⇒ use a specialized recognizer

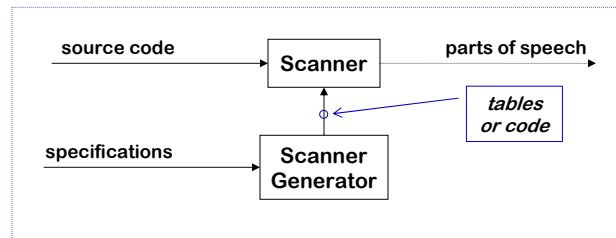
CMSC430 Spring 2007

2

The Big Picture

Why study lexical analysis?

- We want to avoid writing scanners by hand
- We want to harness automata theory



Goals:

- > To simplify specification & implementation of scanners
- > To understand the underlying techniques and technologies

Review of Scanners

Lexical Analysis Strategy: **Simulation of Finite Automaton**

- > States, characters, actions
- > State transition $\delta(\text{state}, \text{charclass})$ determines next state

Next character function

- > Reads next character into buffer
- > Computes *character class* by fast table lookup

Transitions from state to state

- > Current state and next character determine (via δ)
 - Next state and action to be performed
 - Some actions *preload* next character

Identifiers distinguished from keywords by hashed lookup

- > This differs from EAC advice (discussion later)
- > Permits translation of identifiers into **<type, symbol_index>**
 - Keywords each get their own type

Examples of Regular Expressions

Identifiers:

Letter → (a|b|c|...|z|A|B|C|...|Z)

Digit → (0|1|2|...|9)

Identifier → *Letter* (*Letter* | *Digit*)*

Numbers:

Integer → (+ | - | ε) (0 | (1 | 2 | 3 | ... | 9) (*Digit*)*)

Decimal → *Integer* . *Digit**

Real → (*Integer* | *Decimal*) E (+ | - | ε) *Digit**

Complex → (*Real* . *Real*)

Numbers can get much more complicated!

Regular Expressions

(the point)

Regular expressions can be used to specify the words to be translated to parts of speech by a lexical analyzer

Using results from automata theory and theory of algorithms, we can automatically build recognizers from regular expressions

⇒ We study REs and associated theory to automate scanner construction !

Example

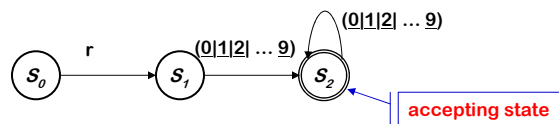
(from Lab 1)

Consider the problem of recognizing register names

$Register \rightarrow r (0|1|2| \dots | 9) (0|1|2| \dots | 9)^*$

- Allows registers of arbitrary number
- Requires at least one digit

RE corresponds to a recognizer (or DFA)



Recognizer for *Register*

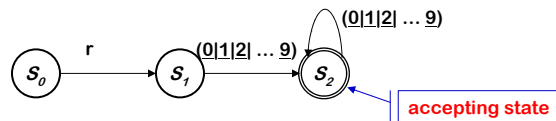
Transitions on other inputs go to an error state, s_e

Example

(continued)

DFA operation

- Start in state S_0 & take transitions on each input character
- DFA accepts a word \underline{x} iff \underline{x} leaves it in a final state (S_2)



Recognizer for *Register*

So,

- r17 takes it through s_0, s_1, s_2 and accepts
- r takes it through s_0, s_1 and fails
- a takes it straight to s_e

Example

(continued)

δ / action	r	0,1,2,3, 4,5,6, 7,8,9	other
0	1 / <u>start</u>	e / <u>error</u>	e / <u>error</u>
1	e / <u>error</u>	2 / <u>add</u>	e / <u>error</u>
2	e / <u>error</u>	2 / <u>add</u>	x / <u>exit</u>
e	e / <u>error</u>	e / <u>error</u>	e / <u>error</u>

- The recognizer translates directly into code
- To change DFAs, just change the tables

What if we need a tighter specification?

r *Digit Digit** allows arbitrary numbers

- Accepts r00000
- Accepts r99999
- What if we want to limit it to r0 through r31 ?

Write a tighter regular expression

- > *Register* \rightarrow r ((0|1|2) (*Digit* | ϵ) | (4|5|6|7|8|9) | (3|30|31))
- > *Register* \rightarrow r0|r1|r2| ... |r31|r00|r01|r02| ... |r09

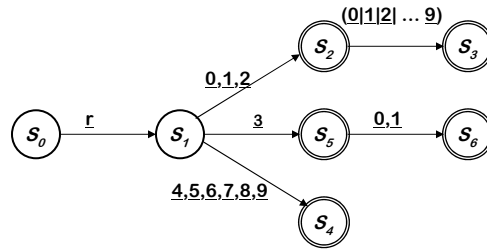
Produces a more complex DFA

- Has more states
- Same cost per transition
- Same basic implementation

Tighter register specification (continued)

The DFA for

$Register \rightarrow r ((0|1|2) (Digit | \epsilon) | (4|5|6|7|8|9) | (3|30|31))$



- Accepts a more constrained set of registers
- Same set of actions, more states

Tighter register specification (continued)

state action	r	0,1	2	3	4,5,6 7,8,9	other
0	1 start	e	e	e	e	e
1	e	2 add	2 add	5 add	4 add	e
2	e	3 add	3 add	3 add	3 add	x exit
3,4	e	e	e	e	e	x exit
5	e	6 add	e	e	e	x exit
6	e	e	e	e	e	x exit
e	e	e	e	e	e	e

Automating Scanner Construction

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 *(define all parts of speech)*
- You could build one in a weekend!

CMSC430 Spring 2007

13

Automating Scanner Construction

RE → NFA (*Thompson's construction*)

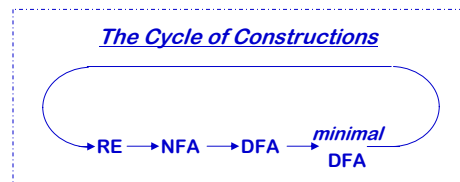
- Build an NFA for each term
- Combine them with ϵ -moves

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
- (Not a particularly obvious algorithm)

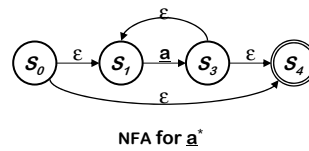
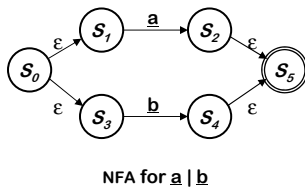
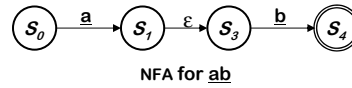
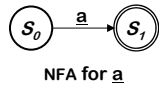
CMSC430 Spring 2007

14

RE \rightarrow NFA using Thompson's Construction

Key idea

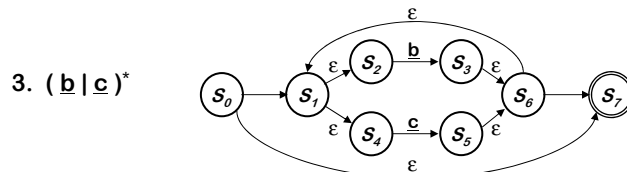
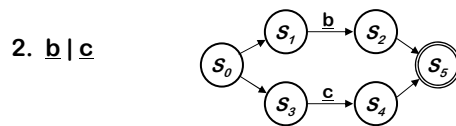
- NFA pattern for each symbol & each operator
- Join them with ϵ moves in precedence order



Ken Thompson, CACM, 1968

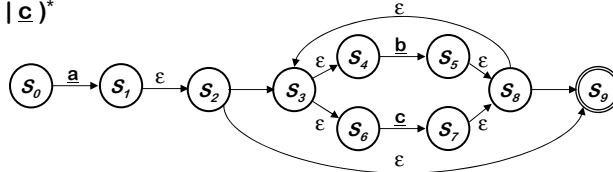
Example of Thompson's Construction

Let's try a (b | c)*

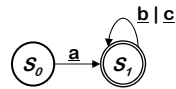


Example of Thompson's Construction (continued)

4. $a(b|c)^*$



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



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

NFA → DFA with Subset Construction

Need to build a simulation of the NFA

Two key functions

- $Move(s_i, a)$ is set of states reachable from s_i by a
- ϵ -closure(s_i) is set of states reachable from s_i by ϵ

The algorithm:

- Start state derived from s_0 of the NFA
- Take its ϵ -closure $S_0 = \epsilon$ -closure(s_0)
- Take the image of S_0 $Move(S_0, \alpha)$ for each $\alpha \in \Sigma$ and take its ϵ -closure
- Iterate until no more states are added

Sounds more complex than it is...

NFA \rightarrow DFA with Subset Construction

The algorithm:

```
 $s_0 \leftarrow \varepsilon\text{-closure}(q_{0n})$ 
while (  $S$  is still changing )
  for each  $s_j \in S$ 
    for each  $\alpha \in \Sigma$ 
       $s_\gamma \leftarrow \varepsilon\text{-closure}(\text{Move}(s_j, \alpha))$ 
      if (  $s_\gamma \notin S$  ) then
        add  $s_\gamma$  to  $S$  as  $s_j$ 
         $T[s_j, \alpha] \leftarrow s_j$ 
```

Let's think about why this works

The algorithm halts:

1. S contains no duplicates (test before adding)
2. 2^{Q_n} 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 character in each s_j
It builds every possible NFA configuration.

$\Rightarrow S$ and T form the DFA

NFA \rightarrow DFA with Subset Construction

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

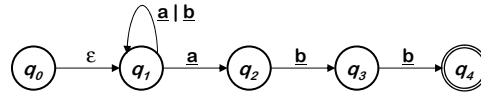
Other fixed-point computations

- Canonical construction of sets of LR(1) items
 - > Quite similar to the subset construction
- Classic data-flow analysis (& Gaussian Elimination)
 - > Solving sets of simultaneous set equations

We will see many more fixed-point computations

NFA → DFA with Subset Construction

Remember $(a | b)^* abb$?



Applying the subset construction:

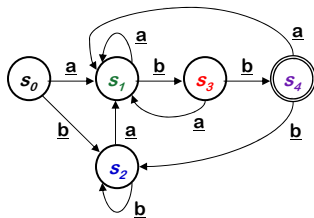
Iter.	State	Contains	ϵ -closure(move(s, <u>a</u>))	ϵ -closure(move(s, <u>b</u>))
0	s_0	q_0, q_1	q_1, q_2	q_1
1	s_1	q_1, q_2	q_1, q_2	q_1, q_3
	s_2	q_1	q_1, q_2	q_1
2	s_3	q_1, q_3	q_1, q_2	q_1, q_4
3	s_4	q_1, q_2	q_1, q_2	q_1

Iteration 3 adds nothing to S , so the algorithm halts

contains q_4
(final state)

NFA → DFA with Subset Construction

The DFA for $(a | b)^* abb$



δ	<u>a</u>	<u>b</u>
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

- Not much bigger than the original
- All transitions are deterministic
- Use same code skeleton as before

Final Task: Generate the Scanner

- How do we specify different **token types**, etc?
 - > One rule per token type
 - RE on right hand side
 - > What about ambiguity?
- How do we provide **actions** in a scanner specified by REs?
 - > One simple action: add character to current token
 - > More complex actions on token end
 - Part of specification
- How do we generate a scanner for **multiple tokens**?
 - > Combine rules for each token
 - What about ambiguity?
 - > How do we identify the tokens (in the input, etc)
 - Legal token if error (or end) transition taken from accepting state
 - Leave error char in input buffer

CMSC430 Spring 2007

23

A Lex Specification, Part I

```
%{
    /* definition of constants BEGIN, END, NAME,
       NUM, STRNG, SPCL, PLUS, MINUS, LT, LE */
%}
/* Regular Definitions */
blank      [ ]
lb         [{]
rb         [}]
comment    {lb}.*{rb} /* . = any but newline */
ws         ({blank})+
letter     [A-Za-z]
digit      [0-9]
name       {letter}({letter})+
numb       {digit}+
quote      ["]
string     {quote}.*{quote}
%%
```

CMSC430 Spring 2007

24

Lex Specification, Part II

```
/* Translation Rules */
{ws}          { /* no action and no return */ }
begin         { return(BEGIN); }
end           { return(END); }
{name}       { yylval = install_name(); return(NAME); }
{number}     { yylval = install_num(); return(NUM); }
{string}     { yylval = install_str(); return(STRNG); }
"+"          { yylval = PLUS; return(SPCL); }
"_"          { yylval = MINUS; return(SPCL); }
"<"          { yylval = LT; return(SPCL); }
"<="        { yylval = LE; return(SPCL); }
%%
install_name() {
    /* procedure to install the lexeme
       whose first character is pointed to by yytext
       and whose length is yyleng into symbol table
       and return pointer to entry */
}
```

CMSC430 Spring 2007

25

Automating Scanner Construction

RE → NFA (Thompson's construction) 

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

NFA → DFA (subset construction) 

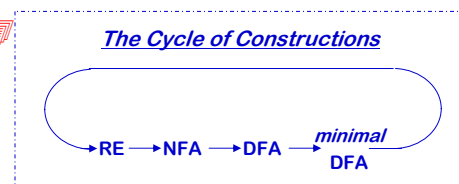
- Build the simulation

DFA → Minimal DFA

- Hopcroft's algorithm

DFA → RE (not really part of scanner construction)

- All pairs, all paths problem
- Union together paths from s_0 to a final state



CMSC430 Spring 2007

26

DFA Minimization

The Big Picture

- Discover sets of equivalent states
- Represent each such set with just one state

Two states are equivalent if and only if:

- The set of paths leading to them are equivalent
- $\forall \alpha \in \Sigma$, transitions on α lead to equivalent states (DFA)
- α -transitions to distinct sets \Rightarrow states must be in distinct sets

A partition P of S

- Each $s \in S$ is in exactly one set $p_i \in P$
- The algorithm iteratively partitions the DFA's states

DFA Minimization

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: $\{F\}$ & $\{Q-F\}$ ($D=(Q,\Sigma,\delta,q_0,F)$)

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

- Assume q_x & $q_y \in s$, and
- $\delta(q_x, \underline{a}) = q_x$ & $\delta(q_y, \underline{a}) = q_y$
- If q_x & q_y are not in the same set, then s must be split
- One state in the final DFA cannot have two transitions on \underline{a}

DFA Minimization

The algorithm

```
P ← {F, {Q-F}}
while ( P is still changing)
  T ← {}
  for each set s ∈ P
    for each α ∈ Σ
      partition s by α
        into s1, s2, ..., sk
      T ← T ∪ s1, s2, ..., sk
  if T ≠ P then
    P ← T
```

This is a fixed-point algorithm!

Why does this work?

- Partition $P \in 2^Q$
- Start off with 2 subsets of Q $\{F\}$ and $\{Q-F\}$
- *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 $|Q|$ sets
- Maximum of $|Q|$ splits

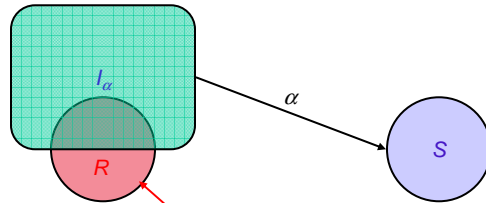
Note that

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

Hopcroft's Algorithm

```
W ← {F, Q-F}; P ← {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 α in Σ do begin
    let Iα ← δα-1( S );
    for each R in P such that R ∩ Iα is not empty
      and R is not contained in Iα do begin
        partition R into R1 and R2 such that R1 ← R ∩ Iα; R2 ← R - R1;
        replace R in P with R1 and R2;
        if R ∈ W then replace R with R1 in W and add R2 to W;
        else if || R1 || ≤ || R2 ||
          then add R1 to W;
          else add R2 to W;
      end
    end
  end
end
```

Key Idea



This part must have an α -transition to some other state!

CMSC430 Spring 2007

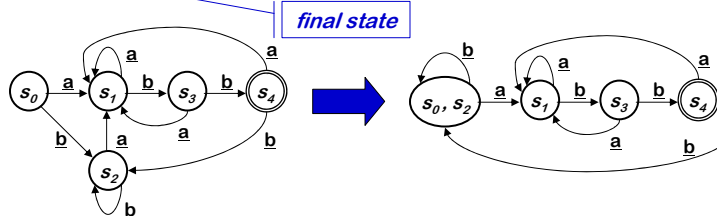
31

DFA Minimization

Enough theory, does this stuff work?

> Recall our example: $(a|b)^*abb$

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

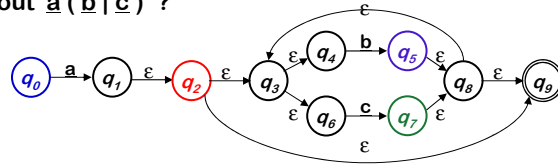


CMSC430 Spring 2007

32

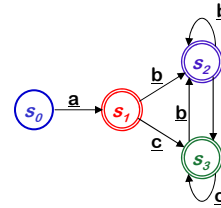
DFA Minimization

What about $a(b|c)^*$?



First, the subset construction:

		ϵ -closure(move(s,*))		
	NFA states	a	b	c
s_0	q_0	$q_1, q_2, q_3, q_4, q_6, q_8$	none	none
s_1	$q_1, q_2, q_3, q_4, q_6, q_8$	none	$q_5, q_8, q_9, q_{3,1}, q_{4,1}, q_6$	$q_7, q_8, q_9, q_{3,2}, q_{4,2}, q_6$
s_2	$q_5, q_8, q_9, q_{3,1}, q_{4,1}, q_6$	none	s_2	s_3
s_3	$q_7, q_8, q_9, q_{3,2}, q_{4,2}, q_6$	none	s_2	s_3

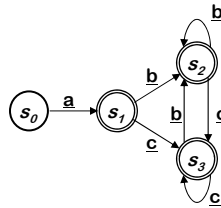


Final states

DFA Minimization

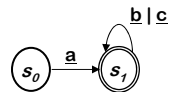
Then, apply the minimization algorithm

		Split on		
	Current Partition	a	b	c
P_0	$\{s_1, s_2, s_3\}, \{s_0\}$	none	none	none



To produce the minimal DFA

final states



In lecture 5, we observed that a human would design a simpler automaton than Thompson's construction did.

The algorithms produce that same 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

$Term \rightarrow [a-zA-Z] ([a-zA-Z] | [0-9])^*$

$Op \rightarrow + | - | * | /$

$Expr \rightarrow (Term Op)^* Term$

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

- $L = \{p^k q^k\}$ *(parenthesis languages)*
- $L = \{wcw^r \mid w \in \Sigma^*\}$

Neither of these is a regular language *(nor an RE)*

But, this is a little subtle. You can construct DFA's for

- Strings with alternating 0's and 1's
 $(\epsilon | 1)(01)^*(\epsilon | 0)$
- Strings with an even number of 0's and 1's
[See Homework 1!](#)

RE's can count bounded sets and bounded differences

What can be so hard?

Poor language design can complicate scanning

- Reserved words are important
if then then then = else; else else = then (PL/I)
- Significant blanks (Fortran & Algol68)
do 10 i = 1,25
do 10 i = 1.25
- String constants with special characters (C, others)
newline, tab, quote, comment delimiters, ...
- Finite closures
 - > Limited identifier length
 - > Adds states to count length

What can be so hard?

(Fortran 66/77)

```
INTEGERFUNCTIONA
PARAMETER(A=6,B=2)
IMPLICIT CHARACTER*(A-B)(A-B)
INTEGER FORMAT(10), IF(10), DO9E1
100 FORMAT(4H)=(3)
200 FORMAT(4 )=(3)
   DO9E1=1
   DO9E1=1,2
   9 IF(X)=1
     IF(X)H=1
     IF(X)300,200
300 CONTINUE
   END
C THIS IS A "COMMENT CARD"
$ FILE(1)
   END
```

How does a compiler do this?

- First pass finds & inserts blanks
- Can add extra words or tags to create a scannable language
- Second pass is normal scanner

Example due to Dr. F.K. Zadeck

Building Faster Scanners from the DFA

Table-driven recognizers waste a lot of 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

```
char ← next character;
state ← s0;
call action(state, char);
while (char ≠ eof)
  state ← δ(state, char);
  call action(state, char);
  char ← next character;
```

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

```
if T(state) = final then
  report acceptance;
else
  report failure;
```

Building Faster Scanners from the DFA

A direct-coded recognizer for *r Digit Digit**

```
goto s0;
s0: word ← ∅;
   char ← next character;
   if (char = 'r')
     then goto s1;
     else goto se;
s1: word ← word + char;
   char ← next character;
   if ('0' ≤ char ≤ '9')
     then goto s2;
     else goto se;
s2: word ← word + char;
   char ← next character;
   if ('0' ≤ char ≤ '9')
     then goto s2;
     else if (char = eof)
       then report acceptance;
       else goto se;
se: print error message;
   return failure;
```

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

Building Faster Scanners

Hashing keywords versus encoding them directly

- Some compilers recognize keywords as identifiers and check them in a hash table *(some well-known compilers do this!)*
- Encoding it in the DFA is a better idea
 - > O(1) cost per transition
 - > Avoids hash lookup on each identifier

It is hard to beat a well-implemented DFA scanner