

CS143: Lexical Analysis

David L. Dill

Stanford University

Introduction

- PA1 (or PA1J) is assigned today.
 - Due in one week
 - I recommend starting immediately.
- Review notes on regular languages based on CS103
 - Includes DFA minimization – optional
- Notes on the material of lectures 1, 3, 4.

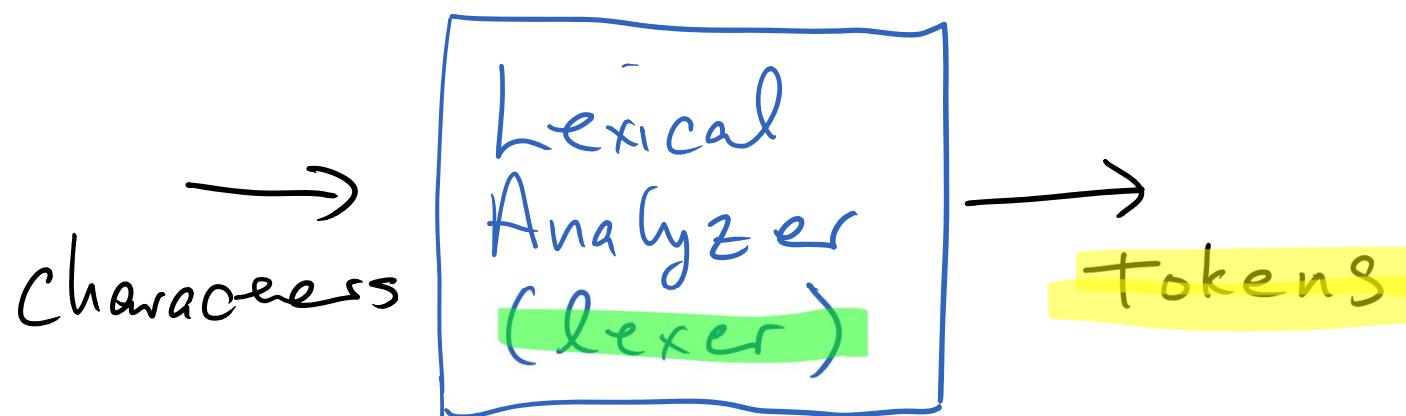
Lexical Analysis

- Interface to Rest of Compiler
- Formal Languages Concepts
- Regular Expressions
- Deterministic Finite Automata
- Theory vs. Practice

Interfaces to Rest of Compiler

Lexical analysis

characters → words

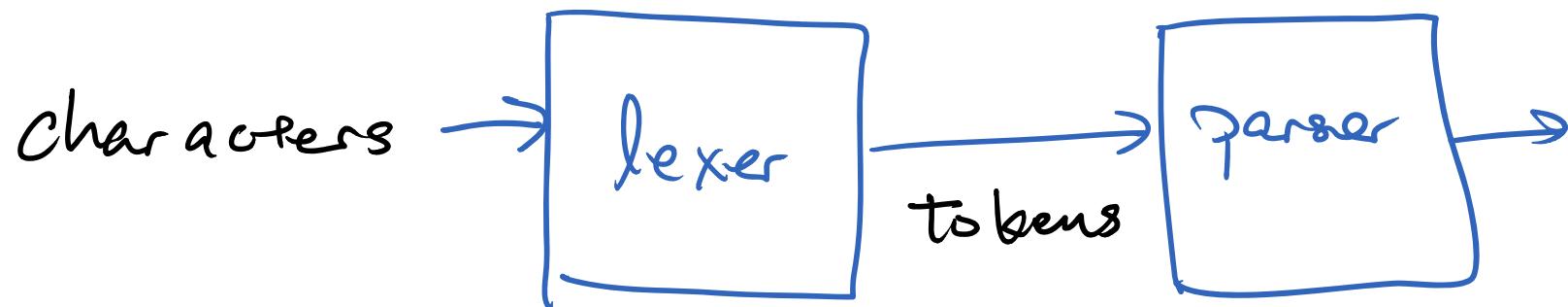


Lexical analysis

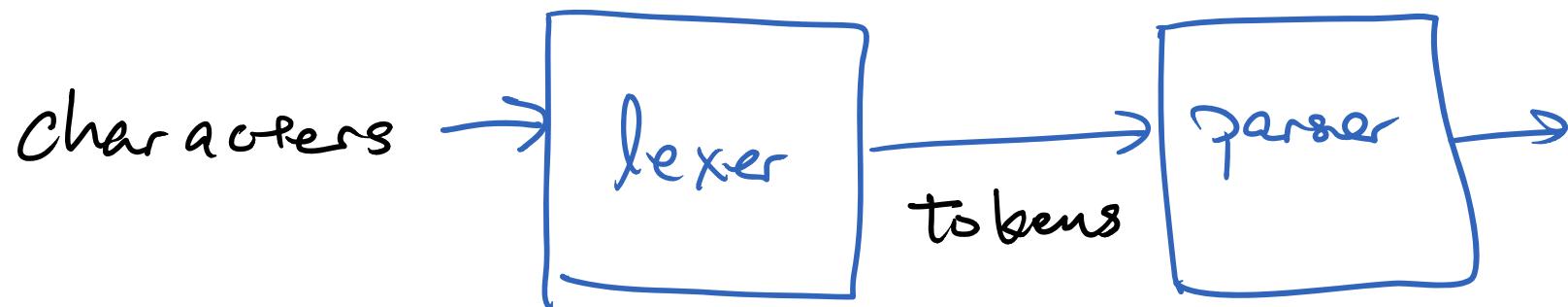
if $x == y$ then $z = 1$, else $z = 2$;

if x == then z = 1 ; else
 | | | | | | |
 z = 2 ;

Interface

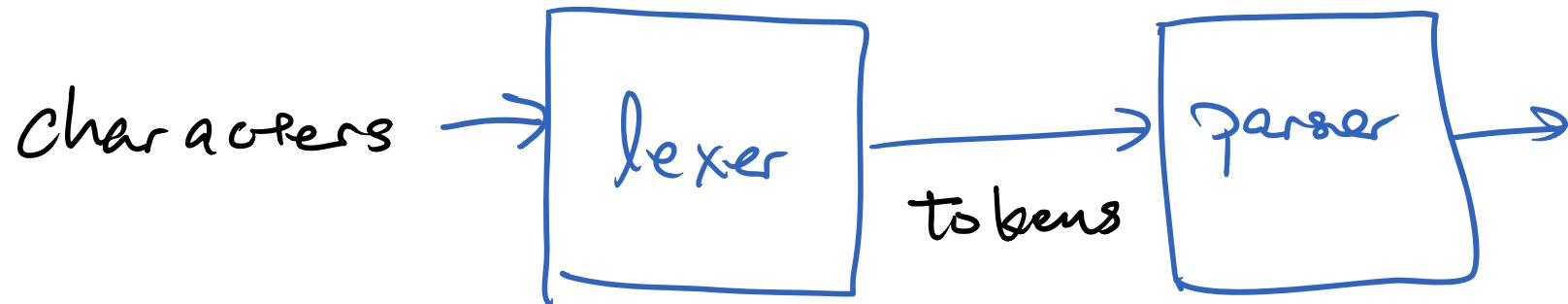


Interface



tokens are generated "on demand"
for parser.

Interface



tokens are generated "on demand" for parser.

getlex() → returns next token

if $x == y$ then $z = 1$, else $z = 2$;
↑
next input

if $x == y$ then $z = 1$, else $z = 2$;



next input

getlex() returns if token

if $x == y$ then $z = 1$, else $z = 2$;

↑

next input

getlex() returns if token

^{next} getlex() returns x,

if $x == y$ then $z = 1$, else $z = 2$;

next input

getlex() returns if token

^{next} getlex() returns x

(discards space after "if" and keeps
going until it can return a token.)

Terminology

If $x = y$ then $a \circ a$

Terminology

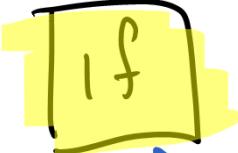
If $x = y$ then $a \circ a$



lexeme - input string for next token

Terminology

If $x == y$ then $a \circ a$

getlex() →  token

data structure returned
by getlex
(goes to parser)

Performance

Lexer used to be performance
-critical
e.g. before 1990?

of characters vs. # of tokens

lexer had to deal with all of these

later stages just had tokens & trees.

Performance

hexer used to be performance
-critical

So we know how to do it **FAST**

Q: Why is the North American
Antelope so speedy?

Formal Languages Concepts

Lexical Analysis Theory

Precise descriptions - Regular expressions

$$[a-zA-Z][a-zA-Z0-9-]^*$$

Automatic generation

Regular expressions \rightarrow NFA \rightarrow DFA

\rightarrow lexical analyzer program

Formal language Concepts

Alphabet - A finite set of **symbols** (a.k.a **characters**).

Σ - standard symbol for an alphabet.

Formal language Concepts

Alphabet - A finite set of **symbols** (a.k.a **characters**).

Σ - standard symbol for an alphabet.

For compilers, a character set like
ASCII, UTF8

Formal language Concepts

String - finite sequence of characters

E.g. abc

Formal language Concepts

$|x|$ - length of the string x

ϵ - empty string $|\epsilon| = \emptyset$

Formal language Concepts

$x \cdot y$ - concatenation of strings
(also written xy)

$$abc \cdot def = abcdef$$

Formal language Concepts

Formal language - set of strings
over some alphabet Σ .

Formal language Concepts

If X, Y are languages over Σ ,
 $X \cdot Y = \{xy \mid x \in X \wedge y \in Y\}$ is
their concatenation.

Formal language Concepts

If X, Y are languages over Σ ,
 $X \cdot Y = \{xy \mid x \in X \wedge y \in Y\}$ is
their concatenation.

$$\begin{aligned} \text{Ex: } & \{a, ab\} \cdot \{b, bb\} \\ &= \{ab, abb, abb\} \end{aligned}$$

Formal language Concepts

Let X be a language. Then X^i
is X concatenated with itself i times

$$X^0 = \{ \epsilon \}$$

$$X^1 = X$$

$$X^2 = XX$$

etc.

Formal language Concepts

$$X^* = X^0 \cup X^1 \cup X^2 \cup \dots$$

is the Kleene closure of X .

Regular Expressions

Regular Expressions

Expression

ϵ

$a \ (a \in \Sigma)$

$R_1 | R_2$

$R_1 \cdot R_2$

R_1^*

Language

$$L(\epsilon) = \{ \epsilon \}$$

$$L(a) = \{ a \}$$

$$L(R_1 | R_2) = L(R_1) \cup L(R_2)$$

$$L(R_1 \cdot R_2) = L(R_1) \cdot L(R_2)$$

$$L(R_1^*) = L(R_1)^*$$

Regular Expressions

Also: Precedence

* , - , |
↑

"stickiest"
like exponentiation

Regular Expressions

Also: Precedence

* , ?, |
↑

next stickiest

like multiplication

Regular Expressions

Also: Precedence

* , - , |
 ↑

least sticky
like addition

Regular Expressions

Also: Precedence

* , - , |

Parentheses to enforce grouping

$((a|b)(c|d))^*$

Notation used in tools is somewhat different
(and more complex).

Abbreviations

$R^?$ — "optional R "

abbrev for $(R|\varepsilon)$

R^+ — positive closure (one or more)

abbrev for RR^* or R^*R .

Regular Definitions

Allow new named abbreviations for regular expressions

Example: Java floating point (partial)

Expt: $((e|E)(+|-)? \text{digit}^+)$

Suffix: $(f|F|d|D)$

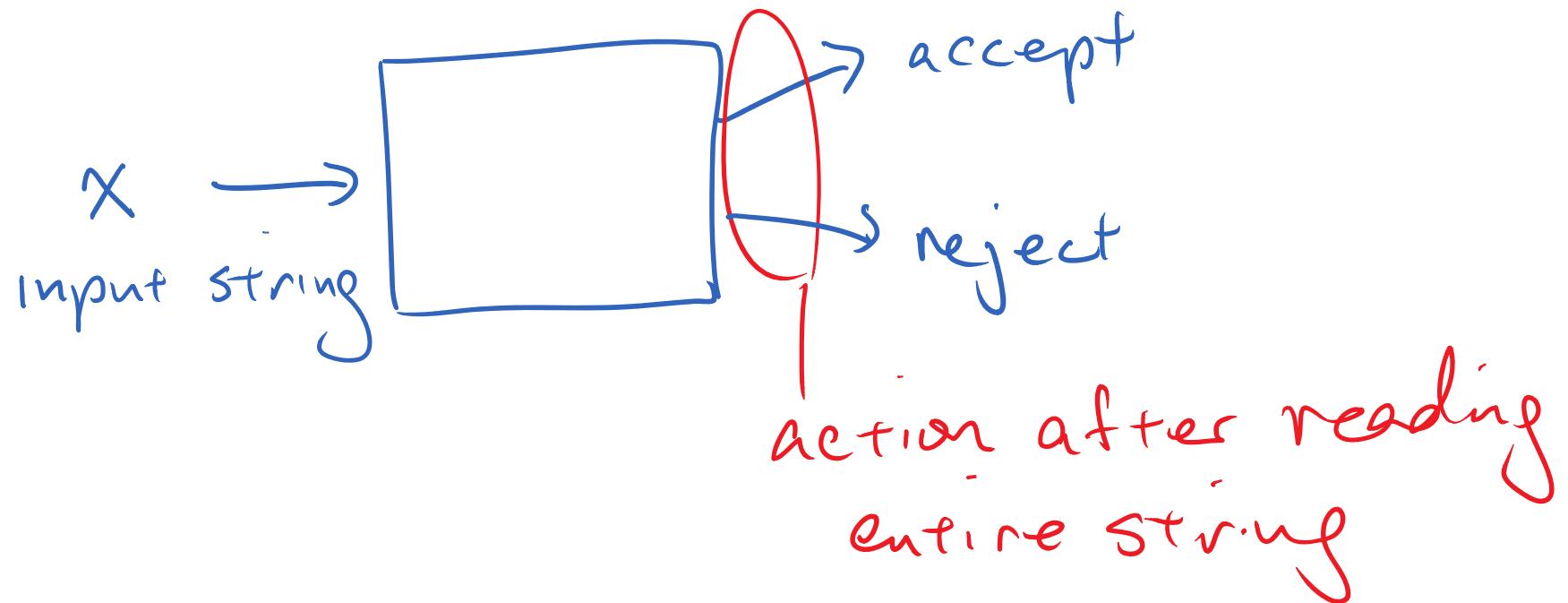
$\text{digit}^+ . \text{digit}^* \text{Expt? Suffix?}$
| $\cdot \text{digit}^+ \text{Expt? Suffix?}$
| $\text{digit}^+ \text{Expt Suffix?}$
| $\text{digit}^+ \text{Expt? Suffix}$

Reg expr
for Java
FP numbs

Deterministic Finite Automata

Deterministic Finite Automata

Good for implementing lexers.



DFA

Q - finite set of states

Σ - alphabet

$\delta: Q \times \Sigma \rightarrow Q$ next state function

$q_0 \in Q$ start state

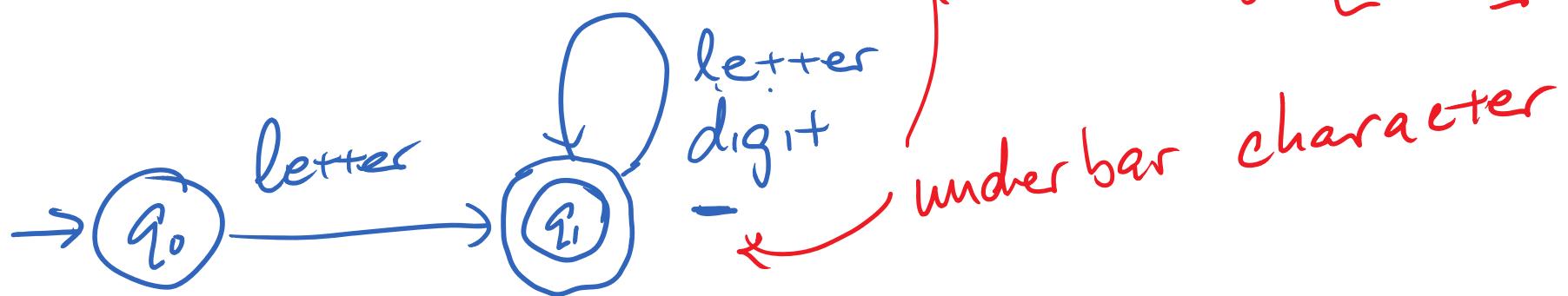
$F \subseteq Q$ accepting states.

partial function, unlike
CS103

Example: Identifiers

letter (letter | digit | _)^{*}

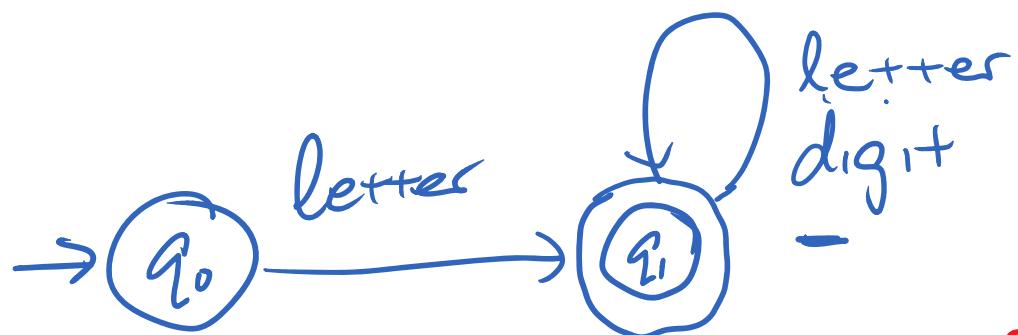
letter [a-zA-Z]
digit [0-9]



Example: Identifiers

letter (letter | digit | -)^{*}

letter [a-zA-Z]
digit [0-9]



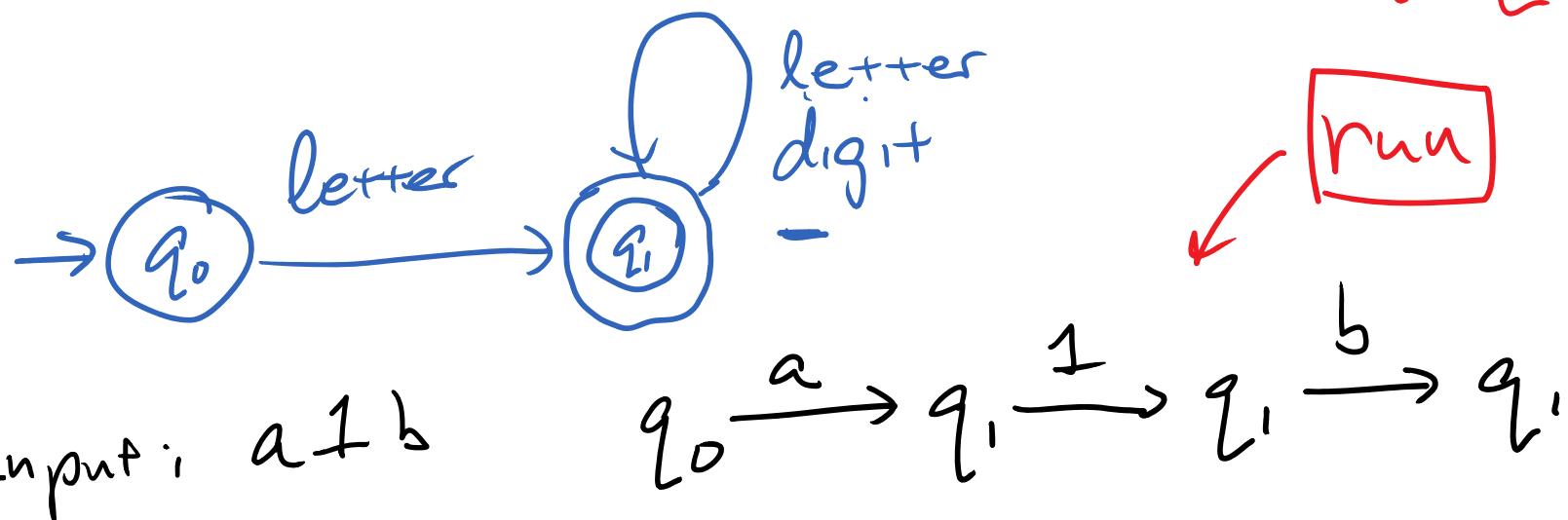
At most one arrow from each state
on each symbol ("deterministic")

Ok to have no arrow on symbol (S partial function)

Example: Identifiers

letter (letter | digit | -) *

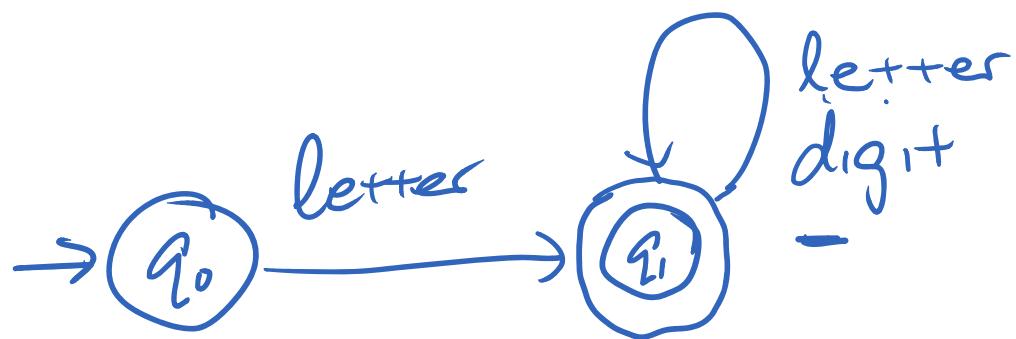
letter [a-zA-Z]
digit [0-9]



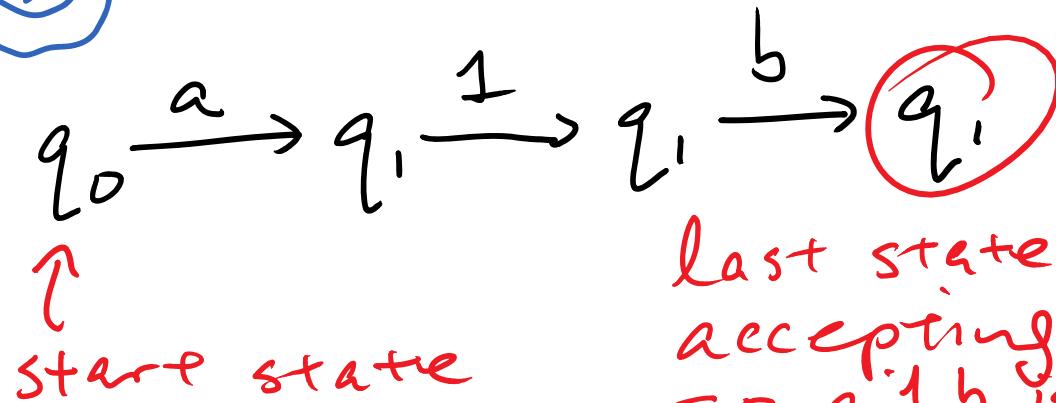
Example: Identifiers

letter (letter | digit | -) *

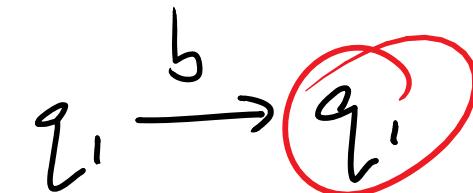
letter [a-zA-Z]
digit [0-9]



Input: a1b



run



last state is
accepting -
so a1b is ident

Automatic Conversion

Regular expression \rightarrow NFA with ϵ transitions

\rightarrow DFA

\rightarrow minimized DFA

\uparrow not in CS103

Review notes posted

Example: Modula-2 Real Numbers

digit⁺ . digit^{*} (E (+|-)? digit⁺)?

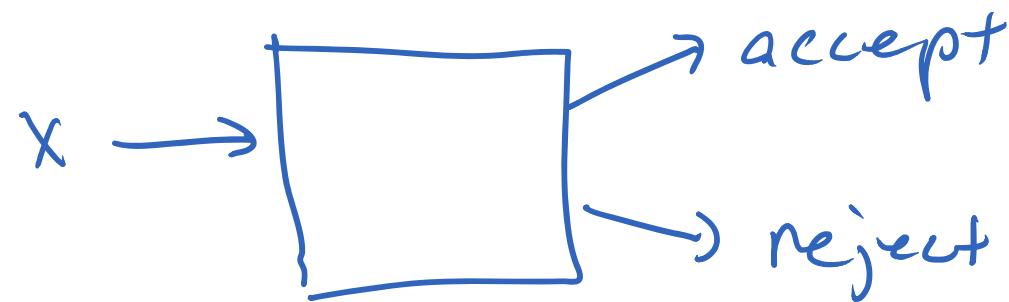
DFA?

Example: Modula-2 Real Numbers

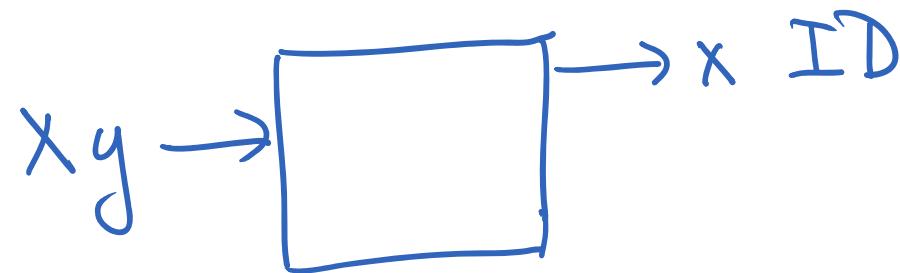
digit⁺ . digit* (E (+|-)? digit⁺)?

Theory vs. Practice

Theory



Practice

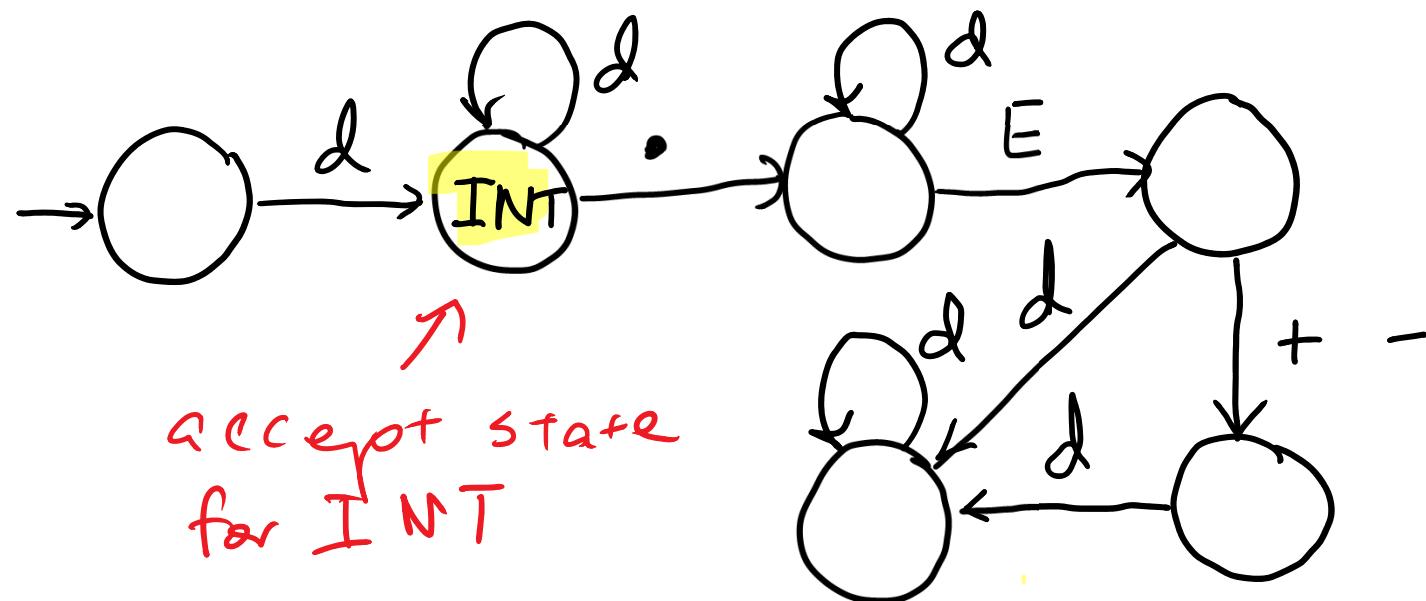


- ① Lexeme is prefix of input
- ② Multiple token types

Combined Automaton

Real: $\text{digit}^+ \cdot \text{digit}^* (E (+|-)? \text{digit}^+) \cdot$

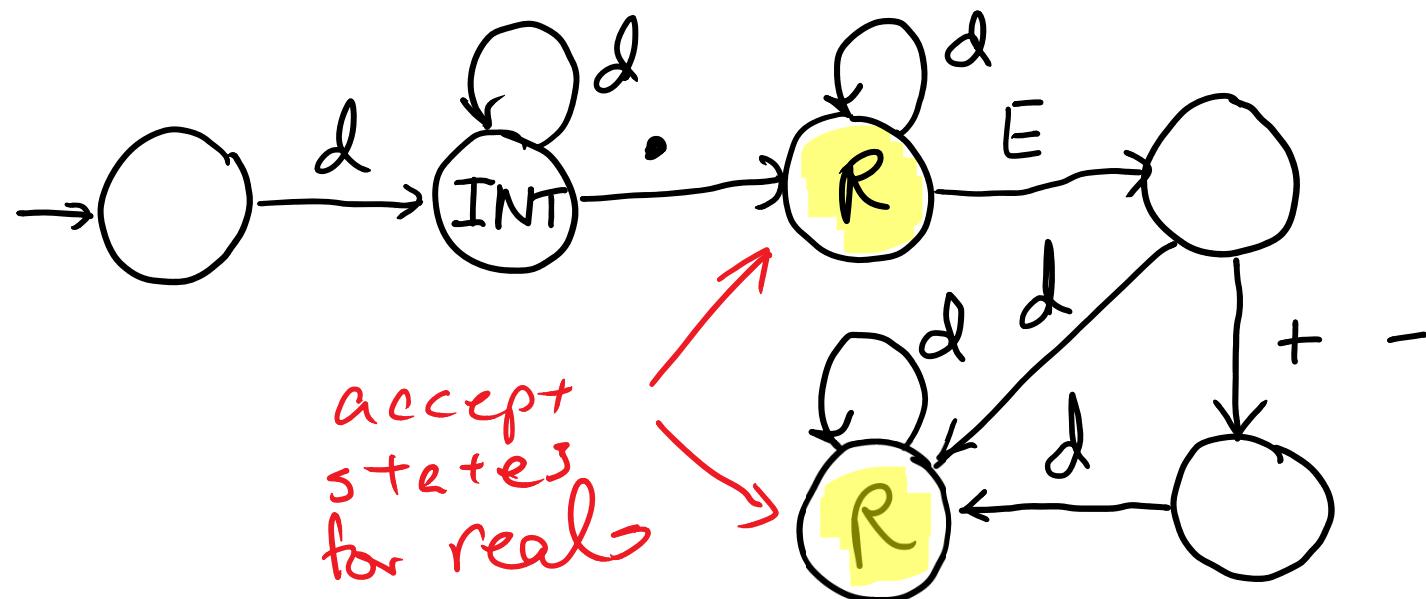
Int: digit^+



Combined Automaton

Real: $\text{digit}^+ \cdot \text{digit}^* (E (+|-)? \text{digit}^+)$.

Int: digit^+

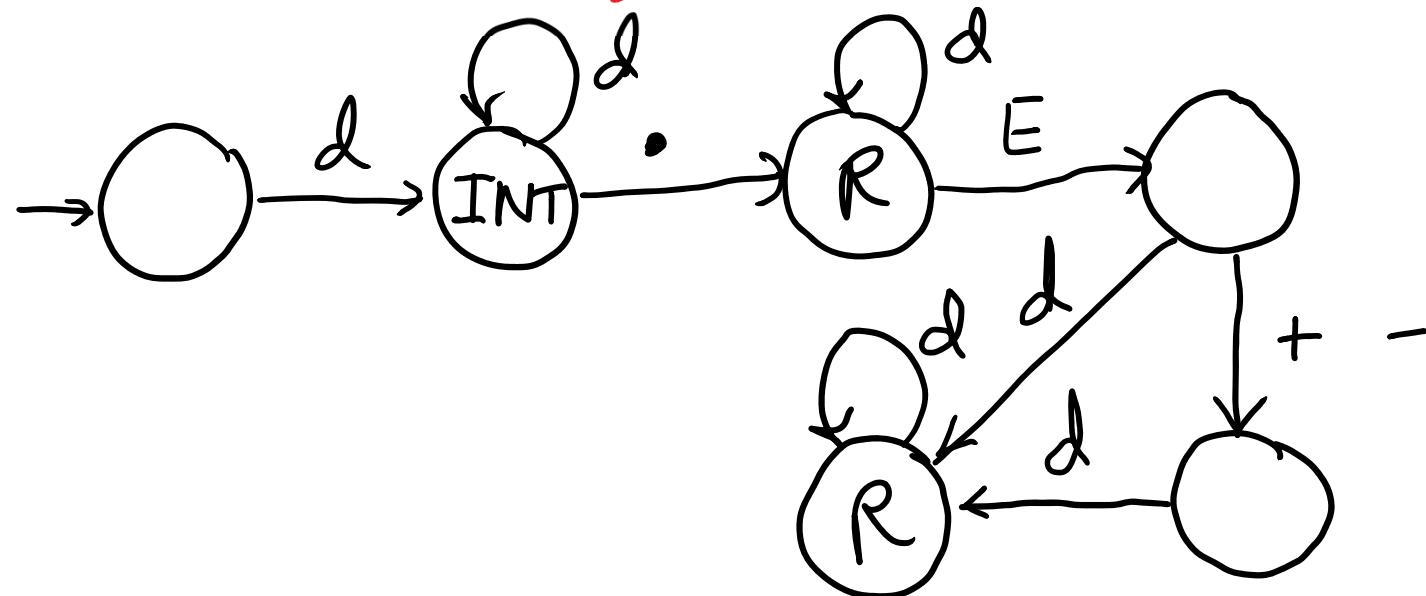


Input: 1 2.3 4+6

Tokens 1 2.3 4 ?_b

Longest lexeme rule
("maximal munch")

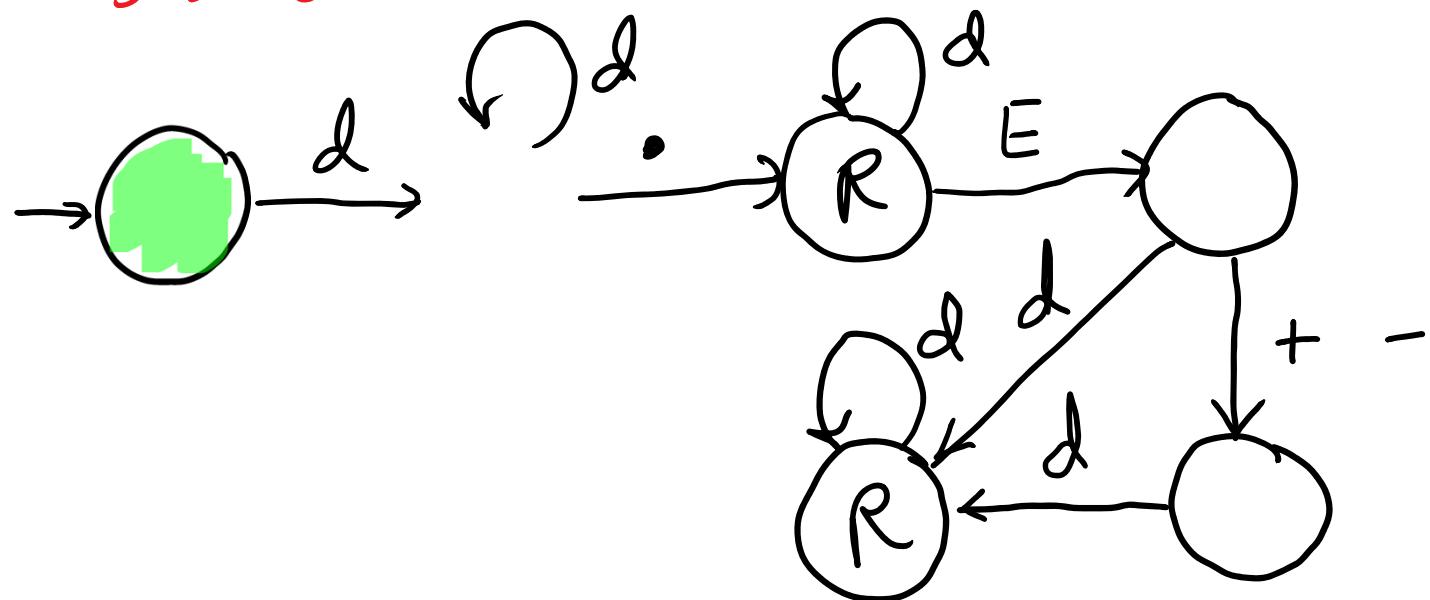
Tokens 12.34 ?_a



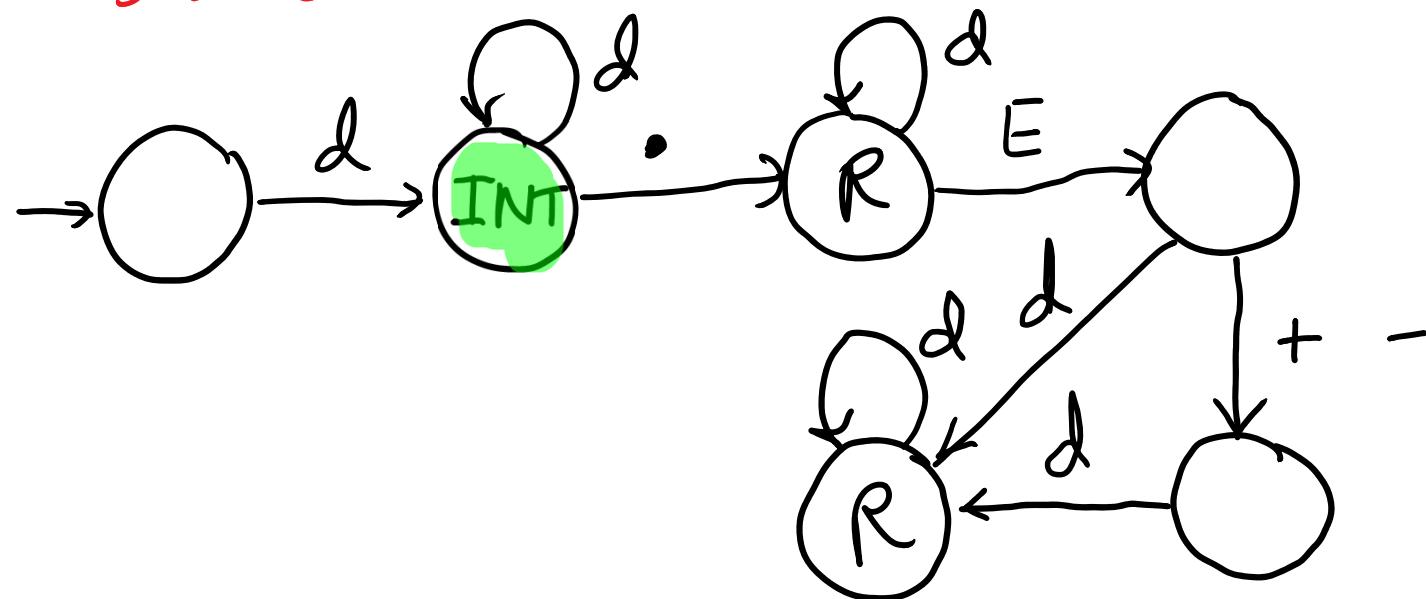
Input: 1 2.3 4+6

startchar ↗
nextchar

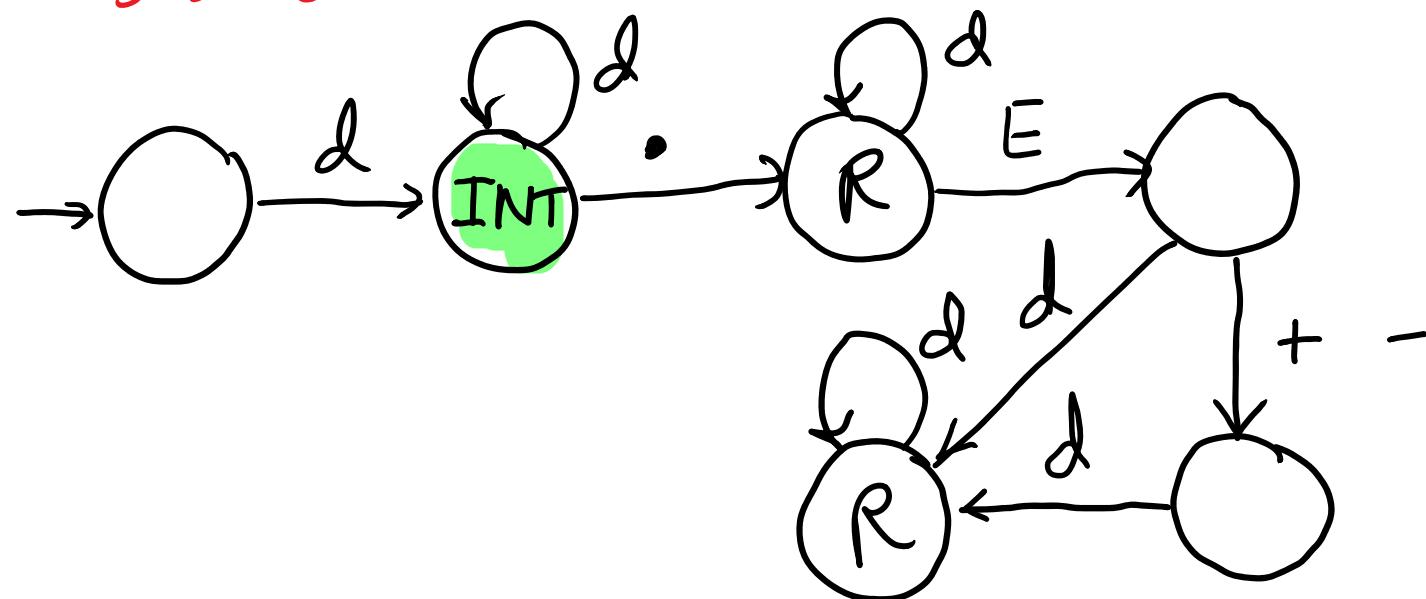
(next char to read)



Input: 1 2.3 4+6
 startchar nextchar (advanced)

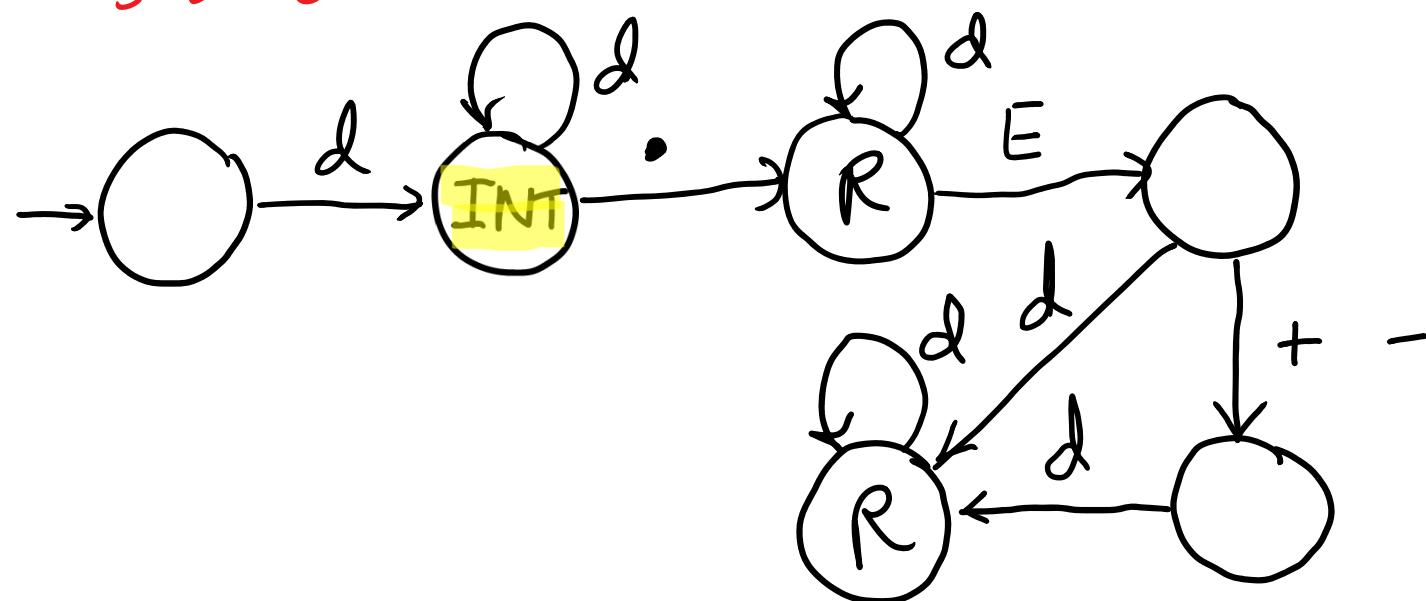


Input: $1 \overset{\text{startchar}}{\nearrow} 2.34+6 \overset{\text{nextchar}}{\nearrow}$ last char (end of longest lexeme
 so far - "1" in this case)



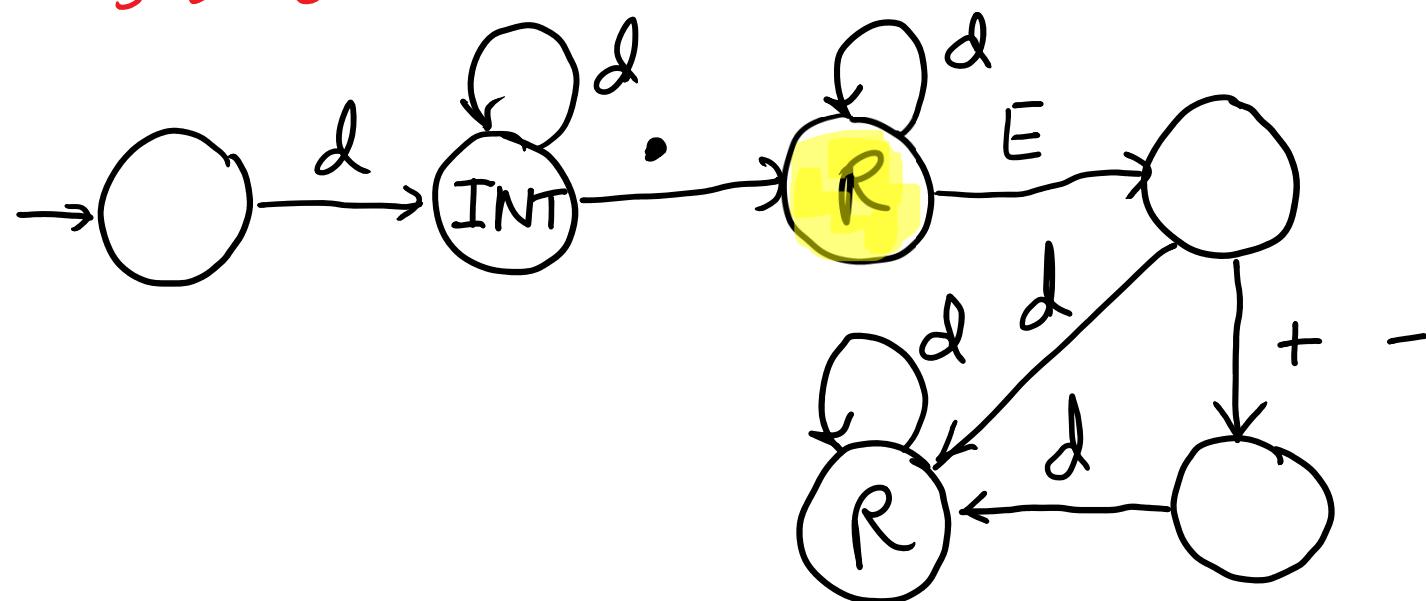
Input: 1 2 . 3 4 + 6
 startchar last char
 nextchar

longest so far:
"12"



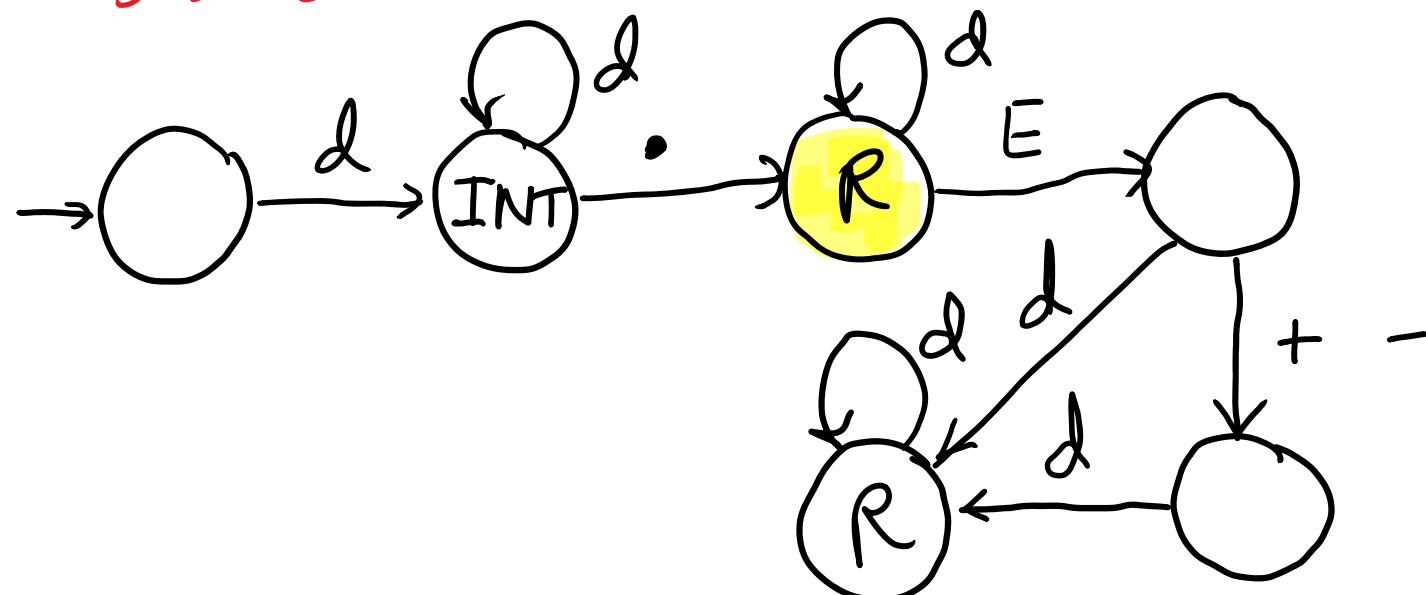
Input: 1 2 . 3 4 + 6
 startchar last char
 nextchar

longest so far:
"12."



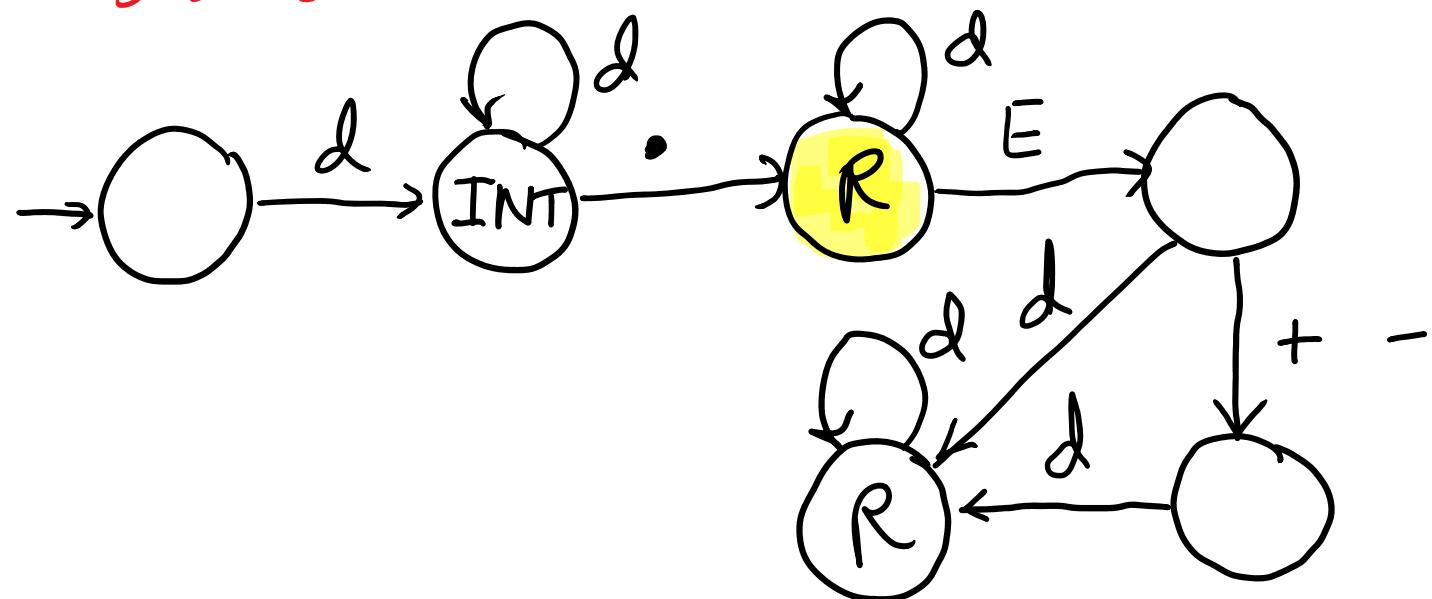
Input: 1 2 . 3 4 + 6
 startchar last char
 nextchar

longest so far:
"12.3"



Input: 1 2 . 3 4 + 6
 startchar →
 lastchar ↓
 nextchar ←

longest so far:
 "12.34"



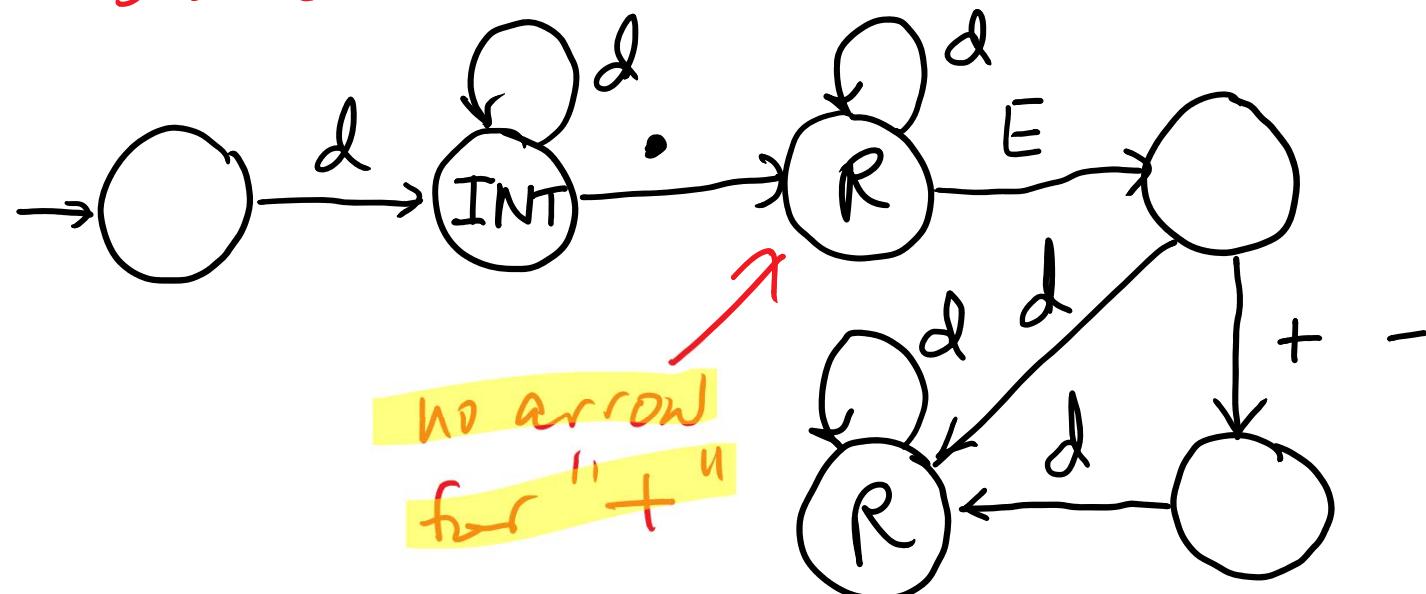
Input: 1 2 . 3 4 + 6

last char

longest so far:
"12.34"

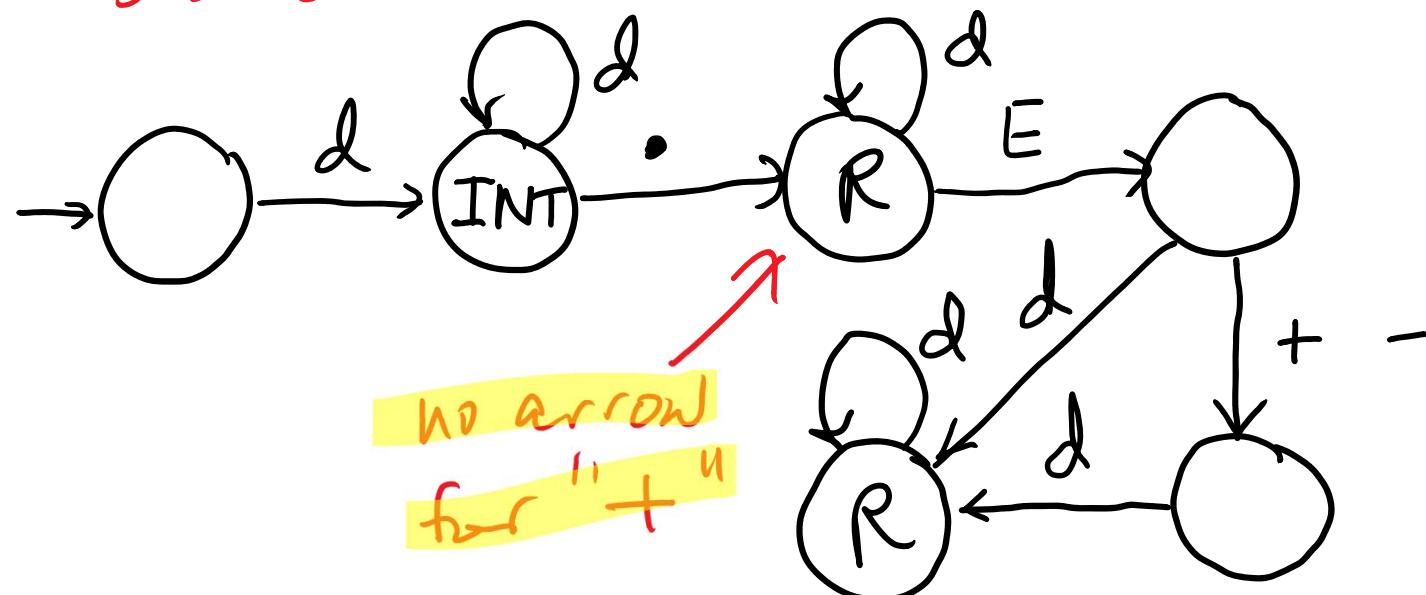
start char

next char



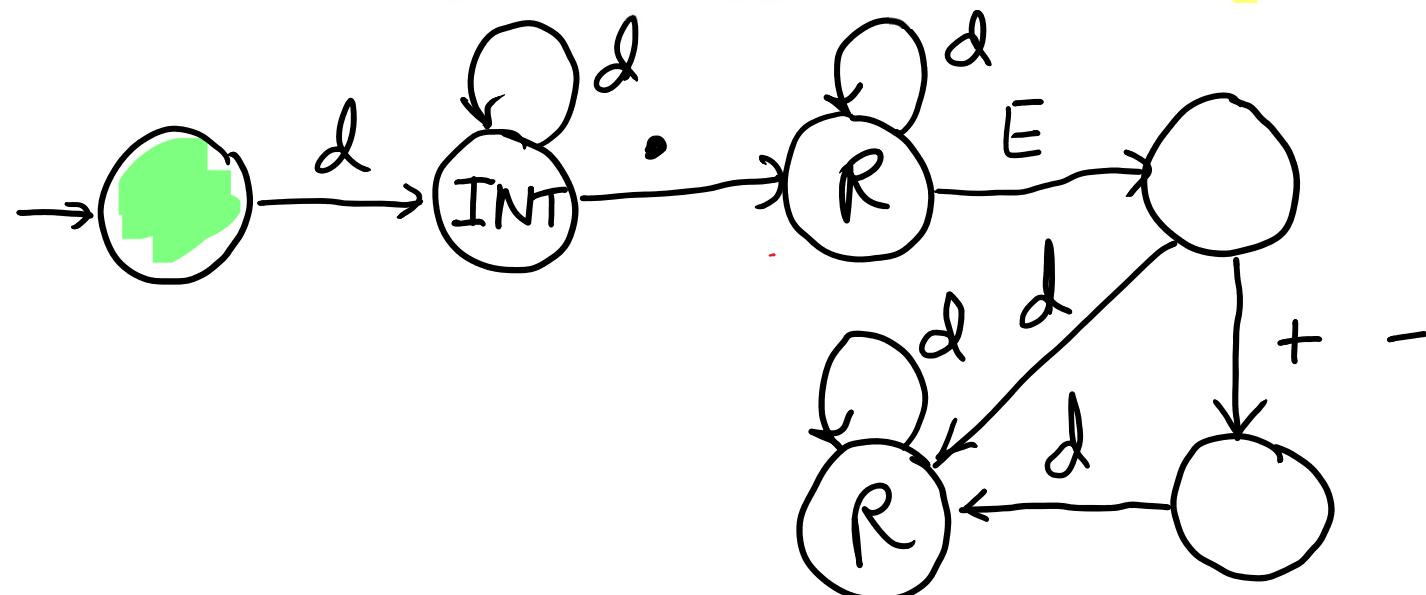
Input: 1 2 . 3 4 + 6
 startchar →
 lastchar ↓
 nextchar ←

longest ~~so far~~:
 "12.34"



Input: 1 2.3 4+6
startchar nextchar

Ready for next lexeme



Input: 12.34+6

startchar \nearrow

lastchar \downarrow

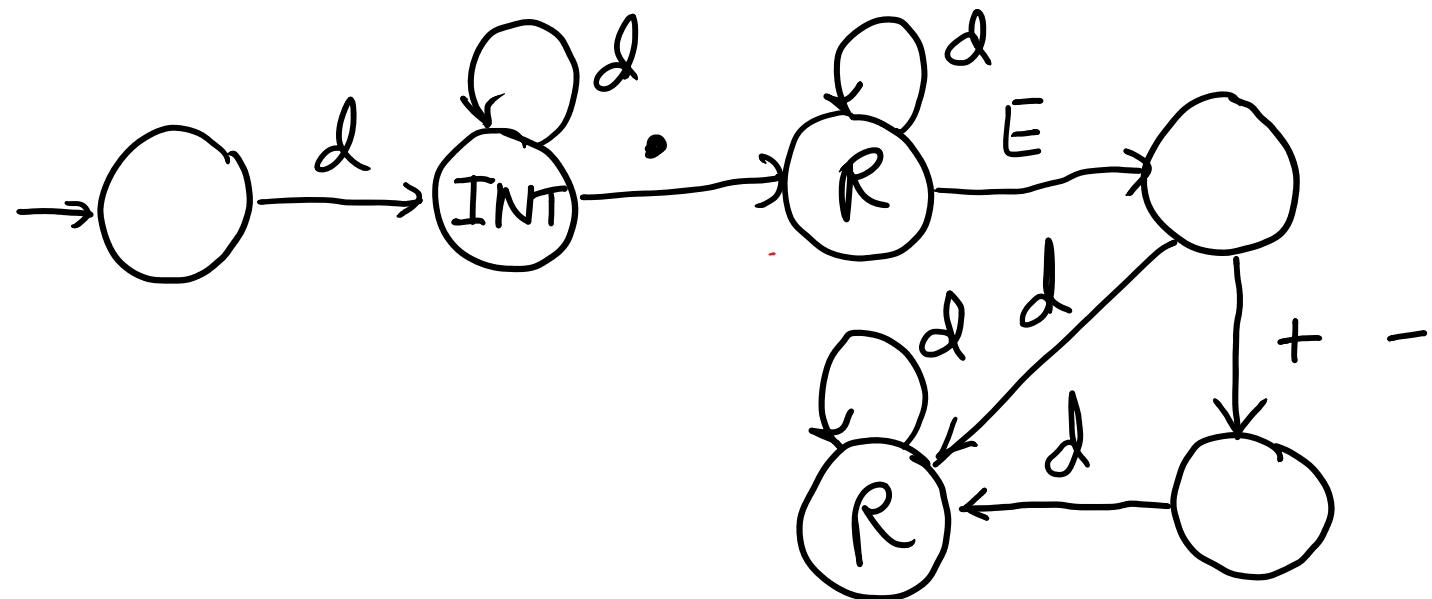
nextchar \nwarrow

longest ~~so far~~: "12.34"

"Look ahead" – had to read
"+" to see that end of lexeme
was previous char.

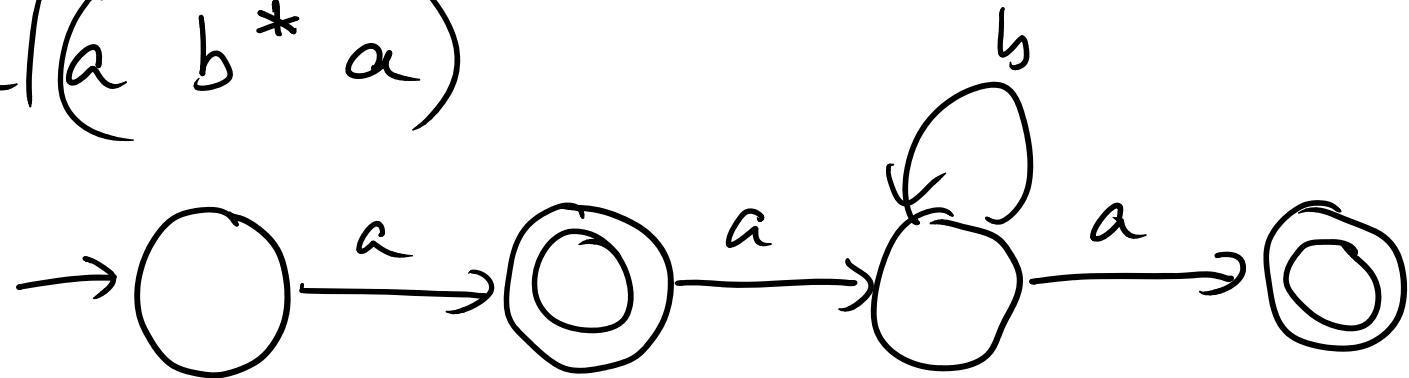
Lookahead = max(nextchar - lastchar)

Lookahead for this example is at least 1. Is it more?



Unbounded lookahead?

$a(a\ b^* a)$



Input: aa b b b b b b b ... b

↑

last char

↑

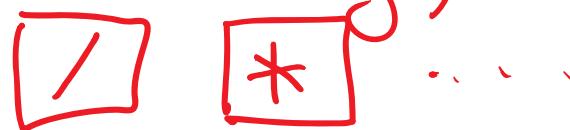
next char

May need to read whole file to find first lexeme!

C Comments

```
/* ... */
```

Has unbounded lookahead if handled
with a single pattern.

If closing "*/" is missing, tokenizes
the comment : 

Saves string unnecessarily

Start Conditions

A better way to handle comments

When you see "/*"

Use new DFA

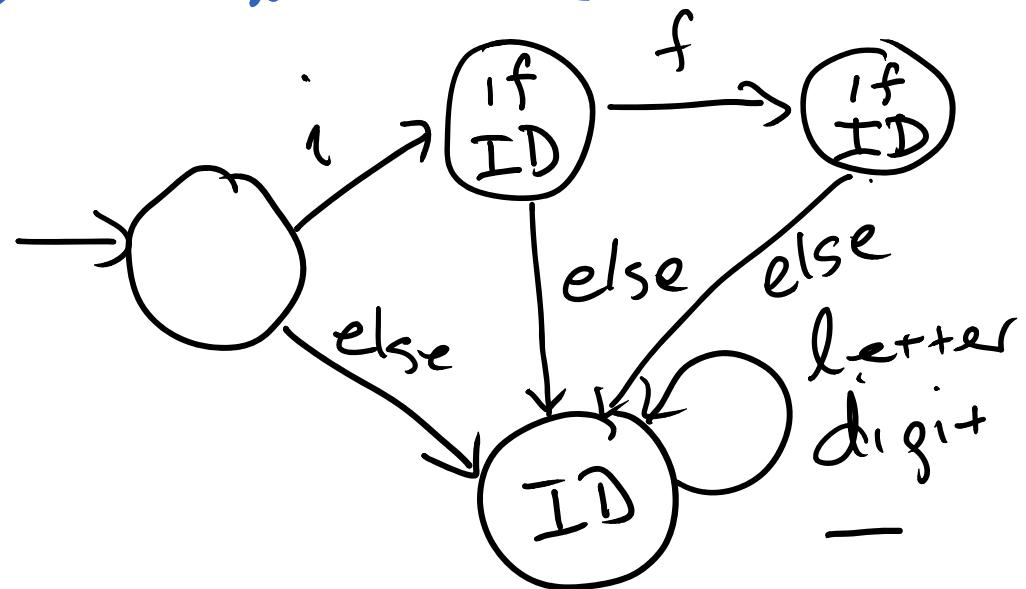
Patterns:

- any single char - discard
- "*/" - discard, return to main DFA

Rule Ordering

if if

ID letter(letter | digit | -)

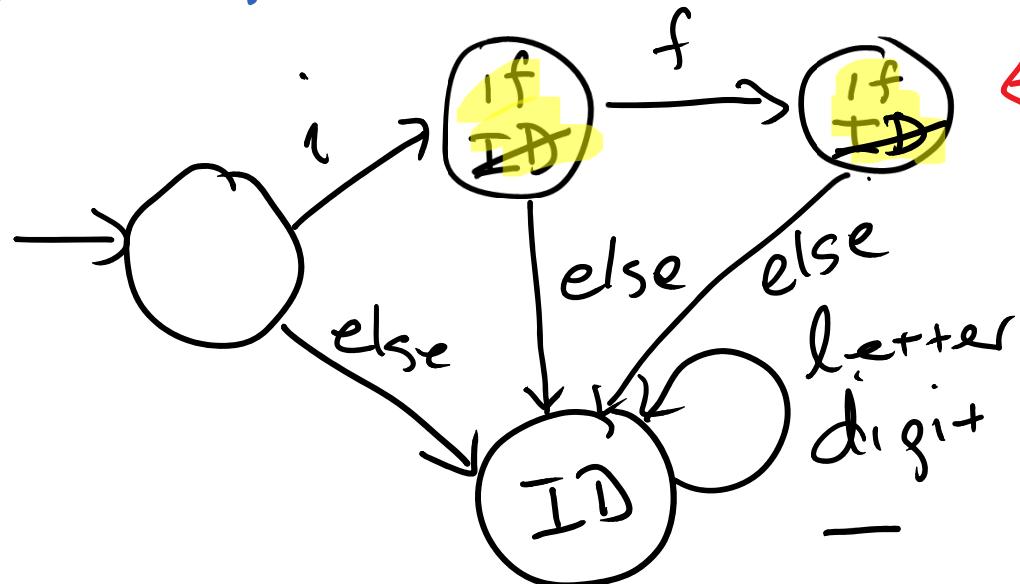


Rule Ordering

if if

ID

letter(letter | digit | -)



← label for
first rule in
file wins.

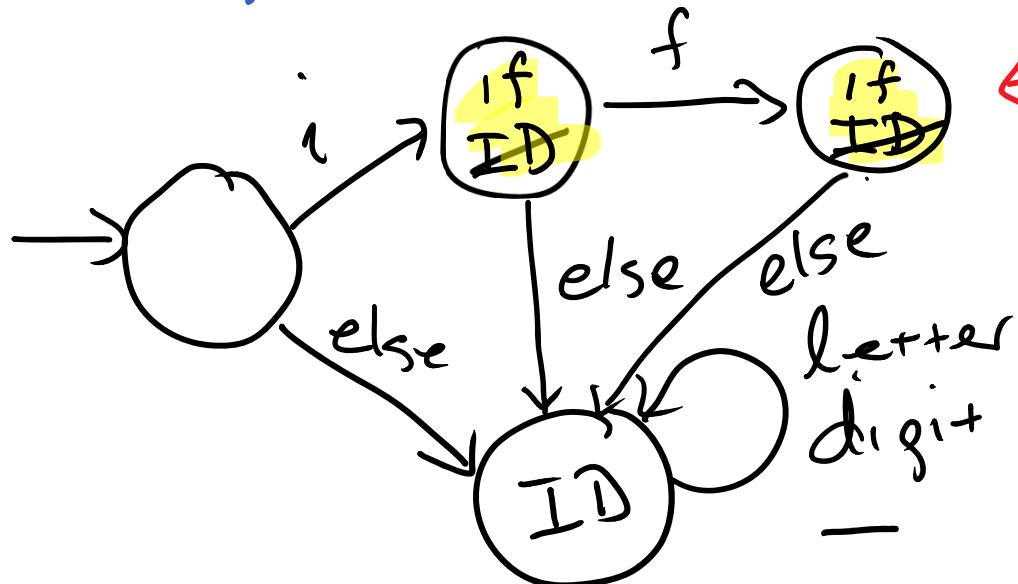
Rule Ordering

if

if ← special case first

ID

letter (letter | digit | -) ← default



← label for
first rule in
file wins.