

# Compiler Construction

## Lecture 2: Lexing

Jeremy Yallop

`jeremy.yallop@cl.cam.ac.uk`

Lent 2023

## Lexing



Regexes

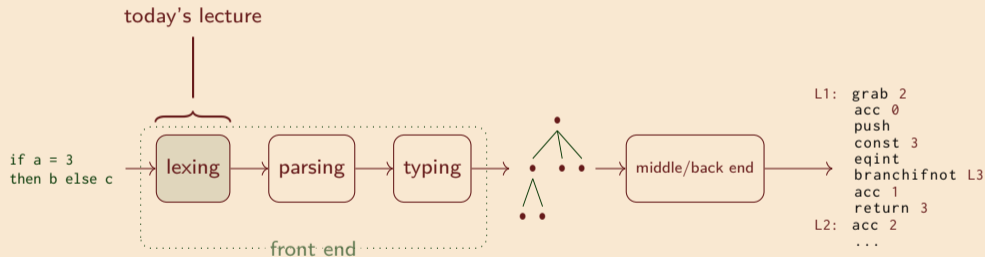
NFA, DFA

RE → NFA

NFA → DFA

Lexing  
(reprise)

∂



## Lexing



Regexes

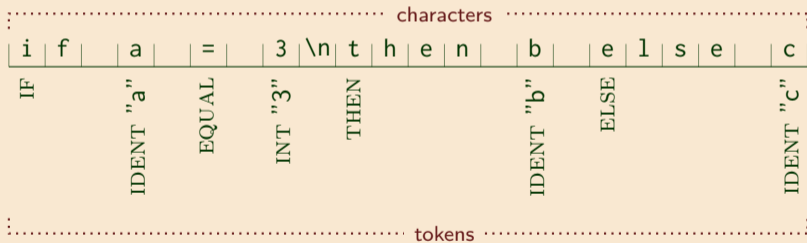
NFA, DFA

RE → NFA

NFA → DFA

Lexing  
(reprise)

Lexing converts a sequence of characters into a sequence of tokens.



# What do lexers look like?

Lexing



Regexes

NFA, DFA

RE → NFA

NFA → DFA

Lexing  
(reprise)

A *lexer* is typically specified as a sequence mapping regexes to tokens:

regular expressions	if	⇒	IF	tokens
	then	⇒	THEN	
	else	⇒	ELSE	
	=	⇒	EQUAL	
	[a-zA-Z]+ as s	⇒	IDENT s	
	[0-9]+ as i	⇒	INT i	
	[ \t\n]	⇒	<i>skip</i>	

Token data type:

```
type token =  
  INT of int  
  | IDENT of string  
  | EQUAL  
  | IF  
  | THEN  
  | ELSE  
  | ...
```

Today's Q: how can we turn this declarative specification into a program?

# Regular expressions

(“regexes”)

Lexing

Regexes



NFA, DFA

RE → NFA

NFA → DFA

Lexing  
(reprise)

∂

Regular expressions  $e$  over alphabet  $\Sigma$  are written:

$$e \rightarrow \emptyset \mid \epsilon \mid a \mid e \vee e \mid ee \mid e^* \quad (a \in \Sigma)$$

A regular expression  $e$  denotes a language (set of strings)  $L(e)$ . For example,

$$L((a \vee b)^* abb) = \{abb, \\ aabb, \\ babb, \\ aaabb, \\ ababb, \\ baabb, \\ bbabb, \\ aaaabb, \\ \dots\}$$

# The regular language problem

Lexing

Regexes



NFA, DFA

RE  $\rightarrow$  NFA

NFA  $\rightarrow$  DFA

Lexing  
(reprise)

The  $L(-)$  function can be defined inductively:

$$L(e) \subseteq \Sigma^*$$

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

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

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

$$L(e_1 \vee e_2) = L(e_1) \cup L(e_2)$$

$$L(e_1 e_2) = \{w_1 w_2 \mid w_1 \in L(e_1), w_2 \in L(e_2)\}$$

$$L(e^0) = \{\epsilon\}$$

$$L(e^{n+1}) = L(ee^n)$$

$$L(e^*) = \bigcup_{n \geq 0} L(e^n)$$

The **regular language problem**: is  $w \in L(e)$ ? This is insufficient for lexing.

# Finite-state automata



# An NFA example

Lexing

Regexes

**NFA, DFA**

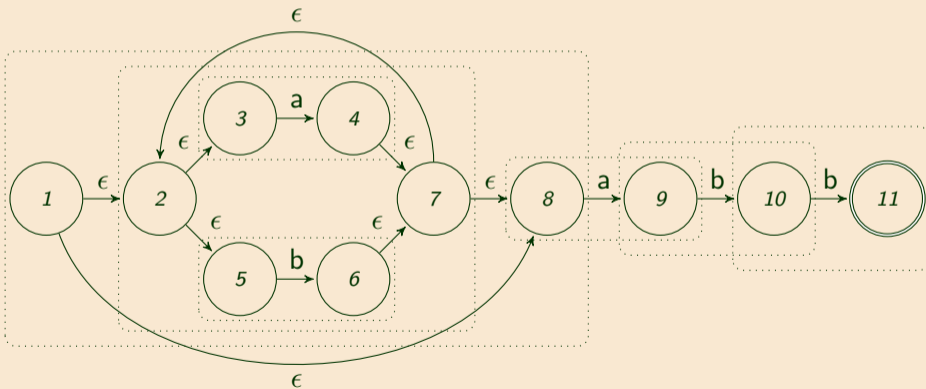


RE  $\rightarrow$  NFA

NFA  $\rightarrow$  DFA

Lexing  
(reprise)

A nondeterministic finite-state automaton for recognising  $L((a \vee b)^* abb)$ :



# Review of Finite Automata (FA)

Lexing

Regexes

NFA, DFA

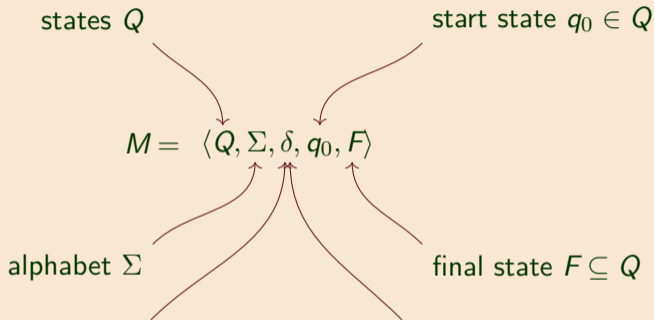


RE  $\rightarrow$  NFA

NFA  $\rightarrow$  DFA

Lexing  
(reprise)

$\partial$



For NFAs:

**(nondeterministic)**

$\forall q \in Q$

$\forall a \in (\Sigma \cup \{\epsilon\})$

$\delta(q, a) \subseteq Q$

For DFAs:

**(deterministic)**

$\forall q \in Q$

$\forall a \in \Sigma$

$\delta(q, a) \in Q$

Lexing

Regexes

**NFA, DFA**



RE  $\rightarrow$  NFA

NFA  $\rightarrow$  DFA

Lexing  
(reprise)

Notation for DFAs:

$$q \xrightarrow{\epsilon} q$$

$$q_1 \xrightarrow{aw} q_3 \text{ if } \delta(q_1, a) = q_2 \text{ and } q_2 \xrightarrow{w} q_3$$

$$L(M) = \{w \mid \exists q \in F, q_0 \xrightarrow{w} q\}$$

Notation for NFAs:

$$q \xrightarrow{\epsilon} q$$

$$q_1 \xrightarrow{w} q_3 \text{ if } q_2 \in \delta(q_1, \epsilon) \text{ and } q_2 \xrightarrow{w} q_3$$

$$q_1 \xrightarrow{aw} q_3 \text{ if } q_2 \in \delta(q_1, a) \text{ and } q_2 \xrightarrow{w} q_3$$

$$L(M) = \{w \mid \exists q \in F, q_0 \xrightarrow{w} q\}$$

Regular expressions  $\longrightarrow$  NFAs

Lexing

$N(-)$  takes a regex  $e$  to an NFA  $N(e)$  accepting  $L(e)$  with a single final state.



Regexes

NFA, DFA

$N(-)$  is defined by induction on  $e$ .

Regex  $\rightarrow$  NFANFA  $\rightarrow$  DFALexing  
(reprise)

Lexing

Regexes

NFA, DFA

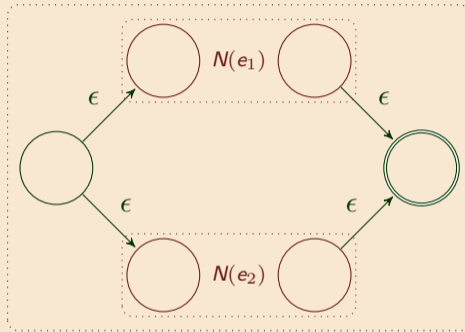
Regex  $\rightarrow$  NFA



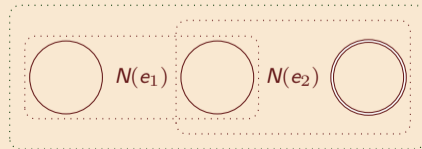
NFA  $\rightarrow$  DFA

Lexing  
(reprise)

$$N(e_1 \vee e_2) =$$



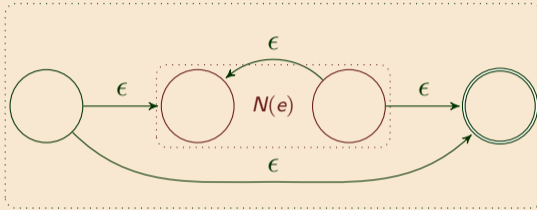
$$N(e_1 e_2) =$$



Lexing

Regexes

NFA, DFA

Regex  $\rightarrow$  NFANFA  $\rightarrow$  DFALexing  
(reprise) $N(e^*) =$ 

*Note: an alternative to this simple construction is Glushkov's (1961) algorithm, which produces an equivalent automaton without the  $\epsilon$  transitions.*

NFAs  $\longrightarrow$  DFAs



Lexing

The *powerset construction* takes a NFA

$$M = \langle Q, \Sigma, \delta, q_0, F \rangle$$

and constructs an DFA

$$M' = \langle Q', \Sigma', \delta', q'_0, F' \rangle$$

where the components of  $M'$  are calculated as follows:

$$Q' = \{S \mid S \subseteq Q\}$$

$$\delta'(S, a) = \epsilon\text{-closure}(\{q' \in \delta(q, a) \mid q \in S\})$$

$$q'_0 = \epsilon\text{-closure}\{q_0\}$$

$$F' = \{S \subseteq Q \mid S \cap F \neq \emptyset\}$$

and the  $\epsilon$ -closure is:

$$\epsilon\text{-closure}(S) = \{q' \in Q \mid \exists q \in S, q \xrightarrow{\epsilon} q'\}$$

Regexes

NFA, DFA

RE  $\rightarrow$  NFANFA  $\rightarrow$  DFALexing  
(reprise)

# How do we compute $\epsilon$ -closure(S)?

Lexing

Regexes

NFA, DFA

RE  $\rightarrow$  NFA

**NFA  $\rightarrow$  DFA**



Lexing  
(reprise)

$\epsilon$ -closure:

push all elements of S onto a stack

result := S

while stack not empty

pop  $q$  off the stack

for each  $u \in \delta(q, \epsilon)$

if  $u \notin$  result

then result :=  $\{u\} \cup$  result

push  $u$  on stack

return result

(NB: just an instance of transitive closure)

# DFA(N((a ∨ b) \* abb))

Lexing

Regexes

NFA, DFA

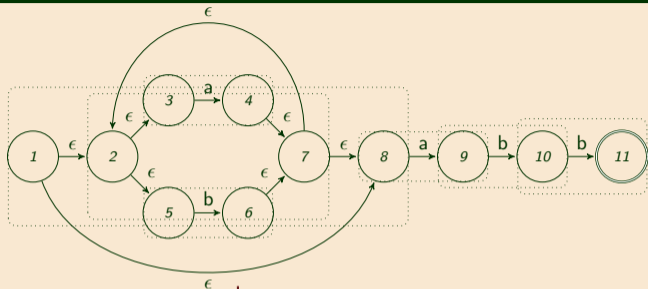
RE → NFA

NFA → DFA

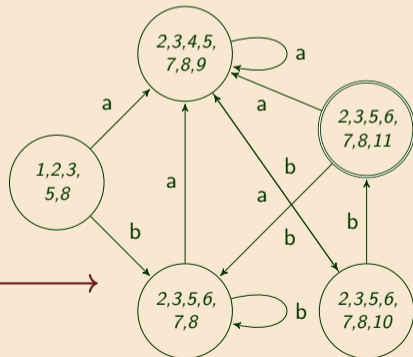
● ● ●

Lexing  
(reprise)

∅



powerset construction

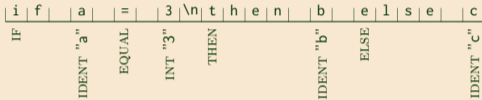


# The lexing problem

# The lexing problem

Lexing

The **regular language problem** (i.e. “is  $w \in L(e)$ ?”) is **insufficient** for lexing.  
We need to tokenize a string using a lexer specification



if  $\Rightarrow$  IF  
...  
[a-zA-Z]+ as s  $\Rightarrow$  IDENT s  
[0-9]+ as i  $\Rightarrow$  INT i  
[ \t\n]  $\Rightarrow$  skip

taking into account that

**Expressions are ordered by priority.**

(treat if as a keyword because the IF rule comes before the IDENT rule)

We should find the **longest match.**

(treat ifif as a variable, not two keywords)

We should **skip whitespace.**

(because whitespace is irrelevant to the parser)

Lexing  
(reprise)



# Define tokens with regexes (automata)

Lexing

Regexes

NFA, DFA

RE → NFA

NFA → DFA

Lexing  
(reprise)

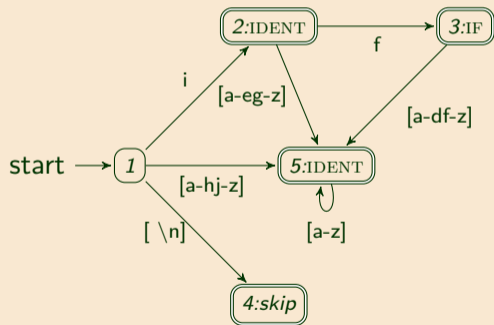
Regex	Finite automaton	Token
if	<pre> graph LR     1((1)) -- i --&gt; 2((2))     2 -- f --&gt; 3(((3)))             </pre>	IF
then	<pre> graph TD     1((1)) -- t --&gt; 2((2))     2 -- h --&gt; 3((3))     3 -- e --&gt; 4((4))     4 -- n --&gt; 5(((5)))     2 --&gt; 2             </pre>	THEN
$[a-zA-Z][a-zA-Z0-9]^*$	<pre> graph LR     1((1)) -- "[a-zA-Z]" --&gt; 2(((2)))     2 -- "[a-zA-Z0-9]" --&gt; 2             </pre>	IDENT s
$[0-9][0-9]^*$	<pre> graph LR     1((1)) -- "[0-9]" --&gt; 2(((2)))     2 -- "[0-9]" --&gt; 2             </pre>	INT n
$[\ \backslash t \backslash n ]$	<pre> graph LR     1((1)) -- "[ \t\n]" --&gt; 2(((2)))             </pre>	<i>skip</i> (not really a token)

# Constructing a Lexer

Lexing

Input:  $e_1 \Rightarrow t_1, e_2 \Rightarrow t_2, \dots, e_k \Rightarrow t_k$ , priority-ordered lexing rules, highest first  
 $\Rightarrow$  Tagged NFA for  $e = e_1 \vee e_2 \vee \dots \vee e_k$   
 $\Rightarrow$  DFA with each accepting state tagged for the  $e_i$  of highest priority.

Regexes



lexer rules		
if	$\Rightarrow$	IF
$[a-z]^+ as s$	$\Rightarrow$	IDENT s
$[\ \backslash n]$	$\Rightarrow$	skip

State 3 could be either an IDENT or the keyword IF. The priority rule eliminates this ambiguity, associating state 3 with the keyword.

Lexing  
(reprise)



# What about longest match?

Lexing

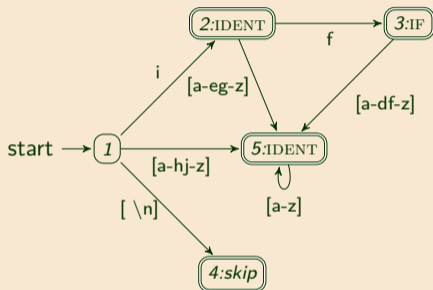
Regexes

NFA, DFA

RE → NFA

NFA → DFA

Lexing  
(reprise)



Start in initial state, and repeatedly:

1. read input until failure (no transition).  
Emit tag for last accepting state.
2. reset state to start state

(| = current position, \$ = EOF)

Input	Current state	Last accepting state	Action
if ifx\$	1	⊥	
i f ifx\$	2	2	
if i ifx\$	3	3	
if ifx\$	⊥	3	emit IF
if i ifx\$	1	⊥	reset
if iifx\$	4	4	
if i fx\$	⊥	4	skip
if iifx\$	1	⊥	reset
if i fx\$	2	2	
if if x\$	3	3	
if ifx \$	5	5	
if ifx \$	⊥	5	emit IDENT "ifx"



Lexing with derivatives

# Matching with derivatives

Lexing

Brzozowski (1964)'s formulation of regex matching, based on *derivatives*.

Regexes

**Derivative of regex  $r$  w.r.t. character  $c$**  is another regex  $\partial_c r$  that matches  $s$  iff  $r$  matches  $cs$ .

NFA, DFA

E.g.: consider  $(b \vee c)^+$ . After matching  $c$ , can accept either  $\epsilon$  or more  $b/c$ , so:

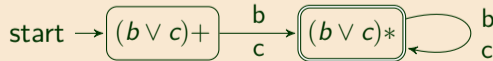
RE  $\rightarrow$  NFA

$$\partial_c (b \vee c)^+ = \epsilon \vee (b \vee c)^+ = (b \vee c)^*$$

NFA  $\rightarrow$  DFA

Construct DFA for  $r$ , taking regexes  $r$  as states, adding transition  $r_i \xrightarrow{c} r_j$  whenever  $\partial_c r_i = r_j$ . For example, for  $(b \vee c)^+$ :

Lexing  
(reprise)



NB:  $\partial_c (b \vee c)^* = (b \vee c)^*$ . (Can you see why?) Also:  $\epsilon$ -matching states are accepting.

Lexing

$\partial_c$  is defined inductively over regexes.

Regexes

Can you see the similarities with derivatives of numerical functions?

(Hint: read  $r_1 r_2$  as  $r_1 \times r_2$  and  $r_1 \vee r_2$  as  $r_1 + r_2$ .)

NFA, DFA

$$\partial_c \emptyset = \emptyset$$

RE  $\rightarrow$  NFA

$$\partial_c \epsilon = \emptyset$$

$$\partial_c b = \emptyset$$

$$\partial_c c = \epsilon$$

NFA  $\rightarrow$  DFA

$$\partial_c (rs) = (\partial_c r)s \mid \nu(r)(\partial_c s) \qquad \nu(r) = \epsilon \text{ if } \epsilon \in L(r)$$

$$\partial_c (r \vee s) = \partial_c r \vee \partial_c s \qquad = \emptyset \text{ if } \epsilon \notin L(r)$$

Lexing  
(reprise)

$$\partial_c r^* = (\partial_c r)r^*$$

More information: *Regular-expression derivatives re-examined* (Owens et al, 2009).

# Lexing with derivatives

Lexing

Lexers match input string against multiple regexes in parallel. Automaton for matching one token; each state corresponds to vector of regexes, one per lexer rule.  $\partial_c$  acts pointwise on the regex vector.

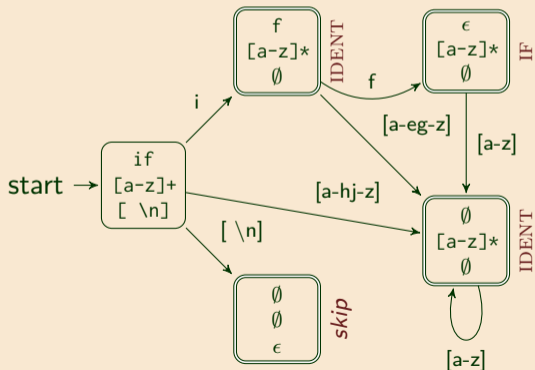
Regexes

NFA, DFA

RE  $\rightarrow$  NFA

NFA  $\rightarrow$  DFA

Lexing  
(reprise)



$\partial$



Next time: context-free grammars