

CSE 443  
Compilers

Dr. Carl Alphonse  
alphonse@buffalo.edu  
343 Davis Hall

# Announcements

- Teams should be formed
- Starting Friday, sit with your teammates
- Attendance sheets will be by team
- Team meetings will start next week

# Phases of a compiler

## Lexical structure

Symbol Table

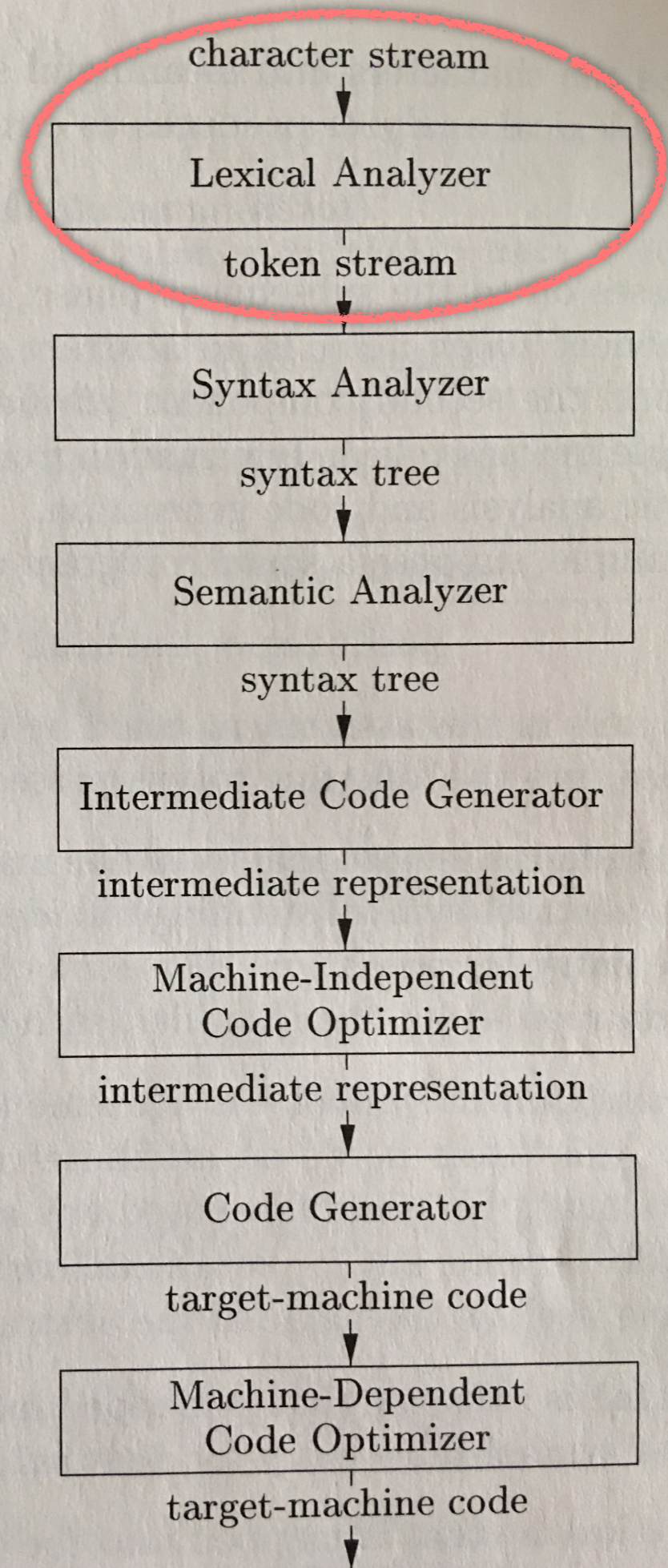


Figure 1.6,  
page 5 of text

# Phases of a compiler

## Lexical structure

Symbol Table

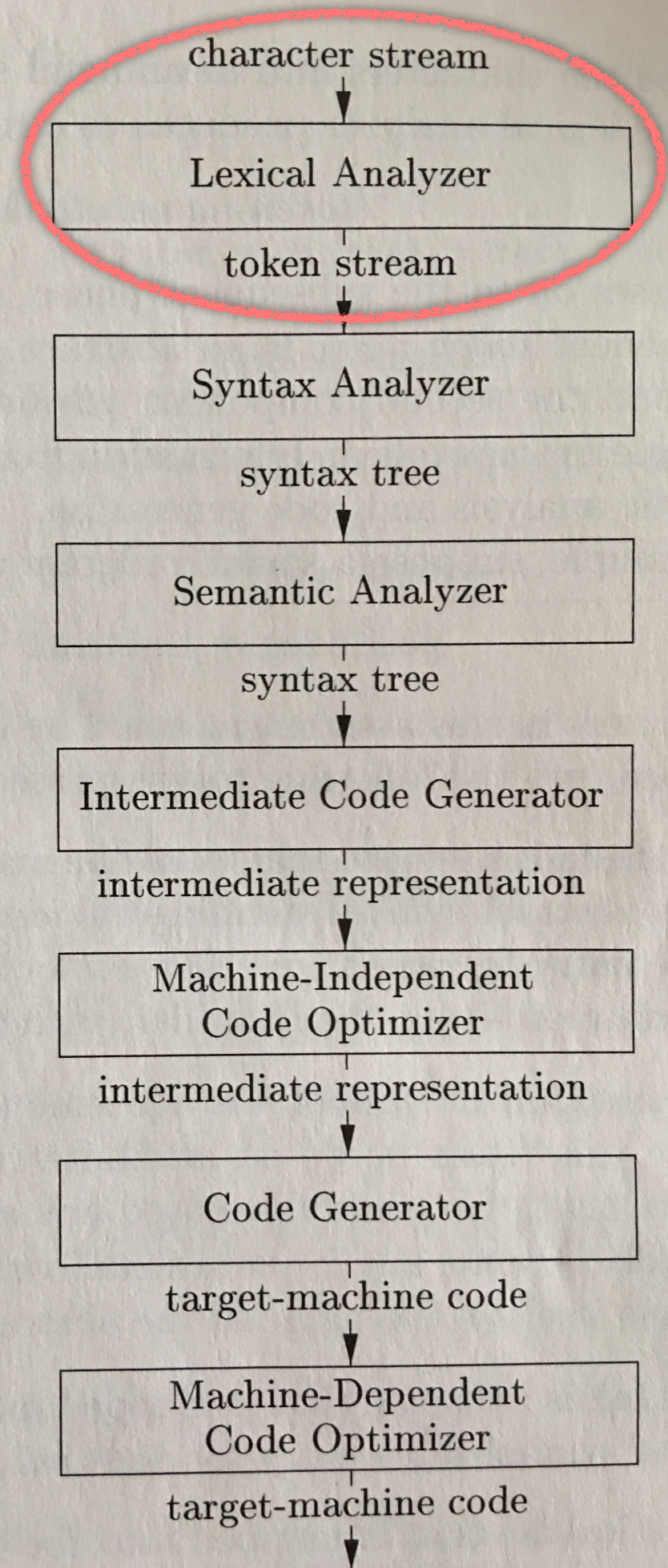


Figure 1.6,  
page 5 of text

# Languages & grammars

$\mathcal{L}(G)$  is the set of all strings derivable from  $G$  starting with the start symbol; i.e. it denotes the language of  $G$ .

# Languages & grammars

Given a grammar  $G$  the language it generates,  $\mathcal{L}(G)$ , is unique.

Given a language  $L$  there are many grammars  $H$  such that  $\mathcal{L}(H) = L$ .

# Languages & grammars

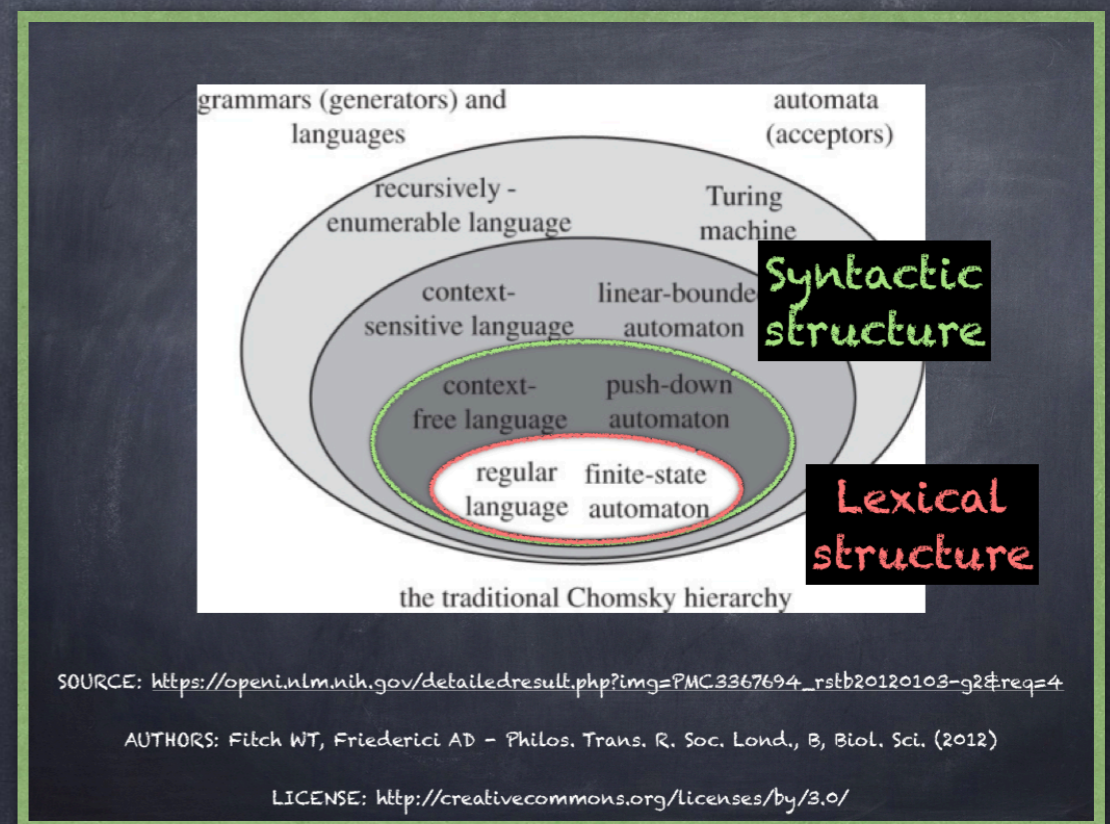
Think about what this means for us: there is no single "correct" grammar for a language.

In fact, grammars for users vs. tool writers vs. compiler writers can all be different.

Given a language  $L$  there are many grammars  $H$  such that  $\mathcal{L}(H) = L$ .

# Lexical Analysis

- Lexical structure described by regular grammar
- Deterministic finite state machine performs analysis





# How is a regular language defined?

- Recall that a language is a set of strings. This set can be finite or infinite.
- The possible regular languages over a given alphabet are defined inductively - construction given on next two slides.

# LANGUAGE operations

## base cases

- $\{\varepsilon\}$  is a regular language
- $\forall a \in \Sigma, \{a\}$  is a regular language

$\varepsilon$  is the empty string

# LANGUAGE operations

If  $L$  and  $M$  are regular, so are:

- $L \cup M = \{ s \mid s \in L \text{ or } s \in M \}$  UNION
- $LM = \{ st \mid s \in L \text{ and } t \in M \}$  CONCATENATION
- $L^* = \bigcup_{i=0, \infty} L^i$  KLEENE CLOSURE
- No other languages are regular

$L^i$  is  $L$  concatenated with itself  $i$  times:

$L^0 = \{\epsilon\}$ , by definition

$L^1 = L$

$L^2 = LL$

$L^3 = LLL$ , etc.

$L^*$  is the union of all these sets!

# Example of $L^*$

Suppose  $L$  is  $\{a, bb\}$

$L^0 = \{\epsilon\}$ , by definition

$L^1 = L = \{a, bb\}$

$L^2 = LL = \{aa, abb, bba, bbbb\}$

$L^3 = LLL = \{aaa, aabb, abba, abbbb, bbaa, bbabb, bbbba, bbbbbb\}$

$L^4 =$

...and so so...

$L^* = \bigcup_{i=0, \infty} L^i = \{\epsilon, a, bb, aa, abb, bba, bbbb, aaa, aabb, abba, abbbb, bbaa, bbbba, bbaa, bbabb, bbbba, bbbbbb, abbbb, bbabb, \dots\}$

Some regular languages over  $\Sigma = \{0,1\}$

The base cases yield these regular languages:

$\{\epsilon\}, \{0\}, \{1\}$

The inductive cases yield many more. Some are:

$\{0, 1\}, \{01\}, \{10\}, \{01, 10\}, \{0, 01\}, \{1, 01\}, \{0, 10\}, \{1, 10\}, \{0, 1, 01\}, \{0, 1, 10\}, \{0, 01, 10\}, \{1, 01, 10\}, \{00\}, \{000\}, \{0000\}, \{11\}, \{111\}, \{1111\},$  and many many more.

Can you demonstrate how each of these is regular?

# Why use grammars?

- Recall that a language is a possibly infinite set of strings.
- A grammar gives us a way to describe, using finite means, an infinite set.
- Regular expressions are equivalent to regular grammars in expressive power: both regular grammars and regular expressions describe regular languages.
- If  $X$  is a regular expression,  $\mathcal{L}(X)$  denotes the set of strings recognized by  $X$ .

# Inductive definition of REGular EXPRESSIONS (regex) over a given alphabet $\Sigma$

$\epsilon$  is a regex

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

For each  $a \in \Sigma$ ,  $a$  is a regex

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

# Regular expressions (regex)

## Inductive definition

Assume  $r$  and  $s$  are regexes.

$r|s$  is a regex denoting  $\mathcal{L}(r) \cup \mathcal{L}(s)$

$rs$  is a regex denoting  $\mathcal{L}(r)\mathcal{L}(s)$

$r^*$  is a regex denoting  $(\mathcal{L}(r))^*$

$(r)$  is a regex denoting  $\mathcal{L}(r)$

Precedence: KLEENE CLOSURE  $>$  CONCATENATION  $>$  UNION

Associativity: all left-associative (minimize use of

parentheses:  $(r|s)|t = r|s|t$  )



# Algebraic Laws

Assume  $r$  and  $s$  are regexes.

COMMUTATIVITY  $r|s = s|r$

ASSOCIATIVITY  $r|(s|t) = (r|s)|t$  and  $r(st) = (rs)t$

DISRIBUTIVITY  $r(s|t) = rs|rt$  and  $(s|t)r = sr|tr$

IDENTITY  $\epsilon r = r\epsilon = r$

IDEMPOTENCY  $r^{**} = r^*$

We can describe a  
regular language  
using a  
regular expression

# Why do we care?

- We will be using a tool called FLEX to construct a lexical analyzer (a lexer) for the programming language we're constructing a compiler for.
- If we give FLEX a regular expression describing the lexical structure of our language, FLEX will produce a C program which acts as our lexer.
- The next step for us to understand (at a high level) how FLEX converts a regex to a C program.

A regular expression can be implemented using a finite state machine.

Finite state machines can be deterministic or non-deterministic:

DFA

deterministic finite automaton

NFA

non-deterministic finite automaton

# Process of building lexical analyzer

1) spell out the language

Language

# Process of building lexical analyzer

2) formulate a regular expression



# Process of building lexical analyzer

3) build an NFA



# Process of building lexical analyzer

4) transform NFA to DFA





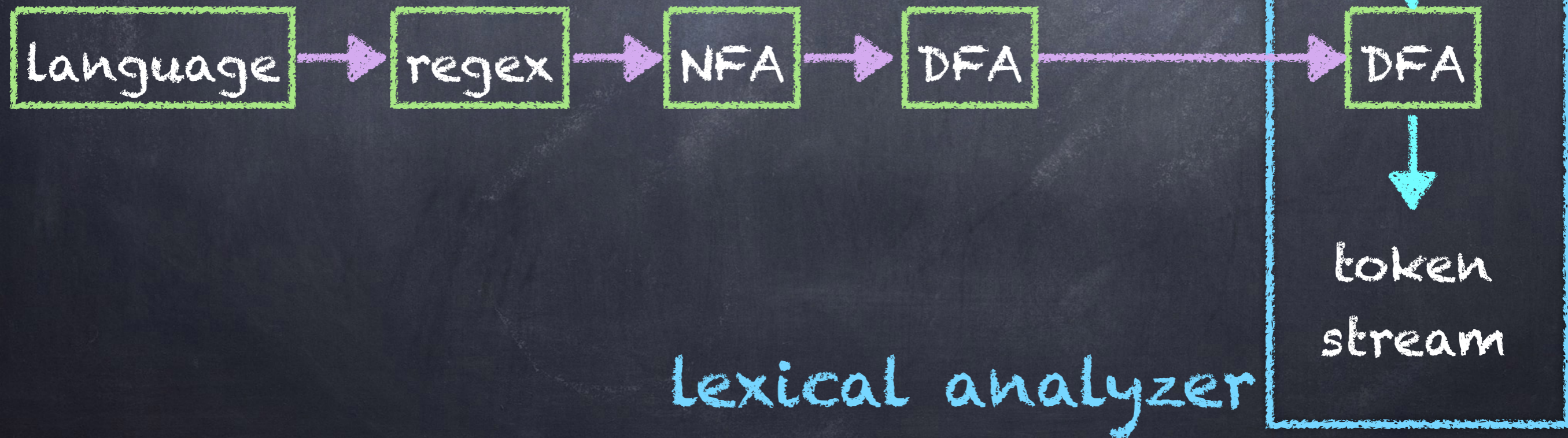
# Process of building lexical analyzer

5) transform DFA to a minimal DFA



# Process of building lexical analyzer

5) The minimal DFA is our lexical analyzer



Step 1:

Construct NFA from regex



# Nondeterministic Finite Automata (NFA)

- A finite set of states  $S$
- An alphabet  $\Sigma$ ,  $\epsilon \notin \Sigma$
- $\delta \subseteq S \times (\Sigma \cup \{\epsilon\}) \times \mathcal{P}(S)$  (transition function)
- $s_0 \in S$  (a single start state)
- $F \subseteq S$  (a set of final or accepting states)

# Deterministic Finite Automata (DFA)

- A finite set of states  $S$
- An alphabet  $\Sigma$ ,  $\epsilon \notin \Sigma$
- $\delta \subseteq S \times \Sigma \times S$  (transition function)
- $s_0 \in S$  (a single start state)
- $F \subseteq S$  (a set of final or accepting states)

# NFA vs DFA

transition function

- $\delta \subseteq S \times (\Sigma \cup \{\epsilon\}) \times \mathcal{P}(S)$
- $\delta \subseteq S \times \Sigma \times S$

no  $\epsilon$ -transitions

no multiple transitions

A state is a circle with its  
state number written  
inside.



Initial state has an arrow from nowhere pointing in. State 0 is often the initial state.





A final (accepting) state is drawn with a double circle.

①

Arrows are labeled  
with  $\epsilon \dots$

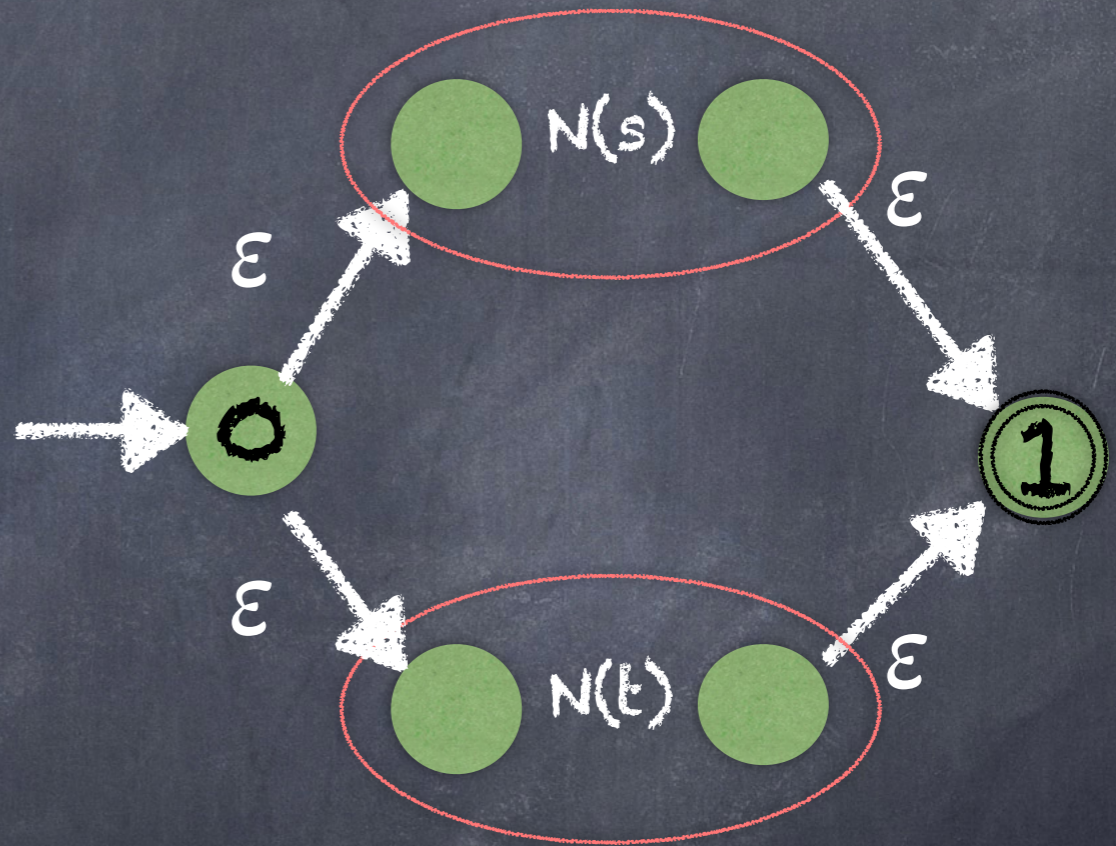


... or  $a \in \Sigma$ .



for each  $a \in \Sigma$

# Regex $\rightarrow$ NFA

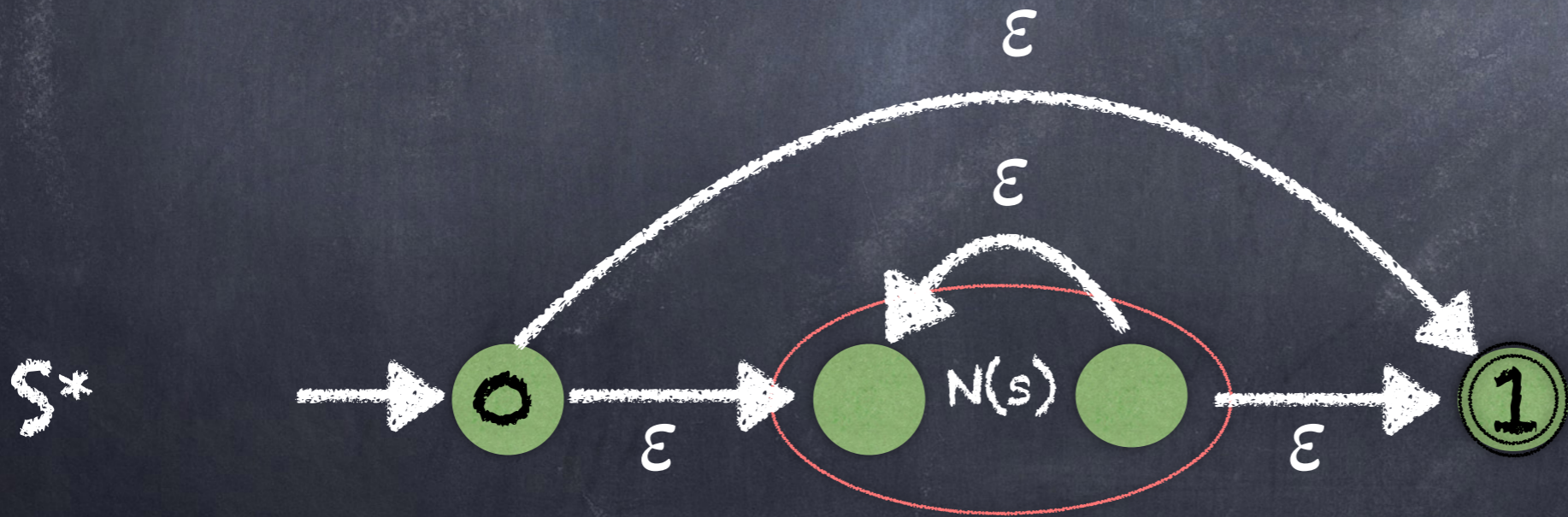
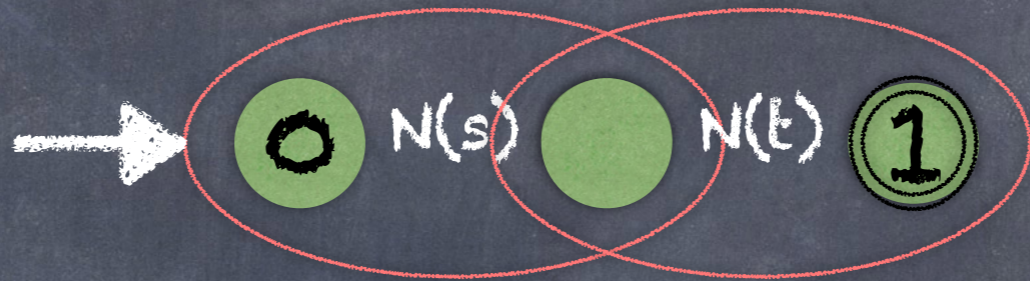


for each  $a \in \Sigma$

$s | t$

# Regex $\rightarrow$ NFA

st



Simple example

static

# Simple example

static



# Simple example

static  
struct

