

Lecture 12 - June 12

Lexical Analysis, Syntactic Analysis

***Regular Expression to Minimized DFA
Implementing a Scanner
Parser in Context
Context Free Grammar (CFG): Intro***

Announcements/Reminders

- **Assignment 1** due next Monday (June 16)
- **Assignment 2** to be released (June 16)
- Required for A2 and Project: ANTLR4 tutorial videos
- Review Slides on Math posted

Pre-study

↳

EECS7030:

Polymorphism/dynamic binding

↳ Composite/visitor design patterns
(ANTLR4)

Tracing Minimizing DFA

$P = T$ means no further opt. possible.
 $T \rightarrow$ last refined partition

$\checkmark \rightarrow$ maximal splittable proper subset of states

start of 1st it.

①

$P \rightarrow$ refined partitions so far

$\{ \{s_0, s_1, s_2, s_4\}, \{s_3, s_5\} \}$

\emptyset

\perp

$\{ \{s_0, s_1, s_2, s_4\}, \{s_3, s_5\} \}$ P

$\{ \{s_2, s_4\}^S, \{s_0, s_1\}^P, \{s_3, s_5\}^S \}$
 $\rightarrow P-S$
 \rightarrow not splittable

$\{ \{s_2, s_4\}^S \}$
 S

$\{ \{s_2, s_4\}, \{s_0, s_1\}, \{s_3, s_5\} \}$

$P \neq T$
 \rightarrow more opt. might be possible!

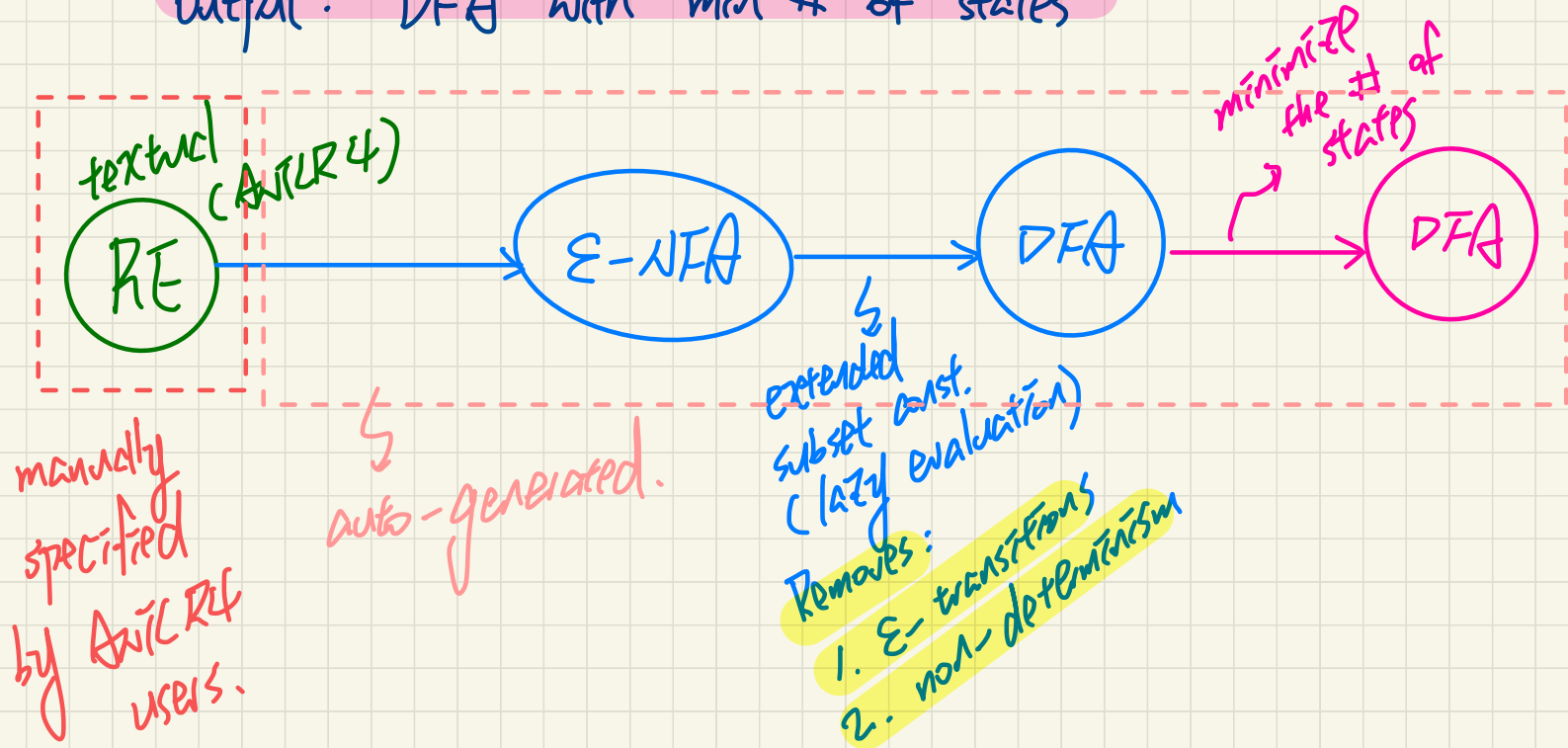
\emptyset

\perp

②

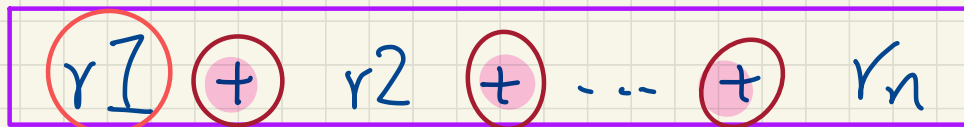
Input: Regular Expression

Output: DFA with min # of states



Scanner

n syntactic categories



$r.e.$
for a
particular
category of tokens

a word
recognized/accepted by
either of the patterns
is returned as a **token**

From RE to Scanner (1)

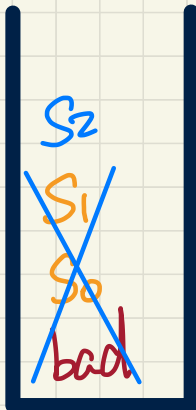
Token Type (CharCat)

r	0, 1, 2, ..., 9	EOF	Other
Register	Digit	Other	Other

Transition

δ of min. DFA.

	Register	Digit	Other
S₀	S₁	S _e	S _e
S₁	S _e	S₂	S _e
S₂	S _e	S ₂	S _e
S_e	S _e	S _e	S _e



Token Type (Type)

S ₀	S ₁	S ₂	S _e
invalid	invalid	register	invalid

Regular Expression: r[0..9]+

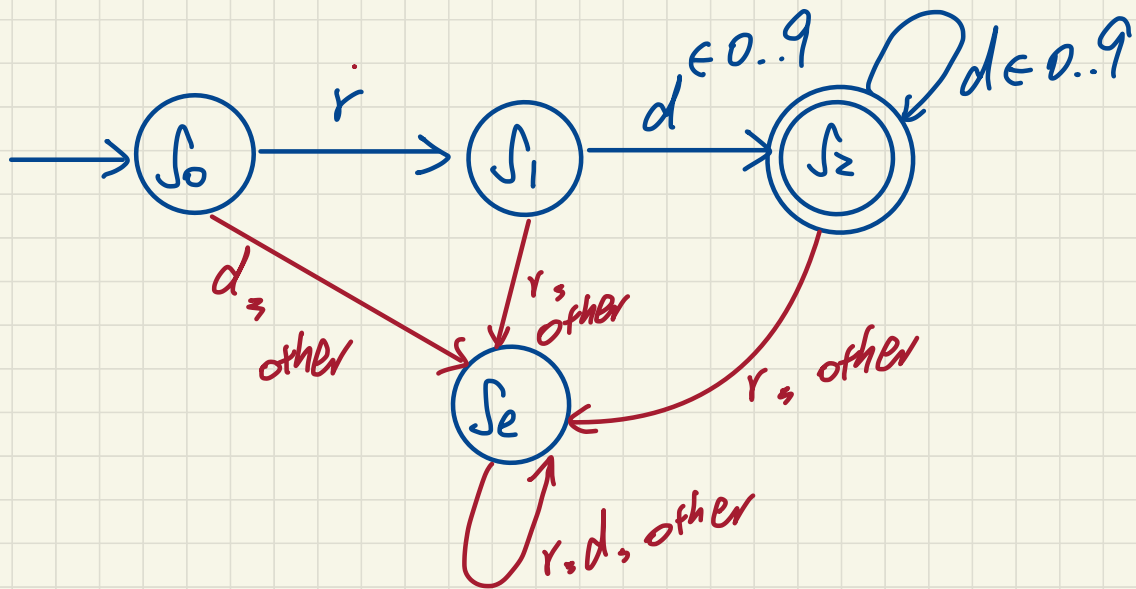
```

NextWord()
-- Stage 1: Initialization
state := S0 ; word := ε
initialize an empty stack S : S.push(bad)
-- Stage 2: Scanning Loop
while (state ≠ Se)
    NextChar(char) ; word := word + char
    if state ∈ F then reset stack S end
    S.push(state)
    cat := CharCat[char]
    state := δ[state, cat]
-- Stage 3: Rollback Loop
while (state ∉ F ∧ state ≠ bad)
    state := S.pop()
    truncate word
-- Stage 4: Interpret and Report
if state ∈ F then return Type[state]
else return invalid
end
    
```

Example input: r2_{EOF}

eventually return r2 as a register

word: r2
 state: S₀ S₁ S₂
 cat: Register Digit



From RE to Scanner (2)

Token Type (CharCat)

r	0, 1, 2, ..., 9	EOF	Other
Register	Digit	Other	Other

Transition

	Register	Digit	Other
s ₀	s ₁	s _e	s _e
s ₁	s _e	s ₂	s _e
s ₂	s _e	s ₂	s _e
s _e	s _e	s _e	s _e

Token Type (Type)

s ₀	s ₁	s ₂	s _e
invalid	invalid	register	invalid

Regular Expression: r[0..9]+

```
NextWord()
-- Stage 1: Initialization
state := s0 ; word := ε
initialize an empty stack s ; s.push(bad)
-- Stage 2: Scanning Loop
while (state ≠ se)
  NextChar(char) ; word := word + char
  if state ∈ F then reset stack s end
  s.push(state)
  cat := CharCat[char]
  state := δ[state, cat]
-- Stage 3: Rollback Loop
while (state ∉ F ∧ state ≠ bad)
  state := s.pop()
  truncate word
-- Stage 4: Interpret and Report
if state ∈ F then return Type[state]
else return invalid
end
```

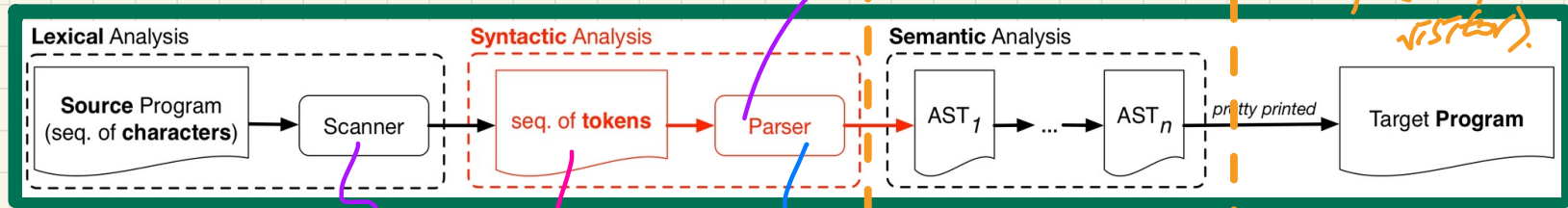
Example input: r24*3

word:

state:

cat:

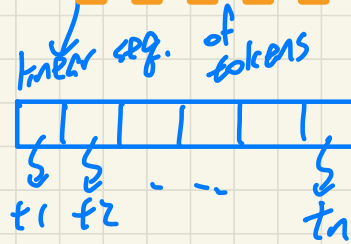
Parser in Context



$r1 + r2 + \dots + rn$

each token has a category (e.g. Register, Digit)

non-linear ASTs



top-down parsing (or) bottom-up parsing

1. CFG
2. derivations
3. PTs

focus: tokens rather than characters.

design patterns (composers, visitor).

Context-Free Grammar (CFG): Terminology

The following language that is **non-regular**

$$\{0^n \# 1^n \mid n > 0\}$$

cannot be expressed/recognized by any RE/NFA/E-NFA/DFA

can be described using a **context-free grammar (CFG)**:

A	→	0A1
A	→	B
B	→	#

Context-free grammar

↳ focus only on recursive pattern of strings

↳ no semantic analysis is done.

non-terminals (variables)

production rule.

start variable

terminals (tokens)

recursive case (more derivation steps expected)

A → 0A1

B → #

base case (terminals only)