

어휘분석기

String \rightarrow tokens



Token

1. 어휘분석 결과는 어떻게 출력되는가?
2. 토큰이 어떤 유형이 있는가?

1. ex) if (a \geq 10)

	if	(a	\geq	10)
	↓	↓	↓	↓	↓	↓
Token number	32	7	4	25	5	8
Token value	0	0	'a'	0	10	0

2.

special form

keyword if while for else

operator symbol + - * < >

Delimiter , ! () []

정규표현 필요 X, 변형할 없음

General form - Programmer

Identifier 변수명

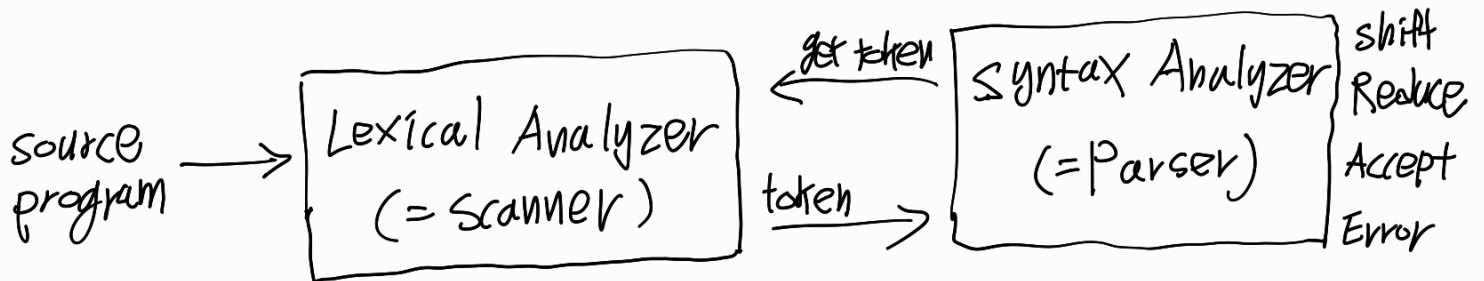
Constant 상수

정규표현으로 기술

ex) $id = (L + _)(L + d + _)^*$

그래서 어휘 분석은 왜 하는가?

구문 분석기 (Syntax Analyzer, parser) 에 token을 넘겨주기 위해서이다



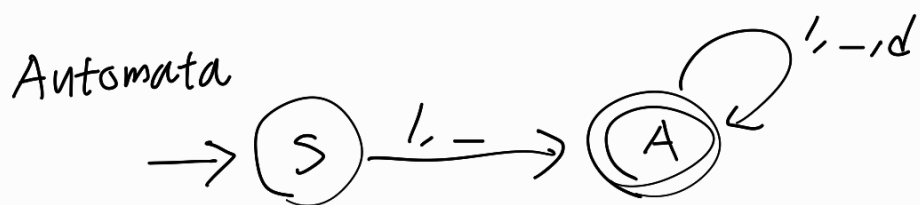
어휘 분석과 파싱 단계는 더욱 쉽고 효율적이기 위해
서로 구분한다

general token의 인식

1. Token description \rightarrow R.E , F.A
2. develop a scanner
3. verify throug regular Language theory

목표는 정규표현을 얻는 것이다

명칭



Regular grammar

$$S \rightarrow 1A \mid _A$$

$$A \rightarrow 1A \mid _A \mid dA \mid \varepsilon$$

Regular expression

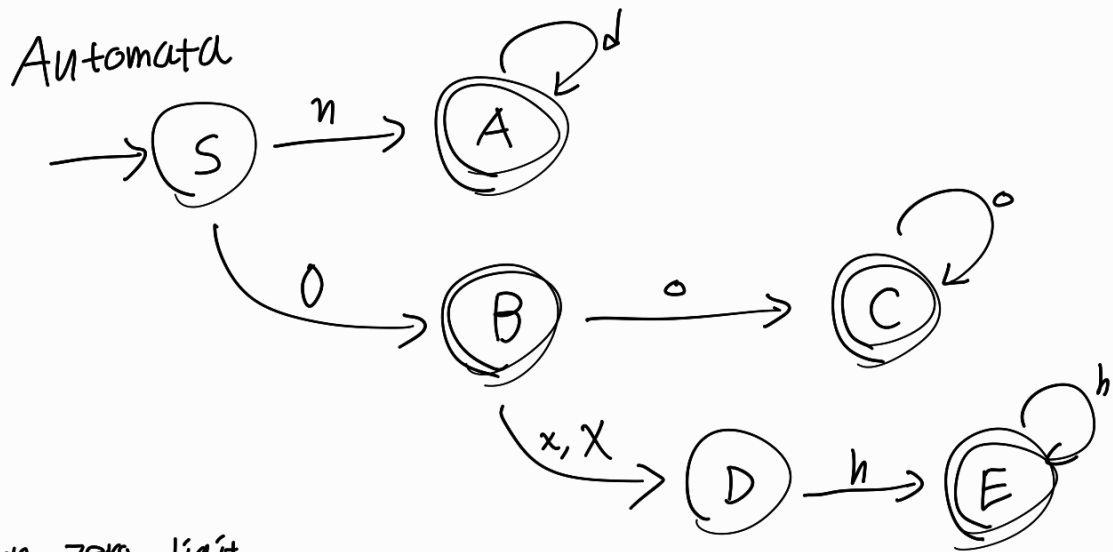
$$S = 1A + _A = (1 + _)A$$

$$\begin{aligned} A &= 1A + _A + dA + \varepsilon = (1 + _ + d)A + \varepsilon \\ &= (1 + _ + d)^* \end{aligned}$$

$$\therefore S = (1 + _)(1 + _ + d)^*$$

정수

10진수 8진수 16진수를 모두 표현해야 한다



n : non-zero digit

o : $0 \sim 7$

h : $0 \sim F$

Regular grammar

$$S \rightarrow nA \mid 0B$$

$$A \rightarrow dA \mid \epsilon$$

$$B \rightarrow 0C \mid xD \mid \epsilon$$

$$C \rightarrow 0C \mid \epsilon$$

$$D \rightarrow hE$$

$$E \rightarrow hE \mid \epsilon$$

Regular expression

$$S = nd^* + 0 + 0o^+ + 0(x+X)h^*$$

실수

Automata

regular grammar

regular expression

String 상수

Automata

regular grammar

regular expression

Comment

Automata

regular grammar

regular expression

이제 얻은 정규표현을 통해서 어휘분석기를 개발할 수 있다

1. 직접 개발

2. 도구 사용 (lex)

도구사용



Lex 입력 파일 구조

정의부

%%

* 규칙부

%%

사용자 서브루틴

정의부

format

name

translation

ex)

D

[0-9]

L

[a-zA-Z]

%%

또다른 기능

%[... %]

^
C언어의 선언문들

include

전역변수 등...

구식부

format

regular expression + actions

ex) integer

color

C언어 코드

[0-9]⁺

사용자 서브루틴

input file에 사용되는 subprogram 정의