

CS344 - Compilers: Scanning

Note Title

1/20/2012

Announcements

- We have moved!
(Surprised me too...)
- Essay is due.
- Look for HW2 tomorrow.

Compilers

The process by which programming languages are turned into assembly or machine code is important in programming languages.

We'll spend some time on these compilers, although it isn't a focus of this class.

Compilers

Compilers are essentially translators, so must semantically understand the code

Output: either assembly, machine code
some other output

Java → bytecode

C++ → assembler

Compilers begin by preprocessing:

- remove white space + comments

- include macros or libraries

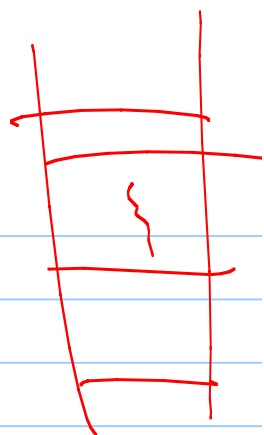
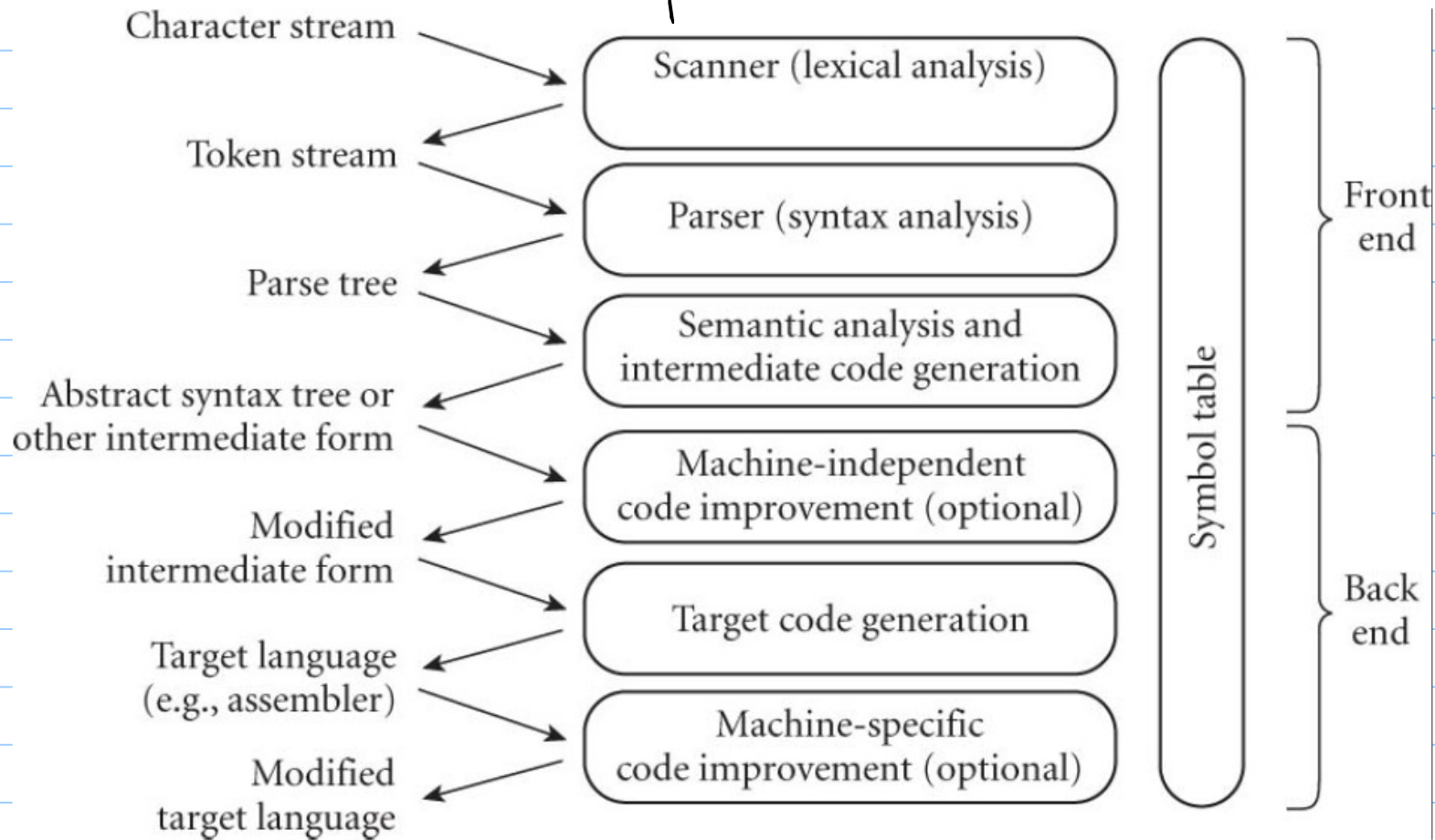
- group characters into tokens

ex:

for (int i = 0 ; i < n ; i ++)

- identify ex: high-level syntactical structures

Overview of Compilation



The steps:

Front end:

Ⓐ Scanner

Ⓑ Parser

Ⓒ Semantic Analysis

Let's dive into these first...

Scanning (lexical analysis)

- Divide program into tokens, or smallest meaningful units

Ex: keywords, {, }, +, variables, etc.

- Scanning + tokenizing makes parsing much simpler.

- While parsers can work character by character, it is slow.

- Note: Scanning is recognizing a regular language, eg via DFA

Parsing

- Recognizing a context-free language,
e.g. via PDA
- Finds the structure of the program
(or the syntax)

Ex: iteration-statement \rightarrow
while (expression) statement

statement \rightarrow compound-statement

Outputs a parse tree.

Semantic Analysis (after parsing)

This discovers the meaning of the commands.

Actually only does static semantic analysis, consisting of all that is known at compile time.

(Some things - eg array out of bounds - are unknown until run time.)

→ error generation

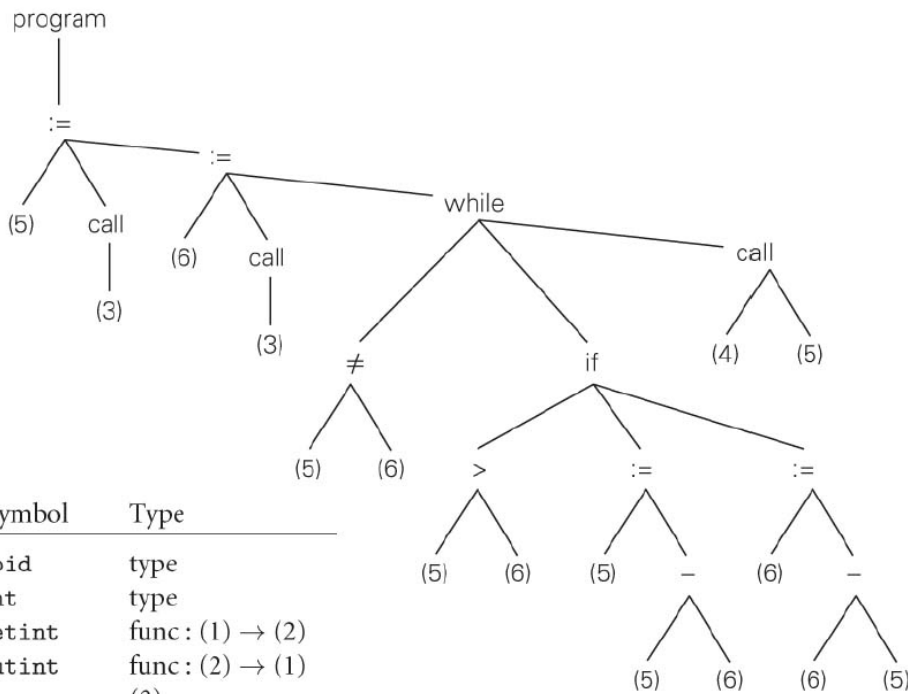
Ex: (semantic analysis)

- Variables can't be used before being declared. (C-like)
- Type checking.
- Identifiers are used in proper context.
- Functions have correct inputs & returns.

etc... (very language dependent)

Intermediate Form

This is the output of the "front end"

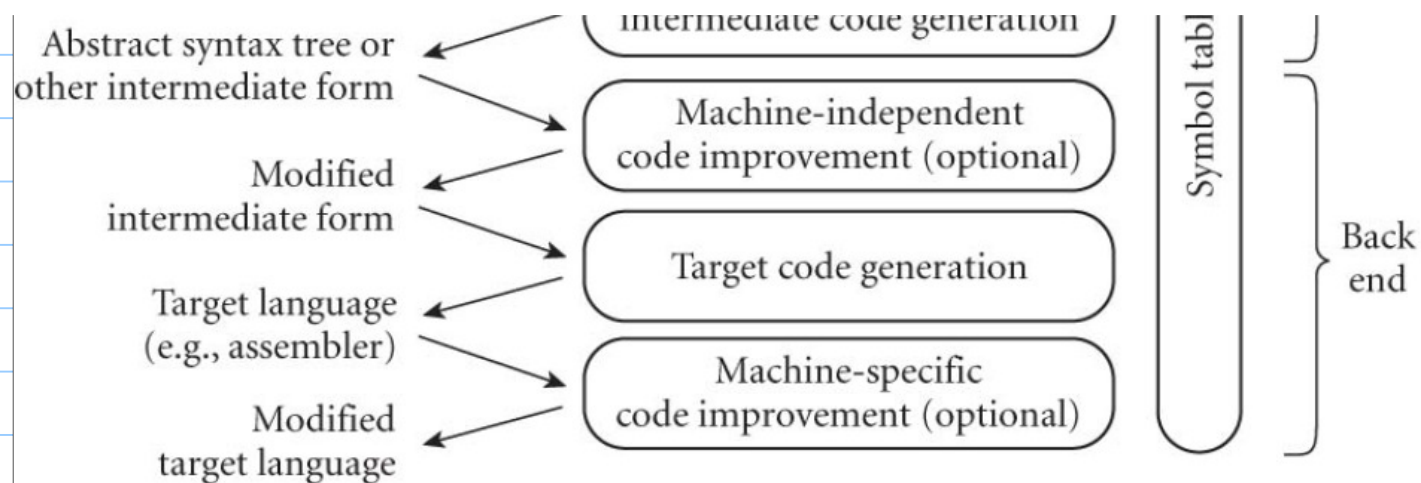


Index	Symbol	Type
1	void	type
2	int	type
3	getint	func: (1) → (2)
4	putint	func: (2) → (1)
5	i	(2)
6	j	(2)

- Often, this is an abstract syntax tree - a simplified version of a parse tree

- May also be a type of assembly-like code

Back end: (Actual code generation)



Creating correct code is generally not difficult.

Optimization of that code is.

Back to front end:

① How is this actually done?

Input is actually a string of
ASCII.

Need to find a way to scan letter
by letter + decide what is
a token.

Then pass the tokens on to
the parser.

Regular Expressions : some theory

Defined as: Language generated as follows:

- A character : 0 or 1
- The empty string, ϵ
- 2 regular expressions concatenated
- 2 regular expressions separated by an or (written |)
- A regular expression followed by * (Kleene star - 0 or more occurrences)

Regular Languages

The class of languages described by a regular expression.

Ex: $0^*10^* = L$

Any number of 0's, followed by a single 1, followed by any # of 0's.

$$1 \in L$$

$$0 \notin L$$

$$001 \in L$$

$$\epsilon \notin L$$

Ex: Give the regular expression for
 $\{w \mid w \text{ begins with a } 1 \text{ and ends with a } 0\}$

$$1 (1|0)^* 0$$

$$(1|0)(1|0)(1|0) \leftarrow \text{any \# of times}$$

Ex: $\{w \mid w \text{ starts with } 0 \text{ and has an odd length}\}$

$$= 0 (11|10|01|00)^*$$

$$= 0 ((1|0)(1|0))^*$$

Example: Numbers in Pascal

digit \rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

unsigned_int \rightarrow digit digit*

unsigned_number \rightarrow
unsigned_int (ϵ | . unsigned_int)
(ϵ | V ((e | E) (+ | -) ϵ) unsigned_int)

\rightarrow digit \rightarrow [0-9]

unsigned_num.
 \rightarrow unsigned_int ϵ ϵ
 \rightarrow digit digit*
 \rightarrow 1

Deterministic Finite Automate (DFA)

Regular languages are precisely the things recognized by DFAs.

- A set of states

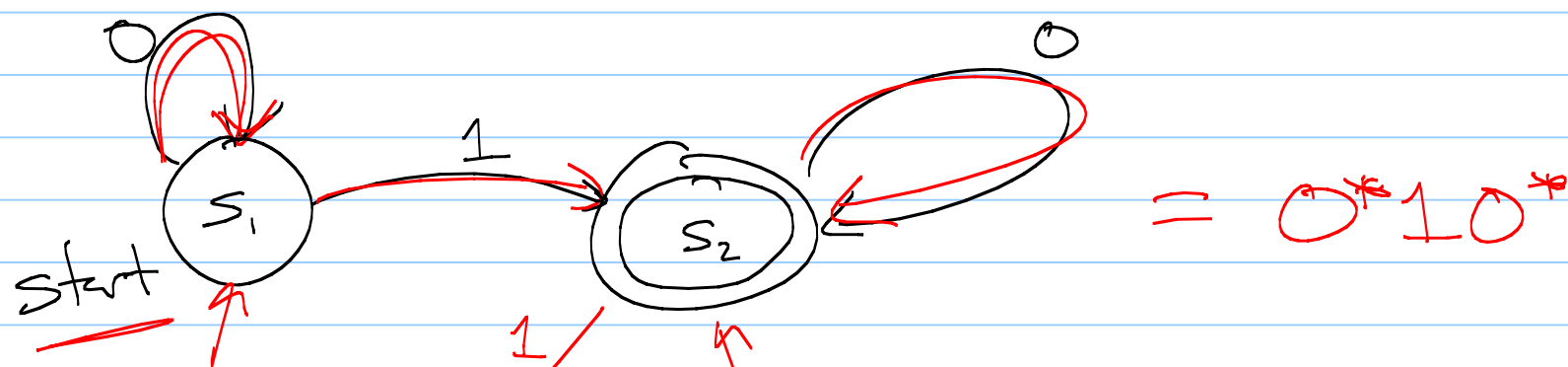
- input alphabet

- A start state

- A set of accept states

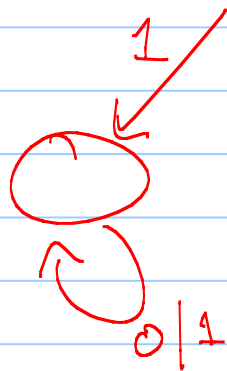
- A transition function: given a state & an input, output a new state

Example: word: 0010 $\in L$ recognized by DFA



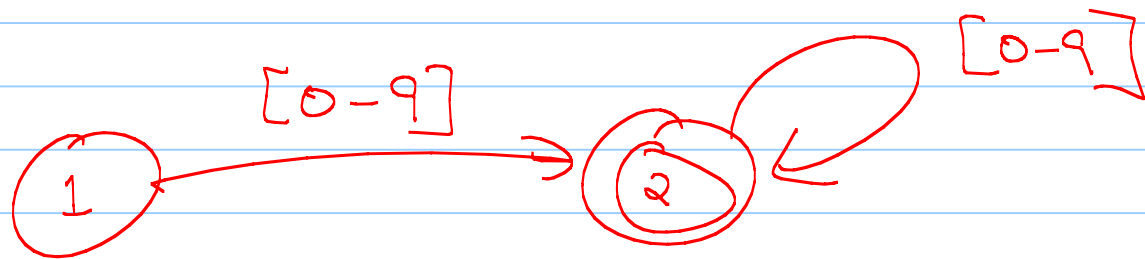
double circle indicates accept state

$00 \notin L$



Ex: unsigned_int \rightarrow digit digit*

digit \rightarrow [0-9]



DFA vs NFA

Non-deterministic Finite Automata

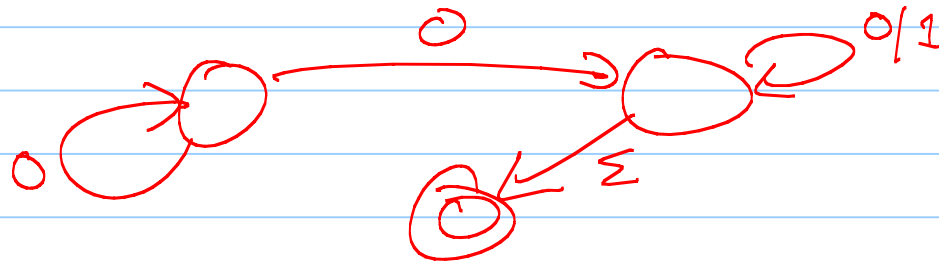
Note: No ambiguity is allowed in DFA's.

So given a state & input, can't be multiple options.

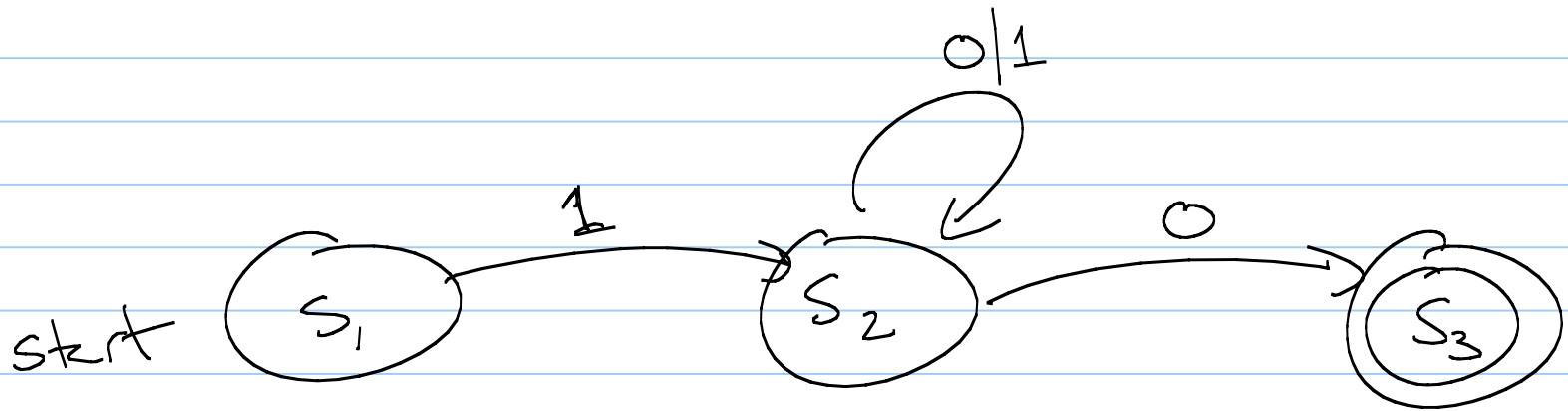
Also - no ϵ -transitions. (in DFA)

If we allow several choices to exist, this is called an NFA.

Ex:



Ex: $L = 1(0/1)^*0$

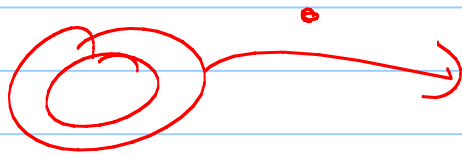


Ex: Some things are easier with NFA!

unsigned_number \rightarrow unsigned_int (ϵ | unsigned_int)

unsigned_int \rightarrow [0-9]

DFA:



Essentially, we can think of an NFA as modeling a parallel set of possibilities (or a tree of them).

Thm: Every NFA has an equivalent DFA.
(size of DFA is much bigger - 2^n)

So: Both recognize reg. languages!

Limitations of Regular Expressions

Certain languages are not regular.

Ex: $\{w \mid w \text{ has an equal number of 0's and 1's}\}$

Somehow, this needs a type of memory, which regular expressions do not have. \cup

Ex: $0^n 1^n$

reg? $00^p 11^p$

\nwarrow (Pumping lemma)

Why do we need this?

Need to "nest" expressions.

Ex: $\text{expr} \rightarrow \text{id} \mid \text{number} \mid \text{-(expr)} \mid \text{expr op expr}$
 $\text{op} \rightarrow + \mid - \mid / \mid *$

non-terminal (pointing to expr)
variable (pointing to -(expr))

Regular expressions can't quite do this.

Context Free Languages

Described in terms of productions
(Called Backus-Naur Form, or BNF)

- A set of terminals T : id, ϵ , +, ...
- A set of non-terminals N $\begin{matrix} N \rightarrow \text{---} \\ \uparrow \end{matrix}$
- A start symbol S (a non-terminal)
- A set of productions

$$\text{Ex: } \{0^n 1^n \mid \underbrace{n > 0}\}$$

$$\left[\begin{array}{l} S \rightarrow 0T1 \\ T \rightarrow 0T1 \mid \epsilon \end{array} \right]$$

$$S \rightarrow 0S1 \mid 01$$

$n > 1$

Ex: $\{ w \mid w \text{ has an equal number of 0's \& 1's} \}$

Expression grammars: Simple calculator

$\text{expr} \rightarrow \text{term} \mid \text{expr} \text{ add_op} \text{ term}$

$\text{term} \rightarrow \text{factor} \mid \text{term} \text{ mult_op} \text{ factor}$

$\text{factor} \rightarrow \text{id} \mid \text{number} \mid -\text{factor} \mid (\text{expr})$

$\text{add_op} \rightarrow + \mid -$

$\text{mult_op} \rightarrow * \mid /$

Example: Show how rules can
generate $3 + 4 \neq 5$

Parse Tree

Ex: $3 + 4 * 5$

