

Compilers 2021/2022: Test 1 Cheat Sheet

Compiler Overview

The frontend

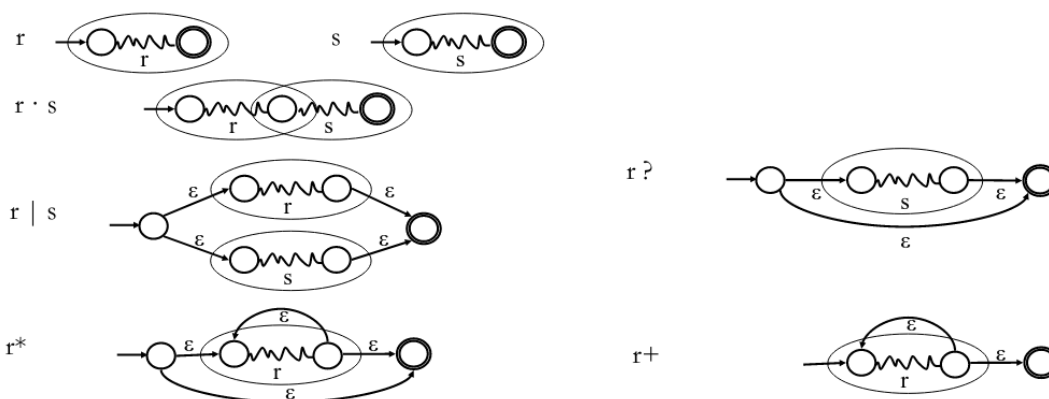
- Scanner
 - Maps character stream into tokens (name and attributes)
- Parser
 - Recognizes context-free syntax and reports error
 - Guides context-free syntax and reports errors
 - Guides context-sensitive ("semantic") analysis (type checking)
 - Builds IR for source program, ie. an AST

Lexical Analysis

Regular expressions

- Lexical patterns form a regular language
- Any finite language is regular
- Recognizable by DFAs

RE to NFA (Thompson Construction)



NFA to DFA (Subset Construction)

DFAedge

Given symbol c and a set of states S , what states can you reach?

$$DFAedge(S, c) = \epsilon\text{-closure}\left(\bigcup_{s \in S} edge(s, c)\right)$$

DFA State	NFA States	ϵ -closure after transition on...		
		0...9	-	.
0	{1, 2}	{2, 3, 4, 8}	{2}	error

DFA State Minimization

- Normalization
 - Assure every state has a transition on every symbol

- Add missing transitions to a trap state
- Algorithm
 - Start with accepting vs non-accepting partitions of states
 - Repartition based on transitions for each symbol: find same partitions for every symbol

DFA to RE (Kleene Construction)

- The sets that take the DFA from state q_i to q_j without going through any state numbered higher than k
- When $k=0$, consider direct transitions
- A dynamic programming approach

$$R_{ij}^k = R_{ik}^{k-1} (R_{kk}^{k-1})^* R_{kj}^{k-1} \mid R_{ij}^{k-1}$$

Syntactical Analysis

Context free grammars

A **context free grammar** $G = (\Sigma, N, S, P)$ is defined by:

- Σ set of *terminal* symbols;
- N set of *non-terminal* symbols;
- $S \in N$ initial symbol;
- P set of *production rules* $X \rightarrow \alpha$ where:
 - ▶ X is non-terminal;
 - ▶ α is a sequence (maybe empty) of terminal or non-terminal symbols

Ambiguity

- A grammar producing same word with different syntax tree
- Eliminate forcing priority and/or associativity

Parsing

Top-down parsing

Recursive descent parsing

- Consume tokens left to right
- Map each non terminal to a function
- Map each production to a different case
- Decide which production to use using the next token

LL Parsing

- Recursive descent parsing technique
- $LL(k)$ means: Left-to-right parse, Leftmost derivation, k -symbols lookahead
- Does not support left recursion

Left recursion removal

$$E \rightarrow E + T$$

$$E \rightarrow T$$

E produces sums of terms, i.e. $E \Rightarrow^* T + T + \dots + T$.

Let us define an equivalent grammar adding a new non-terminal symbol E' :

$$E \rightarrow T E'$$

$$E' \rightarrow + T E'$$

$$E' \rightarrow \varepsilon$$

LL(1): Predictive parsing

- Sufficient for programming languages
- For each non-terminal symbol, the intersection of FIRST sets for each rule must be null (must left factor rules)
- A parsing table maps non-terminals to input and corresponding rule to choose
- Build the table based on NULLABLE, FIRST and FOLLOW
- Rely on the parsing table and an auxiliary stack to parse input

Grammar:

$$S' \rightarrow S\$$$

$$S \rightarrow AB$$

$$A \rightarrow aAb \mid \varepsilon$$

$$B \rightarrow bB \mid \varepsilon$$

Table:

	a	b	$\$$
S'	$S' \rightarrow S\$$	$S' \rightarrow S\$$	$S' \rightarrow S\$$
S	$S \rightarrow AB$	$S \rightarrow AB$	$S \rightarrow AB$
A	$A \rightarrow aAb$	$A \rightarrow \varepsilon$	$A \rightarrow \varepsilon$
B		$B \rightarrow bB$	$B \rightarrow \varepsilon$

We choose a production rule $N \rightarrow \alpha$ on input symbol c if:

1. $c \in \text{FIRST}(\alpha)$, or
2. $\text{Nullable}(\alpha)$ and $c \in \text{FOLLOW}(N)$.

stack	input	action
S'	$aabbb\$$	$S' \rightarrow S\$$
$S\$$	$aabbb\$$	$S \rightarrow AB$
$AB\$$	$aabbb\$$	$A \rightarrow aAb$
$aAbB\$$	$aabbb\$$	consume a
$AbB\$$	$abbb\$$	$A \rightarrow aAb$
$aAbbB\$$	$abbb\$$	consume a
$AbbB\$$	$bbb\$$	$A \rightarrow \varepsilon$
$bbB\$$	$bbb\$$	consume b
$bB\$$	$bb\$$	consume b
$B\$$	$b\$$	$B \rightarrow bB$
$bb\$$	$b\$$	consume b
$B\$$	$\$$	$B \rightarrow \varepsilon$
$\$$	$\$$	consume $\$$
ε	ε	accept

Compute FOLLOW

- FOLLOW does not include ϵ
- $\text{FOLLOW}(S) = \{ \$ \}$
- If $A \rightarrow pBq$ is a production, where p , B and q are any grammar symbols, then everything in $\text{FIRST}(q)$ except ϵ is in $\text{FOLLOW}(B)$.
- If $A \rightarrow pB$ is a production, then everything in $\text{FOLLOW}(A)$ is in $\text{FOLLOW}(B)$.
- If $A \rightarrow pBq$ is a production and $\text{FIRST}(q)$ contains ϵ , then $\text{FOLLOW}(B)$ contains $\{ \text{FIRST}(q) - \epsilon \} \cup \text{FOLLOW}(A)$

Bottom-up parsing

LR Parsing

- LR(k) means: Left-to-right parse, Rightmost derivation (reversed), k symbols looked ahead
- Deals easier with ambiguity and recursion
- Consult the parsing table to parse input using shift, reduce and goto actions
- Read back reductions to get the derivations

	a	b	c	\$	T	R	state	stack	input	action
0	s3	s4	r3	r3	g1	g2	0	ε	aabbbcc\$	shift 3
1				a			3	a	abbbcc\$	shift 3
2			r1	r1			3	aa	bbbcc\$	shift 4
3	s3	s4	r3	r3	g5	g2	4	aab	bbcc\$	shift 4
4		s4	r3	r3		g6	4	aabb	bcc\$	shift 4
5			s7				4	aabbb	cc\$	reduce $R \rightarrow \epsilon$; go 6
6		r4	r4				6	aabbbR	cc\$	reduce $R \rightarrow bR$; go 6
7		r2	r2				6	aabbR	cc\$	reduce $R \rightarrow bR$; go 6
(0)	$T' \rightarrow T \$$						6	aabR	cc\$	reduce $R \rightarrow bR$; go 2
(1)	$T \rightarrow R$						2	aaR	cc\$	reduce $T \rightarrow R$; go 5
(2)	$T \rightarrow aTc$						5	aaT	cc\$	shift 7
(3)	$R \rightarrow \epsilon$						7	aaTc	c\$	reduce $T \rightarrow aTc$; go 5
(4)	$R \rightarrow bR$						5	aT	c\$	shift 7
							7	aTc	\$	reduce $T \rightarrow aTc$; go 1
							1	T	\$	accept

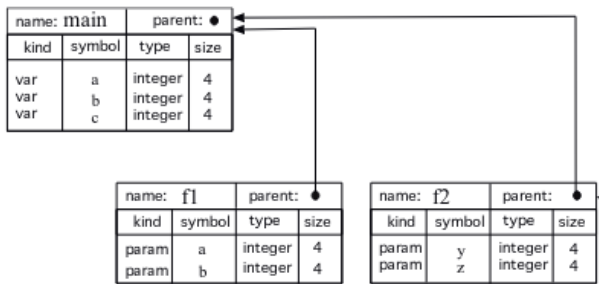
Semantic analysis

Lexical scope (static scope)

- Elements usage correspond to the closest declaration in the AST

Symbol table

- Relates identifiers with semantic information, such as registry location, types and variable values
- Typically implemented with a hash map
- Represent scopes with recursive hash maps



Type checking

- Assert correct function parameter types, variable attribution types
- Generate more efficient code and avoid errors at run-time

Attribute grammars

- Semantic rules for the grammar
- Often implemented with visitor pattern, recursively
- Attributes can be inherited (variable types) or synthesized (types of sub-expressions)

- ▶ Type checking may be made by traversing the AST (one or more times)
- ▶ As the AST is a recursive structure type checking uses recursive functions
- ▶ The compiler builds **node attributes**; examples:
 - ▶ Types;
 - ▶ Symbol Table (**context**)
- ▶ **Synthesized attributes**: bottom-up
- ▶ **Inherited attributes**: top-down

Grammar Rule	Semantic Rules
$decl \rightarrow type\ var\text{-}list$	
$type \rightarrow \text{int}$	$dtype = integer$
$type \rightarrow \text{float}$	$dtype = real$
$var\text{-}list_1 \rightarrow id, var\text{-}list_2$	$insert(id.name, dtype)$
$var\text{-}list \rightarrow id$	$insert(id.name, dtype)$

(Attribute Grammar for Type Declarations)