

WHY LR-PARSING:

1. LR parsers can be constructed to recognize virtually all programming-language constructs for which context-free grammars can be written.
2. The LR parsing method is the most general non-backtracking shift-reduce parsing method known, yet it can be implemented as efficiently as other shift-reduce methods.
3. The class of grammars that can be parsed using LR methods is a proper subset of the class of grammars that can be parsed with predictive parsers.
4. An LR parser can detect a syntactic error as soon as it is possible to do so on a left-to-right scan of the input.

The disadvantage is that it takes too much work to construct an LR parser by hand for a typical programming-language grammar. But there are lots of LR parser generators available to make this task easy.

LR-PARSERS:

LR(k) parsers are most general non-backtracking shift-reduce parsers. Two cases of interest are $k=0$ and $k=1$. LR(1) is of practical relevance.

„L“ stands for “Left-to-right” scan of input.

„R“ stands for “Rightmost derivation (in reverse)”.

K stands for number of input symbols of look-a-head that are used in making parsing decisions. When (K) is omitted, „K“ is assumed to be 1.

LR(1) parsers are table-driven, shift-reduce parsers that use a limited right context (1 token) for handle recognition.

LR(1) parsers recognize languages that have an LR(1) grammar.

A grammar is LR(1) if, given a right-most derivation

$S \Rightarrow r_0 \Rightarrow r_1 \Rightarrow r_2 \dots r_{n-1} \Rightarrow r_n \Rightarrow \text{sentence}$.

We can isolate the handle of each right-sentential form r_i and determine the production by which to reduce, by scanning r_i from left-to-right, going atmost 1 symbol beyond the right end of the handle of r_i .

Parser accepts input when stack contains only the start symbol and no remaining input symbol are left.

LR(0) item: (no lookahead)

Grammar rule combined with a dot that indicates a position in its RHS.

Ex-1: $S^1 \rightarrow .S\$$

$S \rightarrow .$

$x S \rightarrow .(L)$

Ex-2: $A \rightarrow XYZ$ generates 4 LR(0) items

$A \rightarrow .XYZ$

$A \rightarrow X.$

$YZ A \rightarrow XY.$

$Z A \rightarrow XYZ.$

$A \rightarrow XY.Z$ indicates that the parser has seen a string derived from XY and is looking for one derivable from Z.

→ LR(0) items play a key role in the SLR(1) table construction algorithm.

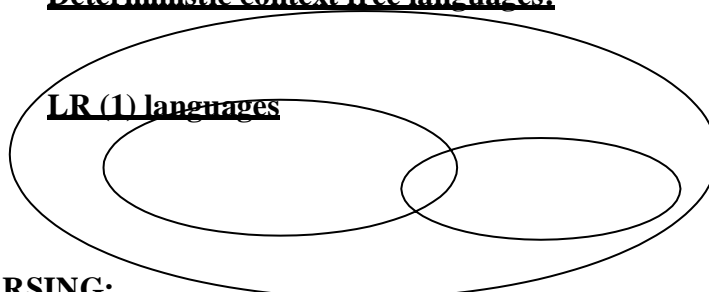
→ LR(1) items play a key role in the LR(1) and LALR(1) table construction algorithms. LR parsers have more information available than LL parsers when choosing a production:

* LR knows everything derived from RHS plus, 'K' lookahead symbols.

* LL just knows, 'K' lookahead symbols into what's derived from RHS.

* **Deterministic context free languages:**

*
*
*
*
*
*
*



LALR PARSING:

Example:

Construct $C = \{I_0, I_1, \dots, I_n\}$ The collection of sets of LR(1) items

For each core present among the set of LR(1) items, find all sets having that core, and replace these sets by their Union# (clus them into a single term)

$I_0 \rightarrow$ same as previous

$I_1 \rightarrow$ “

$I_2 \rightarrow$ “

I_{36} – Clubbing item I_3 and I_6 into one I_{36} item.

$C \rightarrow cC, c/d/\$$

$C \rightarrow cC, c/d/\$$

$C \rightarrow d, c/d/\$$

$I_5 \rightarrow$ same as previous

$I_{47} \rightarrow C \rightarrow d, c/d/\$$

$I_{89} \rightarrow C \rightarrow cC, c/d/\$$

LALR Parsing table construction:

State	Action			Goto	
	c	d	\$	S	C
I_0	S_{36}	S_{47}		1	2
1			Accept		
2	S_{36}	S_{47}			5
36	S_{36}	S_{47}			89
47	r_3	r_3			
5			r_1		
89	r_2	r_2	r_2		

Ambiguous grammar:

A CFG is said to be ambiguous if there exists more than one derivation tree for the given input string i.e., more than one **LeftMost Derivation Tree (LMDT)** or **RightMost Derivation Tree (RMDT)**.

Definition: $G = (V, T, P, S)$ is a CFG is said to be ambiguous if and only if there exist a string in T^* that has more than one parse tree.

where V is a finite set of variables.

T is a finite set of terminals.

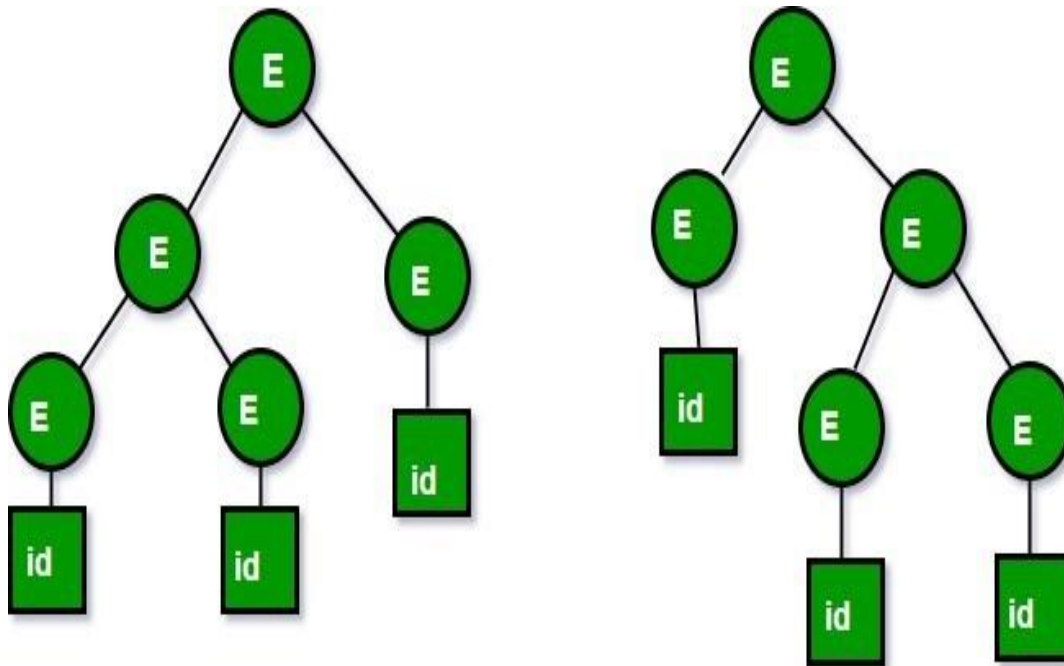
P is a finite set of productions of the form, $A \rightarrow \alpha$, where A is a variable and $\alpha \in (V \cup T)^*$ S is a designated variable called the start symbol.

For Example:

1. Let us consider this grammar : $E \rightarrow E+E|id$

We can create 2 parse tree from this grammar to obtain a string $id+id+id$:

The following are the 2 parse trees generated by left most derivation:



Both the above parse trees are derived from same grammar rules but both parse trees are different. Hence the grammar is ambiguous.

YACC PROGRAMMING

A parser generator is a program that takes as input a specification of a syntax, and produces as output a procedure for recognizing that language. Historically, they are also called compiler-compilers.

YACC (yet another compiler-compiler) is an **LALR(1)** (LookAhead, Left-to-right, Rightmost derivation producer with 1 lookahead token) parser generator. YACC was originally designed for being complemented by Lex.

Input File:

YACC input file is divided in three parts.

```
/* definitions */
....

%%

/* rules */
....

%%

/* auxiliary routines */
....
```

Input File: Definition Part:

- The definition part includes information about the tokens used in the syntax definition:
- %token NUMBER

`%token ID`

- Yacc automatically assigns numbers for tokens, but it can be overridden by

`%token NUMBER 621`

- Yacc also recognizes single characters as tokens. Therefore, assigned token numbers should no overlap ASCII codes.
- The definition part can include C code external to the definition of the parser and variable declarations, within `%{and %}` in the first column.
- It can also include the specification of the starting symbol in the grammar:

`%start nonterminal`

- The rules part contains grammar definition in a modified BNF form.
- Actions is C code in `{ }` and can be embedded inside (Translation schemes).

Input File: Auxiliary Routines Part:

- The auxiliary routines part is only C code.
- It includes function definitions for every function needed in rules part.
- It can also contain the `main()` function definition if the parser is going to be run as a program.
- The `main()` function must call the function `yyparse()`.

Input File:

- If `yylex()` is not defined in the auxiliary routines sections, then it should be included:
`#include "lex.yy.c"`

- YACC input file generally finishes with:

`.y`

Output Files:

- The output of YACC is a file named **y.tab.c**
- If it contains the **main()** definition, it must be compiled to be executable.
- Otherwise, the code can be an external function definition for the function **int yyparse()**
- If called with the **-d** option in the command line, Yacc produces as output a header file **y.tab.h** with all its specific definition (particularly important are token definitions to be included, for example, in a Lex input file).
- If called with the **-v** option, Yacc produces as output a file **y.output** containing a textual description of the LALR(1) parsing table used by the parser. This is useful for tracking down how the parser solves conflicts.

Semantics

Syntax Directed Translation:

- A formalist called as syntax directed definition is used for specifying translations for programming language constructs.
- A syntax directed definition is a generalization of a context free grammar in which each grammar symbol has associated set of attributes and each and each productions is associated with a set of semantic rules

Definition of (syntax Directed definition) SDD :

- SDD is a generalization of CFG in which each grammar productions $X \rightarrow \alpha$ is associated with it a set of semantic rules of the form

$a: = f(b_1, b_2, \dots, b_k)$

Where a is an attributes obtained from the function f .

A syntax-directed definition is a generalization of a context-free grammar in which:

- Each grammar symbol is associated with a set of attributes.
- This set of attributes for a grammar symbol is partitioned into two subsets called synthesized and inherited attributes of that grammar symbol.
- Each production rule is associated with a set of semantic rules.
- Semantic rules set up dependencies between attributes which can be represented by a dependency graph.
- This dependency graph determines the evaluation order of these semantic rules.
- Evaluation of a semantic rule defines the value of an attribute. But a semantic rule may also have some side effects such as printing a value.

The two attributes for non terminal are :

The two attributes for non terminal are :

Synthesized attribute (S-attribute) : (\uparrow)

An attribute is said to be synthesized attribute if its value at a parse tree node is determined from attribute values at the children of the node

Inherited attribute: (\uparrow, \rightarrow)

An inherited attribute is one whose value at parse tree node is determined in terms of attributes at the parent and | or siblings of that node.

- The attribute can be string, a number, a type, a, memory location or anything else.
- The parse tree showing the value of attributes at each node is called an annotated parse tree.

The process of computing the attribute values at the node is called annotating or decorating the parse tree. Terminals can have synthesized attributes, but not inherited attributes.

Annotated Parse Tree

- A parse tree showing the values of attributes at each node is called an Annotated parse tree.
- The process of computing the attributes values at the nodes is called annotating (or decorating) of the parse tree.
- Of course, the order of these computations depends on the dependency graph induced by the semantic rules.

Ex1:1) Synthesized Attributes : Ex: Consider the CFG :

$S \rightarrow EN$

$E \rightarrow E+T$

$E \rightarrow E-T$

$E \rightarrow T$

$T \rightarrow T * F$

$T \rightarrow T / F$

$T \rightarrow F$

$F \rightarrow (E)$

$F \rightarrow \text{digit}$ $N \rightarrow ;$

Solution: The syntax directed definition can be written for the above grammar by using semantic actions for each production

Productionrule

$S \rightarrow EN$
 $E \rightarrow E1 + T$
 $E \rightarrow E1 - T$
 $E \rightarrow T$
 $T \rightarrow T * F$
 $T \rightarrow T | F$
 $F \rightarrow (E)$
 $T \rightarrow F$
 $F \rightarrow \text{digit}$
 $N \rightarrow ;$

Semanticactions

$S.val = E.val$
 $E.val = E1.val + T.val$
 $E.val = E1.val - T.val$
 $E.val = T.val$
 $T.val = T.val * F.val$
 $T.val = T.val | F.val$
 $F.val = E.val$
 $T.val = F.val$
 $F.val = \text{digit.lexval}$
 can be ignored by lexical Analyzer as; I is terminating symbol

For the Non-terminals E, T and F the values can be obtained using the attribute "Val".

The taken digit has synthesized attribute "lexval".

In $S \rightarrow EN$, symbol S is the start symbol. This rule is to print the final answer of expressed.

Following steps are followed to Compute S attributed definition

Write the SDD using the appropriate semantic actions for corresponding production rule of the given Grammar.

The annotated parse tree is generated and attribute values are computed. The Computation is done in bottom up manner.

The value obtained at the node is supposed to be final output.

L-attributed SDT

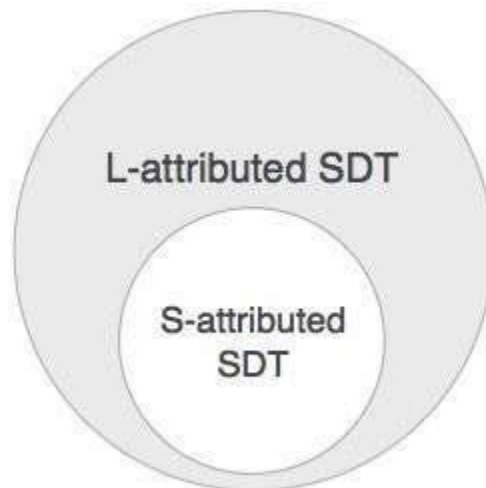
This form of SDT uses both synthesized and inherited attributes with restriction of not taking values from right siblings.

In L-attributed SDTs, a non-terminal can get values from its parent, child, and sibling nodes. As in the following production

$S \rightarrow ABC$

S can take values from A, B, and C (synthesized). A can take values from S only. B can take values from S and A. C can get values from S, A, and B. No non-terminal can get values from the sibling to its right.

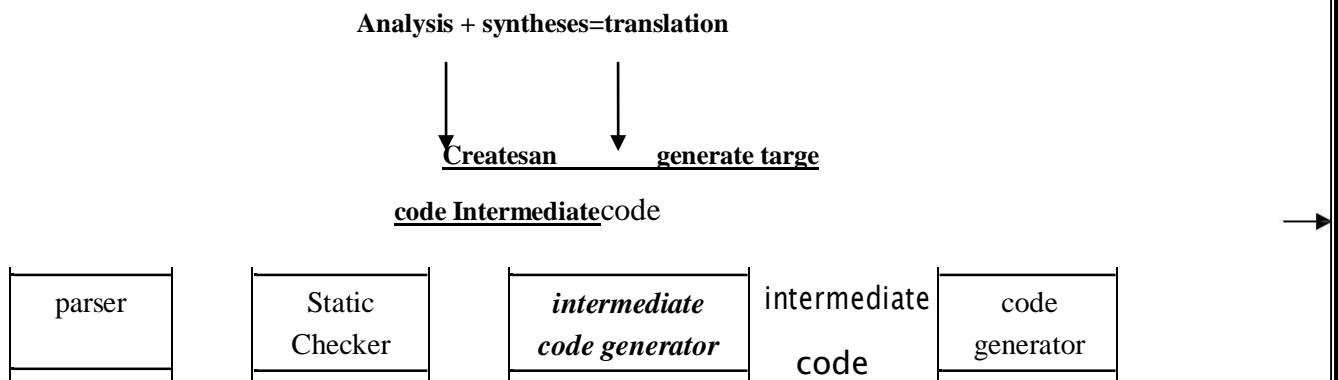
Attributes in L-attributed SDTs are evaluated by depth-first and left-to-right parsing manner.



We may conclude that if a definition is S-attributed, then it is also L-attributed as L-attributed definition encloses S-attributed definitions

Intermediate Code

An intermediate code form of source program is an internal form of a program created by the compiler while translating the program from a high-level language to assembly code(or)object code(machine code).an intermediate source form represents a more attractive form of target code than does assembly. An optimizing Compiler performs optimizations on the intermediate source form and produces an objectmodule.



In the analysis –synthesis model of a compiler, the front-end translates a source program into an intermediate representation from which the back-end generates target code, in many compilers the source code is translated into a language which is intermediate in complexity between a HLL and machine code .the usual intermediate code introduces symbols to stand for various temporary quantities.

We assume that the source program has already been parsed and statically checked..the various intermediate code forms are:

- a) Polishnotation
- b) Abstract syntax trees(or)syntaxtrees
- c) Quadruples
- d) Triples three address code
- e) Indirecttriples
- f) Abstract machinecode(or)pseudocopde

postfix notation:

The ordinary (infix) way of writing the sum of a and b is with the operator in the middle: $a+b$. the postfix (or postfix polish) notation for the same expression places the operator at the right end, $asab+$.

In general, if e_1 and e_2 are any postfix expressions, and \emptyset to the values denoted by e_1 and e_2 is indicated in postfix notation nby $e_1e_2\emptyset$.no parentheses are needed in postfix notation because the position and priority (number of arguments) of the operators permits only one way to decode a postfix expression.

Syntax Directed Translation:

- A formalist called as syntax directed definition is used for specifying translations for programming language constructs.
- A syntax directed definition is a generalization of a context free grammar in which each grammar symbol has associated set of attributes and each and each productions is associated with a set of semantic rules

Definition of (syntax Directed definition) SDD :

SDD is a generalization of CFG in which each grammar productions $X \rightarrow \alpha$ is associated with it a set of semantic rules of the form

$$a := f(b_1, b_2, \dots, b_k)$$

Where a is an attributes obtained from the function f.

- A syntax-directed definition is a generalization of a context-free grammar in which:
- Each grammar symbol is associated with a set of attributes.

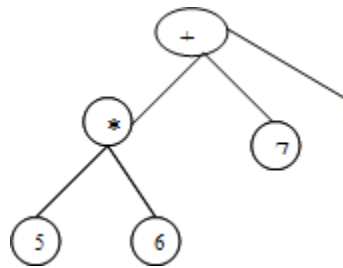
This set of attributes for a grammar symbol is partitioned into two subsets called synthesized and inherited attributes of that grammar symbol.

- Each production rule is associated with a set of semantic rules.
- Semantic rules set up dependencies between attributes which can be represented by a dependency graph.

Annotated Parse Tree

- A parse tree showing the values of attributes at each node is called an Annotated parse tree.
- The process of computing the attributes values at the nodes is called annotating (or decorating) of the parse tree. Of course, the order of these computations depends on the dependency graph induced by the

Syntax tree:



Annotated parse tree :

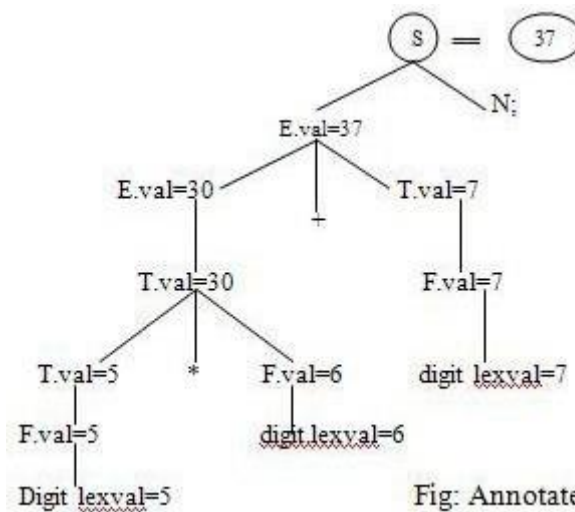


Fig: Annotated parse tree

ASSIGNMENT STATEMENTS

Suppose that the context in which an assignment appears is given by the following grammar. P

□ M D

M □ ε

D □ D ; D | **id** : T | **proc id** ; N D ; S

N □ ε

Nonterminal P becomes the new start symbol when these productions are added to those in the translation scheme shown below.

Translation scheme to produce three-address code for assignments

S → id := E { p := lookup (id.name);
 if p ≠ nil **then**
 emit(p ' := ' E.place)
 else error }

E → E₁ + E₂ { E.place := newtemp;

```

                                emit(E.place ': =' E1.place ' + ' E2.place ) }
E → E1 * E2      { E.place := newtemp;
                  emit(E.place ': =' E1.place ' * ' E2.place ) }
E → - E1         { E.place := newtemp;
                  emit ( E.place ': =' 'uminus' E1.place ) }
E → ( E1 )      { E.place := E1.place }

E → id          { p := lookup ( id.name);
                if p ≠ nil then
                  E.place := p
                else error }

```

Flow-of-Control Statements

We now consider the translation of boolean expressions into three-address code in the context of if-then, if-then-else, and while-do statements such as those generated by the following grammar:

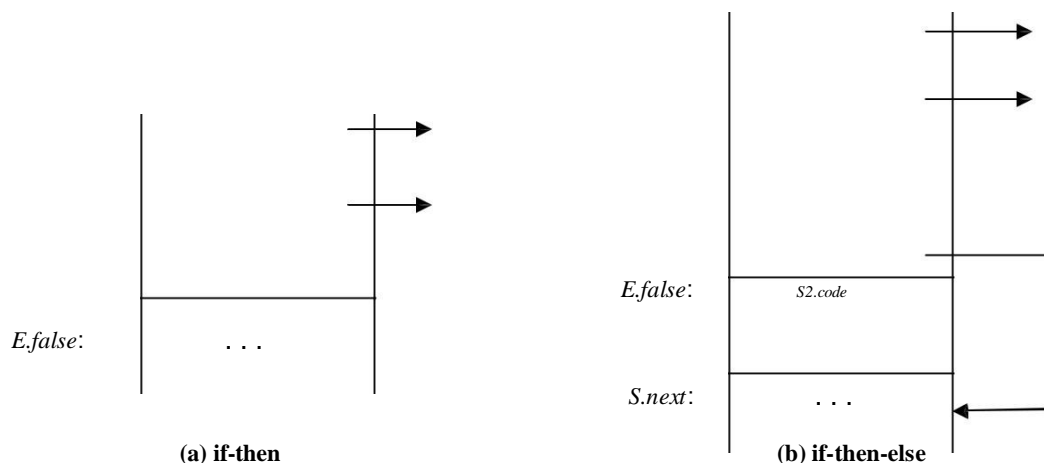
```

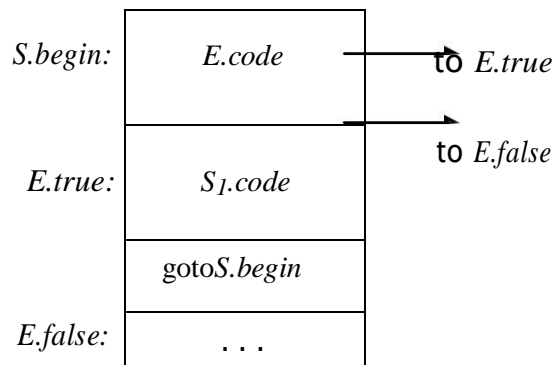
S → if E then S1
   if E then S1 else
   | S2
   | while E do S1

```

In each of these productions, E is the Boolean expression to be translated. In the translation, we assume that a three-address statement can be symbolically labeled, and that the function *newlabel* returns a new symbolic label each time it is called.

- $E.true$ is the label to which control flows if E is true, and $E.false$ is the label to which control flows if E is false.
- The semantic rules for translating a flow-of-control statement S allow control to flow from the translation $S.code$ to the three-address instruction immediately following $S.code$.
- $S.next$ is a label that is attached to the first three-address instruction to be executed after the code for **Code for if-then , if-then-else, and while-do statements**





(c) while-do

PRODUCTION	SEMANTIC RULES
$S \xrightarrow{\quad} \text{if } E \text{ then } S_1$	$E.true := \text{newlabel};$ $E.false := S.next;$ $S_1.next := S.next;$ $S.code := E.code \parallel \text{gen}(E.true, ",:") \parallel S_1.code$
$S \xrightarrow{\quad} \text{if } E \text{ then } S_1 \text{ else } S_2$	$E.true := \text{newlabel};$ $E.false := \text{newlabel};$ $S_1.next := S.next;$ $S_2.next := S.next;$ $S.code := E.code \parallel \text{gen}(E.true, ",:") \parallel S_1.code \parallel$ $\text{gen}(, \text{goto } S.next) \parallel$ $\text{gen}(E.false, ",:") \parallel S_2.code$
$S \xrightarrow{\quad} \text{while } E \text{ do } S_1$	$S.begin := \text{newlabel};$ $E.true := \text{newlabel};$ $E.false := S.next;$ $S_1.next := S.begin;$ $S.code := \text{gen}(S.begin, ",:") \parallel E.code \parallel$ $\text{gen}(E.true, ",:") \parallel S_1.code \parallel$ $\text{gen}(, \text{goto } S.begin)$