

# FORMAL SEMANTICS I

Lectured : T.S. McDermott

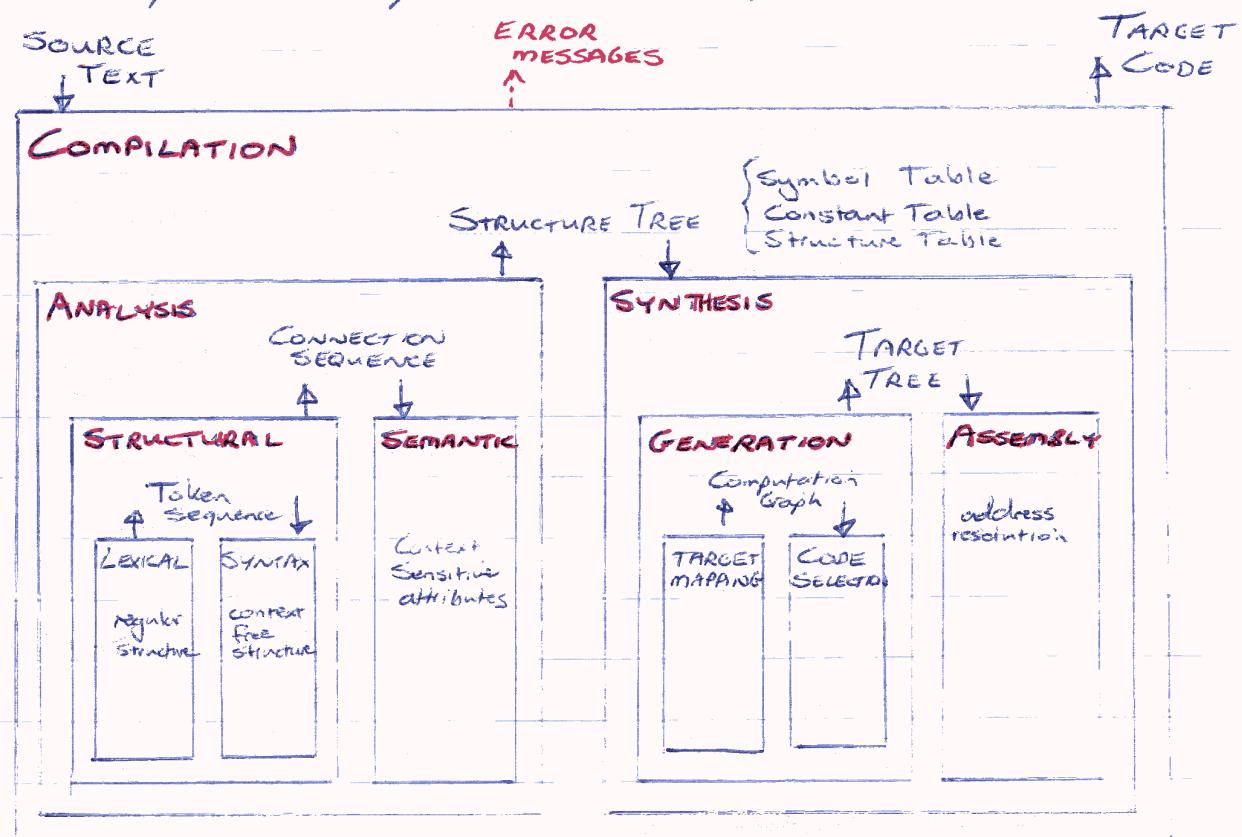
24 Feb 1986 - 1 April 1986

Analysis +  
Synthesis

1.	GENERAL STRUCTURE OF A COMPILER	1
2	CHOMSKY GRAMMARS	1
3	LR PARSING	9
4	ATTRIBUTE GRAMMARS	27
5	W - GRAMMARS	39
6	THE SECD MACHINE	45
7	OPERATIONAL SEMANTICS	61
8	CODE GENERATION	79
9	ERROR HANDLING	81
10	OPTIMISATION	85
11	ADDENDA	

## ① GENERAL STRUCTURE OF A MODULARISED COMPILER

An abstract structure of a modularised compiler is shown below - abstract as could be one-pass (parallel-) or two-pass (sequential module execution).



MODULE := an action without a state (i.e. procedures have no internal state, at least in their mathematical defns?) local variables protected from above routines

## ② PHRASE STRUCTURE GRAMMARS & LANGUAGES

V = Vocabulary (alphabet) set of symbols

L = Language subset of V\* (strings over V)

G = Grammar generator of language (grammar finite but strings produced not necessarily)

$L = \{ X : S \xrightarrow{*} X \}$  S a start symbol  
 $X \in T^*$  terminal strings

$V$  - vocab =  $N \cup T$   
 $P$  - productions  
 $T$  - terminals  
 $Z$  - single start symbol  
 $N$  - non-terminals

$(V, P)$  form a rewriting system

$(T, N, P, Z)$  a phrase structure (Chomsky) grammar  
for  $\{X \in T^*: Z \xrightarrow{*} X\}$  language

Chomsky uses only sets which are grammable or enumerable.

most restrictive Type 3  $\subset$  Type 2  $\subset$  Type 1  $\subset$  Type 0 <sup>most general</sup>

### ACCEPTORS FOR PHRASE STRUCTURE GRAMMARS (AUTOMATA)

Given a string, is it in the language defined by the grammar? Much formal work has been done relating classes of grammars to classes of automata. This relationship is intimate, however for a grammar to be useful to a computer writer, it must be computationally tractable as well as reasonably easy to generate the parse tree for any input string. Type 3 & 2 languages can be thus parsed; Types 1 & 0 present fantastic problems for automatic parsers.

#### Type 3 (Finite State) Grammars. (also Regular Grammars)

Rules of two forms only:  $N \xrightarrow{*} T^* | NT^*$

Exactly one parse tree.

#### Type 2 (Context Free) Grammars

Rules of form  $N \xrightarrow{*} S$  where  $S \in T^*$ . We can generate semantically meaningful trees (eg. creator procedures)

## Type 1 (Context Sensitive) Grammars.

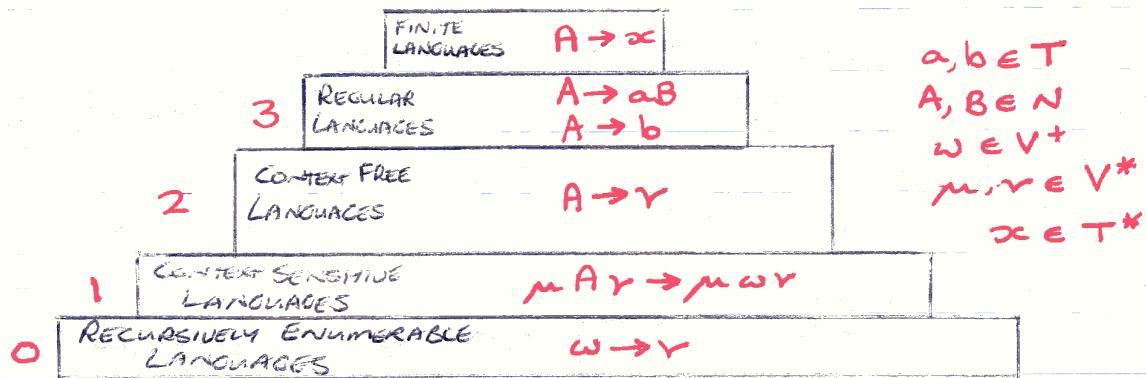
$$\alpha A \gamma \Rightarrow \alpha S \gamma \quad A \in N, S \in V^+, \alpha, \gamma \in V^*$$

As  $S \in V^+$  RHS always longer or same length as LHS.  
 This property means that type 1 grammars are decidable  
 i.e given a string we can decide if it is in the language

## Type 0 (Phrase Structure) Grammars (also Recursively Enumerable)

$$w \rightarrow \mu \quad \text{where } w \in V^+; \mu \in V^*$$

### SUMMARY OF THE CLASSIFICATION OF PHRASE-STRUCTURE G'S.



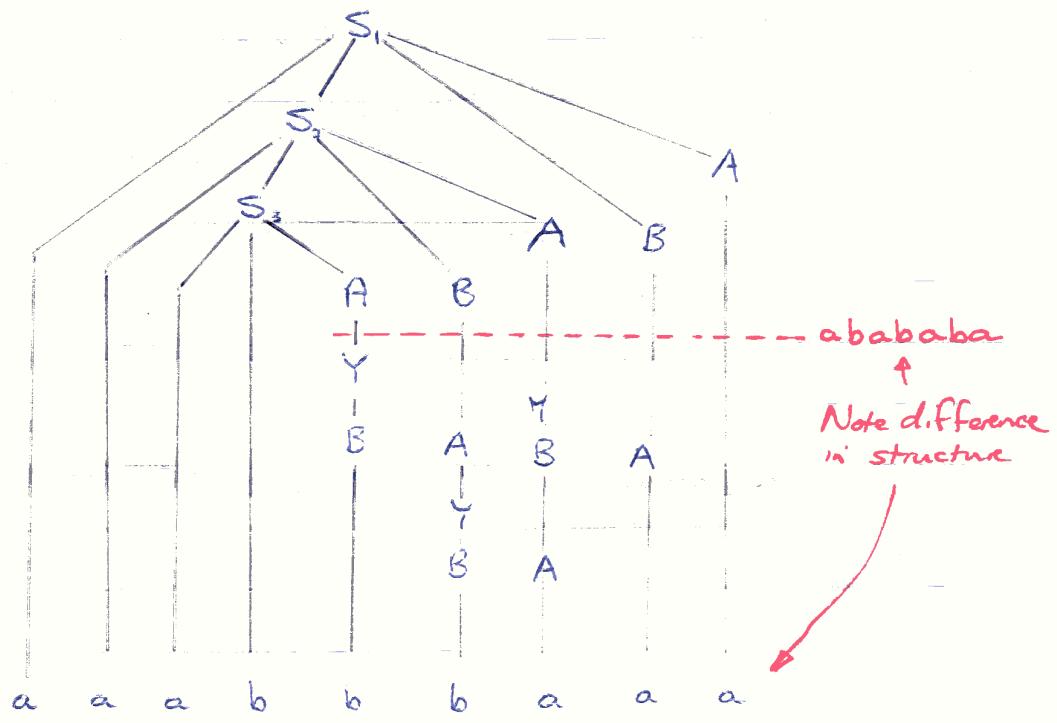
### PHRASE STRUCTURE & MEANING.

A phrase is a substring which derives from a single non-terminal. We interpret a phrase as a semantic unit, if we are not satisfied with any grammar, we choose one upon which we can base a semantic analysis of our sentences. For this reason we use context-free grammars.

We lose this semantic content in context sensitive languages. For example,  $a^n b^n a^n$  (which cannot be produced by context free) can be produced with a context sensitive grammar.

$$\begin{aligned}
 S &\rightarrow aSBA \mid abA \\
 bB &\rightarrow bb \\
 A &\rightarrow a \\
 AB &\rightarrow YB \quad YB \rightarrow YA \quad YA \rightarrow BA \quad [= AB \rightarrow BA]
 \end{aligned}$$

e.g.



Note that the non-terminals no longer represent some 'abstracted' characterisation of the structure of the string they derive.  $S_1 \xrightarrow{*} a a a b b b a a a$   $S_2 \xrightarrow{*} a a b b a$   $S_3 \xrightarrow{*} a b b$

NB Chomsky considered CS to be phrase-structure but for our purposes we consider CF to be, not CS. Chomsky does not deal with semantics but rather with strings which can be analysed mathematically.

## Why Chomsky Grammars?

1. over against BNF unifies the description of all types of grammatical languages (eg unified notion of production)
2. obscures difference of generative analytic syntax
3. obscures difference of constituency & dependency syntax
4. relationship to automata and formal systems

Startsymbol = axioms

productions = rules of inference

derivation = proof

terminal strings = theorems

Points 2 & 3 come from the 2 disjoint alphabets,  
1 startsymbol and 1-place rules.

eg  $\langle \text{subj} \rangle \langle \text{pred} \rangle \langle \text{obj} \rangle$

They made him do it

$\langle \text{subj} \rangle \langle \text{pred} \rangle \langle \text{obj} \rangle$

## (3) PARSING

3.1

We can parse various grammars using automata

③ regular grammars = finite automata

① context sensitive = push-down automata

② context free = limit bounded automata

④ recursively enumerable = Turing Machines

## PARSING REGULAR GRAMMARS WITH FINITE STATE AUTOMATA

$$RG = (T, N, P, Z)$$

$$FSA = (T, N, R, Z, F) \quad R = \text{productions} \quad F = \text{final state}, F \in N.$$

$(T \cup N, R)$  form a rewriting system

If  $A \rightarrow aB \in P$ , then  $Aa \rightarrow B \in R$ .

$$(e.g. R = \{ ch \rightarrow q | F, c \rightarrow I, F \rightarrow I, FE \rightarrow S, Ih \rightarrow q | X, \dots \})$$

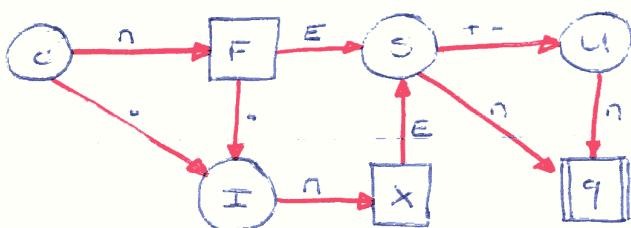
$$P = \{ c \rightarrow n | nF | \cdot I, F \rightarrow I | ES, I \rightarrow n | nx, X \rightarrow ES, S \rightarrow +u | -u | n, u \rightarrow n \}$$

$$Z = C \quad F = \{q\}$$

a regular grammar for real numbers.

All non-terminals are states and  $q$  is the empty non-terminal.

This is a non-deterministic FSA but can be changed into a deterministic FSA (resulting in a state being a set of non-terminals, not just one). Each non-terminal represents the rest of the control string (or at least what it is expected to be).



NDFSA for Recognizing  
real numbers.

## SELF-EMBEDDING.

Nesting (self-embedding) cannot be done by regular grammars, but requires context-free grammars, which can be parsed by pushdown automata using a stack.

$$PDA = (T, Q, S, R, q_0, S_0, F) \quad S = \text{set of states of stack}$$

$q_0$  : start symbol

$S_0 \in$  start stack

$Q = \text{set of states of control string}$

$(T \cup Q \cup S, R)$  form a rewriting system.

Each member of  $R$  has the form  $sqcz \rightarrow s'q'z$



Note that:

FSA has a finite number of states

PDA has a finite number of stack states, but an infinite number of parser states

DPDA has an infinite number of parser states, but these are partitioned into a finite number of equivalence classes. The language of the stack is a finite state automaton

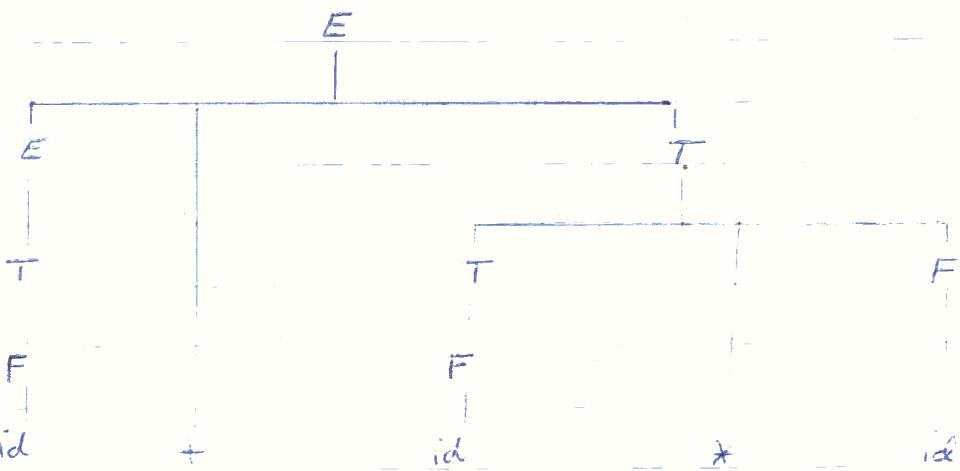
## DETERMINISTIC PARSERS

- Top Down - begin with start symbol and try to derive the null string - produces a left derivation
- Bottom Up - begin with an empty stack and try to produce the start symbol, shifting symbols onto the stack and trying to reduce the contents of the stack (see eg overleaf) produces right derivation

## Why DETERMINISTIC?

To parse without backtracking. For a regular grammar we can always find a deterministic grammar. As PDA's can have an infinite number of parser states we cannot always construct deterministic grammars. All computer languages are designed to be deterministically parseable.

$$E \rightarrow T \mid E + T \quad T \rightarrow F \mid T * F \quad F \rightarrow (e) \mid id$$



### Topdown

Stack      Control

$E$	$i+i*i$
$T+E$	$i+i*i$
$T+T$	$i+i*i$
$T+F$	$i+i*i$
$T+i$	$i+i*i$
$F*T$	$i*i$
$F*F$	$i*i$
$F+i$	$i*i$
$i$	$i$
$\Lambda$	$\Lambda$

### Derivation

Leftmost      Rightmost

$E$	$E$
$E+T$	$E+T$
$T+T$	$E+T+F$
$F+T$	$E+T+i$
$i+T$	$E+F+i$
$i+T+F$	$E+i+i$
$i+F*F$	$T+i+i$
$i+i*F$	$F+i*i$
$i+i*i$	$i+i*i$

### Bottom Up

Stack      Control

$E$	$\Lambda$
$E+T$	$\Lambda$
$E+T+F$	$\Lambda$
$E+T+i$	$\Lambda$
$E+F$	$*i$
$E+i$	$*i$
$T$	$+i*i$
$F$	$+i*i$
$i$	$+i*i$
$\Lambda$	$i+i*i$

### LL As LR

The letters define the directions of control and stack.

### 3.2 LR PARSING

In every LR - parser there is a FSA controlling the stack. E is called the dominant non-terminal of the grammar; G the preliminary start symbol, and cannot appear on the right of a production. The followset is the set of terminals which can occur after a state containing the non-terminal.

#### GRAMMAR

$$G \rightarrow E \quad T \rightarrow F \mid T \times F$$

$$E \rightarrow E + T \mid T \quad F \rightarrow (E) \mid id$$

(refer bottom up down opposite)

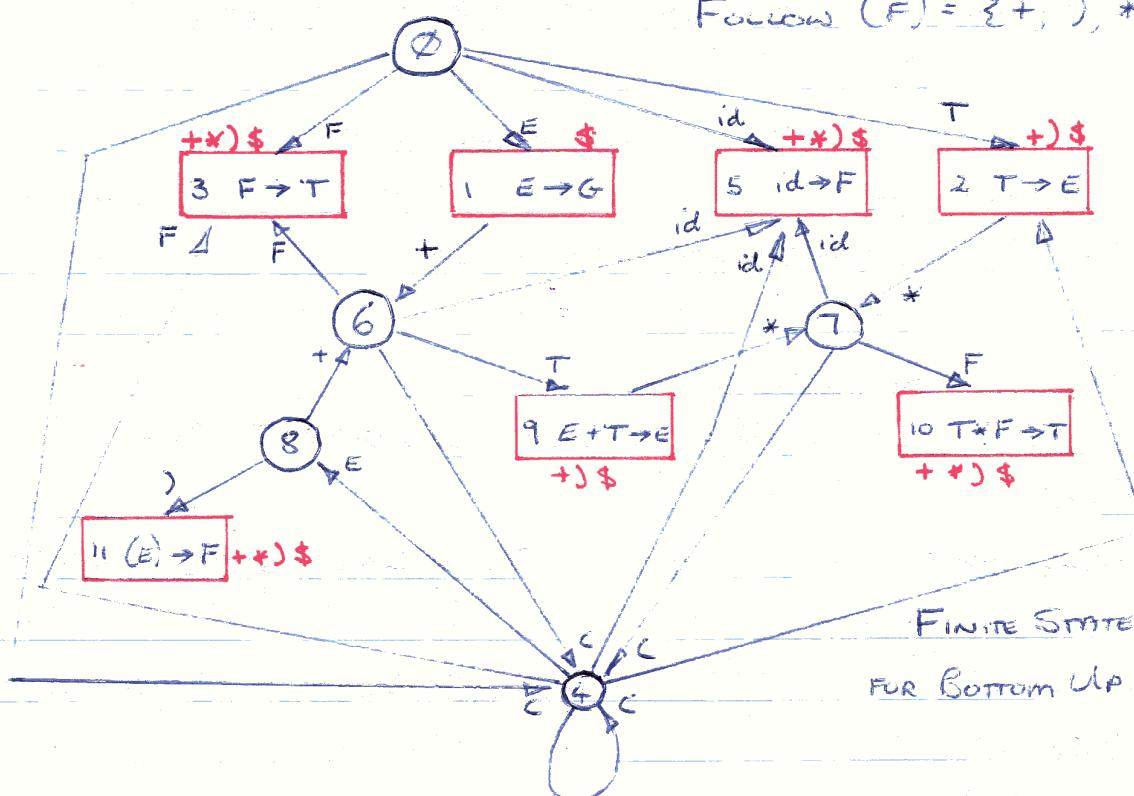
#### FOLLOW SETS

$$\text{Follow}(G) = \{\$, +\}$$

$$\text{Follow}(E) = \{+, ), \$\}$$

$$\text{Follow}(T) = \{+, ), *, \$\}$$

$$\text{Follow}(F) = \{+, ), *, \$\}$$



FINITE STATE CONTROL  
FOR BOTTOM UP LR PARSER.

Note: State 5 is a reduce state only, reducing id to F, then going back to state 0. At each reduce state, if a shift can be done it is also back to state 0. The red letters show the follow sets \ the shift symbols (i.e. those symbols which force a return to state 0).

○ shift state

■ reduce state

## CONSTRUCTING THE FSA

### DEFNS

Item - an item ... contains a non-terminal followed by a stack and control in the form.

(non-terminal)  $\rightarrow$  <stack> . <control>

item.set - a set of items associated with a state.  
(continued later)

## LR PARSER IMPLEMENTATION

LR Parsers generally consist of a driver and table  
(usually machine generated). We will consider

3 types : SLR (simple LR)

LR (canonical LR)

LALR (look-ahead LR)

## LR PARSERS

An LR parser has an input, a stack, and a parsing table. The input is read from left to right, one symbol at a time. The stack contains a string of the form  $S_0 X_1 S_1 X_2 \dots X_m S_m$ , where  $S_m$  is on top. Each  $X_i$  is a grammar symbol and each  $S_i$  a state. States summarize the information on the stack and guide the shift/reduce action. The parsing table consists of two parts, a parsing function action and a control function goto.

The driver consults entry action  $[S_m, a_i]$  where  $a_i$  is the current input symbol. This entry can have one of four values:

- shift s
- reduce  $A \rightarrow B$
- accept (parsing complete)
- error

`goto (<state>, <grammar symbol>)` produces the new state.

All LR parsers behave in this way: only the table construction methods vary. A parsing table for our example grammar is:

### ■ Grammar (productions numbered)

0  $G \rightarrow E$     1,2  $E \rightarrow E + T / E$     3,4  $T \rightarrow T * F / F$   
 5,6  $F \rightarrow (E) / id$

State	Action (Shift Reduce)					Goto (state)
	id	+	*	(	)	
0	s4		s4			
1		s6				accept
2	r2	s7		r2	r2	
3	r4	r4		r4	r4	
4	s5		s4			8 2 3
5	r6	r6		r6	r6	
6	s5		s4			9 3
7	s5		s4			10
8	s6		s11			
9	r1	s7		r1	r1	
10	r3	r3		r3	r3	
11	r5	r5		r5	r5	

All empty actions are errors.  $s_k$  shift & goto state  $k$

NB refer to Work & Goos  
 pp 176 for implementation of sparse table.

$r_k$  reduce according to production  $k$

A grammar for which we can construct an LR parsing table in which every entry is uniquely defined is an LR grammar. There are context-free grammars which are not LR but <sup>(eg any ambiguous grammar)\*</sup> these are avoided. The derive routine of an LR parser is essentially a finite automaton.

Note that the LR requirement that we be able to recognise the occurrence of the right side of a production having seen what is derived from that right side, is far less stringent than the requirement for a predictive parser, namely that we be able to recognise the apparent use of the production seeing only the first symbol it derives.

### 3.3 CONSTRUCTING THE DFA ( $\stackrel{FS}{A}$ ( $LR(0)$ ) (i.e. CANONICAL) $LR(0)$ ITEMSETS)

We show how to construct a DFA from the grammar, which can then be converted into a parsing table. The DFA recognises viable prefixes of the grammar, i.e. prefixes of right-sentential forms that do not contain any symbols to the right of the handle. A viable prefix is so called because it is always possible to add terminal symbols to the right of a viable prefix to obtain a right-sentential form. There is therefore no apparent error provided the input to a given point can be reduced to a viable prefix.

We define an LR(0) item of a grammar  $G$  to be a production of  $G$  with a dot at some position of the right

\* NB an ambiguous grammar can be parsed if shifting is given precedence over reducing.

side Eq  $A \rightarrow XYZ$  generates four items:  $A \rightarrow \cdot X Y Z$ ,  $A \rightarrow X \cdot Y Z$ ,  $A \rightarrow X Y \cdot Z$  and  $A \rightarrow X Y Z$ . The production  $A \rightarrow E$  generates only one item  $A \rightarrow \cdot$ . In an implementation we can represent an item by a pair ( $\langle$ production # $\rangle$ , position>) We group items together into sets which give rise to the states of an LR parser.

One collection of sets of items, which we call the canonical LR( $\infty$ ) collection, provides the basis for constructing SLR (simple LR) parsers. To construct the canonical LR( $\infty$ ) collection for a grammar, we need to define an augmented grammar and two functions CLOSURE and CUTOFF. If  $G$  is a grammar with a start symbol  $S$ , then  $G'$  the augmented grammar for  $G$ , is  $G$  with a new start symbol  $S'$  and production  $S' \rightarrow S$ . This is so that the parser knows that when it reduces  $S$  to  $S'$  it is going into the accept state.

### Closure of the Itemset.

If  $I$  is a set of items for a grammar  $G$ , then the set of items CLOSURE( $I$ ) is constructed from  $I$  by:

- every item in  $I$  is in CLOSURE( $I$ )
- if  $A \rightarrow \alpha \cdot B\beta \in \text{CLOSURE}(I)$  and  $B \rightarrow \gamma$  is a production, then  $B \rightarrow \cdot \gamma$  is in CLOSURE( $I$ )

Intuitively  $A \rightarrow \alpha \cdot B\beta \in \text{CLOSURE}(I)$  indicates that at some point in the parsing process we need expect to see a string derivable from  $B\beta$  as input. If  $B \rightarrow \gamma$  is a production, we would also expect to see a string derivable from  $\gamma$  at this point.

## GOTO

The second useful junction is GOTO ( $I, X$ ) where  $I$  is an itemset and  $X$  a grammar symbol. GOTO ( $I, X$ ) is defined to be the closure of the set of all items

$A \rightarrow \alpha X \cdot \beta$  such that  $A \rightarrow \alpha \cdot X \beta$  is in  $I$ .

Intuitively, if  $I$  is the set of items that are valid for some viable prefix  $\gamma$ , then GOTO ( $I, X$ ) is the set of items that are valid for the viable prefix  $\gamma X$ .

- Example : Consider the augmented grammar :

$$\begin{array}{ll} E' \rightarrow E & T \rightarrow T + F \mid F \\ E \rightarrow E + T \mid T & F \rightarrow id \mid (E) \end{array}$$

If  $I$  is the set of one item  $\{[E' \rightarrow \cdot E]\}$  then closure ( $I$ ) contains the items :

$$\begin{array}{lll} E' \rightarrow \cdot E & E \rightarrow \cdot E + T & E \rightarrow \cdot T \\ T \rightarrow \cdot T + F & T \rightarrow \cdot F & F \rightarrow \cdot (E) \quad F \rightarrow \cdot id \end{array}$$

If  $I$  is the set of items  $\{[E' \rightarrow \cdot E], [E \rightarrow E \cdot + T]\}$  then GOTO ( $I, +$ ) is the closure of  $\{[E \rightarrow E + T]\}$ .

We can now give the algorithm to construct  $C$ , the canonical collection of sets of LR(0) items for an augmented grammar  $G'$  (see opposite).

We shall now construct an LR(0) parser for our grammar. We begin with  $E' \rightarrow \cdot E$  and take the closure of this - this gives the itemset for state 0, which is

$$I_0 = \{[E' \rightarrow \cdot E], [E \rightarrow \cdot E + T], [E \rightarrow \cdot T], [E \rightarrow \cdot T + F], [T \rightarrow \cdot F], [F \rightarrow \cdot (E)], [F \rightarrow \cdot id]\}$$

```

procedure ITEMS ( $G'$ );
begin
   $C := \{\text{closure } (\{S' \rightarrow .S\})\}$ ;
  repeat
    for each itemset  $I$  and grammar symbol  $X$ 
      if  $\text{GOTO}(I, X)$  is non-empty
        then  $C := C \cup \text{GOTO}(I, X)$ 
  until all itemsets and grammar symbols have been done
end

```

Algorithm to construct the canonical collection of  $LR(0)$  itemsets.  
 (above - general alg ; below FSA alg)

begin with state  $\emptyset \} \quad (* \text{ } C := \{\text{closure } (\{S' \rightarrow .S\})\})$   
 close item set  
 work out edges  $(x^i)$  and produce their item sets  
 for all new item sets, we have new states  
 Repeat for each state until there are no more states

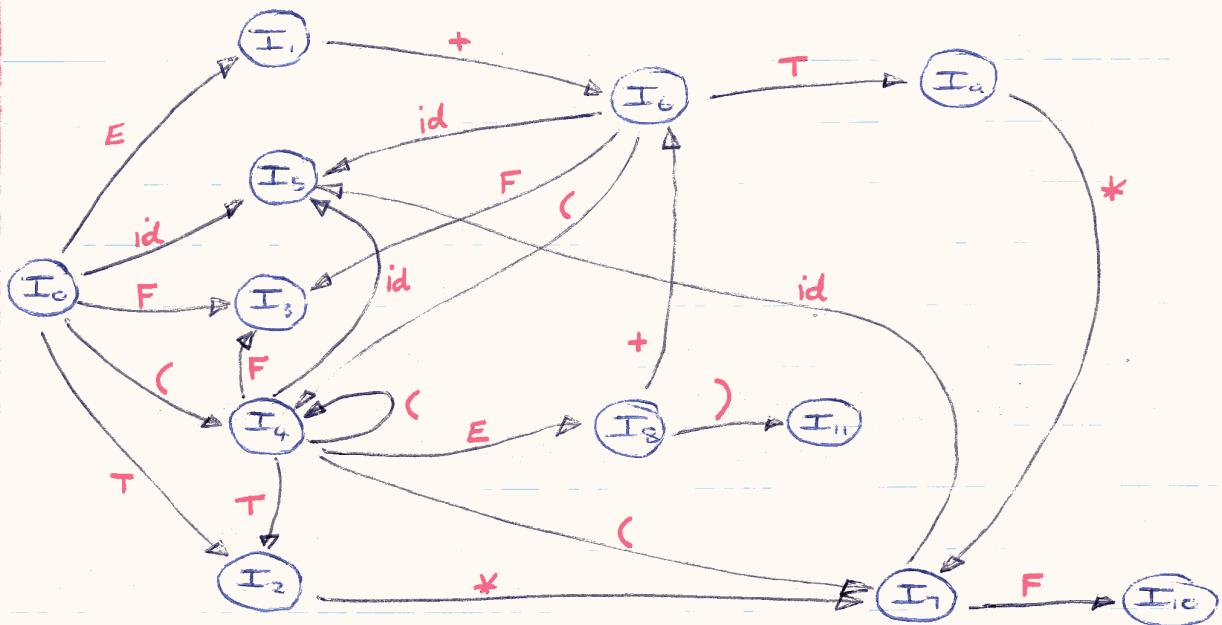
State  $\emptyset$  has edges  $(E, T, F, ( \text{ and id })$ . Consider edge  $E$ . Then  $\text{GOTO}(I_0, E)$  is closure  $\{[E' \rightarrow E.], [E \rightarrow E.+T]\}$ .  
 Let this be state  $I_1$ . Then  $I_1 = \{[E' \rightarrow E.], [E \rightarrow E.+T]\}$ .  
 In a similar fashion we obtain:  
 $I_2 = \{[E \rightarrow T.], [T \rightarrow T.+F]\}$  where edge 2 is  $T$ .  
 $I_3 = \{[T \rightarrow F.]\}$  where edge 3 is  $F$ .  
 $I_4 = \{ \dots \}$  where edge 4 is  $($ .  
 $I_5 = \{[F \rightarrow \text{id}.]\}$  where edge 5 is  $\text{id}$ .

We then work out  $I_6$  from  $\text{GOTO}[I_1, +]$ ,  $I_7$  from  $\text{GOTO}[I_1, *]$ ,  $I_8$  from  $\text{GOTO}[I_4, E]$ ,  $\text{GOTO}[I_4, T]$ ,  $\text{GOTO}[I_4, F]$ ,  $\text{GOTO}[I_4, ()]$  and  $\text{GOTO}[I_4, \text{id}]$ .  $I_9$  comes from  $\text{GOTO}[I_0, n]$  where  $n \in \{T, F, (, \text{id}\}$ ;  $I_{10}$  from  $I_7$  and  $F, (, \text{id}$ ; and  $I_{11}$  from  $I_3$  with  $) +$  and  $*$ .

This gives us the pushdown parser:

edge	state	item	closure
	0	$\{[E \rightarrow . E]\}$	$[F \rightarrow .(E)] [E \rightarrow . id]$
$\phi E$	1	$\{[E \rightarrow E.] , [E \rightarrow E.+T]\}$	$\{[E \rightarrow . E] , [E \rightarrow . E+T] , [E \rightarrow . T] , [T \rightarrow . T+F] , [T \rightarrow F]\}$
$\phi T$	2	$\{[E \rightarrow T.] , [T \rightarrow T.+F]\}$	$\{[E \rightarrow . T] , [T \rightarrow T.+F]\}$
$\phi F$	3	$\{[T \rightarrow F.] \}$	$\{[T \rightarrow F.] \}$
$\phi ($	4	$\{[F \rightarrow (.) E]\}$	$\{[F \rightarrow (.) E] , [F \rightarrow id]\}$
$\phi id$	5	$\{[F \rightarrow id.] \}$	$\{[F \rightarrow id.] \}$
$1 +$	6	$\{[E \rightarrow E.+T]\}$	$\{[E \rightarrow E.+T] , [T \rightarrow T+F] , [T \rightarrow F] , [F \rightarrow (.)]\}$
$2 *$	7	$\{[T \rightarrow T.*F]\}$	$\{[T \rightarrow T.*F] , [F \rightarrow (.)] , [F \rightarrow . id]\}$
$4 E$	8	$\{[F \rightarrow (E.)] , [E \rightarrow E.+T]\}$	$\{[F \rightarrow (E.)] , [E \rightarrow E.+T]\}$
$4 T$		= 2	
$4 F$		= 3	
$4 ($		= 4	
$4 id$		= 5	
$6 +$	9	$\{[E \rightarrow E+T.] , [T \rightarrow T.+F]\}$	$\{[E \rightarrow E+T.] , [T \rightarrow T.+F]\}$
$6 F$		= 3	
$6 ($		= 4	
$6 id$		= 5	
$7 F$	10	$\{[T \rightarrow T+F.] \}$	$\{[T \rightarrow T+F.] \}$
$7 ($		= 4	
$7 id$		= 5	
$8 )$	11	$\{[E \rightarrow (.)]\}$	$\{[F \rightarrow (.)]\}$
$8 +$		= 6	
$9 *$		= 7	

which becomes the DPDA shown opposite. Note that, for example, state 1 is ambiguous, as one item specifies a shift and the other a reduce.



DSFA for Grammar  $G'$ .

### 3.4 CONSTRUCTING SLR PARSING TABLES.

We want to construct action and goto tables from our DPDA. The parsing actions for state  $i$  are determined as follows:

- if  $[A \rightarrow \alpha \cdot aB]$  is in  $I_i$  and  $\text{core}(I_i, a) = I_j$  then action  $[i, a] = \text{"shift } j\text{"}$
- if  $[A \rightarrow \alpha \cdot]$  is in  $I_i$  then for all  $a \in \text{Follow}(A)$  set action  $[i, a] = \text{"reduce } A \rightarrow \alpha\text{"}$
- if  $[S' \rightarrow S \cdot]$  is in  $I_i$ , then set action  $[i, \$] = \text{"accept"}$

If any conflicting actions are generated, the grammar is not SLR(1), and the algorithm fails to produce a valid parser. We consult follow set to decide whether to reduce.

The goto transitions for state  $i$  are constructed using:

- if  $\text{core}(I_i, A) = I_j$  then goto  $[i, A] = j$

The tables for our example are given on page 11.

A problem with SLR is that one can end up with a non-deterministic machine even if the grammar is unambiguous.

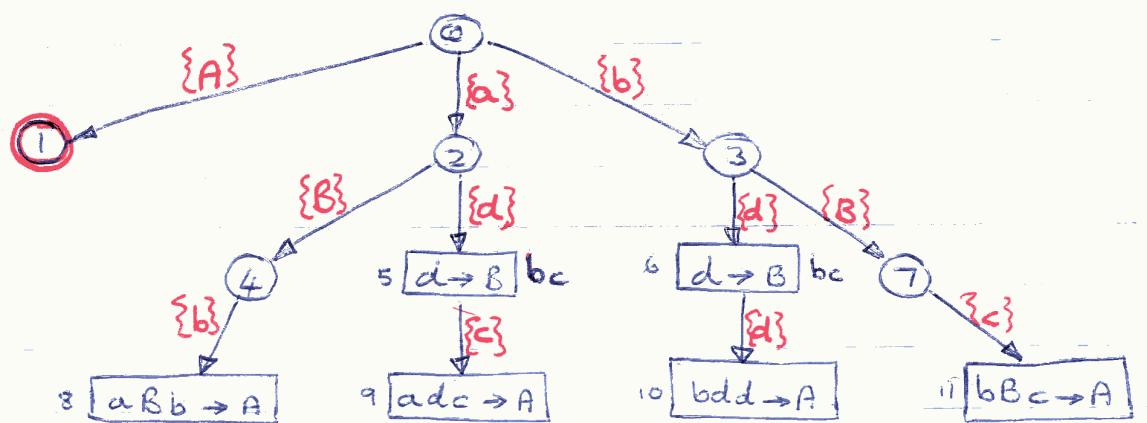
Eg:  $Z \rightarrow A$

$$A \rightarrow aBb \mid adc \mid bBc \mid bdd$$

$$B \rightarrow d$$

$$\text{Follow}(B) = \{b, c\}$$

This gives the FSA:



The problem occurs when the transition from state 2 to state 5 is made. The followset of B in this case is not  $\{b, c\}$  but only  $\{b\}$ . If we use  $\{b, c\}$ , we could reduce d to B (state 5) but this prevents us from recognising the correct string in state 9 as we now have  $aBb$  rather than  $adc$ . This problem is solved in LALR where we consult the follow set for a particular state, eg  $\text{Follow}_5(B) = \{b\}$  (from  $A \rightarrow aBb$ ) and  $\text{Follow}_6(B) = \{c\}$  (from  $A \rightarrow bBc$ )

### 3.5 CONSTRUCTING CANONICAL LR PARSING TABLES.

In the SLR method, state  $i$  calls for a reduction by  $A \rightarrow \alpha$  if some set of items  $I_i$  contains item  $[A \rightarrow \alpha \cdot]$  and  $\alpha$  is in follow( $A$ ). In some situations, however, when state  $i$  appears on top of the stack, the viable prefix (cf p12)  $B\alpha$  on the stack is such that  $\beta A$  cannot be followed by  $\alpha$  in a right-sentential form (cf opposite). Thus the reduction would be invalid on input  $\alpha$ .

It is possible to carry more information in the state that will allow us to rule out some of these invalid reductions by  $A \rightarrow \alpha$ . By splitting states when necessary, we can arrange to have each state of an LR parser indicate exactly which input symbols can follow a handle  $\alpha$  for which there is a possible reduction to  $A$ . The extra information is incorporated into the state by redefining items to include a terminal symbol  $a$  as a second component. The general form of an item becomes  $[A \rightarrow \alpha \cdot \beta, a]$ ; we call such an item an LR(1) item. The lookahead has no effect if  $\beta$  is not  $\epsilon$ , but in an item of the form  $[A \rightarrow \alpha \cdot, a]$  we reduce by  $A \rightarrow \alpha$  only if the next input symbol is  $a$ . Thus, we reduce by  $A \rightarrow \alpha$  only on those input symbols  $a$  for which  $[A \rightarrow \alpha \cdot, a]$  is an LR(1) item in the state on top of the stack.

Formally, we say LR(1) item  $[A \rightarrow \alpha \cdot \beta, a]$  is valid for viable prefix  $\gamma$  if there is a derivation  $S \xrightarrow{*} \gamma S \alpha w \xrightarrow{*} \gamma S \alpha w$  where:  $\gamma = S\delta$  and either  $a$  is the first symbol of  $w$ , or  $w = \epsilon$  and  $a = \$$

The method for constructing the collection of sets of valid LR(1)

items is essentially the same as for canonical LR(0) items (cf 3.3).

To appreciate the definition of the closure operation, consider an item of the form  $[A \rightarrow \alpha \cdot B\beta, a]$  in the set of items valid for some viable prefix  $\gamma$ . Then there is a rightmost derivation  $S \xrightarrow{*} S\alpha x \xrightarrow{*} S\alpha B\beta x$  where  $\gamma = S\alpha$ . Suppose  $B\beta x$  derives terminal string  $b\gamma$ . Then for each production of the form  $B \rightarrow \eta$  for some  $\eta$ , we have derivation  $S \xrightarrow{*} \gamma B\beta x \xrightarrow{*} \gamma\eta B\beta x$ . Thus  $[B \rightarrow \cdot\eta, b]$  is valid for  $\gamma$ . Note that  $b$  can be the first terminal derived from  $B$ , or if  $B$  derives  $E$  in the derivation  $B\beta x \Rightarrow b\gamma$ , in which case  $b = a$ . To summarise both possibilities we say that  $b$  can be any terminal in  $\text{FIRST}(B\beta x)$ . As  $\alpha$  cannot contain the first terminal of  $b\gamma$ ,  $\text{FIRST}(B\beta x) \subseteq \text{FIRST}(B\alpha)$ . We can now give the algorithm.

procedure CLOSURE (I)

begin

repeat

for each item  $[A \rightarrow \alpha \cdot B\beta, a]$  in I, each production  $B \rightarrow \gamma$  and each terminal  $b$  in  $\text{FIRST}(B\beta)$  such that  $[B \rightarrow \cdot\gamma, b]$  is not in I  
do add  $[B \rightarrow \cdot\gamma, b]$  to I

until no more items can be added to I

return (I)

end

procedure GOTO (I, x)

begin

let J be the set of items  $[A \rightarrow \alpha x \cdot B, a]$  such that  $[A \rightarrow \alpha \cdot X\beta, a]$  is in I

return (closure(J))

end

begin

$$C := \text{closure}(\{[S' \rightarrow .S, \$]\})$$

repeat

for each set of items  $I$  in  $C$  and each grammar symbol  $X$  such that  $\text{Goto}(I, X)$  is not empty

$$\text{do } C := C \cup \text{Goto}(I, X)$$

until no more sets of items can be added to  $C$

end

Algorithm to construct  $LR(1)$  itemsets.

The algorithm to construct  $LR(1)$  parse table follows:

- construct  $C = \{I_0, I_1, \dots, I_n\}$  the collection of itemsets for  $G'$
- state  $i$  of the parser is constructed from  $I_i$ . The parsing actions are determined as follows:
  - if  $[A \rightarrow a \cdot a\beta, b] \in I_i \wedge \text{Goto}(I_i, a) = I_j$   
then action  $[i, a] := \text{"shift } j"$
  - if  $[A \rightarrow a \cdot, a] \in I_i$  then action  $[i, a] := \text{"reduce } A \rightarrow a"$
  - if  $[S' \rightarrow S \cdot, \$] \in I_i$ , then action  $[i, \$] := \text{"accept"}$ .

If a conflict results from the above rules, the grammar is said not to be  $LR(1)$  and the algorithm fails.

- the goto transitions are determined by:

$$\text{if } \text{Goto}(I_i, A) = I_j \text{ then goto } [i, aA] = j$$

The table formed by this algorithm is called the canonical  $LR(1)$  parsing table. An LR parser using this table is a canonical LR parser. If the parsing action function has no multiply defined entries, then the grammar is an  $LR(1)$  grammar. Every  $SLR(1)$  grammar is an  $LR(1)$  grammar, but the canonical parser may have more states than the  $SLR$  parser for the same grammar.

## CONSTRUCTING LALR PARSING TABLES

This is, in practice, the method of choice as the tables it produces are considerably smaller than the canonical LR tables, yet most common PL syntactic constructs can be expressed conveniently by an LALR grammar.

We look for sets of  $LR(1)$  items having the same core, that is, set of first components, and we may merge these sets with common cores into one set of items.

e.g.:  $C \rightarrow d \cdot, C/d$  have core  $\{ C \rightarrow d \cdot \}$

$C \rightarrow d \cdot, \$$

Note that, in general, a core is a set of  $LR(0)$  items for the grammar at hand, and that an  $LR(1)$  grammar may produce more than two sets of items with the same core.

Since the core of  $GOTO(I, x)$  depends only on the core of  $I$ , the goto's of merged sets can themselves be merged. Suppose we have an  $LR(1)$  grammar (i.e no parsing conflicts). If we replace all states having the same core with their union, it is possible that there is a conflict, but unlikely because

Suppose in the union there is a conflict on lookahead  $a$  because there is an item  $[A \rightarrow \alpha \cdot, a]$  calling for reduction by  $A \rightarrow \alpha$ , and there is another item  $[B \rightarrow B \cdot a, b]$  calling for a shift. Then some set of items from which the union was formed has item  $[A \rightarrow \alpha \cdot, a]$ , and since the cores of all those states are the same, it must have an item  $[B \rightarrow B \cdot a, c]$  for some  $c$ . But then this state has the same shift-reduce conflict on  $a$ , and the grammar was not  $LR(1)$  as we assumed. Thus the merging of states with a common core can never produce a SR conflict because shift actions depend only

on the core, not the lookahead.

It is possible, however, that a merger will produce a reduce-reduce conflict, as the following example shows:

Consider the grammar  $S \rightarrow S$

$$S \rightarrow aAd \mid bBd \mid aBe \mid bAe$$

$$A \rightarrow c \quad B \rightarrow c$$

generating the four strings acd, ace, bcd and bce.

This grammar is LR(1). If we construct the sets of items, we find  $\{[A \rightarrow c, d], [B \rightarrow c, e]\}$  valid for viable prefix ac and  $\{[A \rightarrow c, e], [B \rightarrow c, d]\}$  valid for bc. Neither of these sets has a conflict, and their cores are the same. The union  $\{[A \rightarrow c, d/e], [B \rightarrow c, d/e]\}$  does have a conflict, since reductions by both  $A \rightarrow c$  and  $B \rightarrow c$  are called for on inputs d and e.

An easy, but space consuming LALR(1) table construction method follows, which serves primarily as a definition of LALR(1) grammars. Constructing the entire collection of LR(1) sets of items requires too much space & time to be used in practice.

- Construct  $C = \{I_0, I_1, \dots, I_n\}$ , the collection of sets of LR(1) items.
- For each core present amongst the sets of LR(1) items, find all sets having that core, and replace these sets by their union.
- Let  $C' = \{J_0, J_1, \dots, J_m\}$  be the resulting sets of LR(1) items. The parsing actions for state i are constructed as before for CLR.
- The goto table is constructed as follows: if  $J$  is the union of one or more sets of LR(1) items, e.g.  $J = I_1 \cup I_2 \cup I_k$ , then the cores of  $Goto(I_1, x)$ ,  $Goto(I_k, x)$  are the same as  $I_1$ ,  $I_k$  all have the same core. Then  $Goto(J, x)$  is their union.

## CONSTRUCTION OF LR(1) PARSER

Grammar  $G \rightarrow E$        $T \rightarrow T * F \mid E$   
 $E \rightarrow E + T \mid T$        $F \rightarrow (E) \mid id$

NOTE (FOR LR(1)): An itemset is closed if (assuming there is a production  $B \rightarrow x$ ) it contains  $[B \rightarrow \cdot x, t]$  whenever it contains either  
 $[A \rightarrow \alpha \cdot B \beta r]$   
or  $[A \rightarrow \alpha \cdot B, t]$   
or  $[A \rightarrow \alpha \cdot B C \gamma, v]$ , and  $t \in FIRST(c)$ .

In the next state o moves past B without change of right context. The final state corresponds to the closure of  $[G \rightarrow \cdot E, \$]$

An itemset is inadequate if it contains an item  $[A \rightarrow B \cdot, t]$  and either another item  $[C \rightarrow D \cdot, t]$  (reduce/reduce conflict) or an item  $[C \rightarrow D \cdot \in E, t]$  (shift/reduce conflict) eg state 2 is inadequate on the LR(0)

edge	state	item	closure	etc
$\emptyset$	0	$\{[G \rightarrow \cdot, \epsilon, \$]\}$	$\{[G \rightarrow \cdot E, \$], [E \rightarrow \cdot E + T, \$], [E \rightarrow \cdot T, \$]\} [E \rightarrow \cdot E + T, \$]$ or write $\{[G \rightarrow \cdot E, \$], [E \rightarrow \cdot E + T, \$/\+], [E \rightarrow \cdot T, \$/\+],$ $[T \rightarrow \cdot T * F, \$/\+/\*], [T \rightarrow \cdot F, \$/\+/\*],$ $[F \rightarrow \cdot (E), \$/\+/\*], [F \rightarrow \cdot id, \$/\+/\*]\}$	$[E \rightarrow \cdot T, \$/\+/\*]$
$\emptyset E$	1	$\{[G \rightarrow E \cdot, \$], [E \rightarrow E \cdot + T, \$], [E \rightarrow E \cdot + T, \$/\+]\}$	closed	
$\emptyset T$	2	$\{[E \rightarrow T \cdot, \$/\+], [T \rightarrow T \cdot + F, \$/\+/\*]\}$	closed	
$\emptyset F$	3	$\{[T \rightarrow F \cdot, \$/\+/\*]\}$	closed	
$\emptyset ($	4	$\{[F \rightarrow ( \cdot E), \$/\+/\*]\}$	$\{[F \rightarrow ( \cdot E), \$/\+/\* /\)], [E \rightarrow E \cdot + T, \$/\+], [E \rightarrow T, \$/\+],$	
$\emptyset .id$	5	$\{[F \rightarrow id \cdot, \$/\+/\*]\}$	closed $[T \rightarrow \cdot T * F, \$/\+/\*], [T \rightarrow \cdot F, \$/\+/\*], [F \rightarrow id, \$/\+/\*]$	
$+ E$	6	$\{[E \rightarrow E \cdot + T, \$/\+]\}$	$\cup \{[T \rightarrow \cdot T * F, \$/\+/\*], [T \rightarrow \cdot F, \$/\+/\*], [F \rightarrow (E), \$/\+/\*], [F \rightarrow id, \$/\+/\*]\}$	
$* T$	7	$\{[T \rightarrow T \cdot * F, \$/\+/\*]\}$	$\cup \{[F \rightarrow \cdot (E), \$/\+/\*], [F \rightarrow \cdot id, \$/\+/\*]\}$	
$\emptyset E$	8	$\{[F \rightarrow (E \cdot), \$/\+/\*], [E \rightarrow E \cdot + T, \$/\+]\}$	closed	
$\emptyset T$		$\{[E \rightarrow T \cdot, \$/\+], [T \rightarrow T \cdot * F, \$/\+/\*]\}$	which can be combined with 2... and so on	

3.7

## GENERAL REMARKS.

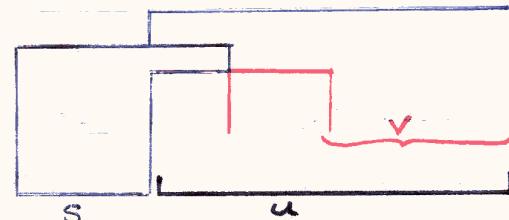
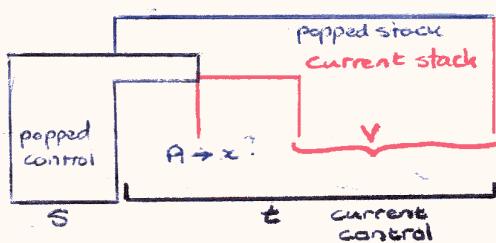
LR parsers are more powerful than LL, and can be automatically produced; in terms of power,

$$LR(1) > LALR(1) > SLR(1) > LR(0).$$

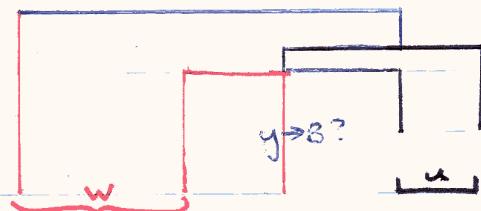
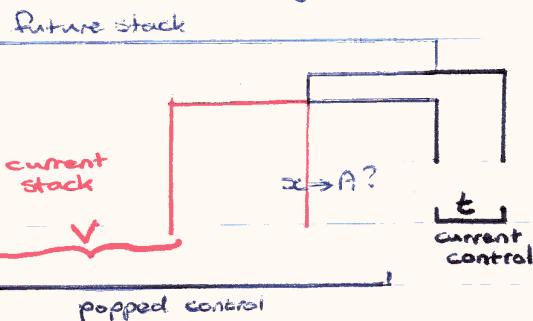
Given  $\alpha \in k$ , it is decidable whether a grammar is  $LR(k)$ , but if  $k$  is not given it is undecidable whether a grammar is  $LR(k)$ .

Under what conditions is a grammar LL or LR? To consider this, imagine a 'snapshot' of the parser at some moment in time:

$LL(k)$  condition for determinism: if  $t[1..k] = u[1..k]$  then  $x = y$   
 (i.e. the first  $k$  symbols of the current control determine which production to choose)

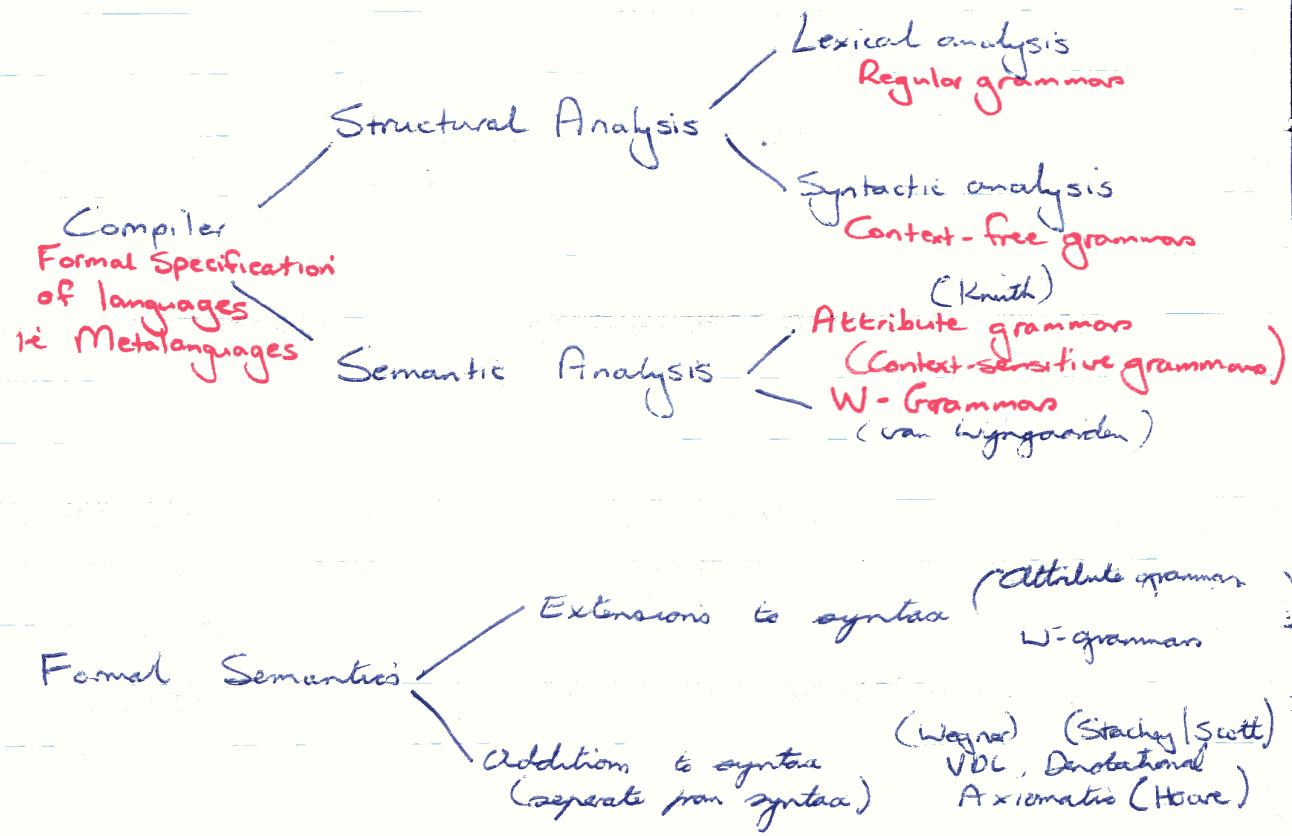


$LR(k)$  condition for determinism: if  $Vxt[1..k] = Wyu[1..k]$  then  
 $V = W$ ,  $x = y$ ,  $A = B$



5/4/23

## GENERAL OVERVIEW



"Semantics" are what happens when we interpret. (We consider compiling as the first step of interpreting, namely interpreting the static semantics, while execution is concerned with the dynamic semantics.)

4

## ATTRIBUTE GRAMMARS. (SEMANTIC ANALYSIS)

Pros:   
 - some analysis  
 - type checking  
 - easier  
 - computable  
 - readability

Semantic analysis cannot be done using a context-free grammar as we need to check across long distances. Using context-free grammars, however, enables us to label our nodes meaningfully.

One way to get around this problem is using attribute grammars. Every non-terminal is given a set of attributes. This can deal with all grammars\*. This can be compared to a compiler using a symbol table and displays (although this would be a useless formalization of semantics). The syntax grammar is extended to contain the attributes of each node (non-terminal) in the parse tree. These attributes are passed up and down the tree in a fashion determined by export and import lists.

Pagan

With each distinct symbol of the context-free grammar there is associated a finite set of attributes, which notationally, are just names. With each distinct attribute there is associated a domain of values. A given attribute may be associated with any number of grammatical symbols. Each node of the syntax tree of a valid program is labeled not only with its grammatical symbol but also by a set of attribute-value pairs, one for each attribute associated with the symbol, and possibly by a logical condition expressing a constraint that must be satisfied by the attribute values involved. The value associated with an attribute occurrence in the tree is determined by various evaluation rules associated with the grammar's production rules.

- \* Eg Suppose we are defining a machine-dependent language more powerful than CS; can handle Turing Machines)

for use on a computer with 32-bit words. An unsigned integer constant is thus syntactically invalid if it exceeds  $2^{31}-1$ . While the set of all unsigned numbers is easily defined in BNF as:

$$\langle \text{Numeral} \rangle ::= \langle \text{digit} \rangle \mid \langle \text{Numeral} \rangle \langle \text{digit} \rangle$$

$$\langle \text{digit} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid \dots \mid 9$$

it is extremely difficult to define the set of numbers  $\{0, 1, \dots, 2147483647\}$  concisely in BNF or its variants.

Instead we associate an attribute *Val*, corresponding to the domain of integers, and write the following specification:

$$\langle \text{Numeral} \rangle ::= \langle \text{digit} \rangle \quad \text{Val}(\langle \text{Numeral} \rangle_1) \leftarrow \text{Val}(\langle \text{digit} \rangle)$$

$$|\langle \text{Numeral} \rangle_2 \langle \text{digit} \rangle$$

$$\text{Val}(\langle \text{Numeral} \rangle_1) \leftarrow 10 * \text{Val}(\langle \text{Numeral} \rangle_2) + \text{Val}(\langle \text{digit} \rangle)$$

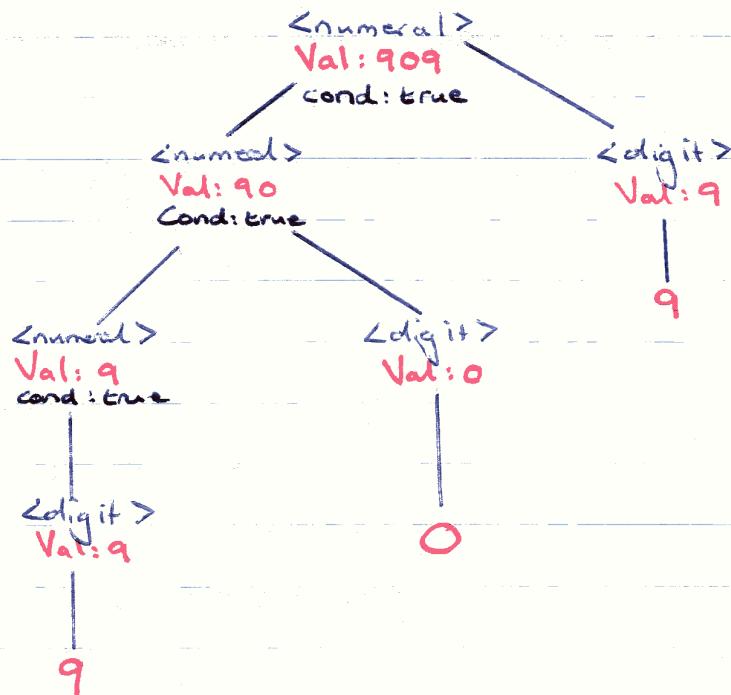
Condition:  $\text{Val}(\langle \text{Numeral} \rangle_1) \leq 2^{31}-1$

$$\langle \text{digit} \rangle ::= 0 \quad \text{Val}(\langle \text{digit} \rangle) \leftarrow 0$$

$$|\dots$$

$$|\ 9 \quad \text{Val}(\langle \text{digit} \rangle) \leftarrow 9.$$

Eg parse tree for 909



In this example, we have information moving up the tree from the leaves to the root. We say that Val is a synthesised attribute of `<numeral>` and `<digit>`. It is possible for an attribute value at a node to be obtained from the node's parent; in this case we call it an inherited attribute. In general, a given grammatical symbol may have both synthesised and inherited attributes, and a given attribute may be synthesised with respect to one symbol and inherited with respect to another. Intuitively, a synthesised attribute at a node corresponds to information arising from the internal constituents of that construct, while an inherited attribute corresponds to information arising from the external context of the construct. Under each alternative of a production rule, there must be an evaluation rule for each synthesised attribute of the symbol on the left (the symbol being defined) and for each inherited attribute of each symbol on the right (each symbol in the alternative).

In general we will have to deal with structured, non-numeric domains of attribute values. For our purposes we will make use of value domains that are enumerations, sets, tuples or sequences, according to the following conventions:

- the constants of an enumeration domain are arbitrarily chosen names enclosed in single quotes, eg 'proc'.
- for sets, the evaluation rules will employ the usual notation  $\{\}$ ,  $\cup$ ,  $\in$  etc.
- tuples belonging to a given domain consist of a fixed number of values, possibly of different types, in a certain order. The field selector operation is `.`, tuples are enclosed in parentheses () with fields separated by commas.

- a sequence is an ordered collection of any number of values of the same type. The notation  $\langle \dots \rangle$  denotes a sequence, with  $\langle \rangle$  the empty sequence. The following primitive functions apply to sequences:
 

append ( $s, v$ )	First ( $s$ ) (of car)
concat ( $s_1, s_2, \dots, s_n$ )	Last ( $s$ )
length ( $s$ )	Tail ( $s$ ) (of cdr)
	allbutlast ( $s$ )

For the sake of uniformity, we further decree that the specification of a subject language by means of an attribute grammar will consist of four parts:

- attributes and values (a list of attributes & domains)
- attributes associated with non-terminals
- production and attribute evaluation rules
- definition of auxiliary evaluation functions.

McD: Refer to the attribute grammar for Eva on note.

An attribute grammar is a context-free grammar  $G = (T, N, S, P)$   $V = T \cup N$ , augmented by attributes  $A = \bigcup_{x \in V} A_x$  where  $A_x = I_x \cup S_x$  and  $I_x \cap S_x = \emptyset$ . (i.e.  $A_x$  a disjoint union of inserted and synthesised attributes).

Note that if  $G$  is reduced,  $S$  appears only on the lhs of a production, so  $I_S = \emptyset$ . Also,  $\forall x \in T$ ,  $I_x = \emptyset$ .

We write a production  $p \in P$  in the form  $X_0 \rightarrow X_1 X_2 \dots X_{n_p}$  where  $n_p \geq 0$ ,  $X_i \in N$ ,  $X_k \in V$ .

For each occurrence  $X_k$  of a symbol  $X$  in  $p$ , for each attribute  $a \in A_X$ , there is an attribute occurrence  $a_k$  (of  $X_k$ ) occurrences with depths / calls of procedures of attribute occurrences with params (formal / actual), input of inherited attribute domains + types.

For each occurrence  $s_0$  of a synthesised attribute on the LHS of a production  $P$  there is a semantic function  $f_{s_0}^P$  (evaluation function) which yields the value of  $s_0$  as a function of other attribute occurrences in  $P$ . (values move up the tree ; formal output)

For each occurrence  $i_k$  of an inherited attribute on the RHS of a production  $P$  ( $1 \leq k \leq n_p$ ) there is a semantic (evaluation) function which yields the value of  $i_k$  as a function of other attribute occurrences in  $P$ . (values thus move down the tree ; actual input)

We must avoid circular definition of attribute - occurrence values, i.e., we need a well-defined attribute grammar (WAG). Ideally we want to go from left to right in each production; inheritance in a down traverse of the tree followed by synthesis in an up traverse.

Algorithm DESIRED:

Processnode ( $X_0$ ) {for any  $X$  examine its  $P$ ?

```

begin
    for  $k = 1$  to  $n_p$  do
        if  $X_k \in N$  then begin
            update  $\{i_k\}$ 
            processnode ( $X_k$ )
        end
        update  $\{s_0\}$ 
    and

```

Initiate by calling Processnode ( $S$ )

An attribute grammar is :

- COMPLETE if  $\forall x \in V$   $S_x$  is defined in all production  $X \rightarrow \dots$   
 $I_x$  is defined in all production  $\dots \rightarrow \dots x \dots$
- WELL-DEFINED (WAG) if complete, and with acyclic dependency graphs for all generable parse trees  
( $\Leftrightarrow$  all attributes are effectively computable)
- PARTITIONED if  $\forall x \in V$ , the attribute set  $A_x$  can be evaluated in the order  $A_x^1, A_x^2, \dots, A_x^m$  where:  
all disjoint  $\begin{cases} A_x^i \subseteq S_x & \text{for } i = m, m-2, m-4, \dots \\ A_x^i \subseteq I_x & \text{for } i = m-1, m-3, m-5, \dots \end{cases}$   
and the dependency graphs for all productions are acyclic. (Every production has an associated evaluation algorithm 'attached' to every node where that production operates.)
- ORDERED (OAG) if it partitions as follows:  
 $A_x^i = T_x^{m-i+1} - T_x^{m-i}$  where  $m$  is the smallest number such that  $T_x^{m+1} \cup T_x^m = A_x$ , where  
 $T_x^1 = T_x^0 = \emptyset$ .  
 $T_x^k = \{a_x : b_x \text{ indirectly (presimistically) depends on } a_x \text{ implies } b_x \in T_x^j \text{ where } j \leq k\}$   
where  $a$  is synthesised / inherited for  $k$  odd / even.

Ordered attributes can be partitioned deterministically (look at actual dependencies, and produce all potential dependencies, and thus construct the total (presimistic) dependencies)

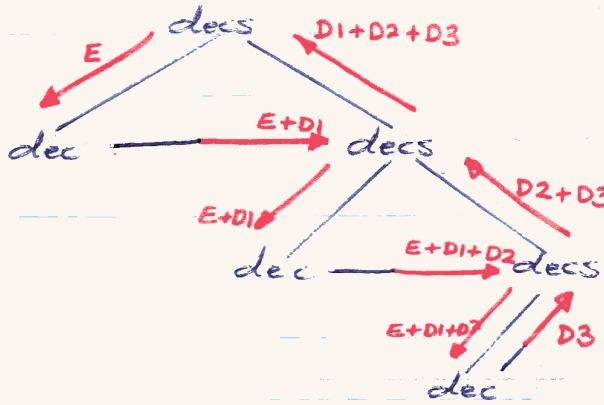
We now consider some aspects of AG's.

### EVALUATION ORDER

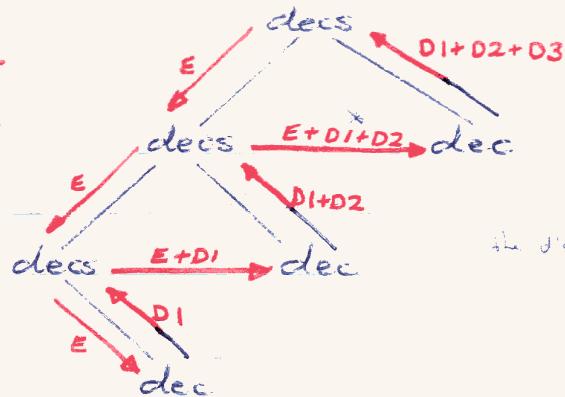
In the AG for Eva, we have the following production  
 $\text{decs} - E - d_1 + d_2 \rightarrow \text{dec} - e - D_1$      $\text{decs} - e + d_1 - D_2 \xrightarrow{\text{inherited}}$   
 synthesised

We can evaluate this in different ways. e.g:

RIGHT RECURSIVE     $\text{decs} \rightarrow \text{dec decs}$



LEFT RECURSIVE     $\text{decs} \rightarrow \text{decs dec}$

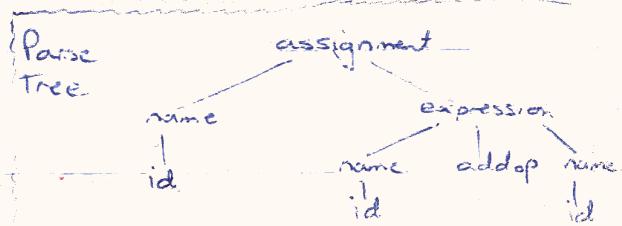


— Given, the local attributes are synthesized.  
 — Also, there is no half at \*, instead the whole takes place in the special node — remember, the D's are sets,  $\{d_1, d_2, d_3\}$

STEPWISE     $\text{decs} - E_1 - e_2 + d_2 - d_1 + d_2 \rightarrow \text{decs} - e_1 - E_2 - D_1 \text{ dec} - e_2 - D_2$

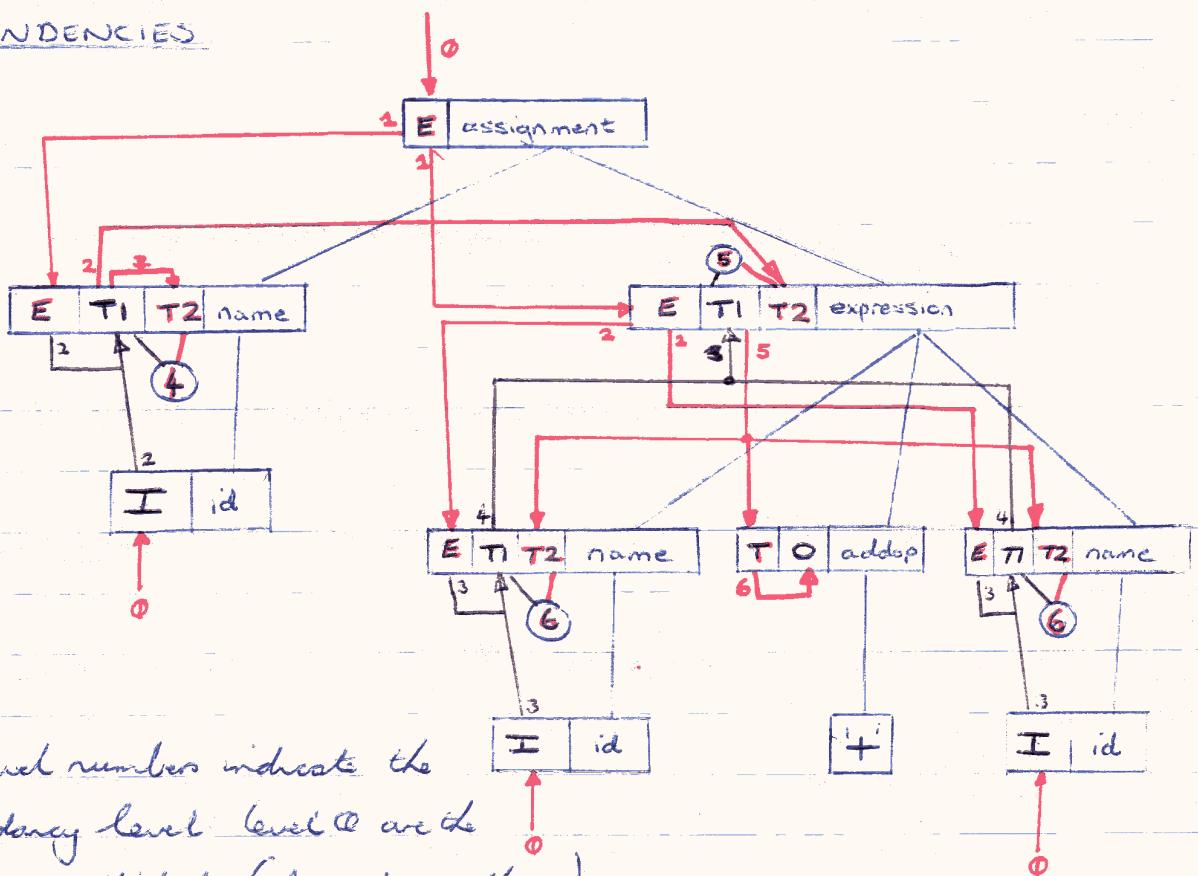
Consider the following attribute grammar:

- assignment -  $E \rightarrow \text{name-}e\text{-INT-int expression-}e\text{-T-int}$   
 $\quad | \quad \text{name-}e\text{-T}_1\text{-t}_1 \quad \text{expression-}e\text{-T}_2\text{-real}$
- expression -  $E = \text{cvt}(T_2, \text{int}) \wedge \text{cvt}(T_3, \text{int}) \rightarrow \text{int, real : t-T}_1$   
 $\rightarrow \text{name-}e\text{-T}_2\text{-t addop-}t\text{-O name-}e\text{-T}_3\text{-t cvt(t,t)}$
- name -  $E = \text{tp}(i, e) : t-T_1 \rightarrow \text{id-I cvt(t,t)}$
- addop - INT - addint  $\rightarrow '+'$
- addop - REAL - addreal  $\rightarrow '+'$



\* If  $T_3$  and  $T_2$  can be converted to int, answer is int, else real

### DEPENDENCIES

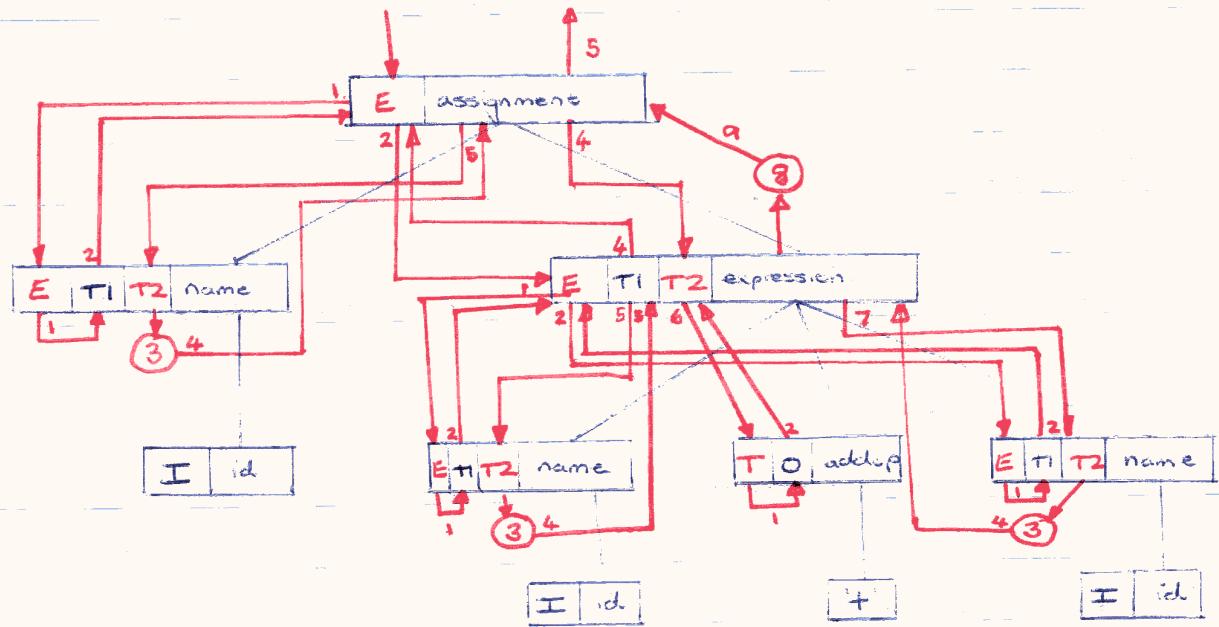


The level numbers indicate the dependency level: level 0 are the intrinsic attributes (depend on nothing).

This dependency diagram shows us two things: (how do we do this for all trees?)

- our grammar (tree) is well-defined: no circular dependencies
- certain attribute calculations can be done in parallel - we have partitioned the evaluation into 7 stages (0-6). By finding these partitions we can in general detect our conflicts.

## EVALUATION ORDER.



## Partition Sets (cf p32)

Eval order	Expression	Name	Addup	
1	{E}	{E}	{T}	inher
2	{T <sub>1</sub> }	{T <sub>1</sub> }	{O}	synth
3	{T <sub>2</sub> }	{T <sub>2</sub> }		inher
4	{O}	{O}		synth

conditions

The partitions always alternate between inherited and synthesised, and always start on inherited and end on synthesised attribute.

## ORDERED ATTRIBUTE GRAMMARS

In order to check that the grammar as a whole has no cycles we need to be totally pessimistic. If we can partition the attributes in a grammar into subsets which are evaluated sequentially in order, with the odd partitions inherited and the even ones synthesised the grammar is ordered (not all well-defined AG's are ordered). We can consider these partitions

to be sets of co-routines, and can say that if we can create co-routines for the grammar, then it is ordered.

## IMPLEMENTATION CONSIDERATIONS (figs refs are to Date & Yone)

Evaluation functions can be co-routines attached to each node through pointers pointing to the routine's re-entrant code.

- simple translation of code in Simul (Fig 8.16)
- recursive procedures using case statements and parameter passing (Figs 8.17 & 8.18)
- using tables and a table driven evaluation function

Attribute storage:

Intermediate attributes:

- overlaying of attributes (i.e. store only those which are necessary; only set aside fixed storage for the final attributes)
- use of local variables of the evaluation procedure

Final attributes:

- storage only at specified nodes (pointers at other nodes referring to it)
- global storage (e.g. at root)

Refs SIGPLAN V14N2? 1979 David A. Watt - "AN EXTENDED AG FOR PASCAL"

Acta Informatica V13 P229-256 Ulwe Kastens -  
"Ordered Attributed Grammars"

## TYPES OF ATTRIBUTE GRAMMARS.

A list of <sup>classes</sup> of AG's is given in descending order of expressive power (i.e. increasing restrictions on dependencies)

AG		Knuth 1968
WAG	well-defined	"
ANCAG	absolutely non-circular	Kennedy & Warren 1976
OAG	ordered	Kostens 1980
m-PAG	m-alternating passes	Jazayeri & Walter 1975
n - PAG	n - left to right passes	} Bochmann 1976
1 - PAG	1 - left to right pass	
= L - AG	L - attributed	} Lewis, Rosenkrantz & Stearns 1974
S - AG	S - attributed	

## 5. W-GRAMMARS (2-level Grammars)

Pagan W-grammars were introduced by Van Wijngaarden for the Algol-68 report. They are capable with dealing with context-sensitive aspects of a language.

We first describe a metalanguage for context-free grammars, and then extend these to W-grammars.

A protonotion is any sequence of lowercase letters and spaces where the spaces are ignored (i.e. something and some thing are the same protonotion). A protonotion is a notion of a particular grammar if there is a production rule that defines it, so that "notion" is basically a new word for "nonterminal". A protonotion ending in symbol corresponds to a terminal symbol of the subject language; in order to eliminate all possible conflicts between metalanguage and subject language, such protonotions are always used in productions, while the correspondence between those and the actual symbols of the language are listed in a separate representation table.

There are four punctuation metasymbols; the left and right sides of a rule are separated by a colon; alternative definitions are separated by semicolons; the protonotions are separated by commas, and the entire rule is terminated by a period.

A hyper-rule is an abstraction of a number of different production rules sharing a common pattern. For example, the BNF grammar for EVA contains three production rules which we may now rewrite as:

declaration sequence : declaration; declaration sequence, declaration.

statement sequence : statement; statement sequence, statement.

letter sequence : letter; letter sequence, letter.

All these rules can be replaced by the hyper-rule

SEQITEM Sequence : SEQITEM; SEQITEM sequence, SEQITEM.

where the use of upper-case letters signifies that SEQITEM is a metanotion which stands for any number of protonotations. The latter are specified by a metaproduction rule or metarule, with the metanotion as its left side. This particular metanotion could be defined by the metarule  
SEQITEM :: declaration ; statement ; letter.

The production rules represented by a given hyper-rule are obtained by consistent substitution of protonotations (specified by the metarules) for the contained metanotations.

The RHS of a metarule may also involve metanotations, possibly the metanotion being defined. By convention, metanotations are not broken up by spaces, spaces are used to separate meta/protonotations in a metarule.

Apart from the representation table, a two level grammar consists of a set of metarules and a set of hyperrules. Each hyperrule stands for a number of production rules, possibly an infinite number, but only one if it involves no metanotations. For each metanotion used, there must be an associated metarule to specify what protonotations it represents. The metarules may be thought of as a higher-level grammar that (with the aid of hyperrules) generates a grammar, where the generated grammar may involve an infinite number of notions and have an infinite number of production rules (or alternatives thereof). This is basically the reason why 2-level grammars are fundamentally more powerful than BNF variants which are finite.

To illustrate the ability of 2-level grammars to handle context sensitive aspects of languages, consider the

grammars "def'ee", which cannot be defined by a context-free grammar.

We use two metanotations:

LETTER :: letter x; letter y; letter z.

TALLY :: i; TALLY ;.

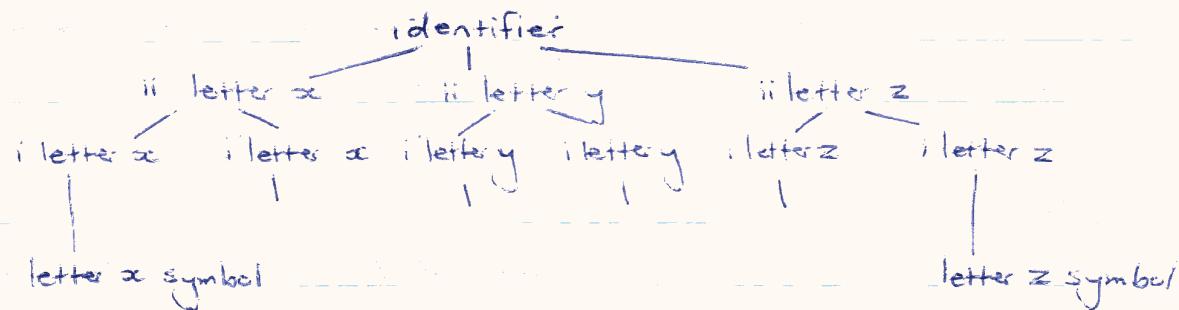
and the hypotheses:

identifier :: TALLY letter x, TALLY letter y, TALLY letter z.

TALLY ; LETTER :: TALLY LETTER, ; LETTER .

i LETTER :: LETTER symbol .

Tree for  $\alpha\alpha\gamma\gamma\alpha\alpha$



A 2-level grammar can be conceived of as a 2-level macro builder which builds up a tree containing expansions of the non-terminals. Parsing such grammars requires a pattern-matching algorithm. W-grammars are very precise but not very practical.

We will now consider aspects of a 2-level grammar for Eva.

program:: new block containing DECS, where DECS consistent.

• A program is a new block containing consistent DECS

DECS :: DEC ; DECS DEC.

DEC :: TYPE type TAG.

TYPE :: VALTYPE ; proc ; proc with PARAMETERS.

TAG :: LETTER ; TAG LETTER . LETTER :: letter ALPHA

ALPHA :: a; b; c; ... x; y; z. VALTYPE :: char; string

PARAMETERS :: PARAMETER ; PARAMETERS PARAMETER.

PARAMETER :: VALTYPE type TAG.

- We arrange for most of the grammars notions to each begin with a protonotion represented by the metanotation NEST, defined by:

$\text{NEST} :: \text{new} ; \text{new DECS} ; \text{NEST new DECS}$ .

A notion of the form NEST NOTION will characterize a syntactic unit of an Eva program iff the NEST protonotion records the name and type information specified by just those declaration sequences and/or parameter lists which bear upon that syntactic unit. Thus the outer block will have as dummy nest new. The notion for each principal construct inside the main block will have a nest of the form new new DECS; the next inner block a nest of the form new new DECS1 new DECS2, etc.

- Note that complete information on the parameters of a procedure, if any, is included in the protonotion corresponding to a procedure name. For example, the protonotion
- proc with string type letter s char type letter c letter c  
type letter s letter u letter b ...  
corresponds to proc sub (string s, char cc) = ...

- In the case of the block and declaration constructs, we embed this syntactic notions still further by supplying a DECS protonotion corresponding to the new set of declarations involved. Where we would simply have the notion block in a context-free grammar, we will now be dealing with notions of the form NEST block containing DECS, where the DECS part records the information given by the declarations inside the block. For example, the "topmost" hyperrule in the

43

grammar is:

program: new block containing DECS, where DECS consistent.

The predicate part where DECS consistent corresponds to the context condition that states that a name cannot be declared more than once at the same level in a program. The formalisation of this condition requires a number of auxiliary hypotheses. First, if the DECS is actually a single DEC, there is no problem:

where DEC consistent : EMPTY.

Otherwise, according to the metarules, the DECS must be of the form DECS DEC, which is the same as DECS TYPE type TAG. Since the type information is not relevant here, we merely stipulate that the final TAG must not appear earlier in the protonotion, as well as (recursively) regaining the earlier part to be itself consistent.

where DECS TYPE type TAG consistent :

where DECS consistent, where TAG not in DECS.

where TAG not in DECS DEC : where TAG not in DECS,  
where TAG not in DEC.

where TAG1 not in TYPE TAG2 : where TAG1 is not TAG2.

EMPTY ::

- if the second level grammar is finite then metarules and hypotheses simply abbreviate many production rules
- if the second level grammar is itself context free then the generated grammar can have an infinite number of notions and productions and can be of type 0.  
(see notes)

Eg

LETTER :: x; y; z .

finite choice  $\Rightarrow$  abbreviation

N :: 1, N1 .

infinite language

identifier :: Nx, Ny, Nz .

infinite alternatives on RHS

N1 LETTER : N LETTER, 1 LETTER . infinite number of rules

LETTER :: LETTER symbol . finite number of rules  $\Rightarrow$  abbreviation

## 6 THE SECD MACHINE (an operational semantics)

LANDIN Computer Journal 66 p308

Operational semantics defines semantics by setting up a abstract machine for interpreting the language; meaning is therefore considered to be analogous to value.

McCarthy's LISP language runs on the abstract SECD machine. As LISP is sufficiently powerful to allow other languages to be written in it, the SECD machine can be considered to provide a general programming semantics method of description. The SECD machine is based on the lambda calculus, a calculus of functions (all values are functions).

### Syntax of Lambda Expressions

$\langle \text{exp} \rangle ::= \langle \text{id} \rangle \mid \langle \text{application} \rangle \mid \langle \text{abstraction} \rangle$

$\langle \text{application} \rangle ::= \lambda \langle \text{operator} \rangle \langle \text{operand} \rangle$

$\langle \text{operator} \rangle ::= \langle \text{exp} \rangle$

$\langle \text{operand} \rangle ::= \langle \text{exp} \rangle$

$\langle \text{abstraction} \rangle ::= \lambda \langle \text{bound var} \rangle \langle \text{body} \rangle$

$\langle \text{bound var} \rangle ::= \langle \text{letter} \rangle$

$\langle \text{body} \rangle ::= \langle \text{exp} \rangle$

### McCARTHY Conditions

We are used to conditions of the form

if  $\langle \text{boolean-exp}_1 \rangle$  then  $\langle \text{exp}_1 \rangle$

else if  $\langle \text{boolean-exp}_2 \rangle$  then  $\langle \text{exp}_2 \rangle$

In this formalism, we replace these by:

$\langle \text{boolean-exp}_1 \rangle \rightarrow \langle \text{exp}_1 \rangle, \langle \text{boolean-exp}_2 \rangle \rightarrow \langle \text{exp}_2 \rangle$

e.g.  $\text{zero}(x) \rightarrow 1, t \rightarrow x * f(x-1)$

(effectively just  $\text{zero}(x) \rightarrow 1, x * f(x-1)$  cf cond in LISP)

The SECD machine consists of four components : the Stack (evaluations), Environment (store), Control (current expression) and Lamp (procedure calls).

The evaluate function could be defined as :

$$\text{evaluate } (E, x) = \begin{cases} \text{is-id}(x) \rightarrow E(x), \\ \text{is-application}(x) \rightarrow \text{operator}(x)(\text{operand}(x)), \\ \text{is-abstraction}(x) \rightarrow \\ \quad \lambda u. \text{evaluate}(\text{extend}(E, u, \text{var}(x)), \text{body}(x)) \\ \text{where } \text{extend}(E, u, x) = \lambda y. y = x \rightarrow u, E(y) \end{cases}$$

In the description of the machine, the letters h will be used to stand for head, and t for tail, so, for example hS is the head of the stack (cf car, cdr).

The Sexp interpreter can be considered as a machine which takes symbolic expressions as input, evaluates them, and outputs the values of the Sexp's. The output function could be described as follows

$$\text{output } (E, x) = \text{iterate}(\text{transform}, \langle \text{empty}, E, x, \text{empty} \rangle)$$

where

$$\text{iterate } (F, \langle S, E, C, D \rangle) = \text{null}(C) \wedge \text{null}(D) \rightarrow S,$$

where  $F, F(\langle S, E, C, D \rangle)$

$$\text{transform } (\langle S, E, C, D \rangle) =$$

$$\text{null}(C) \rightarrow \langle hS - S', E, C, D \rangle \quad \text{where } \langle S', E, C, D' \rangle = hD,$$

$$\text{is-id}(hC) \rightarrow \langle E(hC) - S, E, EC, D \rangle,$$

$$\text{is-application}(hC) \rightarrow \langle \text{rand}(hC) - \text{rate}(hC) - Y - EC, D \rangle,$$

$$\text{is-abstraction}(hC) \rightarrow \langle \text{closure}(hC, E) - S, E, EC, D \rangle,$$

$$hC = 'Y' \rightarrow$$

$$\beta\text{-closure}(hS) \rightarrow \langle \text{empty}, \text{extend}(\text{env}(hS), hS, \text{var}(hS)) \rangle$$

$$\text{body}(hS), \langle tS, E, C, D \rangle,$$

$$\langle hS(hS) - tS, E, EC, D \rangle$$

Note the use of functionals to give the meaning of functionals.  
 The program is part of the data structure (state vector).  
 The interpreter is effectively a function from string syntax to interpretation.

### McCarthy's Work in the Field of Computability

McCarthy formulated a list of problems with which the emergent theory of computation should concern itself, as well as a set of methods and a formalism for dealing with those, which is capable of dealing with symbols rather than numeric data. Since his formalism necessitated the use of recursion, he investigated the extent to which recursive programs and iterative programs were equivalent. He also developed the technique of recursive induction for proving pairs of recursive programs equivalent.

McCarthy applied his formalism to semantics with the idea that the meaning of a program can be expressed in terms of the change it effects on the state vector of an idealised interpreter.

McCarthy's formalism is particularly attractive as it has the form of a high level language, and can thus be used to explain its own semantics as well as being far more useful practically than any other such formalism (eg FSM's).

M.'s formalism is based on an idea similar to that of general recursive functions: there is a set  $F$  of primitive functions and a set  $C$  of functionals for building new functions out of the old, the claim being that the closure  $C(F)$  comprises all computable functions.

Actually, McCarthy just suggested a set  $C$  of strategies, leaving the user to define a base set  $F$  appropriate to the

data structures. In the following section, we define S-expressions and a corresponding set  $F_{S\text{-exp}}$  of functions. Core concepts of generalised composition, conditional expressions, and recursion. Composition generalises function nesting and provides a form of sequencing.

### McCarthy Conditionals

General form  $(p_1 \rightarrow e_1, p_2 \rightarrow e_2, \dots, p_n \rightarrow e_n, e_m)$   
which corresponds to the Algol 60 expression

if  $p_1$  then  $e_1$

else if  $p_2$  then  $e_2$

else if ...

else if  $p_n$  then  $e_n$  else  $e_m$

(McCarthy actually suggested the incorporation of conditionals into Algol in 1959.)

He introduced conditional expressions as a piece of formal mathematical notation, and then gave a process interpretation. However, in contemporary mathematics an expression merely has a value; it is not part of mathematics to say how that value is computed. Purists could object to an evaluation process cast in the role of a definition of a piece of mathematical notation. Note that there cannot be a simple function  $f(p_1, e_1, \dots, p_n, e_n)$  for computing the value of conditionals, since  $(1 > 2 \rightarrow 3/0, t \rightarrow 4)$  has the value 4, but  $f(1 > 2, 3/0, t, 4)$  is undefined.

We can use lambda abstraction to 'delay' the evaluation of the  $e_i$ 's to get around this. In particular,

$(p_1 \rightarrow e_1, t \rightarrow e_2)$  can be modelled by:

$$[\text{ifasfn } (p_1) (\lambda().e_1, \lambda().e_2)]()$$

where  $\text{ifasfn}$  ("if as a function") is a function that returns a second argument selected when applied to false, and a first argument selected when applied to true.

Nevertheless, the conditional expression fundamentally involves understanding a process, and this could be a suggestion of the inadequacy of modern mathematics which does not incorporate the idea of process.

NB: lambda notation (Church) should not be confused with lambda calculus (Church). Lambda notation is a way of giving a systematically complicated default name to any function (eg  $\lambda(x,y).x^2+y$  which is the name of the function which computes  $x^2+y$ ). Lambda calculus is a set of conversion rules allowing  $\lambda$ -names to be changed. The conversion rules reflect the operations of renaming bound variables, applying a function to a set of arguments, etc. Church argued that  $\lambda$ -calculus embodied a precise definition of 'computable'; it has the same power as Turing machines.

We can use  $\lambda$ -notation to write our functions. Consider a function:

$$\text{abs}(x) = (x \geq 0 \rightarrow x, t \rightarrow -x).$$

We apply  $\lambda$ -abstraction, giving the default name:

$$\text{abs} = \lambda x. (x \geq 0 \rightarrow x, t \rightarrow -x)$$

There is a problem in applying  $\lambda$ -abstraction to recursive definitions such as:

$$\text{factorial}(x) = (x=0 \rightarrow 1, x * \text{factorial}(x-1))$$

because of the definition containing its own name (where in  $\lambda$ -notation, the definition is the name). However, considering the general form of a recursive definition

$$f(x) = e(x, f)$$

By  $\lambda$ -abstraction,  $f = \lambda x. e(x, f)$ , where there may be any number of  $f$ 's on the RHS. We reduce this to one by:

$$f = \{\lambda g. (\lambda x. e(x, g))\}(f)$$

Now,  $\lambda g. \lambda x. e(x, g)$  is a functional (say  $F$ ), which takes  $f$  as argument. We have:  $f = F(F)$ , so that  $f$  is a fixed point of  $F$ . We now show that there is a function  $Y$  (with a complex  $\lambda$ -name) which, when applied to  $F$ , yields one of its fixed points  $f$ , i.e.  $Y(F) = f$ , so that  $Y(\lambda g. \lambda x. e(x, g))$  is the  $\lambda$ -name for  $f$ .

We want to find  $f$  so that  $F(f) = f$ . This is solved by:

$$\begin{aligned} f &= \text{self apply } (h) \text{ where } h = \lambda y. F(y(y)) \\ \text{since } F &= h(h) \quad (\text{defn of selfapply}) \\ &= \{\lambda y. F(y(y))\}(h) \quad (\text{defn of } \lambda) \\ &= F(h(h)) \quad (\text{subst. for } y) \\ &= F(f) \quad (\text{since } f = h(h)) \end{aligned}$$

Now we can derive a  $\lambda$ -name for  $Y$ . Selfapply has a name  $\lambda g. g(g)$ , so the pure  $\lambda$ -name for  $f$  is  $\{\lambda g. g(g)\}(\lambda y. F(y(y)))$

Now we have  $Y(F) = f$ , so that finally:

$$Y = \lambda G. \{\lambda g. g(g)\}(\lambda y. G(y(y)))$$

$Y$  was invented by Curry and Feys<sup>1958</sup> and is called the paradoxical combinator (since selfapply generates a number of paradoxes). The development of  $Y$  has a number of interesting theoretical consequences; it is not suggested as a practical mechanism for evaluating recursively defined function applications!

## S-EXPRESSIONS

- S-exps are essentially binary trees with identifiers (atoms) for their leaves. von Neumann architecture, using contiguous sequential memory, is inappropriate for computing with symbolic structures. In particular:
- symbolic data structures can become unpredictably large  
(storage must be allocated incrementally)
  - symbolic processing general calls for the replacement of one structure with another of a different size.

These problems resulted in a set of techniques called (collectively) list processing.

The syntax of S-expressions is

$$\begin{aligned} \langle \text{sexpr} \rangle &::= \text{atom} \mid (\langle \text{sexpr} \rangle \cdot \langle \text{sexpr} \rangle) \mid \langle \text{list} \rangle \mid \langle \text{liststructure} \rangle \\ \langle \text{list} \rangle &::= (\{\text{atom}\}^*) \\ \langle \text{liststructure} \rangle &::= (\{\langle \text{sexpr} \rangle\}^*) \end{aligned}$$

(dotted pair)

$\langle \text{list} \rangle$ s and  $\langle \text{liststructure} \rangle$ s are in fact a special case of dotted pairs, useful as shorthand. (see later)

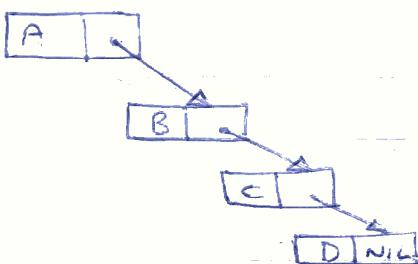
We now describe Flisp. Its set of computable functions of S-exps is thus  $\text{Cmc(Flisp)}$  (the set of LISP computable functions).

- predicate atom :  $\text{Sexpr} \rightarrow \{\text{t}, \text{nil}\}$
- predicate eq :  $\text{Sexpr} \times \text{Sexpr} \rightarrow \{\text{t}, \text{nil}, \perp\}$  (defined for atoms only, cf nodes of trees)
- selectors car :  $\text{Sexpr} \rightarrow \text{Sexpr}$  left subtree } undefined for atoms  
cdr :  $\text{Sexpr} \rightarrow \text{Sexpr}$  right subtree }
- constructor cons :  $\text{Sexpr} \times \text{Sexpr} \rightarrow \text{Sexpr}$  builds tree

## List Notation

To describe binary trees we need an ordered data structure. Lists describe 'linear' binary trees. We introduce the empty tree NIL to make relationships generally applicable even to the last\_elt of a list (as  $(\text{car } (\text{atom}))$  is undefined, but  $(\text{cdr } (\text{atom}, \text{nil}))$  is not)

Eg.  $(A. (B. (C. (D. \text{NIL})))) \quad (= (A \ B \ C \ D))$



The following recursive function list-to-pair describes a process for converting list notation into the corresponding longhand dotted-pair as follows.

list-to-pair  $((e_1 \dots e_n)) = (n=0 \rightarrow \text{nil},$   
 $t \rightarrow e_1. \text{list-to-pair} ((e_2 \dots e_n)))$

## APPLICATION OF THE FORMALISM : PROGRAMMING SEMANTICS

In order to give the semantics of a programming language  $L$ , we need a precise way of specifying well-formed programs PROGL in  $L$  (ie syntax). Secondly, we need a set  $M_L$  of possible meanings of programs in  $L$ , so that if  $P \in \text{PROGL}$ , then the meaning of  $P \in M_L$ . Finally, we need some way of associating with any program  $P$  its meaning. We would like our programs to be deterministic, which suggests modelling the meaning of  $L$

by a function  $L\text{-MEANS} : \text{PROG} \rightarrow M_n$

We presumably need to insist that  $L\text{-MEANS}$  should be computable; however, our determinism requirement is not completely general, and we must consider  $L\text{-MEANS}$  being a function as a special case.

### PROG

The meaning of a program is to a large extent independent of the particular (strings of) symbols chosen to represent its syntactic constructs. However, there is generally a strong relation between grammatical structure and semantics. McCarthy observed that an alternative formal specification for syntax of programs more appropriate for semantics can be created from selectors and predicates. The predicates correspond to the distinct syntactic structures (e.g. is-id?), the selector 'parse' these by selecting components of non-terminals. This method of specifying syntax has been taken up in VDL. McCarthy called such a syntactic specification abstract and analytic in contrast with concrete (since it does not focus on a particular external representation) and synthetic (since it details the analysis of programs rather than their construction or synthesis).

$$\begin{aligned} \text{eg } \text{is-prog}(P) = & \text{ is-const}(P) \vee \text{is-vble}(P) \\ & \vee (\text{is-sum}(P) \wedge \text{is-prog}(\text{left}(P)) \\ & \quad \wedge \text{is-prog}(\text{right}(P))) \\ & \vee (\text{is-product}(P) \wedge \text{is-prog}(\text{left}(P)) \\ & \quad \wedge \text{is-prog}(\text{right}(P))) \end{aligned}$$

M<sub>L</sub>

$M_L$  is considered to be a set of functions from state to state

Eg Consider the language  $C$  of a simple one-address computer with a single accumulator. The commands are:

LOAD <reg>	$\text{acc} \leftarrow \text{reg}$
STO <reg>	$(\text{reg}) \leftarrow (\text{acc})$
ADD <reg>	$(\text{acc}) + (\text{reg})$
LI <const>	Load immediate $\text{acc} \leftarrow \text{const}$

Formal syntax

IS-command ( $c$ ) = (is-load ( $c$ ) and is-reg (sel-reg ( $c$ )))  
 or (is-sto ( $c$ ) and is-reg (sel-reg ( $c$ )))  
 or (is-add ( $c$ ) and is-reg (sel-reg ( $c$ )))  
 or (is-lim ( $c$ ) and is-val (sel-reg ( $c$ )))

A state is a pairing of registers  $R$  with values  $V$ . We define the function:

$$\text{update} : R \times V \times S \rightarrow S$$

Here  $\text{update}(r, v, s)$  is the state resulting from updating  $s$  by associating value  $v$  with register  $r$ . The function  $\text{lookup}$  associates registers with their values

$$\text{lookup} : R \times S \rightarrow V$$

We define  $C$ -means in terms of the function step, which describes the result of executing one command

$$\begin{aligned} \text{step}(c, s) &= (\text{is-load } c) \rightarrow \text{update}(\text{acc}, \text{lookup}(\text{sel-reg } c, s), s) \\ &\quad \text{is-sto } c \rightarrow \text{update}(\text{sel-reg } c, \text{lookup}(\text{acc}, s), s), \\ &\quad \text{is-add } c \rightarrow \text{update}(\text{acc}, \text{lookup}(\text{acc}, s) + \text{lookup}(\text{sel-reg } c, s), s) \\ &\quad \text{is-lim } c \rightarrow \text{update}(\text{acc}, \text{sel-reg } c, s) \end{aligned}$$

$$C\text{-means}(P) = \lambda s. \text{null}(P) \rightarrow s, C\text{-means}(\text{rest}(P))(s, \text{step}(\text{first}(P), s))$$

McCarthy's approach to semantics has been substantially developed, particularly by VD Lab. McCarthy and Painter (1967) prove the correctness of a compiler for a simple expression language E (cf bottom p 53) into the machine language C. This proof introduces the technique of structural induction, later developed by Burstall (1969).

### THE METACIRCULAR EVALUATOR.

The LISP evaluator embodies the environment model of evaluation, which has two basic parts:

- to evaluate a compound expression (other than a special form (eg define)) evaluate the subexpressions and then apply the value of the operator subexpression to the value of the operand subexpressions.
- to apply a compound procedure to a set of arguments, evaluate the body of the procedure in a new environment, constructed by extending the environment part of the procedure object by a frame in which the formal parameters of the procedure are bound to the actual arguments to which the procedure is to be applied.

These two rules describe the essence of the evaluation process - a basic cycle in which expressions to be evaluated in environments are reduced to procedures to be applied to arguments, which in turn are reduced to new expressions to be evaluated in new environments, and so on, until expressions are reduced to symbols (whose values are looked up in the environment) and primitive procedures (which are applied directly).

Given the ability to apply primitives, we may question the need for the evaluator. The job of the evaluator is not so much to specify the primitives of the language as to

provide the means of combination and abstraction which bind the primitives to form a language. Specifically the evaluator :

- allows us to deal with nested expressions
- allows us to use variables
- allows us to define compound procedures
- provides the special forms which must be evaluated differently from procedure calls

The core of the evaluator consists of two procedures eval and apply. eval takes as arguments an expression and an environment. It classifies the expression and directs its evaluation. eval is structured as a case analysis of the syntactic type of the expression to be evaluated. Each type of expression has a predicate that tests for it and selectors for accessing its parts. Eval is defined as :

(define (eval exp env))

(cond

- ((self-evaluating? exp) exp) (eg numbers)
- ((quoted? exp) (text-of-quotation exp))
- ((variable? exp) (lookup-variable-value exp env))
- ((defn? exp) (eval-defn exp env))
- ((assignment? exp) (eval-assignment exp env))
- ((lambda? exp) (make-procedure exp env))
- ((conditional? exp) (eval-cond (clauses exp) env))
- ((application? exp) (apply (eval (operator exp) env) (list-of-values (operands exp) env)))

(else (error "Unknown-exp-type --" exp)))

)

)

apply takes two inputs, a procedure and a list of arguments to which the procedure is to be applied.

```
(define (apply proc args)
  (cond ((primitive-proc? proc) (apply-primitive proc args))
        ((compound-proc? proc)
         (eval-sequence (proc-body proc)
                       (extend-environment
                        (parameters proc)
                        arguments
                        (proc-environment proc)))
         ))
        (else (error "unknown proc type" "proc"))))
```

The evaluator uses data abstraction to decouple the general rules of operation from the details of how expressions are represented. The syntax of the language being evaluated is thus solely determined by the procedures which classify and extract pieces of expressions (the predicates and selectors).

### ENVIRONMENTS.

An environment is a sequence of frames, each of which is a table of bindings associating variables with values. A variable in an environment is looked up by finding the binding for the variable in the environment, and returning the value part of the binding. To extend an environment

by a new frame associating variables with values, we construct a new frame of bindings and adjust this frame to the environment. To set a variable to a new value in a specified environment, we alter the value part of the binding. To define a variable, we add a new binding to the head of the first frame in the environment, if the variable is already bound in that frame we can do exception processing. (eg bind new value; general error, etc)

An environment can be represented as a list of frames, where a frame is a list of bindings, where a binding is a pair containing a variable and its value. (This is only one of many plausible representations for an environment)

### TREATING EXPRESSIONS AS PROGRAMS

One operational view of the meaning of programs is that a program is a description of an abstract (perhaps infinitely large) machine. For example, consider the following program to compute factorials:

```
(define (factorial n)
  (if (= n 1)
      1
      (* (factorial (- n 1)) n)))
```

We can regard this program as a description of a machine containing parts that decrement, multiply, and test for equality, together with a two position switch and another factorial machine (the factorial machine is infinite as it contains another factorial machine within it). In a similar way, we can regard the evaluator as a

machine which takes as input a description of a machine. Given this input, the evaluator configures itself to emulate the machine described. For example, if we fed the evaluator with a description of the factorial machine, it would be able to compute factorials. The evaluator can thus be regarded as a universal machine capable of mimicing other machines when given their descriptions. (it thus regards the user programs as its 'data', one of the reasons why this distinction is unclear in LISP)

## OVERVIEW OF LISP

Lisp is:

- an expression language (S-exps)
- a functional language (expression, command & binding structures implemented by functions)
- a symbol-processing language (symbolic atoms)
- a list-processing language (lists are S-exps)

The LISP machine is an already 'impluit' function EVAL to which successive LISP expressions are submitted, and interpreted as function applications to be evaluated.

DATA: linked lists of linked lists of... of atoms (symbolic or numeric) all referenceable through a central object list (OB-list). Each atom has a property list (P-list) attached, structured in ( $\langle$ attribute name $\rangle$ ,  $\langle$ attribute value $\rangle$ ) pairs.

PROGRAM: data list with prepsic structure and atoms playing the role of identifiers. Sequencing is by functional composition, alternative by cons function, looping by recursive function calls.

ENVIRONMENT - most recent association, stored at beginning of association list (a-list); with global associations (stored as 'APVAL') on p-list taking precedence. The a-list is accessible to programmers through setq. Parameters can be passed evaluated or unevaluated (quoted). User functions are of type EXPR or FEXPR, primitive functions are of type SUBLR or FSUBLR.

### Other Readings

"A Formal Description of a Subset of Algol" J. McCarthy.  
Proc. IFIP, North-Holland, 1966.

## OPERATIONAL Semantics & SECD

We generalise Church's lambda expressions (in which even numbers are  $\lambda$ -exps) to get Tandis's applicative expressions (AE's) in which we include numbers, (lists of) identifiers,  $\lambda$  functions, etc. Many aspects of high level languages can be directly carried over into AE's (including recursive functions using Curry's paradoxical combinator). Consider the conditional expression if  $p$  then  $a$  else  $b$ . We could describe this by  $\{\text{ifAsfn}(p)\}(a, b)$  where ifAsfn is a  $\lambda$ -function which returns a function first which selects  $a$ , or a function second which returns  $b$ . More precisely:

$$\text{ifAsfn}(t) = \begin{cases} \text{first where } \text{first} = \lambda(x,y).x & \text{if } p \text{ is true} \\ \text{second where } \text{second} = \lambda(x,y).y & \text{if } p \text{ is false} \end{cases}$$

This satisfactorily handles many conditionals, but not those with the form if  $a=0$  then  $b$  else  $(1/a)$ . We have to delay the evaluation of both the conclusion and alternative.

$\lambda()$ .  $1/0$  is a well-defined function, even if an attempt to apply it results in an error. We can thus render one conditional more adequately by:

$$\{\text{ifAsfn}(p)\}(\lambda().a, \lambda().b)( )$$

The same technique can be used to give meaning to for commands, and call by name. P.Z. Ingerman (1961) invented a technique for implementing call by name based on procedures known as 'thunks'. This idea is useful in using AE's to accomplish the same. If we have the declaration: <sup>Arcs</sup> integer procedure simple( $x$ ) ; integer  $x$ ;

$$\text{Simple} := x + 2;$$

the body is transformed into  $\text{Simple} := x( ) + 2$  and a call  $\text{Simple}(a+b)$  transformed into  $\text{Simple}(\lambda().a+b)$

## EVALUATING AE'S

(This section is in essence equivalent to the SEC0 section).  
The most direct approach to defining the values of AE's is to define a function evaluate( $\alpha$ ) giving the value of the AE  $\alpha$ . We consider the environment to be a function  $E$  giving the values of the free variables (identifiers) in  $\alpha$ .  
Thus we define:

$$\begin{aligned} \text{val}(\alpha) = & \text{is-ident } (\alpha) \rightarrow E(\alpha), \\ & \text{is-const } (\alpha) \rightarrow \{\text{val}(E, \text{sel-opr}(\alpha))\} [\text{val}(E, \text{sel-optk})] \\ & \text{is-lambda } (\alpha) \rightarrow f \\ & \quad \text{where } f(u) = \text{val}(E_{\text{new}}, \text{sel-body}(\alpha)) \\ & \quad \text{where } E_{\text{new}} = \text{extend}(E, u, \text{sel-br}(\alpha)) \end{aligned}$$

The function extend updates the environment of  $E$  by binding  $u$  to the bound variable part of  $\alpha$ .

$$\text{extend}(E, u, \alpha) = \lambda y. y = \alpha \rightarrow u, E(y)$$

Whatever the merits of evaluate as a model for AE's, it is certainly not obviously related to the way AE's can be evaluated by machine. - how are AE's 'run'? Tandem (1964) defined a set of states, and a function transform : state  $\rightarrow$  state, such that repeated application of transform described mathematically a sequence of steps to evaluate an AE. A state  $(S, E, C, D)$  consists of a stack, environment, control and dump. (In terms of ALGOL60 runtime storage, the control is the code being executed, the stack is the current frame, the dump is the frame stack up to and including the previous frame, and the environment is the static link or display.)

transform  $\langle S, E, C, D \rangle =$

$\text{null}(C) \rightarrow \langle \text{hd}(S) :: S', E', C', D' \rangle$  where  $D = \langle S', E', C', D' \rangle$

(let  $x = \text{hd}(C)$ )

$\text{is-ident}(x) \rightarrow \langle \text{lookup}(E, x) :: S, E, C, D \rangle,$

$\text{is-comb}(x) \rightarrow$

where  $C_{\text{new}} = \text{sel-apd}(x) :: \text{sel-apr}(x) :: (\text{ap} :: \text{El}(C))$

$\text{is-lambda}(x) \rightarrow \langle \text{mk-closure}(x, E) :: S, E, C, D \rangle,$

$x = \text{ap} \rightarrow$

(let  $f = \text{hd}(S)$ )

$\text{not is-closure}(f) \rightarrow \langle S_{\text{new}}, E, C, D \rangle$

where  $S_{\text{new}} = f(2^{\text{nd}}(S)) :: \text{El}(\text{El}(S))$

(let  $\text{cl-body} = \text{sel-body}(f)$  /\* body of closure \*/)

let  $\text{cl-bv} = \text{sel-body}(f)$

let  $\text{cl-env} = \text{sel-env}(f)$

$\langle [], \text{extend}(\text{cl-env}, 2^{\text{nd}}(S), \text{cl-bv}), \text{cl-body}, D_{\text{new}} \rangle$

where  $D_{\text{new}} = \langle \text{El}(\text{El}(S)), E, C, D \rangle$

)))

To operate transform, we start with an empty stack and dump, the AE to be evaluated on control, and the initial environment as the environment component of the state. The function transform is repeatedly applied until the control and dump are empty. (see Landin, 1964). Most of the complexity of transform is in the application of closures (corresponding to starting a new runtime <sup>stack</sup> frame).

Eg evaluate  $\text{sq}(a)$  in environment [a.3]. Successive snapshots of the state are:

$\langle [], [a.3], \text{sq}(a), [] \rangle$

$\langle [], [a.3], [\text{a sq ap}], [] \rangle$

$\langle [3], [a.3], [\text{sq ap}], [] \rangle$

$\langle [\text{sq } 3], [a.3], [\text{ap}], [] \rangle$

$\langle [9], [a.3], [], [] \rangle$

Unlike evaluate, transform specifies the evaluation order:  
first the operand, then the operator & then apply the latter  
to the former.

Eg2 evaluate  $u(\lambda x.y)$  in environment  $[y.27]$   
where  $u = \lambda f. \{\lambda y. f(y)\}(z)$

$\langle [], [y.27], u(\lambda x.y), [] \rangle$

$\langle [], [y.27], [\lambda x.y \text{ up}], [] \rangle$

$\langle [\text{Closure}], [y.27], [u \text{ up}], [] \rangle$

where Closure has fm  $\lambda x.y$  and env  $[y.27]$

$\langle [\text{Closure Closure}], [y.27], [\text{up}], [] \rangle$

where Closure is fm  $\lambda f. \{\lambda y. f(y)\}(z)$  with env  $[y.27]$

$\langle [], [f. \text{ Closure } y.27], \{\lambda y. f(y)\}(z), D_1 \rangle$

where  $D_1$  is  $\langle [], [y.27], [], [] \rangle$

$\langle [2], [f. \text{ Closure } y.27], [\lambda y. f(y) \text{ up}], D_1 \rangle$

$\langle [\text{Closure } 2], [f. \text{ Closure } y.27], [\text{up}], D_1 \rangle$

where Closure has fm  $\lambda y. f(y)$  and env  $[f. \text{ Closure } y.27]$

$\langle [], [y.2f. \text{ Closure } y.27], f(y), D_2 \rangle$

where  $D_2$  is  $\langle [], [f. \text{ Closure } y.27], [], D_1 \rangle$

$\langle [2], [y.2f. \text{ Closure } y.27], [f \text{ up}], D_2 \rangle$

$\langle [\text{Closure } 2], [y.2f. \text{ Closure } y.27], [\text{up}], D_2 \rangle$

$\langle [], [\lambda. 2y.27], y, D_3 \rangle$

where  $D_3 = \langle [], [y.2f. \text{ Closure } y.27], [], D_2 \rangle$

Note how the env binding  $y$  to 27 has been reinstated from  
Closure. Now we stack 27 and unwind...

$\langle [27], [\lambda. 2y.27], [], D_3 \rangle$

$\langle [27], [y.2f. \text{ Closure } y.27], [], D_2 \rangle$

$\langle [27], [f. \text{ Closure } y.27], [], D_1 \rangle$

$\langle [27], [y.27], [], [] \rangle$

We now show how the transform function of the SECD machine can be used to define a function mechanical-val which evaluates an AE  $\alpha$  in an environment  $E$ . First we define a function iterate, which repeatedly applies a function  $f$  to a state  $S$ .

$$\begin{aligned} \text{iterate } (f, S) = \text{null } (\text{sel-control } (S)) \rightarrow \\ (\text{null } (\text{sel-dump } (S)) \rightarrow \text{sel-stack } (S), \\ \text{iterate } (f, \text{reinstate-dump } (S)), \\ \text{iterate } (f, f(S))) \\ \text{where reinstate-dump } (S) = \text{consent-to-stack} \\ (\text{sel-stack } (S), \text{sel-dump } (S)) \end{aligned}$$

Next we define a function load-AE  $(E, \alpha)$  which loads an AE  $\alpha$  in an environment  $E$  into an SECD state:

$$\text{load-AE } (E, \alpha) = \langle [], E, \alpha, [] \rangle$$

Finally:

$$\text{mechanical-val } (E, \alpha) = \text{iterate } (\text{transform}, \text{load-AE } (E, \alpha))$$

Most writers on semantics (including Fandin) consider evaluate to evaluate the 'real' meaning of AE's, and think of mechanical-val as a description of an implementation of AE's. The reason for this is that the  $\langle S, E, C, D \rangle$  model is closely related to runtime storage structures in recursive languages, and that mechanical-val specifies an evaluation order whereas evaluate doesn't (this reflects again the malleability of mathematical formalisms to incorporate the concept of process). Strictly speaking, mechanical-val and evaluate are different descriptions of the same function. The

description of mechanical val is inspired by, and suggests, a particular way to evaluate AE's; but that does not alter its status as a function.

(Algol semantics - cf. Brady p 238 ff.)

NB Lucas (1971) points out that the state of a Sardini; McCarthy or Urraca-style approach to semantics depends crucially on the language being defined, and he discussed a correlation between the complexity of the state needed to define a language using the Urraca machine and the complexity of the language.

## 7 OPERATION SEMANTICS : VDL.

INTRO  
(Body)

One of the outstanding early achievements in formal semantics was the IBM Vienna Laboratories defn of the complex programming language PL/I. This project was based on Zandri's AE defn of Algol 60 and McCarthy's abstract syntax / state transformation semantics. The overall aim of the Vienna definition project was to lay down precise and detailed conditions which had to be met by an implementation.

The over-riding concern for consistent implementations is reflected in the V labs defn of 'meaning' as 'the entire sequence of states corresponding to a program run'. This has the advantage of being very precise, but is extremely difficult to reason with about the behaviour of a program.

Following McCarthy and Zandri, the V defn interprets an abstract syntax representation of a program. One of the contributions of the Vienna work was a calculus of structured objects which is applicable equally to machine states and abstract syntax trees. The cornerstone of the calculus is a powerful replacement operator  $\mu$ . As in the case of the SECO machine, one of the state components is a control, consisting of the piece of program currently being interpreted. Since the order in which expressions are evaluated is deliberately left partially unspecified, the control component is itself a tree, and is represented as a structured object, with instructions at the tree nodes. Any of the instructions at the leaves of the control tree is a candidate to be executed, so that there is generally a set of possible next states. Thus the function analogous to the SECO machine transform returns

a set of successor states, and a program execution is a sequence  $s_1, s_2, \dots$  of states such that  $s_{i+1} \in \text{transform}(s_i)$ . An execution step modifies the state and, in particular, the control tree. Two commonly occurring modifications to the control tree are:

- delete the current node and pass an appropriate value to higher nodes in the tree.
- replace the current node by a tree.

The VDL is unsatisfactory for reasoning about programs. Part of the problem seems to be that the sequence of states is completely unstructured. There is scope for modification, however, as the control tree really is a tree.

7.2  
(Pagan)

## THE OPERATIONAL APPROACH - VDL

We take as a starting point some suitable formulation of the abstract syntax of the subject language and assume that there will never be any attempt to use the semantic specifications for the analysis of syntactically invalid programs. In the operational approach, meaning is described in terms of such devices as abstract machines with discrete 'states' and more-or-less explicit sequences of computational operations. The operational approach is most directly oriented towards language implementors - taken to the extreme, an actual concrete or interpreter for a subject language can be regarded as the language's definitive specification. However, most compilers have features / restrictions which have little or nothing to do with the language definition. Thus most operational definitions specify an abstract rather than a real processor for the

subject language.

A compiler-oriented operational definition emphasises the specification of a process for translating the subject language into some other (presumably simple) representation which is independent of the subject language. An interpreter-oriented defn, on the other hand, is largely concerned with modelling the process of executing subject-language programs (usually in abstract form); the specifications for any necessary translations of programs from textual to abstract form are deemphasised or omitted. VDL is an interpreter-oriented method.

VDL is one of the oldest and best-known metalinguages for writing defns of formal systems that model the interpretation of (usually abstract) programs. The central aim is to define an abstract machine for interpreting abstract programs of the subject language by passing through a sequence of discrete states, and to this extent its mode of operation is independent of the subject language. Actually, however, it is highly tailored to the subject language by virtue of the fact that the structure of its states and the allowable transitions from one state to another must be specifically defined. The allowable state transitions are defined by a set of instruction definitions written in a special notation; these typically make up the bulk of a VDL lang. spec.

### 7.3 NOTATION FOR OBJECTS & ABSTRACT SYNTAX.

VDL incorporates a general notation for describing and manipulating hierarchical structures; these structures may be pictured as trees with labelled edges but no left-right ordering.

An elementary object is an object with no internal structure of its own, so that it corresponds to a leaf node of a tree representing some structured object that contains it. A structured or composite object is denoted by a set of one or more pairs of the form  $\langle s, a \rangle$ , where  $s$  is a selector and  $a$  denotes another object (either elementary or structured). A set of  $n$  selector-object pairs may be represented as a tree with  $n$  branches, each of which is labeled with one of the selectors and leads to a subtree representing the object selected. The presence of the selectors make left-right ordering unnecessary.

The null object  $\perp$  is a special elementary object which serves to indicate the "emptiness" of an object or absence of some component of an object; the null object (and its selector) have no corresponding branch in the tree diagram, so for example  $\langle \langle s\_p1, A \rangle \rangle$  is the same object as  $\langle \langle s\_p1, A \rangle, \langle s\_p2, \perp \rangle \rangle$ .

A selector\* may be used as a function that takes an object as an argument and yields the corresponding component of the object, considering the object as a tree, the function yields the subtree below the branch labeled by the selector. If the object does not have an object corresponding to the selector, the null object results. The selectors for the immediate components\*\* of any object must be all different.

The mutation operator  $\mu$  builds new objects which are modifications of old ones. An expression of the form  $\mu(a; \langle s, b \rangle)$  returns the object obtained by replacing the  $s$  component of  $a$  with the object  $b$ .

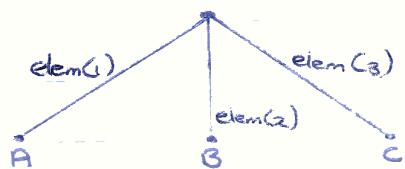
There is often a need for objects which are to be regarded as lists of other objects, such that there may be any number

\* \* immediate descendants in the tree.

\* by convention, selector names begin with 's'

of objects in a list. The convention is to use  $\text{elem}(1)$ ,  $\text{elem}(2)$ ... as the selectors of the list components (elements). The list  
 $f = \{\langle \text{elem}(1) : A \rangle, \langle \text{elem}(2) : B \rangle, \langle \text{elem}(3) : C \rangle\}$

can be pictured as:



and can be abbreviated to  $\langle A, B, C \rangle$ .  $\text{elem}(2)(f) = B$ , etc.

### Standard functions:

- $\text{head} \equiv \text{elem}(1)$  so  $\text{head}(f) = A$
- $\text{tail}$  as cdr in Lisp, so  $\text{tail}(f) = \langle B, C \rangle$
- $\text{length} = \# \text{ elts in list}$ , so  $\text{length}(f) = 3$

The concatenation operator for lists consists of a caret  $\wedge$ .

$$\text{eg } f \wedge \text{tail}(f) = \langle A, B, C, B, C \rangle$$

$f \wedge \text{head}(f)$  undefined, as  $\text{head}(f)$  is not a list.

$$f \wedge \langle \text{head}(f) \rangle = \langle A, B, C, A \rangle$$

The empty list is  $\langle \rangle$

A predicate is a function  $p: \text{object} \rightarrow \{T, F\}$ . By convention, names of predicates begin with 'is-'. Predicates which test for <sup>specific</sup> elementary objects are formed by appending the symbol for the object, eg 'is-'+'(opr)', 'is- $\pi_2(z)$ '. Predicates which test objects for membership in various classes must be explicitly defined by equations, eg

$$\text{is-asmt-stmt} = (\langle s-\text{lhs} : \text{is-var} \rangle, \langle s-\text{rhs} : \text{is-exp} \rangle)$$

states that an object will satisfy 'is-asmt-stmt' iff it consists of a component satisfying is-var and selected by s-lhs, and a component satisfying is-exp and selected by s-rhs.

Predicates may also be defined as disjunctions of other predicates - eg

$$\text{is\_abcd} = \text{is\_A} \vee \text{is\_B} \vee \text{is\_C} \vee \text{is\_D}.$$

By convention, predicates that test for lists end in '-list' and are not explicitly defined by equations, eg the predicate `is-aest-stmt-list` is satisfied by any list of objects that all satisfy `is-exp`.

A VDL semantics specification includes a set of predicates describing the structure of a state  $\mathcal{S}$  of the abstract machine. All VDL abstract machines have a control component in the state, a special object called its control tree.

M.D. Both the program syntax and the state structure are described in the same syntactic notation. The advantages of the tree representation used are:

- closed under substitution of subtrees for subtrees
- unique access paths to components (of DOAG's)
- finite access paths to components (of DAG's)
- can represent DOAG's and DAG's by allowing composite selection operators at terminal vertices.

The VDL approach is based on abstract syntax (semantics of the program as a whole is defined by associating semantics with each syntactically-specified program component) and McCarthy's notion (1962) that "the meaning of a program is defined by its effect on the state-vector" (i.e. in terms of runtime transformations).

## THE FORMALISM (NB slightly different notation to Pagan)

A Vienna Definition System  $V$  is a sextuple  $(E, C, \Omega, S, \delta, \mu)$  where

- $E$  is a set of elementary objects  $E \cap C = \emptyset$
- $C$  is a set of composite objects  $E \cup C = \Omega$  (object)
- $\Omega$  is the empty object  $\Omega \in E$
- $S$  is a set of selectors  $S \subseteq E$
- $\delta$  is a general selector operator  $\delta: S \times E \rightarrow E$
- $\mu$  is the construction (mutation) operator  $\mu: \Omega \times S \times E \rightarrow C \subseteq \Omega$

- Axioms:
1.  $s(e) = \Omega \quad (s(e) = \delta(s, e))$  where  $s \in S, e \in E$
  2.  $(e_1 = e_2)$  is decidable
  3.  $c_1 = c_2 \rightarrow \forall s, s(c_1) = s(c_2)$
  4.  $\delta'(e, \mu(\Omega, s, \delta)) = (s = s_1 \rightarrow \delta_1, \delta(e, s_1))$

Wegner: "To obtain ... understanding of a class of computational structures, it is necessary to consider both representation independent axiomatic characterisations and specific models that are realisations of the axiom system."

" $(E, C, \Omega, S, \delta)$  with axioms 1-3 define a purely syntactic system ... for which equality of objects implies identity of structure. The operator  $\mu$  with axiom 4 introduces semantic content ... in the sense that two different expressions containing ...  $\mu$  may define the same ... object".

"Composite objects may be regarded as tree-structured computer memories. Composite selectors may be regarded as addresses that select locations in memory, where locations may have components that are locations\*. The object selected may be regarded as the value contained in the location"

\* cf tree-structured filing systems.

The tree has selectors as edges, E's objects as leaves, and C objects as branches.

A programming language is characterised by a set of programs. Programs and program components are represented as Vanna trees. The composition of sets from simpler sets is represented by composition of predicates (cf V, 1). e.g:

$P = (\langle S_1 : p_1 \rangle, \langle S_2 : p_2 \rangle, \dots, \langle S_n : p_n \rangle)$  defines the predicate satisfied by all objects that can be constructed from objects  $x_1, x_2, \dots, x_n$  satisfying the predicates  $p_1, p_2, \dots, p_n$  by the construction operator  $\text{pro}(\langle S_1 : x_1 \rangle, \langle S_2 : x_2 \rangle, \dots, \langle S_n : x_n \rangle)$  where  $\text{pro}(\langle S_1 : o_1 \rangle, \langle S_2 : o_2 \rangle, \dots, \langle S_n : o_n \rangle)$   
=  $\text{pr}(\text{pr}(\dots(\text{pr}(\langle S_1 : o_1 \rangle, S_2, o_2), \dots), S_n, o_n))$

(which constructs tree



## EXPRESSION EVALUATION

is-state =  $(\langle s\_c : \text{is-control} \rangle, \langle s\_env : \text{is-env} \rangle, \langle s\_o : \text{is-output} \rangle)$

The control tree has instructions at its nodes of the form

$$\text{instruction name } (\bar{x}) = \begin{array}{c} \text{pred. } \text{action} \\ p_1 \rightarrow a_1 \\ \vdots \\ p_n \rightarrow a_n \end{array} \{ \text{OR} = a$$

Instructions can be:

- macro (self-replacing) : syntactic transformations replacing the instruction with a subtree of instructions (expands tree, modifies only control)

- value-returning (assignments) : semantic transformations +  
delete instruction, pass value to predecessor vertices  
(reduces tree)

(Only) any terminal vertex can be executed at any time

Example eval - expr ( $t$ ) = print ( $a$ ),  
 $a$ : value ( $t$ ) (order ok as a first  
line is undefined until  
 $a$  has been valid)

value ( $t$ ) = is-binary ( $t$ )  $\rightarrow$  apply ( $a, b, s\text{-rator} (t)$ )  
 $a$ : value ( $s\text{-rand}_1 (t)$ )

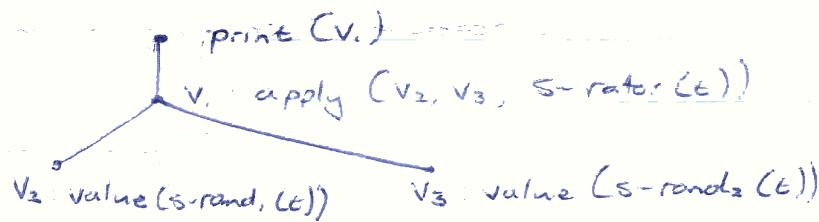
$b$ : value ( $s\text{-rand}_2 (t)$ ),

is-var ( $t$ )  $\rightarrow$  PASS  $\leftarrow t, s\text{-env} (\xi)$

is-const ( $t$ )  $\rightarrow$  PASS  $\leftarrow t$

apply ( $val_1, val_2, op$ ) =  $op = '+' \rightarrow$  PASS  $\leftarrow val_1 + val_2$   
 $op = '*' \rightarrow$  PASS  $\leftarrow val_1 * val_2$

### Control Tree



$V_1, V_2$  and  $V_3$  are used as  $a, b$  and  $t$  are bound variables  
and hence will be different actual variables at different  
parts of the tree.

## BLOCK-STRUCTURE ENVIRONMENTS

We add denotation (store) and attribute tree to the state, giving:

ids  $\rightarrow$  cells      cells  $\rightarrow$  values  
 is-state = ( $\langle s\text{-c} : \text{is-control} \rangle$ ,  $\langle s\text{-env} : \text{is-env} \rangle$ ,  $\langle s\text{-den} : \text{is-den} \rangle$ ,  
 $\langle s\text{-at} : \text{is-attribute} \rangle$ ,  $\langle s\text{-d} : \text{is-dump} \rangle$ ,  $\langle s\text{-n} : \text{is-integer} \rangle$ )

cells  $\rightarrow$  types

such that

is-environment = ( $\{\langle id : \text{is-number} \rangle \parallel \text{is-id}(id)\}$ )

is-attribute = ( $\{\langle n : \text{is-type} \rangle \parallel \text{is-n}(n)\}$ )

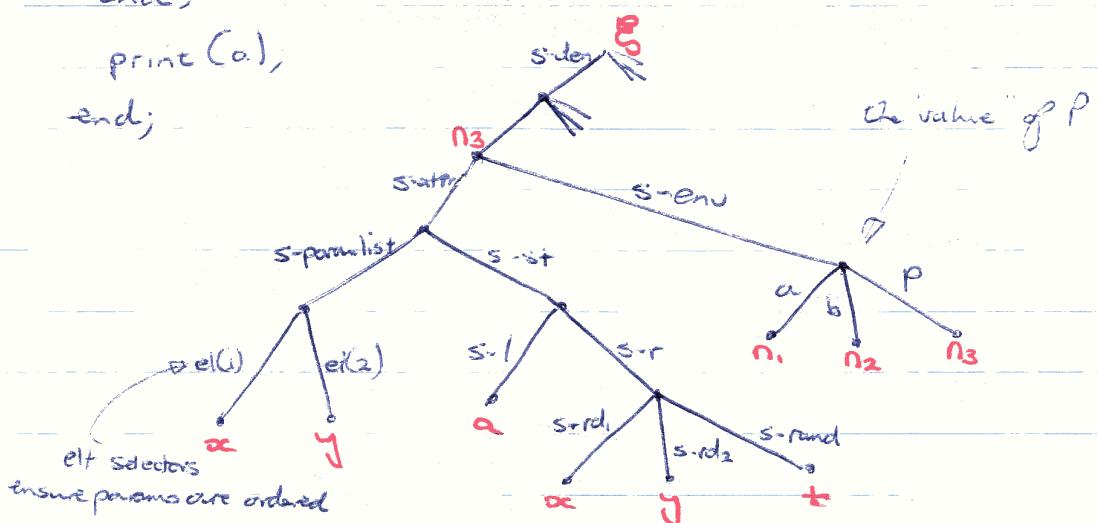
is-denotation = ( $\{\langle n : \text{is-value} \vee \text{is-procden} \vee \text{is-funcden} \rangle \parallel \text{is-n}(n)\}$ )

{ is-procden = ( $\langle s\text{-attr} : \text{is-procattr} \rangle$ ,  $\langle s\text{-env} : \text{is-env} \rangle$ )  
 is-procattr = ( $\langle s\text{-paramlist} : \text{is-id-list} \rangle$ ,  $\langle s\text{-st} : \text{is-st} \rangle$ )

→ (the value of a procedure is its carried environment, formal parameters, and body statements)

The initial state is ( $\text{interpret}(t)$ ,  $\langle \rangle$ ,  $\langle \rangle$ ,  $\langle \rangle$ ,  $\langle \rangle$ , 1)

Example 1: begin integer  $a, b$ ;  
 procedure  $P(x, y)$ ;  $a = x + y$ ;  
 begin integer  $a$ ;  
 $a := 1$ ,  $P(a, a)$ , print( $a$ );  
 end;  
 print( $a$ );  
 end;



$$? \quad \text{is-d} = (\langle s\text{-env}: \text{is-env} \rangle \langle s\text{-c}: \text{is-control} \rangle, \langle s\text{-d}: \text{is-dump} \rangle) \vee \text{is-L}$$

Ex2

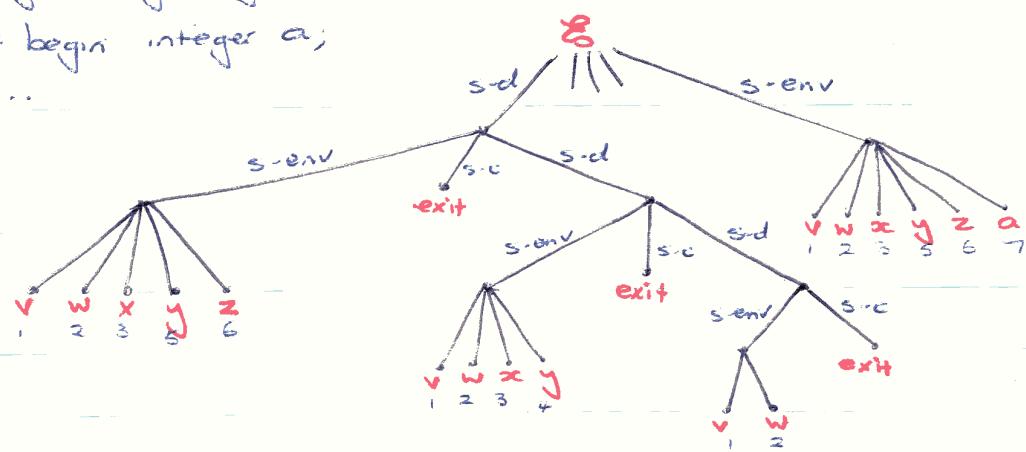
B1: begin integer v, w;

B2: begin integer x, y;

B3: begin integer y, z;

B4: begin integer a;

...



## SYNTHESIS (BRIEF NOTES)

### 8. CODE GENERATION

#### INTERMEDIATE FORMS OF TRANSLATION

- Size and alignment of elementary storage objects  
 ↗ placement of numbers used to build composite objects
- Composition of allocatable storage blocks  
 ↗ heap objects, activation records at compile time  
 (e.g. packed structures smaller than unpacked, etc.)
- Register allocation 1 : bookkeeping for blocks, keeping some values in fast registers.
- Register allocation 2 : temporaries for expressions
- Code reduction
  - stores - reducing need for temporaries (eval order)
  - operators - try to simplify expressions
  - load - check if something which must be loaded is already loaded

We now have a computational graph (see page ①)

- REGISTER-ASSIGNMENT ALGORITHM (simulation of eval order)  
 Registers have the following attributes : free / copy counts / unique copy (i.e. if we need this register we must save its contents first) / locked (in the process of placing a value) A table of registers is kept.

- INSTRUCTION SEQUENCE SELECTION (decision table - leaf conditions and root tables make action sequences)

This produces a target tree (cf assembler)

- Internal address resolution - access paths (i.e. change labels to addresses)
- External address resolution (link-load)

## 9. ERROR HANDLING.

### 9.1 WHAT THE COMPILER SHOULD TELL THE USER.

- negative feedback (errors) for debugging
  - positive feedback (listings) for maintenance
  - reporting should be as comprehensive as possible without interfering with the listings excessively
    - source & target language, compiler name & version, date & time of compilation, options on & off
    - text, co-ordinate system (line, statement, location #, nesting), cost measures (statistics)
    - pretty-printing, target code options
    - cross-reference & symbol-table printing options.
- {   
 true feedback }

### 9.2 ERROR REACTIONS.

- The compiler reacts with: when

Invalid translation

no redundancy detected

Crash or Loop

error sensitivity occurs

Quit

error detection & report used.

Recover & continue reporting

error recovery used

Repair & continue compiling

error repair used

NB Some errors transform valid programs into other valid programs, and only manifest themselves dynamically. LL(1) and LR parsing detects syntax errors at the earliest point possible before accepting the first symbol to offend against validity.

We can change the stack (unpleasant) or change the control, either by inserting tokens or skipping tokens. Insertion is difficult. The most widely used method is:

PANIC MODE: discard control until something "solid" is found (eg statement delimiter); then discard stack until something that can validly precede (succeed) the solid token is found. We build up a pack of solid symbols - the price a correct program must pay for error recovery.

## Alternatives

- error isolation (smallest reducible phrase containing error)
  - error productions  $A \rightarrow \underline{\text{error}}$  ad hoc additions to grammar  
to deal with "most common" errors
  - combination of these

Wants & Yields: Find a continuation, discard to an anchor, and replace with something consistent

$$wt X \in (T^* \setminus L) \rightarrow wt' x' \in L$$

A horizontal line is divided into four segments by vertical arrows pointing upwards. The first segment is labeled "stack". The second segment is labeled "error token". The third segment is labeled "control". The fourth segment is labeled "valid programs". Below the line, the word "all token combinations" is written.

Find  $w_{\mu} \in L$ ; construct the set of anchors  
 $\{ \det : wv\alpha \dots \in L \wedge v \dots = \mu \}$

Find the shortest  $n \in T^*$  such that  $\varepsilon X = n \alpha \mu$ ;

Replace  $\eta$  with the shortest  $v \in T^*$  such that  $wv\alpha \dots \in L$ .

NB: Errors are arbitrary & irrational - we cannot build a rational system to deal with them systematically.

## 9.4 SOME REMARKS ON ERROR TYPES.

Level	Type	Description
-	Notes	e.g. reporting of non-standard constructs
-	Comments	programming style
-	Warnings	possible error
-	Errors	reparable
-	Fatal errors	recoverable
-	Deadly errors	Goodbye!

There should be an option to set the report level, with a cut-off on density (avalanche)

ANALYSIS - lexical  
 'syntactic'  
 'semantic' ↓      NB Recovery based upon redundancy  
 which increases in direction of  
 arrow (so easiest to deal with  
 semantic errors)

SYNTHESIS - detects only "compiler errors", including violation  
 of table limits, etc.

RUNTIME - The compiler can prepare for those by some  
 backmapping from store to source data declarations  
 and / or some backmapping from target code to source  
 text.