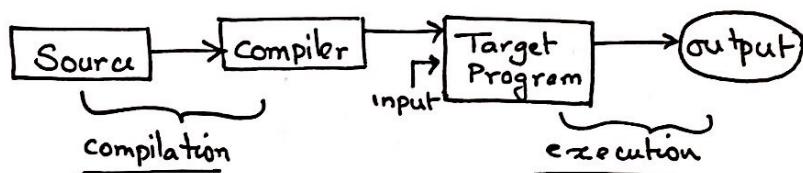


Compiler Design

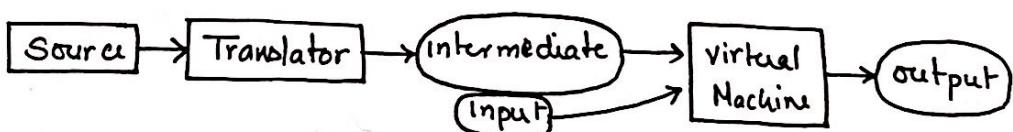
- For a program to be run, it must first be translated into a form in which it can be executed by the machine. This translation is done by a system software called the compiler.
- The compiler reads the program in one language and translates it to another. From source to target language.
- The compiler also reports errors it encounters during the translation.



- Apart from compilers, there are interpreters which instead of translating to a target program directly execute the code in the source itself.

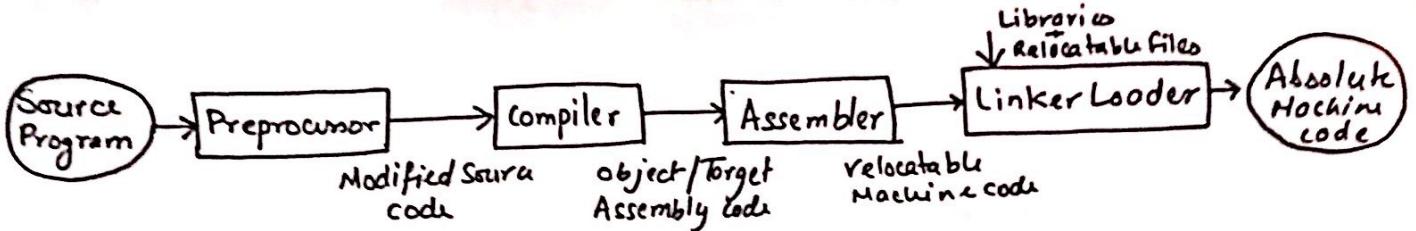


- Languages like Java use a hybrid system which first translates to an intermediate language and then this is executed by an interpreter.



Language Translation System

1. The source is sometimes first collected by the preprocessor. This resolves any preprocessor directives and Macros. Both the input and output are high level languages.
2. The compiler produces assembly code. This is easy to produce and debug. The output is called the object / target assembly code.
3. The assembler produces the relocatable machine code. This is in binary.
4. The linker handles linking of different components and libraries as well as external.
5. The loader loads all the files into memory for execution.



Structure of Compiler

The Compiler has two Functional Parts.

— Analyzer

- > This breaks up the input source into its constituent sequence pieces and imposes a grammatical structure on them. This is used to create an intermediate representation of the source.
- > The analyzer also alerts the user with errors (semantic/syntactic) through informative error messages.
- > This also collects and stores information about the source in a symbol table.

— Synthesizer

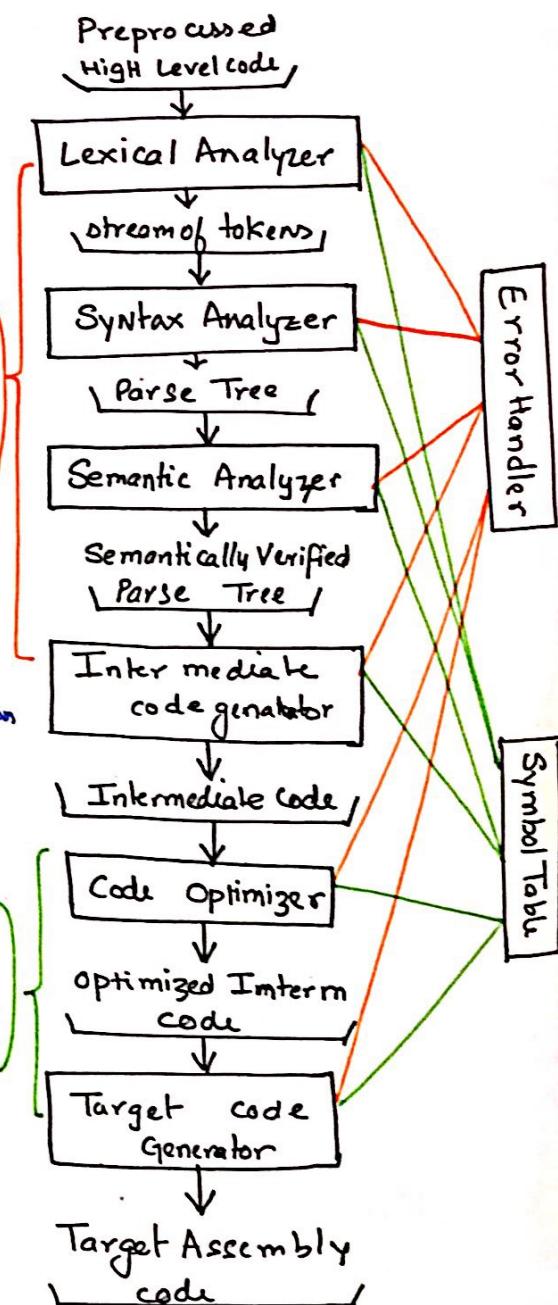
- > This constructs the desired target Assembly code from the intermediate code.
- > This may also include some form of machine dependant/independent optimizers. These are for improving code quality.
- > This uses the shared Error Handler and symbol tables.

The Error Handler and symbol table is shared by all phases.

Phases of Compiler

— Lexical Analyzer

- > Reads the stream of characters making up the source program and groups them into meaningful sequences called lexemes.
- > The output of the L.A is a sequence of tokens. $\langle \text{token-name}, \text{attr} \rangle$



Syntax Analyzer

- > This is considered to be the main part of the compiler. It builds the parse tree using a set of grammatical rules particular to that language.
- > The output of the Syn.A is the parse tree.

Semantic Analyzer

- > This uses the parse tree and the symbol table to check for semantic consistency with language definition.
- > The type information is either stored in the tree or the table.
- > It also performs type checking to check for matching operands and if needed do type conversion i.e coercion.
- > The output is the semantically verified parse tree.

Intermediate Code Generation

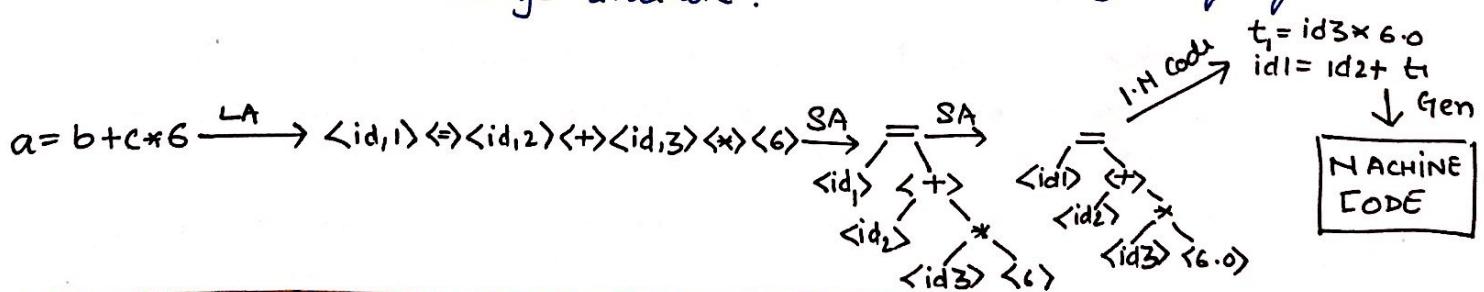
- > This is an explicit low-level machine like code generator which uses the output of the Sem.A and Symbol table.
- > This code must be easy to produce & translate to machine.
- > They can be of the following type : 1) 3 address code, 2 addr.

Code Optimization

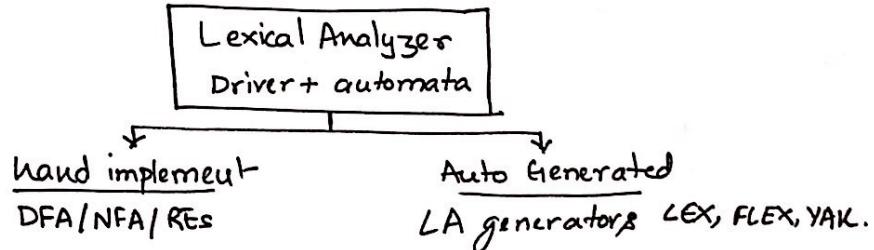
- > This is done to improve the running time of the program without slowing down the compilation much.
- > This removes redundant / duplicate instructions, type conversion etc.
- > Compiler that do this are called optimizing compilers.

Code Generation

- > Maps the intermediate [optimized] code to machine language.
- > The mapping can be to ~~any~~ other non machine codes too.
- > In case of machine code, memory locations, registers etc are selected and defined.
- > The most essential part is to judiciously assign registers and storage allocation.

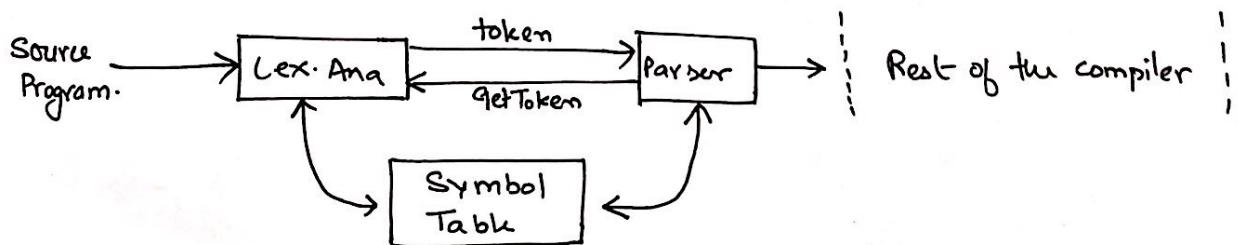


Lexical Analyzer (Scanner)



Role of the L.A

- > Read the input program, group the characters into lexemes and produce a sequence of tokens of each lexem.
When a new lexeme containing an identifier is encountered, an entry must be made into the symboltable. [L A]
 - > Stripping of comments and white spaces. [scanning]
 - > Correlating errors with the source program. Some times a copy of the source is made with the errors messages inserted at the appropriate places.
- # The LA and Syn.A often interact with each other. This is implemented by the parser calling the LA with the getToken command. The LA then reads in the next chars, matches the lexeme & produces a token which it gives to the parser.

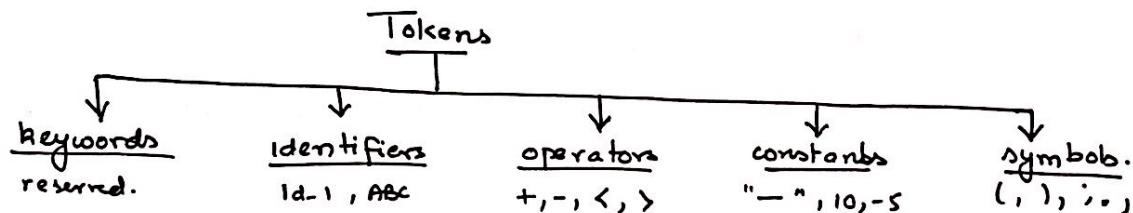


LA as an independant component

1. Simplifies overall design. The LA deals with token structures and scanning while the parser deals with syntactic structure. A parser that does both will be very complex to create & maintain.
2. Improves Compiler Efficiency. Reading is a very time consuming part. Having an independant L.A allows us to apply specialized tasks and buffering strategies.
3. Due to the inherent advantage of modularity, makes changes easier.

4. Improves the portability of the system.
5. Allows for the consolidation of device specific issues to the LA.

- Token : A token is a pair consisting of a token name and an optional attribute value. The name is abstract & reflects the type of lexical unit.
- Pattern : A pattern is the description of the form that the lexeme of a token may take. These are used to match the lexeme with the type
- Lexeme : A sequence of characters in the source program that matches the pattern for a token and is identified by the lex. A as an instance of that token.



Attributes of Tokens

- > When more than one lexeme matches a pattern the LA provides additional info to the rest of the compiler.
 - > The LA provides the token name and attribute, the token name influences the parser while the attribute influences the translation.
 - > This information i.e the lexeme, type, location [for error reporting] is maintained in the symbol table.
- # The handling to which spaces is a serious implementational challenge sometimes a keyword can be mistook as an id or something else e.g.

eg> DO 5 I =1.25
 can be DO5I as an Id
 or DO 5 I as a do while

Token Representation

< token, attributes > eg> < add-op >, < id, pointer to id, >
 < literal, string value "My appriciation" >
 < Keyword, if value keyword >

Lexical Errors

- > The LA on its own cannot tell errors in the code & thus requires some aid in this task.
- > Errors arise when the lexical analyzer is unable to proceed because nof the pattern of the token matches any of the remaining input.
- > The simplest strategy (panic mode) is to delete successive characters from the remaining input until the LA can find a well formed token.
- > other techniques include
 1. delete one char.
 2. insert any missing char.
 3. replace a char.
 4. Transpose two adjacent.
- > The simplest transform strategy is to check if the prefix of the input can be transformed to a valid lexeme in one transform.
- > The general approach is to find the minimum transforms.

Input Buffering

- > To be sure that we have the right lexeme, we often have to look one or more characters beyond the next lexeme.



- Single Buffer

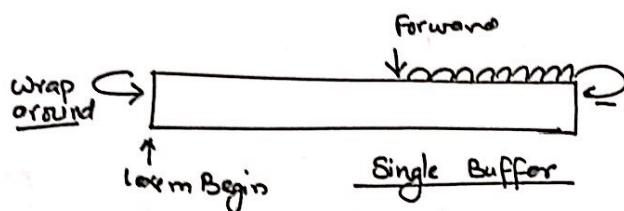
- > The Buffer is usually of the size of one disk block i.e 4096 b.
- > On one system call, N chars are read into the buffer at once and if there's < N chars a special eof is added.

- > Two pointers i.e

lexemeBegin : marks the start of the current lexeme.

forward : scans ahead till the match is found.

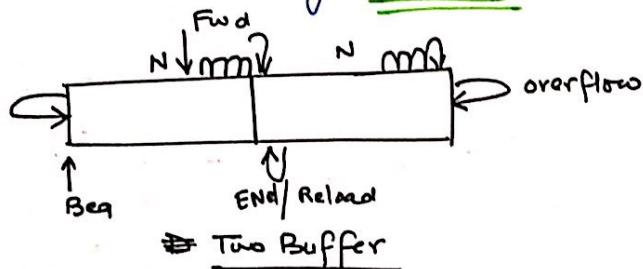
- > Once the lexeme is determined forward is set to its right. Then it is recorded and reflected in the ST. Then the LexemeBegin is set to the immediate next char of forward.
- > One major drawback of single Buffer is if the lexeme is longer than the block, wrap around occur.



Two Buffer

- > This involves two buffers that are alternately ~~reloaded~~ loaded. Each Block is N i.e one char block.
- > When the forward is advanced, we check if we haven't reached the buffer end. If so the other buffer must be reloaded. The fwd is moved to the start of it.
- > This suffers from the same wrap around drawback.
- > Additionally there's two checks each time forward is moved
 - to check buffer end
 - to check char read.

This can be resolved using sentinels.

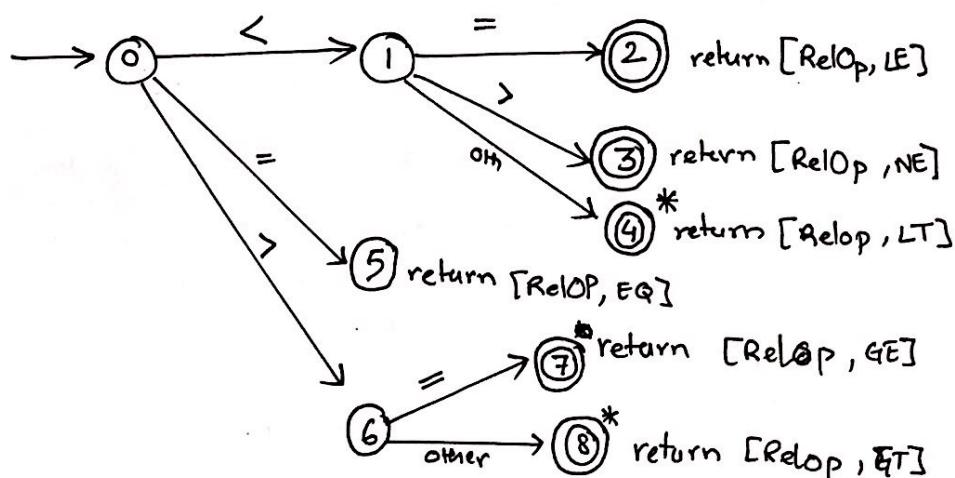


- # For very long sequences like strings, they are broken up and treated as combinations.

Token Recognition

- > There is a set of patterns (defined by regular exp) for each type of token which are used to match the lexeme.
- > These patterns are converted to transition diagrams (DFA/NFA)
The start state can be analogous to Lexeme Begin and each state transition can be analogous to moving the forward.
- > The basic conventions for such transition diagrams are:
 - 1> Some states are marked as accepting / final. If there is an action to be taken, it is attached to it.
 - 2> If it is needed to retract the fwd pointer i.e the lexeme does not include the transition symbol, we place a *.
 - 3> The TD always starts with the start state.
- > We start with one TD and keep moving to the next till we can find a match. If no match is found error recovery is invoked
- > Frequently used tokens have their TD higher in the matching scheme.

Eg > Relational Op



- Keyword and Identifier Recognition

- > The task of Identifier and keyword Recognition is problematic because in terms of pattern they're indistinguishable.
- > There are two possible ways to solve this.

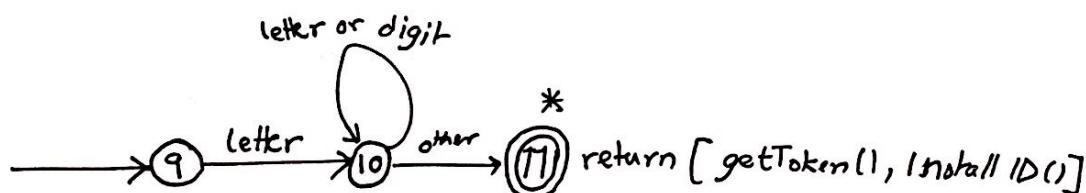
1. Use the symbol table intelligently

Install all key words into the symbol table initially. When an ID is found do a call to InstallID() to place it in the table if its not already there.

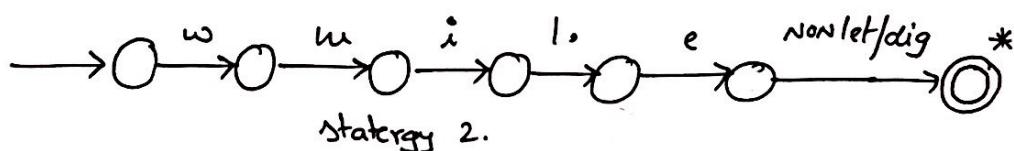
GetToken() examines the symbol table entry for the lexeme and returns its type i.e. KW / ID.

2. Make a separate TD for keywords

Build a separate TD for Each keyword. This TD has higher precedence over ID when the lexeme matches both ID & KW.



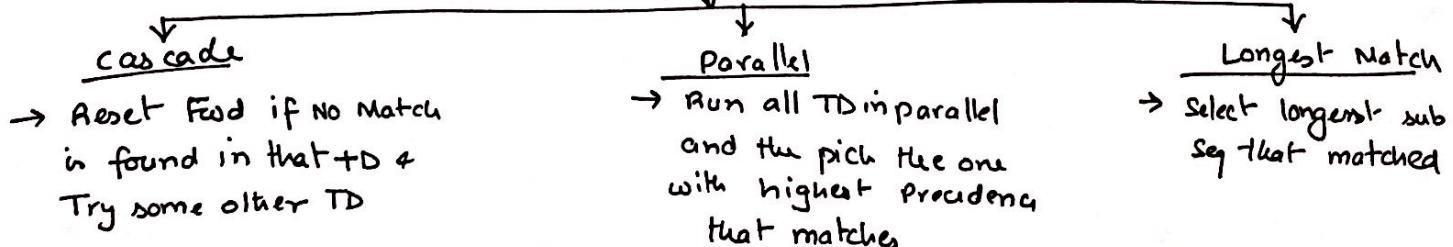
Strategy 1.



Strategy 2.

- > The main benefit of using 2. over 1. is it allows to easily add or modify the keywords instead of having to redefine the symbol table.

TD Based Lex Analyzer

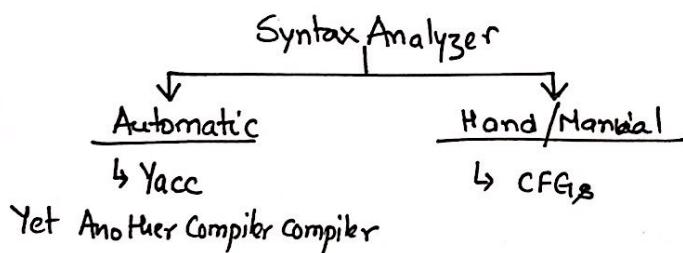


Syntax Analysis [Parsing]

- > The syntax analyzer takes in a sequence of tokens and generates a parse tree.
- > The syntactical constructs of a language are specified using Context Free Grammars.

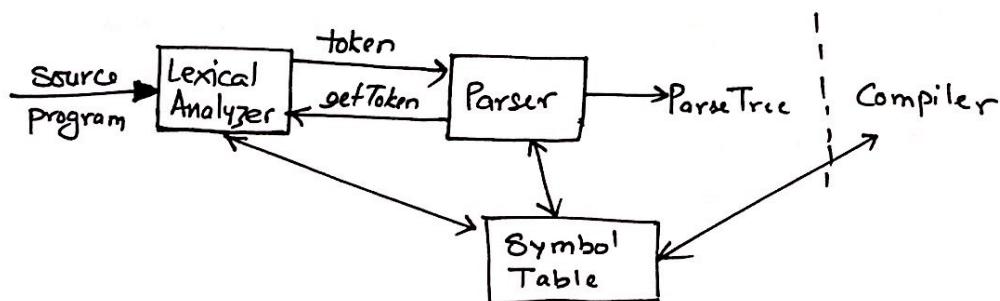
Grammar Based Parsers

1. The grammar provides a specific & easy-to-understand specification.
2. For some grammars, we can automatically generate parsers. This also allows us to uncover missed ambiguities of the Language.
3. The structured method of using grammar helps ease the translating task.
4. Grammars allow the language to evolve & accommodate new constructs.

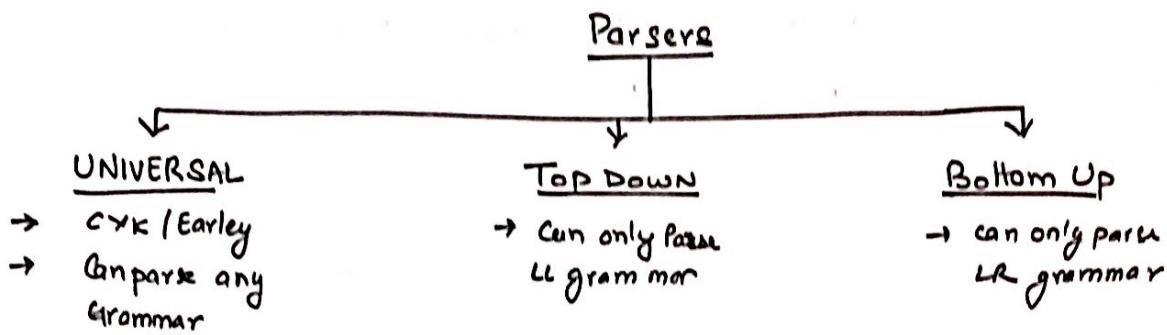


The Role of Parser

- > Obtain a string of tokens from the L.A and verify if it can be generated by the Grammar by generating the parse tree.
- > To report any errors in an intelligible manner.
- > Recover from commonly occurring errors and resumer its task.
- > In essence, this is the most important part of the compiler.



Types of Parsers



UNIVERSAL PARSERS

- These use methods like those of the Cocke-Younger-Kasami algorithm or the Earley's algorithm and can parse any type of grammar.
- These are proven to be highly inefficient in production use and are thus rarely used.

Top Down Parsers

- Just as the name suggests, the parser work from the root to the leaves scanning the input left to Right, one symbol at a time
- These only work (efficiently) with LL grammars i.e L-R read | Left Most derivation.

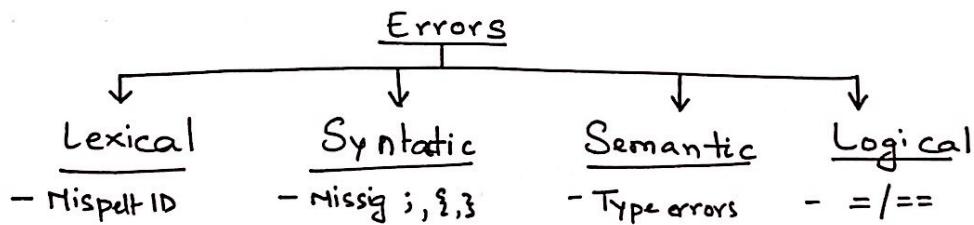
Bottom-Up Parsers

- Similar to Top-Down parser but start with the leaf and work to the top.
- These only work efficiently on LR grammar i.e L-R read / Right Most derivation.

Larger LR grammar are generated automatically in most case.

Apart from this parsers also collect info on tokens, do type checking and other analysis in syntax & semantics.

Syntax Error Handling -



Error Types

— Lexical

- These include mispelt ID, KW or operators as well as missing quotes around strings.

— Syntactic

- These include misplaced, missing, extra ";" or braces.
else without if, case without switch etc.

— Semantic

- These include type mismatch errors, return type matching.

— Logical

- These can be incorrect reasoning by the coder, using assignment instead of comparison, equal op.

> Most LL and LR Parsers detect errors as soon as they occur i.e they exhibit the variable prefix property.

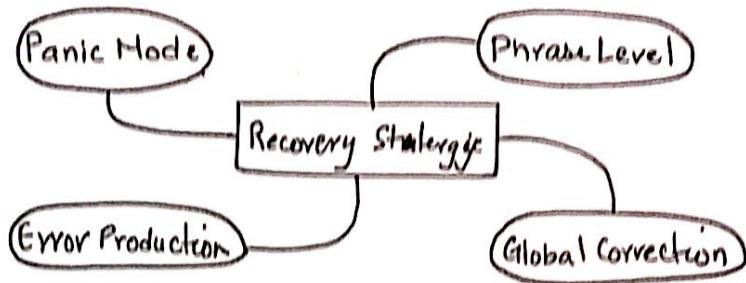
"Report an error if the prefix of the input cannot be completed to form a string in that language."

Role of the Error Handler

1. Report the presence of errors clearly and accurately
2. Recover from error quickly to detect further errors.
3. Add a minimal overhead to the parsing phase.

> The most common error reporting method is to show the offending line with a pointer to the error.

Error Recovery Strategies



- > The simplest method is to quit and display an informative message. If the compiler is rerun, additional errors may pop-up and there may be a case when an avalanche of spurious errors occur in this case the compiler must prevent this.

→ Panic Mode Recovery

- On finding an error, discard the error one at a time until a set of designated synchronizing token is found (';', ';;', '?'). The role of these symbols is always clear and ambiguous.
- This has a tendency to skip a considerable amount of errors. This is simple and is guaranteed never to go into an oo loop.

→ Phrase Level Recovery

- Perform local correction on the ~~rest~~ remaining input i.e replace/ remove / add a ';', ';;'. This choice is left to the designer.
- This can correct any input string but misses errors that occurred before the point of detection.
- Can cause error avalanche.

☰ Error Productions

- Anticipate common errors and add productions that generate those errors.
- Whenever one of these error production is used the parser can generate the appropriate ~~the~~ error diagnostic.
- This is the most commonly used strategy.

四 Global Correction

- Select the minimal no of change to be done to the input for it to be error free.
- These are very costly methods and can cause unintended behavior.
- Input x is changed to least cost y using G.

Context Free Grammar (Review)

- Terminals: They're the basic symbols which make up the strings
They're shown using small letters or in bold.
- NonTerminals: They're the syntactic variables that denote strings
They're shown using capital letters.
- Start State: This is a special NonTerminal and is the starting point of any string.
- IV Productions: They specify the manner in which NT and T are combined

$$S \rightarrow \alpha \quad [\alpha \text{ can be NT, T, or a mix}]$$

$$\text{i)} \quad S^* \rightarrow \alpha$$

S in 0 or more derivations derives the sentential α .

$$\text{ii)} \quad S^+ \rightarrow \beta$$

S in one or more derivation derives the sentential β .

Parse Trees as Derivation

- > A parse tree is the graphical representation of a derivation detailing out the order in which the productions are applied.
- > Each internal node is the application of a Production
- > Each leaf node is a terminal. The tree is read left to right.
- > The main challenge in building the parse tree is selecting the production to expand first and then which production.

Ambiguity in Grammars

"A grammar that produces more than one parse tree for a sentence" is called ambiguous

OR

"An ambiguous grammar is one which produces more than one LMD or RND for the same sentence"

> It is desirable for the grammar to be unambiguous as the compiler has no method for determining which parse tree to select. The only exception to this is the operator precedence parser.

> An ~~unambiguous~~ grammar can be unambiguous using disambiguating rules which discard wrong P.Trees.

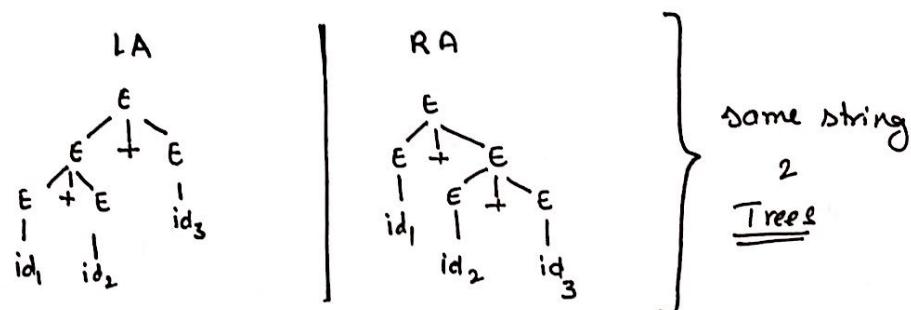
$$\text{AG} \xrightarrow{\text{disambi Rule}} \text{UG}$$

Why does Ambiguity Exist?

→ Ambiguity in associativity -

- > operators can be left/ right associative. The parse tree for both are very different but semantically same.
- > In most languages $+, -, *, \div$ are left associative and \wedge and $=$ are right associative.

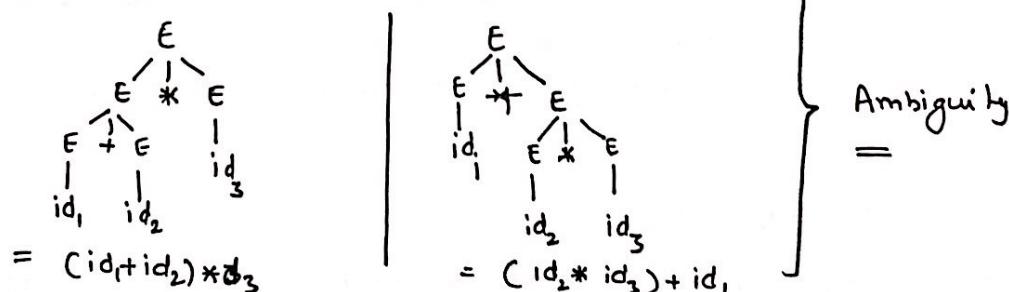
eg $id_1 + id_2 + id_3$



⇒ Ambiguity in precedence

- > operators have a scale of precedence i.e if two are to simultaneously occur which should be evaluated first.
- > Ideally the operators with least precedence must be evaluated first.

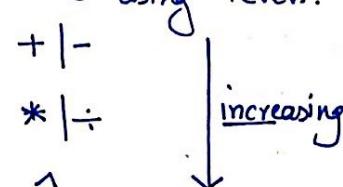
eg $id_1 + id_2 * id_3$



Associativity is handled using recursion. . .

- > Left Asso \Rightarrow LR
- > Right Asso \Rightarrow RR

Precedence is handled using levels.



Writing Grammar for a Parser

- > CFGs are capable of describing most of the syntax but cannot handle task like declaration testing. This is done by consequent Phases.

Why Have a Grammar Based Parser ?

- > As a grammar can also be used to ~~not~~ describe regular expressions. why is there a need of separation the lexical & syntactic unit into RE and CFG based units?
 - Modularizes the Front End.
 - ⇒ Lexical Rules are very simple and using complex grammars to describe them is wasteful.
 - ☰ RE are more consistent and easy-to-understand notation for tokens.
 - ☒ RE based LRAs can be automatically constructed.

"^{IF} Ideal Grammar for Parser is Unambiguous, deterministic and be free " of states which might cause ~~as~~ loops
(esp for top-down parsing)

1 Eliminating Ambiguity

- > An ambiguous grammar can be made unambiguous using a set of disambiguating rules.
- > For operators associativity, it can be resolved using Recursion.
(Top Down parsers go into as loops with LR ∵ they're items changed to RR)
- > For operator precedence, it can be resolved by using levels with low precedence at the top.
- > For conditionals "Match else with closest unmatched then"
- > These unambiguating rules can be built directly into the grammar but are not in practice.

e.g. $E \rightarrow E+E \mid E*E \mid id$

Ambiguous

$$\left. \begin{array}{l} E \rightarrow TE' \\ E' \rightarrow +TE' \mid \epsilon \\ T \rightarrow FT' \\ T' \rightarrow *FT' \mid \epsilon \\ F \rightarrow id \end{array} \right\}$$

unambiguous.

2. Eliminating Non Determinism

- > If a grammar shows the common prefix problem then it suffers from non determinism.

$$A \rightarrow \alpha\beta_1 / \alpha\beta_2 / \alpha\beta_3 \dots / \alpha\beta_n$$

- > In the above production, there is no way of telling which production is to be selected without having to look ahead.
> The common prefix problem is solved by Left Factoring.

$$\begin{aligned} A &\rightarrow \alpha A' \\ A' &\rightarrow \beta_1 / \beta_2 \end{aligned}$$

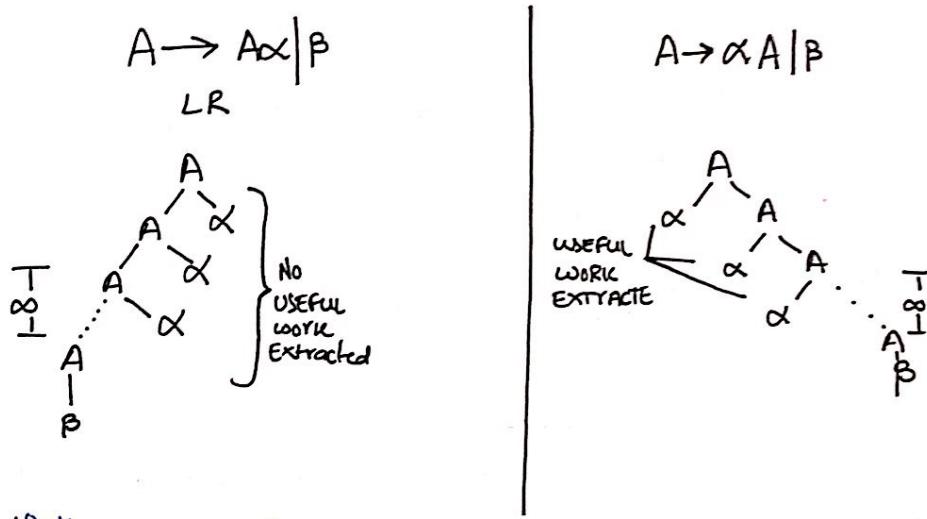
Algorithm : Left Factor (G)

1. For each NT in grammar say A, find the longest common prefix common to 2 or more productions α
2. if $\alpha \neq \epsilon$
Replace $A \rightarrow \alpha\beta_1 / \alpha\beta_2 / \dots / \alpha\beta_n / r$
with
 $A \rightarrow \alpha A' / r$
 $A' \rightarrow \beta_1 / \beta_2 / \dots / \beta_n$
3. Repeat till no NT exist showing the Common Prefix Problem.
4. Return Left factored G.

- > Non-Determinism and ambiguity are very different and solving for one doesn't solve for the other.
- > With ambiguity we know the solution but here we don't.
- > This is only a problem for parsers that cannot back track but ambiguity is a problem for all.

3. Eliminating Left Recursion

- > A grammar shows Left Recursion if for some A in G there is a production of the type $A \xrightarrow{+} A\alpha$.
- > Top down parsers enter an ∞ loop with such grammar.



- > Even if the parser enters an ∞ loop in RR there is still useful work being done thus there'll arise a case when it breaks.
- > This Recursion can be immediate as in

$$A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_n$$

or be non immediate as in

$$\begin{aligned} A &\rightarrow B | C \\ C &\rightarrow A\alpha | id \end{aligned}$$

- > Immediate and non immediate recursion are handle by the following algorithm

Algorithm : Eliminate LR (G)

1. Arrange NT as A_1, A_2, \dots, A_n .
 2. for each i from 1 to n {
 3. for each j from 1 to $i-1$ {
 4. replace all $A_i \rightarrow A_j \delta$ by $A_i \rightarrow \delta_1 r | \delta_2 r | \dots | \delta_k r$
where $A_j \rightarrow \delta_1 | \delta_2 | \dots | \delta_k$ are current A_j Prods.
 5. }
 6. Eliminate all immediate LR among A_i productions.
 7. }
- $\leftarrow RR-G$

> The above algorithm has two main requirements

1. No cycles . $[A \xrightarrow{*} A]$
2. No ϵ productions . $[A \rightarrow \epsilon]$

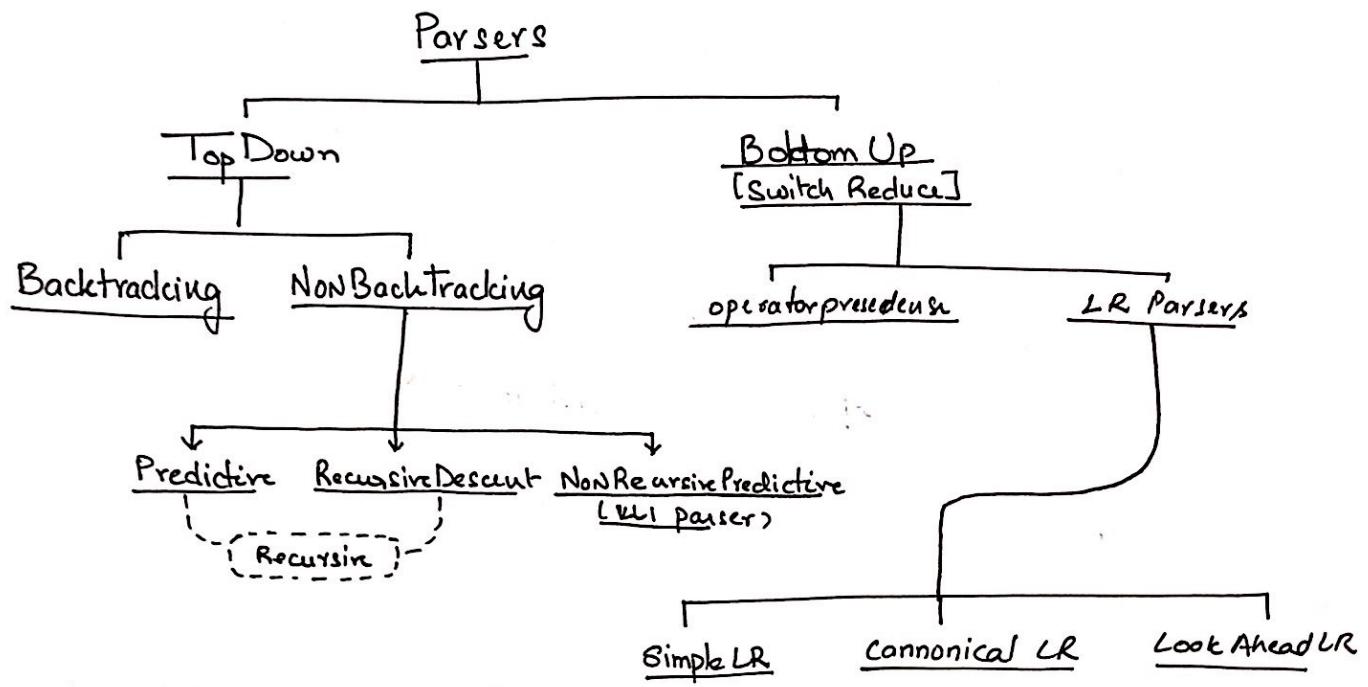
> Immediate recursion is removed as follows

$$\begin{aligned} & A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_n | \beta_1 | \beta_2 | \dots | \beta_m \\ = & A \rightarrow \beta_1 A' | \beta_2 A' | \dots | \beta_m A' \\ & A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_n A' | \epsilon \end{aligned}$$

> There can be ϵ in the grammar if and only if it doesn't induce the recursion.

$$\begin{aligned} & A \rightarrow A\alpha | \beta \\ = & A \rightarrow \beta A' \\ & A' \rightarrow \alpha A' | \epsilon \end{aligned}$$

Classification of Parsers



Top Down

$S \Rightarrow aABe \Rightarrow aAde \xrightarrow{*} abbcde$

Bottom Up

$aabbcd \Rightarrow aA bcd \Rightarrow aAd \xrightarrow{*} S$

Top down Parsers

- > This can be defined as the task of constructing a parse tree for the input starting from the root and creating the nodes in pre order.
- > It is basically finding the leftmost derivation.
- > A predictive parser can derive the correct A production by looking ahead a fixed number of symbols.
- > Top down, predictive parsers only work on unambiguous, deterministic, right recursive grammars i.e $LL(k)$ ~~where k is the look ahead~~.

— Recursive Descent Parsing

- > This contains a set of procedures, one for each NT.
- > The execution begins with the start symbol & halts when the entire string is scanned.
- > General algorithms may require backtracking but normally they're not used.
- > The main point when dealing with recursive descent parsing is that when using LR grammar even with backtracking they can go into an ∞ loop.

First and Follow

- > First and follow allow us to determine which production to apply based on the next input symbol.
- > First and follow are very useful when we're building the parsing table.
- > In panic mode recovery, the tokens defined by follow can be the synchronizing tokens.

— First(α)

- > $\text{First}(\alpha)$ where α is a string of grammar symbols is the set of terminals that begin strings derived from α and if $\alpha \Rightarrow^* \epsilon$ then ϵ is in $\text{first}(\alpha)$.

Compute $\text{First}(\alpha)$:

1. if α is a terminal then $\text{First}(\alpha) = \{\alpha\}$
2. if $\alpha \Rightarrow \epsilon$ is a production add ϵ to $\text{first}(\alpha)$.
3. if X is a NT & $X \Rightarrow Y_1 Y_2 \dots Y_k$ for $k \geq 1$ then place α in $\text{first}(X)$ if α is in $\text{First}(Y_i)$ and ϵ is small of $\text{First}(Y_1), \dots, \text{First}(Y_k)$ i.e. $Y_1 \dots Y_{i-1} \xrightarrow{*} \epsilon$.
if ϵ is in $\text{First}(Y_j) \forall j=1, 2, \dots, k$ add ϵ to $\text{first}(X)$

- > The First of a string like $x_1 x_2 \dots x_n$ is $\text{first}(x_1)$ if $x_1 \Rightarrow^* \epsilon$ is not a production otherwise it'll also have $\text{first}(x_2)$ and so on.

Follow(A)

- > $\text{Follow}(A)$ for a NT A is the set of terminals that can appear immediately to the right of A in some sentential form. i.e a set of terminals a such that there exists $S \xrightarrow{*} \alpha A \beta$.
- > If A is the right most symbol in a sentence then \$ is in $\text{Follow}(A)$.

Compute $\text{Follow}(A)$:

1. Place \$ in $\text{follow}(S)$. S is the start symbol.
2. If there is a production $A \rightarrow \alpha B \beta$ then every item in $\text{first}(\beta)$ except ϵ is in $\text{follow}(B)$.
3. If there is a production $A \rightarrow \alpha B$ or $A \rightarrow \alpha B \beta$ and $\text{first}(\beta)$ contains ϵ then everything in $\text{Follow}(A)$ is in $\text{Follow}(B)$.

LL(1) Grammars

- > The first L stands for Left-to-Right Reading of input.
The second L stands for left most derivation.
The (1) means a look ahead of 1.
- > A grammar is LL(1) iff whenever $A \rightarrow \alpha / \beta$ are two distinct productions then
 - For no T 'a' do both α and β derive string starting with a
 - ⇒ at most one of α or β can derive empty.
 - ⇒ if $B \Rightarrow \epsilon$ then α does not derive any str beginning with a T in $\text{Follow}(A)$. & v.v.
- > Basically in LL(1) grammar with $A \rightarrow \alpha / \beta$ production
 - $\text{First}(\alpha)$ and $\text{First}(\beta)$ are disjoint .
 - ⇒ if $\epsilon \in \text{First}(\beta)$
 $\text{First}(\alpha) \cap \text{Follow}(\alpha) = \emptyset$.
 - ⇒ if $\epsilon \in \text{First}(\alpha)$
 $\text{First}(\beta) \cap \text{Follow}(\beta) = \emptyset$.

Building The Parsing Table

Construct_Parse_Table(G):

1. For each ~~terminal~~ ~~non-terminal~~ $A \rightarrow \alpha$ in G
2. For each terminal a in $\text{First}(\alpha)$ add $A \rightarrow \alpha$ to $M[A, a]$
3. if ϵ is in $\text{First}(\alpha)$, For each terminal b in $\text{Follow}(A)$ add $A \rightarrow \alpha$ to $M[A, b]$.
If ϵ is in $\text{First}(\alpha)$ and $\$$ in $\text{follow}(A)$ add $A \rightarrow \alpha$ to $M[A, \$]$.
4. Mark each undefined entry in M as error.

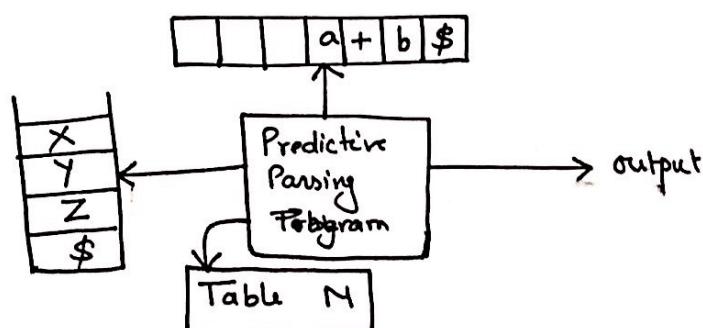
AMBIGUOUS GRAMMAR WILL HAVE MULTIPLE ENTRY IN $M[A, a]$.

Non Recursive Predictive Parser (LL(1) Parser)

- > This type of parser can be built by maintaining a stack explicitly.
- > They are driven by the parsing table M .
- > The parser is controlled by the symbol X on top of the stack
if X is a NT then $M[X, a]$ is used otherwise T is matched.

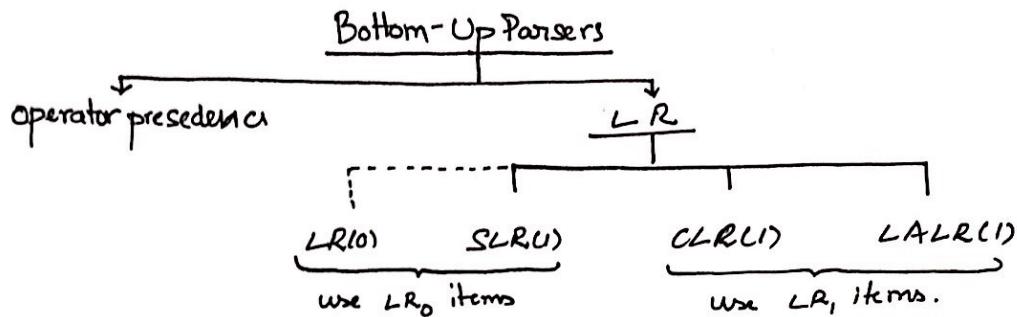
N-R-Pr-Parser(w, M):

1. let a be the starting (first) symbol of w .
2. let X be at the top of stack.
3. while ($X \neq \$$)
 4. if ($X = a$) Pop & let a be the next symbol in w .
 5. else if ($M[X, a]$ is an error) error()
 6. else if ($M[X, a] = X \rightarrow Y_1 Y_2 \dots Y_k$)
 7. output $X \rightarrow Y_1 Y_2 \dots Y_k$
 8. Pop
 9. Push $Y_k Y_{k-1} \dots Y_1$
 10. let X be the top of stack.



Bottom-Up Parsers

- Construction of a parse tree for an input string beginning at the leaves and working up towards the root.
- It follows a left-to-right scan, rightmost derivation in reverse.



Reductions

- Bottom-up parsing is a process of reducing a string w to the start symbol of a grammar.
 - At each reduction step a specific substring matching the body part of a production is replaced by the non-terminal at the head of the production.
 - A reduction is the reverse of a derivation step.
- # Bottom up thus constructs a derivation in reverse.

Handles

" A handle is a substring that matches the body of a production & whose reduction represents one step along the reverse RND. "

Formally if $S \xrightarrow{r_m}^* \alpha Aw \xrightarrow{r_m}^* \alpha Bw$ then $A \rightarrow B$ in the position following α is a handle.

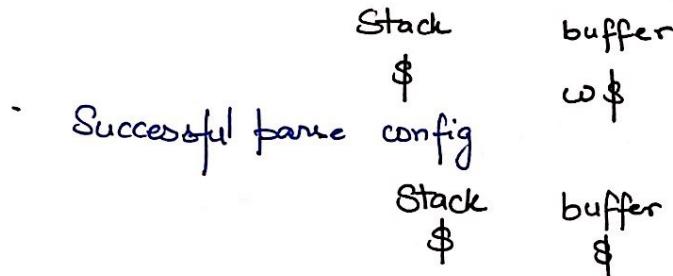
If the grammar is unambiguous then every right sentential form of the grammar has exactly one handle.

The Reverse RND can be obtained using handle pruning i.e start with a string of terminals w , if w is a sentence. Let $w = r_n$; r_n is the n^{th} right sentential form of some deriv.

To get the deriv locate handle B_n in r_n and reduce.

Shift-Reduce Parsing (LR(k) Grammars only.)

- A stack holds the grammar symbols
- The buffer holds the string.
- \$ marks both the bottom of stack and end of an input.
- Initial config.



- Scanning is done left to right and at each step the parser can put none or more symbols into the stack until it can reduce with β .
- This is repeated till there is an error or it is successfully parsed.

Operations

1. Shift : shift/push next input symbol onto the top of stack
2. Reduce : Right end of handle must be at top and then locate the left end and replace with LHS.
3. Accept : Announce successful parse
4. Error : Discover syntax error and call error handler.

"The handle will always eventually appear at the top of the stack, never inside."

The main issues in design are what to reduce with and when to reduce.

Conflicts occur when we know we have a handle but the stack and input are insufficient to produce a decision for what to reduce with.

LR(k) Parsing -

L means Left to right scanning of the input buffer.

R means Right most derivation done in reverse.

K means the look-ahead size for getting a parsing decision.

They are table driven and work on LR grammars

LR Grammar: A grammar for which an LR parser can recognize a handle in Right Sentential form.

Advantages:

- LR parsers can be made to recognize virtually all programming constructs for which a CFG can be made.
- LR parsing is the most general non backtracking Shift Reduce parsing.
- LR parsers can detect errors as soon as they occur.
- The LR grammar class is a proper superset of LL grammar class.

Drawback: Way too much work. use Yacc instead.

Items: An Item is a grammar production which has a dot at some place in it

Eg $A \rightarrow \alpha$ has $A \rightarrow \cdot \alpha$ and $A \rightarrow \alpha \cdot$ as items.

Items indicate how much of a production we have seen at any point.

Augmented Grammar

An augmented grammar G' is a grammar with a production as $S' \rightarrow S$ added to the original grammar G .

The use of it is to show the end of parsing i.e when $S' \rightarrow S$ is about to be reduced.

Closure

If I is an item of G then the closure of I is the item constructed from I by:

1> add every item of I to closure(I).

2> if $A \rightarrow \alpha \cdot B \beta \in \text{closure}(I)$ and $B \rightarrow \gamma$ exists then add $B \rightarrow \cdot \gamma$ to closure(I).

Closure means at some point we expect to see a string derived from it.

LR Parsing Algorithm

- An LR parser consists of an input, an output, a driver program and a parsing table.
- Only the table varies b/w LR parsers.
- The parse driver shifts between states while parsing.

Structure of LR Parsing Table

The table has two parts, A parse Action part and a Goto part.

Action takes in a state i and a ~~non-terminal~~ terminal a and gives one of the following outputs.

1. Shift : shift the input a to top of stack and push in j
2. Reduce $A \rightarrow \beta$: Reduces β and leaves A .
3. Accept : Accept the input and stop.
4. Error : discovery of error.

Goto is for a set of items and if $\text{GOTO}[I_i^0, A] = I_j$ map state i to NFA to a state j .

LR-Parsing(w , LR-Table) :

1. put S_0 onto stack
2. let a be first symbol of $w\$$;
3. while ($1 > 0$)
 4. let s be state on top of stack;
 5. if $\text{Action}[s, a] = \text{shift } t$
 6. push t onto stack;
 7. $a = \text{next_char}()$;
 8. else if ($\text{Action}[s, a] = \text{Reduce } A \rightarrow \beta$)
 9. pop $|\beta|$ symbols;
 10. let t be top of stack;
 11. push $\text{GOTO}(t, A)$ onto stack;
 12. else if ($\text{Action}[s, a] = \text{accept}$) break;
 13. else call errorhandler

Construction of SLR Parse Table

- This begins with LR(0) items and the LR(0) automata.
- Get G' from the given grammar G .
- Get the canonical collection of items of G'
- Get the goto function.

Build SLRTable (G')

1. construct $C = \{ I_0, I_1, \dots, I_n \}$ the collection of set of LR(0) items for G' .
2. state it is constructed from I_i the actions are :
 - 2.1. if $[A \rightarrow \alpha \cdot a \beta] \in I_i$ and $\text{GOTO}(I_j, a) = I_j$ set action $[i, a]$ to shift j .
 - 2.2. if $[A \rightarrow \alpha \cdot] \in I_i$ set action $[i, a]$ to reduce $A \rightarrow \alpha \wedge a \in \text{follow}(A)$.
 - 2.3. if $[S' \rightarrow S \cdot] \in I_i$ set Action $[i, \$]$ to accept.
3. The goto are made for state i and NTA as :

if $\text{GOTO}[I_i, A] = I_j$; $\text{GOTO}[i, A] = j$.
4. Mark all undefined entries as error.
5. initial state is the state for $[S' \rightarrow S \cdot]$.

Drawback of SLR

- The method for filling in reduce steps is stupid and should never be used.
- Eg

$$\begin{aligned} S &\rightarrow L = R / R \\ L &\rightarrow *R / id \\ R &\rightarrow L \end{aligned}$$

$$\begin{aligned} I_0 = S' \rightarrow \cdot S & \quad I_1 = S' \rightarrow S \cdot \\ S \rightarrow \cdot L = R & \\ S \rightarrow \cdot R & \\ L \rightarrow \cdot *R & \\ L \rightarrow \cdot id & \\ R \rightarrow \cdot L & \end{aligned}$$

$$I_2 = \boxed{\begin{array}{l} S \rightarrow L \cdot = R \\ R \rightarrow L \cdot \end{array}}$$

⇒ Action $[2, =] = \text{shift to } 6$
 ⇒ Action $[2, =] = \text{Reduce } R \rightarrow L$] conflict.

Not all unambiguous grammars are SLR(1).

Variable Prefix

- The prefixes of right-sentential forms that can appear on the stack of a LR parser
- A variable Prefix is a prefix of a right-sentential form that does not continue past the right end of the rightmost handle of that sentential.
- It is a central theme of LR parsing that the set of valid items for a variable prefix γ is exactly the set of items reached from the initial state along the path labeled γ in the LR(0) automata.

Canonical LR Parser

- Makes full use of look ahead symbols
- Makes use of a large set of items called LR(1) items.

LR(1) items

- By splitting states we can arrange to have each state of an LR Parser to indicate exactly which input follows a handle α .
- The information is incorporated as lookahead.

$[A \rightarrow \alpha \cdot \beta, a]$

$A \rightarrow \alpha \beta$ is a production

a is a terminal or $\$$.

- Lookahead has no effect when β is not ϵ but $[A \rightarrow \alpha \cdot, a]$ calls for a reduction only when input is a .
- LR(1) item $[A \rightarrow \alpha \cdot \beta, a]$ is a valid variable prefix r if

$$\exists S \xrightarrow{*} SAw \xrightarrow{*} S\alpha\beta w$$
 - $\Rightarrow r = S\alpha$
 - $\Rightarrow a$ is first symbol of w or w is ϵ and $a = \$$.

Generating LR(1) Items

- Use a modified closure logic.

Closure (I):

repeat :

foreach Item $[A \rightarrow \alpha \cdot B \beta, a]$ in I

for each production $B \rightarrow r$ in G'

for each terminal b in $\text{first}(B)$:

add $[B \rightarrow \cdot r, b]$ to I ;

Until no more items are added.

GOTO (I, x):

INIT J as empty set

for each item $[A \rightarrow \alpha \cdot x \beta, a] \in I$:

add $[A \rightarrow \alpha x \cdot \beta, a]$ to J

return closure (J)

Building The CLR(1) Table

Build CLR Table (G'):

1. construct $C' = \{ I_0, I_1, \dots, I_n \}$ the set of LR(1) items.
2. State i is constructed from I_i ; the actions are:
 - 2.1 if $[A \rightarrow \alpha \cdot \beta, b] \in I_i$ and $\text{GOTO}(I_i, a) = I_j$ set Action $[i, a] = \text{shift } j$
 - 2.2 if $[A \rightarrow \alpha \cdot, a] \in I_i$ and $A \neq S'$ then set Action $[i, a]$ to reduce $A \rightarrow \alpha$.
 - 2.3 if $[S' \rightarrow S \cdot, \$] \in I_i$ set Action $[i, \$]$ to accept.
3. Goto are constructed from $i \vee \text{NT } A$ st: if $\text{GOTO}(I_i, A) = I_j$; $\text{GOTO}(i, A) = j$.
4. All blank entries are errors
5. The init state is $[S' \rightarrow \cdot S, \$]$

The number of states in CLR(1) > SLR(1) always.

~~LALR~~

Lookahead LR Parsers

- Formed by merging the states of a CLR(1) parser where no conflicts arise,
- a conflict can be a Reduce-Reduce conflict eg:

$$I_1 = \boxed{\begin{array}{l} A \rightarrow a \cdot, c \\ B \rightarrow b \cdot, d \end{array}} \quad I_2 = \boxed{\begin{array}{l} A \rightarrow a \cdot, d \\ B \rightarrow b \cdot, c \end{array}}$$

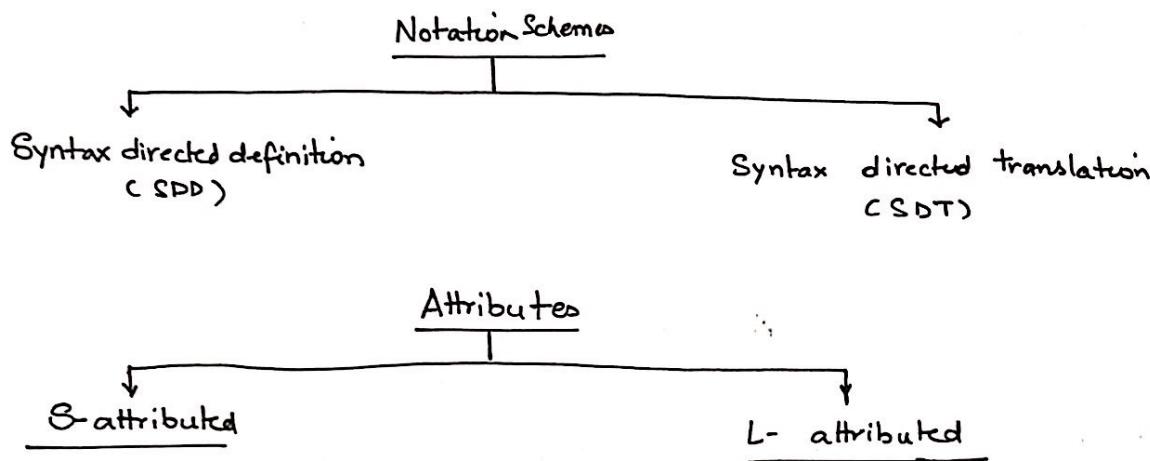
No conflict

$$I_{12} = \boxed{\begin{array}{l} A \rightarrow a \cdot, c/d \\ B \rightarrow b \cdot, c/d \end{array}} \quad R/R \text{ conflict.}$$

- The number of states in an LALR parser are always same as that in SLR parser.
- The idea of creating the LALR table is to construct the set of LR(1) items and if no conflicts arise, merge them, the ones with identical LR(0) cores.
- # All LALR grammars are CLR but all CLR grammars may not be LALR.

Semantic Analysis

- Can associate info with a construct by adding attributes to the grammar.
- Done by augmenting semantic Rules to the CFG.
- The Semantic analyzer ensures the program conforms to a certain semantic.
- Semantic Analysis can capture things like names, scopes, types etc which syntax cannot.
- The semantic rules are attached to the production.



each symbol has at least one associated attribute but depending on implementation exceptions might exist.

S attributes: Value of these depends only on the children (synthesized from).



These are evaluated in a bottom up manner.

L attributes : The value for these depends on the parent of the value or its left sibling.

These can be evaluated with both top down and bottom up parsers.

Semantic Analysis is not an independent process, it is done along with parsing.

Syntax Directed Definition (SDD)

- A SDD is a augmented CFG in which each grammar symbol, terminal or non terminal has an associated set of attributes.
- SDD are abstract specifications i.e It only tells the way in which the attributes relate but does not cover anything about when they are computed.
- The rules set up a dependency graph which can be used to derive an order to execute the rules.
- The parse tree containing values at each node is called the annotated parse tree.

Formally : A SDD is a generalization of a context free grammar where each production $A \rightarrow \alpha$ is associated with a set of semantic rules of the type:

$$X := f(a_1, a_2, \dots, a_k)$$

1. X is a synthesized attribute of A and $a_1, a_2, \dots, a_n \in \alpha$
2. X is an inherited attribute of one of the grammar symbols in α and $a_1, a_2, \dots, a_n \in A$ or α .

- # A SDD is S-attributed if every grammar symbol has synthesized attributes.
- # Inherited values are useful as they express dependence in the context in which they appear.

Eg > 1. Synthesized Values Only

$$\begin{aligned}
 E \rightarrow E_1 + T & \quad \{ E.val := E_1.val + T.val \} \\
 E \rightarrow E_1 - T & \quad \{ E.val := E_1.val - T.val \} \\
 E \rightarrow T & \quad \{ E.val := T.val \} \\
 T \rightarrow T_1 * F & \quad \{ T.val := T_1.val * F.val \} \\
 T \rightarrow T_1 / F & \quad \{ T.val := T_1.val / F.val \} \\
 T \rightarrow F & \quad \{ T.val := F.val \} \\
 F \rightarrow (E) & \quad \{ F.val := E.val \} \\
 F \rightarrow \text{num} & \quad \{ F.val = \text{num.lexval} \}
 \end{aligned}$$

2. Using Both Inherited and Synthesized

$$\begin{aligned}
 D \rightarrow T.list & \quad \{ list.type := T.type \} // sibling \\
 list \rightarrow list, id & \quad \{ list.type := list.type \} // parent \\
 & \quad \{ id.type := list.type \} \\
 list \rightarrow id & \quad \{ id.type := list.type \} \\
 T \rightarrow \text{int} & \quad \{ T.type := \text{integer} \} \\
 T \rightarrow \text{float} & \quad \{ T.type := \text{float} \}
 \end{aligned}$$

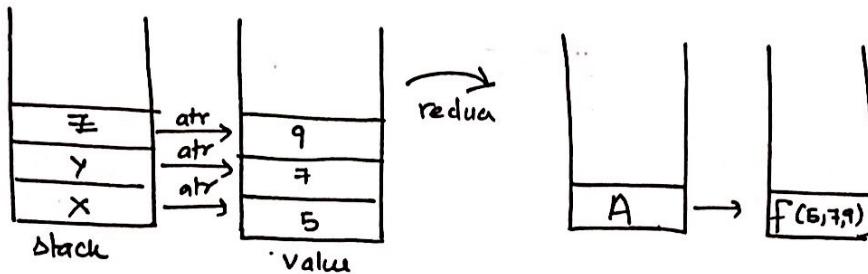
- # Every S-attributed value is always L-attributed.

S-attributed Bottom-Up Parser

- only have synthesized values
- Uses a second stack to hold the attributes. called the value stack.
- if the stack contains the symbol A at i then the value stack at i has its attributes.
- If there is no attribute defined the stack entry is undefined too. (like operators)
- The value of new synthesized attribute can be computed from the attributes appearing on the stack for that grammar symbol.

Eg if there are K terminal RHS of the production then Newtop after reduction of value stack is [Top - K + 1]

$$A \rightarrow XYZ$$

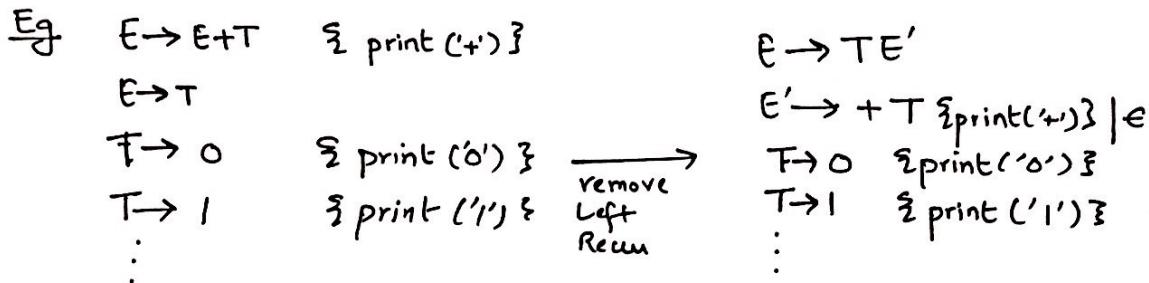


L-attributed Definitions

- S-attributed values don't allow for inherited values and are too restrictive.
- Formally.: A SDD is said to be L-attributed if for every $A \rightarrow x_1 x_2 \dots x_n$ each inherited value of x_j depends only on:
1. The attributes of x_k where $k \leq j-1$
 2. The inherited attributes of A.
- These can be implemented using LL(1) parsers for some grammars and using LR(1) parsers ~~for~~ for a larger grammar set.

Translation Schemes (SDT)

- A Translation scheme is a refinement of a SDD where the semantic actions are embedded within the right side of a production to specify the exact spot / point in time the actions are executed.
- A translation scheme generates an output for each sentence generated by the grammar by executing the actions in the order they appear during the DFS of the parse tree.



The actions are shown using dotted links in the parse tree.

- The restrictions of L-attributed values ensure that the attribute values are available.

Guidelines for designing a Translation Scheme

- if there are only synthesized attributes, we can create an action consisting of an assignment for each semantic rule & place it at the end of it (production)
- if there are both synthesized and inherited values.
 - Compute an inherited attribute for a symbol on the RHS of a prod in an action before that symbol.
 - No action should refer to a synthesized attribute to the right of it.
 - Compute synth attr for non terminals only after computing all the attributes it references. This action can be placed at the end of RHS.

L-Attributed Bottom Up Parser

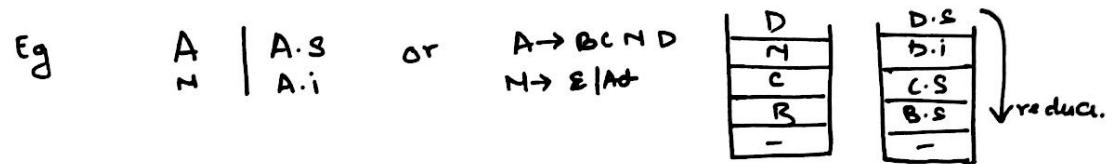
- To accommodate embedded actions we introduce marker symbols in place of the action and add another production for their marker.

Eg $A \rightarrow X_1 \{ \text{ACT} \} X_2$

becomes

$$\begin{aligned} A &\rightarrow X_1 M_1 X_2 \\ M_2 &\rightarrow \epsilon \{ \text{ACT} \} \end{aligned}$$

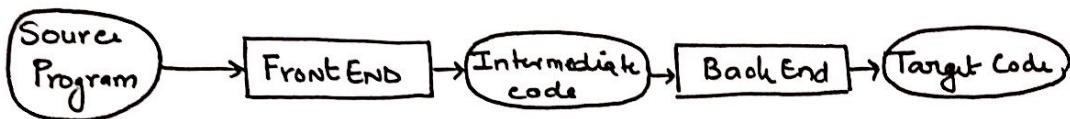
- Now any inherited value will be on the corresponding stack level of the marker symbol on the value stack.



- The remaining parsing follows that of SDD.

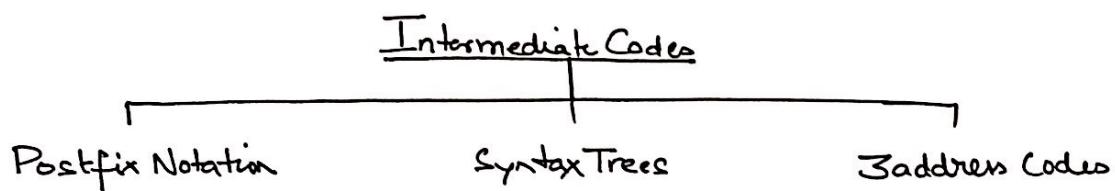
Intermediate Code Generation

- It is not always possible to directly generate target code in one pass of the compiler.
 - Features like delayed declaration may not be possible to handle in just one pass.
 - There may not be enough memory.
 - Single pass may not be able to generate highly efficient code.



Advantages of Machine Independent IC

- We can create compilers for different machines by only creating separate backends.
- We can create compilers for different languages by only making new frontends.
- We can apply machine independent optimisations to the IC and make it easier to translate.



To select a type of IC we need to consider:

1. ease of conversion from source to IC.
2. ease of conversion and optimization to target code.

Syntax Trees

- A syntax tree is a condensed form of the parse tree which is used to represent the syntactic structure of the language constructs. Operators and keywords are at the interior nodes while numbers and identifiers are at the leaves.

Construction of Syntax Trees

1. Create a leaf node for all constants and identifiers.
2. Create subtrees for all subexpressions.
3. operators are at the root of the subtrees with operands as children.

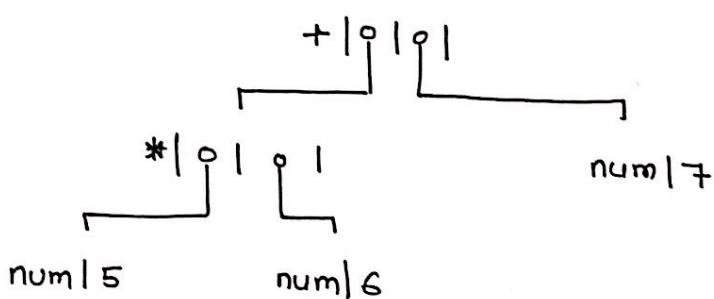
It uses 3 functions

1. CidLeaf (ident, sptr) : This takes in the id and the pointer to symbol table and returns a pointer to the newly made leaf node.
2. CnumLeaf (num, value) : This takes in the token number and its value and return a pointer to the new node.
3. CopNode (oper, first, Second) : This takes in the token oper and the two operands which may be other subtrees. and returns a new node.

Eg:

$E \rightarrow E_1 + T$	{	$E \cdot np := \text{copnode} (+, E_1 \cdot np, T \cdot np)$	}
$E \rightarrow T$	{	$E \cdot np := T \cdot np$	}
$T \rightarrow T_1 * F$	{	$T \cdot np := \text{copnode} (*, T_1 \cdot np, F \cdot np)$	}
$T \rightarrow F$	{	$T \cdot np := F \cdot np$	}
$F \rightarrow (E)$	{	$F \cdot np := E \cdot np$	}
$F \rightarrow id$	{	$F \cdot np := \text{cidleaf} (\text{ident}, \text{sptr})$	}
$F \rightarrow num$	{	$F \cdot np := \text{cnumleaf} (\text{num}, \text{val})$.

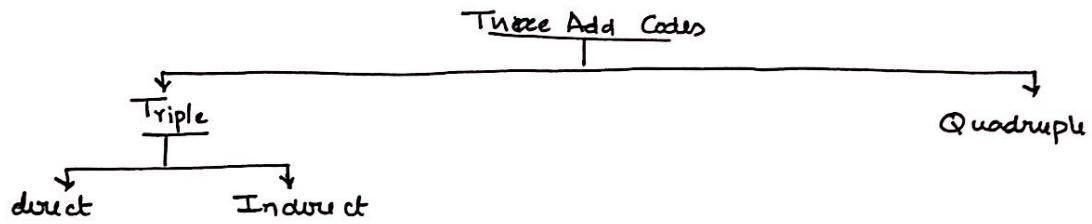
5 * 6 + 7



Add : $S \rightarrow id := E$ {
 id. $np := \text{cidleaf} (\text{ident}, \text{sptr})$
 S. $np := \text{cnode} (:=, id \cdot np, E \cdot np)$

To get SDD for assignment operation.

Three Address Codes



- Each statement at most has 3 addresses ; two for operands and one for the result

$a := b \text{ op } c$

Common 3 address codes

- $a := b \text{ op } c ,$
- $a := \text{op } b$
- $a := b$
- $\text{goto } x$
- if $a \text{ relop } b \text{ Goto } x$
- if $a \text{ goto } x$
- param a call p,n return x
- $a := b[i] \text{ or } b[i]:=a$
- $a := \&a \text{ or } a := *b \text{ or } *a:=b.$
- These are simple records having fields for operators and two operands.

Triple : 3 records i.e operator, operand1, operand2.

Operands are either pointers to the symbol table or pointers back in the table of triples.

Eg $a := a+b*c+c*d *e$

Triple	op	op1	op2	
0	*	b	c	} Table of triples
1	*	c	d	
2	*	(1)	e	
3	+	a	(0)	
4	+	(3)	(2)	
5	:=	x	(4)	

- Triples make identification of identical subexpressions easier.
- Table can be maintained and statements may be represented as a pointer to triples in order. This is Indirect Triple.
- Advantage: makes optimization of cases where statements are moved around better.

Quadruple: 4 records i.e operator, operand1, operand2 and result.

OP1, OP2, and result are references to symbol table entries.

- This use of result allows for greater flexibility for optimizing by making movement of statements even easier.

opr	op1	op2	res	
*	b	c	g	
*	c	d	t ₂	
*	t ₂	e	t ₃	
+	t ₁	a	t ₄	
+				
:=	t ₃	t ₄	t ₅	
	t ₅		x.	

list of temporaries.

Generation of Three Address Code

This is done by writing Semantic Rules ~~not~~.

It uses the following functions:

- 1) GenTemp(): returns a sequence of temporaries.
- 2) Gen(): generates the 3 address statement and add it to the list of output.
- 3) Lookup(id.entry): gives the entry for id if it exists or otherwise gives nil.

Eg) Assignment operation

S → id := E { id.name := lookup(id.entry)
 { if id.name != nil: gen(id.name := E.name) } }

E → E₁ + E₂ { E.name = gettemp
 gen(E.name := E₁.name + E₂.name) }

E → id { id.name = lookup(id.entry)
 { if id.name != nil:
 E.name = id.name } }

Type Conversion

- AKA coercion (done automatically by the compiler).

Eg $X := y * 3 + y * r$

float	int	float	int

$$t_1 = \text{int2float}(y)$$

$$t_2 = 3 * t_1$$

$$t_3 = y * r$$

$$t_4 = \text{INT2FLT}(t_3)$$

$$t_5 = t_4 * t_2$$

$$X = t_5.$$

Assignment OP

$S \rightarrow id := E$

```

    { if lookup(id.entry) = NIL : error
      if id.T = INT and E.T = INT
        gen C id.name := E.name
        S.type := void
      else if id.T = Float and E.T = Float
        gen C id.name := E.name
      else if id.T = F and E.T = int
        t := floatGenTemp()
        gen C t := int2float E.name
        gen C id.name := t
    }
  
```

$E \rightarrow E_1 + E_2$

```

    { E1.T = INT and E2.T = INT
      E.T = INT
      E.name := INTIENTEMP()
      gen C E.name := E1.name 'INT+' E2.name
    }
    elseif E1.T = Float and E2.T = float
      E.T = float
      E.name = floatGenTemp()
      gen C E.name := E1.name 'float+' E2.name
    elseif E1.T = float and E2.T = INT INT
      E.T = float
      E.name = floatGenTemp()
      gen C t := int2float E2.name
      gen C E.name := E1.name 'float+' t.name
    elseif E1.T = INT and E2.T = float
      E.T = float
      E.name := floatGenTemp()
      gen C t := int2float E1.name
      gen C E.name := t.name 'float+' E2.name
    else
      E.type := E-error
    }
  
```

Boolean Expressions

- They can be either `true` / `false` or using numbers `1` and `0`.
- Using Numeric Representation

$E \rightarrow E_1 \parallel E_2$	$E.name := \text{gentemp}$ <code>gen(E.name' := ' E1.name' ' E2.name');</code>
$E \rightarrow E_1 \wedge E_2$	$E.name := \text{gentemp}$ <code>gen(E.name' := ' E1.name' && E2.name');</code>
$E \rightarrow \text{true}$	$E.name := \text{gentemp}$ <code>genname gen(E.name' := ' 1)</code>
$E \rightarrow \text{false}$	$E.name := \text{gentemp}$ <code>gen(E.name' := ' 0)</code>
$E \rightarrow E_1 \text{ relop } E_2$	$E.name := \text{gentemp.}$ <code>gen (if E1.name relop.op E2.name goto nextlabel+3);</code> <code>gen (E.name := 0)</code> <code>gen (goto nextlabel +2)</code> <code>gen (E.name := 1))</code>

- Using Boolean Expression and Control flow

- We use the function `genlabel` which generates a new label.
- each symbol has a `.f` and `.t` values which are jump to labels.

$$E \rightarrow E_1 \parallel E_2 \quad \left\{ \begin{array}{l} E_1.t := E.t \\ E_1.f = \text{genlabel} \\ E_2.t := E.t \\ E_2.f := E.f \\ E.\text{code} := E_1.\text{code} \parallel \text{gen}(E_1.f :') \parallel E_2.\text{code}. \end{array} \right\}$$

$$E \rightarrow E_1 \wedge E_2 \quad \left\{ \begin{array}{l} E_1.t := \text{genlabel} \\ E_1.f = E.f \\ E_2.t = E.t \\ E_2.f = E.f \\ E.\text{code} := E_1.\text{code} \parallel \text{gen}(E_1.t :') \parallel E_2.\text{code}. \end{array} \right\}$$

$$E \rightarrow E_1 \text{ relop } E_2 \quad \left\{ \begin{array}{l} E.\text{code} = \text{gen (if ' E1.name relop.op E2.name 'goto' E.t) //} \\ \text{gen (goto E.f)} \end{array} \right\}$$

Control Flow Constructs

$S \rightarrow \text{if } E \text{ then } S$
| $\text{if } E \text{ then } S \text{ else } S$
| $\text{while } E \text{ do } S$
| Begin SL end | id - E
 $SL \rightarrow SL; S | S$
 $E \rightarrow \text{Arith exp with books.}$

The following instruction is in the label S.follow.

Expression like $S \rightarrow \text{if } E \text{ then } S$; S is done if E is true.

$S \rightarrow \text{if } E \text{ then } S$ $\left\{ \begin{array}{l} E.t := \text{genlabel} \\ E.f := S.\text{follow} \\ S_1.\text{follow} := S.\text{follow} \\ S.\text{code} := E.\text{code} \parallel \text{gen}(E.t':') \parallel S_1.\text{code}. \end{array} \right\}$

$S \rightarrow \text{if } E \text{ then } S_1 \text{ else } S_2$ $\left\{ \begin{array}{l} E.t := \text{genlabel} \\ E.f := \text{genlabel} \\ S_1.\text{follow} := S.\text{follow} \\ S_2.\text{follow} := S.\text{follow} \\ S.\text{code} := E.\text{code} \parallel \text{gen}(E.t':') \parallel S_1.\text{code} \parallel \text{gen}(\text{goto}', S.\text{follow}) \\ \parallel \text{gen}(E.f':') \parallel S_2.\text{code}. \end{array} \right\}$

$S \rightarrow \text{while } E \text{ do } S_1$ $\left\{ \begin{array}{l} S.\text{start} := \text{genlabel} \\ E.t := \text{genlabel} \\ E.f := S.\text{follow} \\ S.\text{follow} := S.\text{start}. \\ S.\text{code} := \text{gen}(S.\text{start}:') \parallel \\ E.\text{code} \parallel \\ \text{gen}(E.t':') \parallel \\ S_1.\text{code} \parallel \\ \text{gen}(\text{goto}', S.\text{start}) \end{array} \right\}$

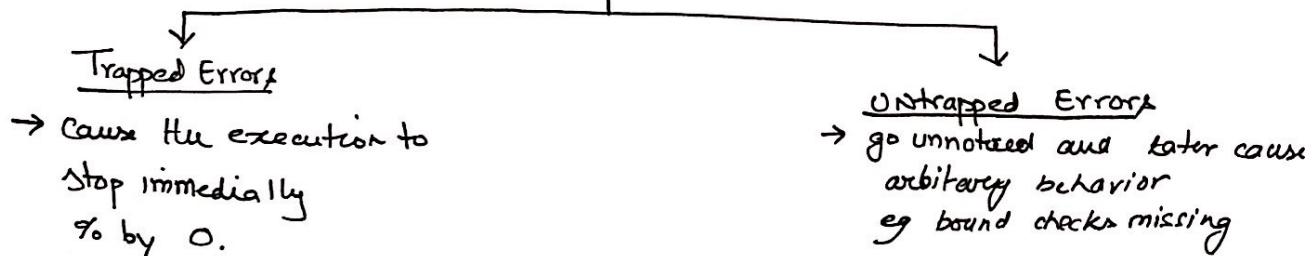
Type Analysis And Type Checking

- Static type checking ensures that certain type of checks like uniqueness, flow and control are performed and errors reported.
- The static type checker needs to detect and report type mismatch errors.

Type System

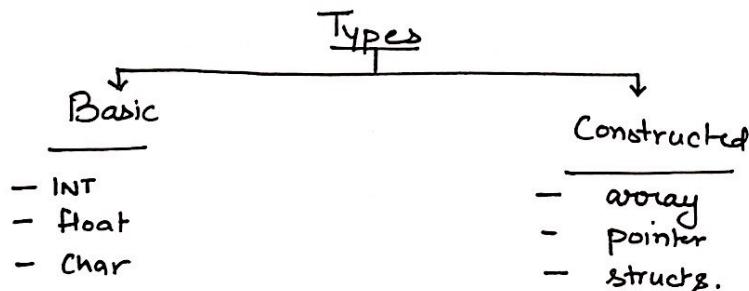
- The type system keeps track of types of variables and of all expressions.
- Type system determines if a program is well behaved.

Execution Errors



A program fragment is safe if there are no trapped errors. Typed languages can force good behavior.

- Type system is a collection of rules for assigning type expressions.
- If a language can guarantee if its compiler ~~can~~ accepts a program with no errors then it is called strongly typed.
- The source language defines the notion of types and the associated rules.



Type Expressions

This is a way to express the type of a language construct.

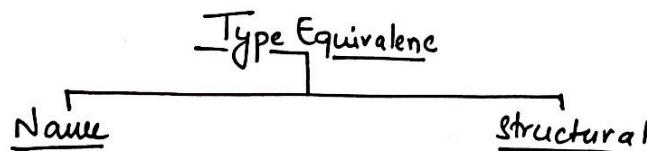
1. Basic : int, char, float
2. error : special basic type to deal with errors.

- 3. Null : special basic type to show no type
- 4. name : type defn and aliases
- 5. Type constructor : arrays, products, structs, pointers, function .
 - 5.1 Arrays: if array elements are type T and T is a type then array (1,T) is a type
 - 5.2 Product: The cartesian product $T_1 \times T_2$
 - 5.3 Structs: The constructor struct applied to the product of the type of field.
 - 5.4 Pointers: Pointer(T) is a type expression if T is.
 - 5.5 function : Function is a mapping from D to R types $D \rightarrow R$.
e.g. array (0...99, float) \times int \rightarrow float

Type expressions may be Trees or DAGs.

Basic Type Checker

- Type checker checks for type compatibility by checking if two expressions are equivalent.



Name Equivalence

Two expressions are name equivalent iff they are identical. i.e. Their type names match and the type names are treated as distinct.

Structural Equivalence

Two expressions are struct eq if they represent structurally equivalent types.

```

set_eq (a, b) {
  if a, b are some basic type
    ← true;
  elif a = array (a1, a2) and b = array (b1, b2)
    return set_eq (a1, b1) and set_eq (a2, b2)
  elif a = a1 × a2 and b = b1 × b2
    return set_eq (a1, b1) and set_eq (a2, b2)
  elif a = pointer(a1) and b = pointer (b1)
    return set_eq (a1, b1)
  elif a = a1 → a2 and b = b1 → b2
    return set_eq (a1, b1) and set_eq (a2, b2)
  else return false
}
  
```

For grammar G :

$P \rightarrow D; S$
 $D \rightarrow D; D \mid TL$
 $T \rightarrow \text{int} \mid \text{float} \mid \text{char}$
 $L \rightarrow L[\text{num}] \mid *L \mid id$
 $S \rightarrow id := E \mid \text{if}(E)S \mid \text{while}(E)S \mid S_1; S_2$
 $E \rightarrow \text{literal} \mid \text{num} \mid id \mid E \% E \mid E \text{ binop } E \mid *E \mid E[E] \mid E(E).$

1> Type express for main record.

$D \rightarrow TL \quad \{ L.T_y = T.T_y \}$
 $T \rightarrow \text{int} \quad \{ T.T_y = \text{int} \}$
 $T \rightarrow \text{float} \quad \{ T.T_y = \text{float} \}$
 $T \rightarrow \text{char} \quad \{ T.T_y = \text{char} \}$
 $L \rightarrow id \quad \{ \text{addtype(id.entry, L.type)} \} \quad // \text{associates type in ST} //$
 $L \rightarrow *L_i \quad \{ L_i.T_y := \text{pointer(L.type)} \} \quad // \text{Inherited}$
 $L \rightarrow L[\text{num}] \quad \{ L.T_y = \text{array(0...NUN, L.type)} \} \quad // \text{Inherited}$

2> Expressions

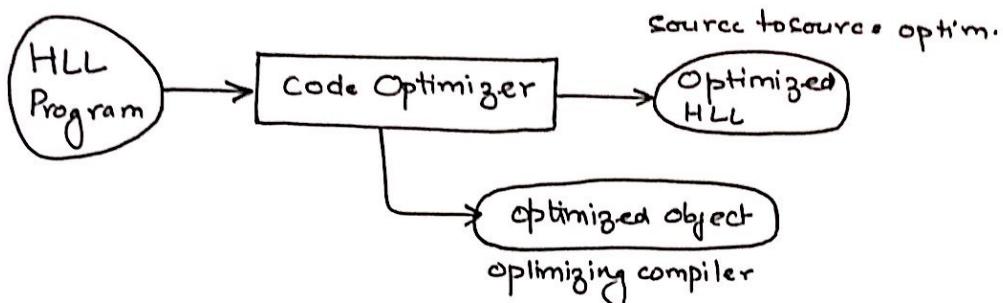
$E \rightarrow \text{literal} \quad \{ E.T = \text{char} \}$
 $E \rightarrow \text{num} \quad \{ E.T = \text{int} \}$
 $E \rightarrow id \quad \{ E.T = \text{check(id.entry)} \}$
 $E \rightarrow E_1 \text{ binop } E_2 \quad \{ \text{if } E_1 \text{ is } \text{if} \text{ then } E_2 \text{ else } \text{error} \}$
 $E \rightarrow E_1 \% E_2 \quad \{ \text{if } E_1 \text{ is } \text{if} \text{ then } E_2 \text{ else } \text{error} \}$
 $E \rightarrow *E_1 \quad \{ \text{if } E_1.T_y = \text{pointer(t)} \}$
 $E \rightarrow E_1[E_2] \quad \{ \text{if } E_2.T_y = \text{int} \text{ and } E_1.T_y = \text{array(int,t)} \}$
 $E \rightarrow E_1(E_2) \quad \{ \text{if } E_2.T_y = S \text{ and } E_1.T > S \rightarrow t \}$

3> Statement Checking

$S \rightarrow id := E \quad \{ S.T_y := \text{void if } id.T_y = E.T_y \text{ else error} \}$
 $S \rightarrow \text{if}(E)S \quad \{ S.T_y := S_1.T_y \text{ if } E.T_y = \text{int} \text{ else error} \}$
 $S \rightarrow \text{while}(E)S \quad \{ S.T_y = S_1.T_y \text{ if } E.T_y = \text{int} \text{ else error} \}$
 $S \rightarrow S_1; S_2 \quad \{ S.T_y = \text{void if } S_1.T_y = \text{void and } S_2.T_y = \text{void else error} \}$

Code Optimization

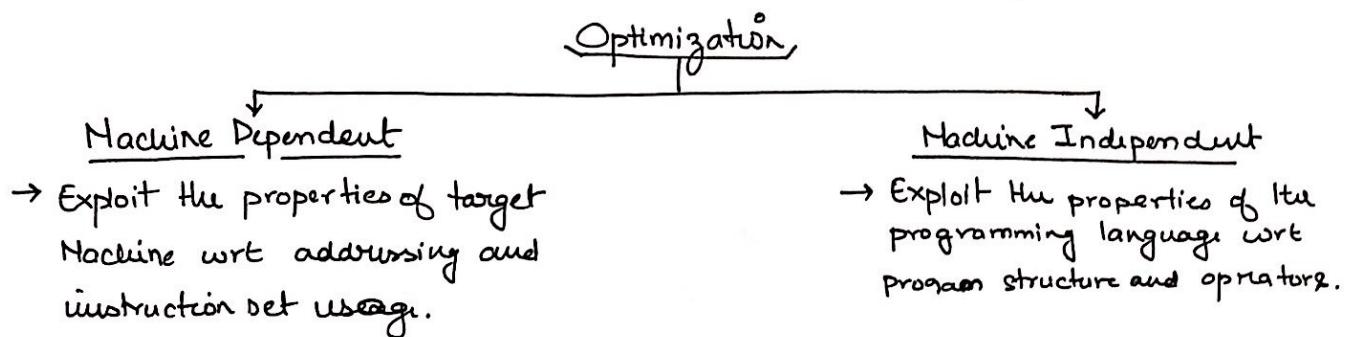
- Techniques to obtain more efficient object code than what would be obtained by the code generator.
Goal: Must ensure the output is semantically same as the input
 - : Must not alter any algorithms used in the source program.



- Source-to-Source produces code that is portable and can be executed on any machine that has that compiler.
- Optimizing compilers have the advantage of being able to do array reference and register optimization.

Levels of Optimization

1. At source level by the coder by changing loops or transforming loops.
2. At the ICG level by the compiler by improving loops and address calculations.
3. At the target code level by keeping heavily used values in registers.



Machine Dependent Optimization

- Allocation of scarce resources to achieve high efficiency in execution.
- Use of intermediate instructions, incrementation, indexing, indirection & vector ops.
- Intermixing data with instructions.

Machine Independent Optimization

- Based on the mathematical properties of a sequence of source statements.
- Optimization essentially means analysing the overall purpose of statement to find an equivalent sequence.

Optimizing Transform Properties

1. Preserve the meaning of the program.
2. Speed up programs by a measurable amount.
3. Reduce size of program.
4. The time and resource spent optimizing must be worthwhile.

1. Common Subexpression Elimination

- Expressions which have already been computed need not be evaluated again if the value in the variables hasn't changed.

e.g. $a := b * c;$
 ⋮
 $d := b * c + x - y;$

We can eliminate the 2nd $b * c$ from the exp if no intervening statement has changed its value.

$t_1 := b * c;$
 $a := t_1;$
 ⋮
 $d = t_1 + x - y;$

Typically; two expressions are common if :

1. They are lexically equivalent.
2. They evaluate to identical values. i.e. No assignment for any operands exists between the two expression.
3. The value of operands must not change even due to a procedure calls.

$$\left. \begin{array}{l} a := b * c; \\ b := 10; \\ \vdots \\ d = b * c + x - y; \end{array} \right\} \text{Not reducible} \quad \left\{ \begin{array}{l} a := b * c; \\ x := b \\ \vdots \\ d = x * c + x - y \end{array} \right.$$

2. Compile Time Evaluation

- Some runtime actions can be shifted to compile time to reduce re-running of code sections.
- Expressions with constant operands can be evaluated at compile time and replaced by a single value. This is called folding.

$$A = 2 * (22.0 / 7.0) * r;$$

can be optimized to

$$A := (6.28...) * r;$$

- Folding is done by most compilers during array addresses.

$$x := 12.4;$$

:

$$y := x / 2.3;$$

- Between x and y no operation changes x . \therefore we can instead have

$$y := 12.4 / 2.3$$

done during compile time, this is called constant propagation.

- Constant propagation means using the constant value assigned rather than variable reference. This enhances the scope of folding.
- This is only possible when the variable is assigned the same value along all paths reaching the point of use.
- Compile Time Evaluation can also eliminate identity operations like $1*a$, $0*b$.

3. Variable Propagation

- Replacing the use of one variable by another if it has been assigned the same value along all paths reaching its use and none of the variable is modified after this use.

$$c = a ** b$$

$$x = a$$

$$\vdots$$

$$d = x ** b + x - y \Rightarrow d = a ** b + x - y \Rightarrow \text{folding possible now.}$$

4. Dead Code Elimination

- If the value assigned to a variable is not used beyond a point in code then the variable is said to be dead at that point.
- Any assignment made to a dead variable are wast and can be eliminated.
- A block of code is said to be dead if it computes values which are dead.
- Assignment to a var results in dead code if another assignment is made before it is used.
- Variable propagation helps assignment operations making them dead.

$$\begin{aligned}
 & c = a * * b \\
 & x = a \\
 & \vdots \quad \text{propagate} \\
 & d = x * * b + x - y \\
 & \Downarrow \\
 & c = a * * b \quad \text{common} \\
 & x = a \quad \text{dead code} \Rightarrow \quad t_1 = a * * b \\
 & \vdots \\
 & d = a * * b + x - y \quad c = t_1 \\
 & \quad \quad \quad d = t_1 + x - y
 \end{aligned}$$

5. Code Movement

- Reduces evaluation frequency of expressions. This is done by moving the evaluation to other parts of the program.

Eg if $a < 10$
 $b = x^2 - y^2$

else
 $b = 5$

$a = x^2 - y^2 * 10;$

- $x^2 - y^2$ is evaluated twice which can be optimized to.

$$\begin{array}{l}
 \text{if } a < 10 \\
 \quad t = x^2 - y^2 \\
 \quad b = t \\
 \text{else} \\
 \quad t = x^2 - y^2 \\
 \quad b = 5 \\
 \quad a = t * 10;
 \end{array}
 \xrightarrow{\text{common + Repeat}}
 \begin{array}{l}
 t = x^2 - y^2 \\
 \text{if } a < 10 \\
 \quad b = t \\
 \text{else} \\
 \quad b = 5 \\
 \quad a = t * 10;
 \end{array}$$

Reduced to partial Redundance \longrightarrow moved code.

This is called Code Hoisting i.e moving code out of loop before in the original code.

This also reduces space usage.

This is significantly useful while the code is hoisted out of loops when they have invariant code in them.

6. Strength Reduction

This is like frequency reduction but instead of reducing the number i.e. frequency of execution we reduce the computational complexity of the operations.

By strength Reduction we mean replacing high strength op with low ones.

Eg replacing $*$ with $+$ when possible.

$i=1$

while ($i < 10$): can be transformed to:

$$y = i * 4$$

$i=1$

$t=4$

while ($i < 10$):

$$y = t;$$

$$t = t + 4;$$

This is normally done for induction variables + light $*$ operators.

• Induction variable : a var of the type $i := i \pm \text{const}$ ".

This reduces overall running time.

It is unfavorable to do Strength Reduction to floats.

7. Loop Test Replacement

- we can replace a loop termination test with one variable with another equivalent variable.
- following S.R the loop induction variable is useless & \therefore replaced by tmp.

$$\begin{array}{l} i=1 \\ t=4 \end{array}$$

loop : $y = t$
 $i = i + 1$
 $t = t + 4$
if $i < 10$ goto loop



$$t=4$$

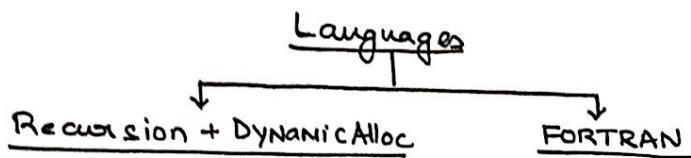
loop $y = t;$
 $t = t + 4$
if $t < 40$ goto loop

Memory Allocation

- To bind data items to memory locations. i.e. To collect data names & info and determine size of instructions.
- To allocate data we must know:
 - length of data item.
 - Type of data item.
 - dimension if any
 - Scope if any.

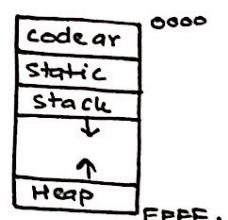
Features Affecting Malloc

- 1> Recursion : need to maintain various states.
- 2> Param Passing method.
- 3> Local Names.
- 4> Non Local Names.
- 5> Dynamic data structure.



Subdivision of run time Storage

1. code area : holds the generated target code
2. static area : holds the absolute addresses of determined items
3. Stack area : data objects or procedures.
4. Heap area : data items created at run time



Allocation Strategies

- Activation Record : a conceptual aggregate of data that contains all the info req for a single execution of procedure.
- The record is kept on the stack and new one is pushed for a call and popped for return

DYNAMIC LINK	STATIC LINK	SAVED Regist	Return value	Return add	Actual Parame	Local Vars	Tempor.
Nonlocal data in some other record.		Storage value of item to be returned		return address in calling	calling Parame		selfexplanitory

Annotations below the table:

- Dynamic Link: Pointer to ACTR of calling funct.
- Static Link: values of program counters
- Saved Registers: return address in calling
- Return Value: selfexplanitory

Static Allocation

- Vars remain permanently allocated irrespective of accessibility.
- Names are bound to storage as the program is compiled.
- Local names are retained across various activation records.
- The amount of storage is determined by
 - 1> elementary datatype → integral no of bytes.
 - 2> contiguous blocks → arrays / strings.
- It is possible to find address at which code can find data it needs.

Limitations

- Memory is bound for subprograms which may be mutually exclusive.
- Memory is bound for subprogram which are never called.
- Cannot recurse.
- No runtime alloc.
- All data item sizes must be known.

Stack Allocation

- Storage is organized as a stack. An ACTR is pushed when a new procedure is run.
- The storage free on popping is available for reuse.
- Mutually exclusive blocks no longer bind space.
- Activation records contain storage for the local vars of that call, which are lost when it is popped.

Advantages:

1. supports recursion
2. supports runtime malloc
3. allows for array declarations like $A(I,J)$.

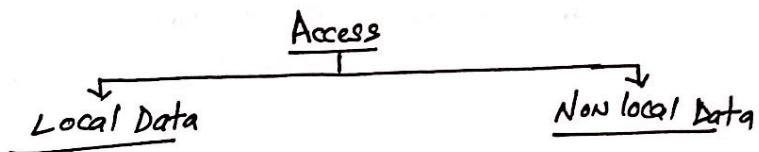
Drawbacks:

Memory addressing is via pointers and index regs which slows stuff down

Heap Allocation

- Allows retention of local values when the ACTR is popped. The programmer explicit control over the alloc / dealloc.
- Alloc is done from a large free area and access is through pointers.
- The space after a dealloc is not reused rather garbage collection is done.
- Heap alloc is useful for cases which are not LIFO or which are asynch.
- Heap alloc is more general but not as efficient as stack.

Memory Allocation In Block Structure Language



Local Data

- Local data may be referred within the record using an offset wrt the base of the record.
- There is a pointer to the start of the record: Base Pointer .
- The BP always sit at the Base of Record in execution .
- When a new procedure is entered, space is allocated and BP is assigned current value of Top. & then Top is incremented.

New Procedure : $*DL = BP$ $*DL \rightarrow$ DynamicLink.
 $BP = Top$
 $Top = Top + |actr|$

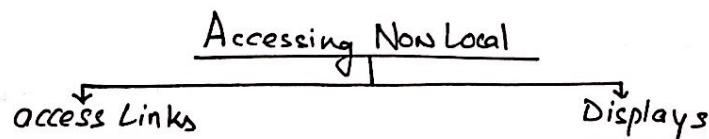
Procedure return : $Top = BP$
 $BP = *DL$.

Non Local Data

- values of non local data cannot be accessed with just offsets .
- " A Block is a statement containing its own data declarations enclosed b/w Delimiters .

The names that can be referred to at any point in Block language using static scoping :

1. Locally declared
2. Name of enclosing procedures
3. Name declared more closely to the block overrides the other .



- To access non locals we use the static link pointer which points to the activation record of the most closely enclosing block.
- The non local vars are accessed through these pointers starting from base of ACTR.

To setup static links:

1. Level of main is 0.
 2. if level of enclosing is l then level of enclosed is $(+1)$.
- Static link is traversed $4 - l_2 + 1$ times starting from ACTR of calling.
 - Static link of called is \neq set to base of new ACTR.
 - Using static links is costly if large static block nesting exists.

To improve efficiency we use an array (called the display) of pointers to the activation records.

The array is such that the act of x at level i is at $display[i]$.

- To setup $display[i]$ to point to the base of current act at level i .
The old $display[i]$ is saved within act of the procedure.
- Current value of $display$ is saved in the new record & $display[i]$ now points to it.
- Display is reset to old value at the end of execution.
- Display is updated whenever new activation occurs & whenever control returns from a new activation.
- Using Display we can reference all non locals at some cost.

Symbol Table Organization

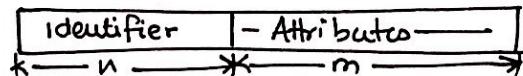
- It is a table containing all names used in the source program and their associated values.
- The table is looked up every time a name is encountered in the source program.
- It has two phases, building where symbols and associated values are inserted and a lookup/reference phase where values/ symbols are search.
- Records are deleted when the symbol goes out of scope.

Desired Capabilities

1. Provide entry and retrieval of symbols.
2. Provide a mechanism for search.
3. Be able to delete elements which are no longer needed.
4. Provide fast lookup by keeping the table on memory.
5. It should be able to grow as per the needs.
6. It must be able to handle duplicate entries. i.e support scopes.
7. Must be designed in a manner that some 3rd party can manage the table data structures and modules.

Elementary Symbol Table Org.

- Entries are stored as consecutive seq of words in memory. The nature of info depends on the name, this can be made uniform by moving some info about the name out.
- When the maximum allowed size of a id is known a fixed length entry can be used, where n unit are reserved for the ID name and m for the attributes.

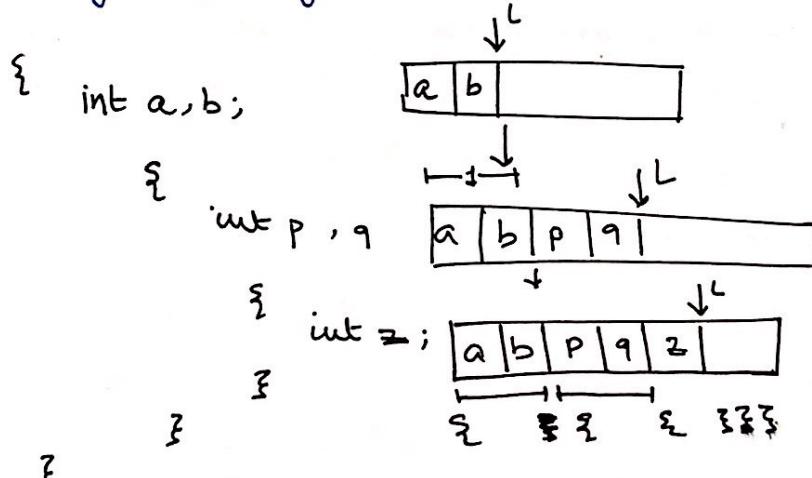


- Using a fixed length strategy has following advantages :
 1. All entries are alike and can be swapped.
 2. Table length can be fixed / estimated
 3. Position of an entry can easily be calculated.
- When the maximum size of identifier is not known, a variable length entry is used.
- l units of the block are kept for the length of the ID.
- Use of variable length entries restricts the table to linear lookup strategies.

length	name	Attribute
$l.$	$n.$	$m.$

1. Linear Search Organization

- Linear list of records with each record corresponding to only one symbol.
- The symbols are stored in order of arrival.
- The search is from last to first entry.
- Scope rules are inherently handled due to the entry scheme, inner levels are closer to the end.
- For very large programs, the search can take a very long time



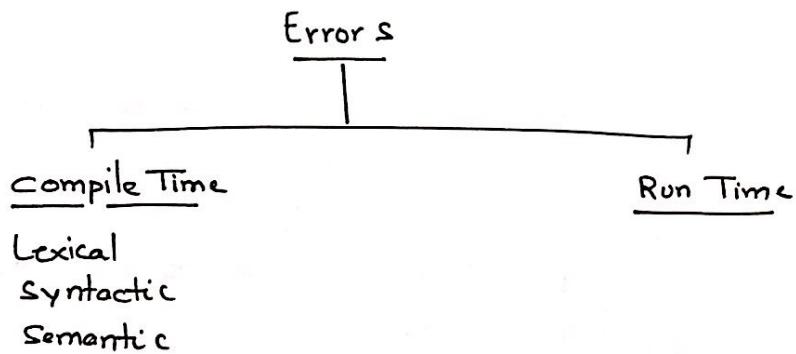
2. Binary Search Organization

- Search time is reduced by storing the entries sorted in a BST.
- The records are stored in groups of starting symbols.
- The lookup is reduced to $O(\log(n))$ time.
- The insert and delete are fairly complex tasks which requires moving of the table.
- This is best for static capacity tables.

Error Handler

Goals of the Error Handler

- To detect errors during various phases of the compiler
- To report the errors clearly providing sufficient info on the cause and nature
- To recover from errors in order to continue processing
- To correct common errors [difficult]
- To Not slow down the compiler.



Lexical Errors

- Difficult to locate due to nature of the LA.
- Possible sources :
 - Extra char.
 - Incorrect char.
 - Missing char.
 - Transposed chars.
 - Max ID-len exceeded.
- One method is to delete successive chars from remaining input till something matches - panic mode.
- Pass the char to the parser and let it decide.
- Error Reporting must be:
 - clear and unambiguous
 - avoiding duplications.
 - highlight the source of error
 - be as close as possible to the error.

- Must be in terms of source program entities
- Error caused by recovery must be different from a genuine error.
- If a correction is made, it should be notified to the user.

Syntax Errors :

- Most common type of error.
- Identified when tokens disobey the grammatical rules of the language.
- Table driven parsers have explicit entries for errors.
- Tables have a viable prefix rule i.e detect error as soon as there's a mismatch.
- Reporting :
 - Info on the token
 - 1. Examine entries in each state till the errors
 - 2. break routine e_i for state i
 - 3. set of valid tokens for i .
 - 4. e_i first tells the line, then token & then the expected token.
 - e_i can be same for 2 states.
- Recovery
 - suspend normal parsing.
 - change the error configuration.
 - resume normal parsing.
- 1. Panic Mode : Attain a state from which the parsing can safely resume.
 - Discard input symbols one-by-one until a sync token is found. The sync need not be the end of the ~~token~~ statement.
 - Syncs are normally delimiters.
 - Examines at the point of error.
 - Guarantees it'll work and not loop ~~only~~
 - Skips a lot of code (which itself may have errors)

2. Phase Level Recovery : perform local correction on remaining input.

- delete a source symbol.
- Insert a synthetic symbol
- Replace source symbol with synthetic symbol.
- present new string which bypass the error.
- can go into ∞ loop if the incorrect symbol is used
- Cannot be used when errors where of error occurs much before detection.

3. Error Productions : Encorporate production for most common type of errors.

when an error occurs, these production are used to create the error prompt & parsing continues.

4. Global Correction : Minimizes the level of modifications at the global level when multiple possible modif exist.
i.e min dist recovery.

Normally Parsers go for simpler and faster logic rather than minimizing cost.

- Top-Down Parsers

- 2 cases, terminal at top if next symbol or NT at Top and input symbol have a blank entry.
- Use sync set. It has symbols from First and Follow.
- do panic mode or parse the stack and create list of syncs till the point of error & then unstack till ~~top~~ some symbol is at top.

- Bottom-Up Parsers

- Use inserting synthetic symbols.
- For panic mode, scan down the stack till until s with goto A. discard input till symbol which matches A is found.
- OR fill blanks with routines.

Semantic Errors

- Due to violation of rules of accessibility i.e type error, arr index.
- Can be local / global in scope.
- Not always possible to detect.
- Can be immediate or delayed.

Immediate : can be detected while processing the erroneous line.

$$A = B^{++}C.$$

- Not that easy to recover/correct.

Delayed : Not detected when the line is processed but when the line takes effect.

Eg: Missing labels, Invalid Transfer (outside of loop to inside)

The errors are detected when the destination is encountered.

- Intersecting DO's / invalid nesting.

Runtime Errors

- Due to violation of semantics of machine computation.
eg underflows / overflows.
- Detected by machine H/W. This passes info to runtime control routine & then it generates the error. Some times recovery is done using assumptions.
- undefined values → Flag bits
- Dimension overrun → Bound check.
- I/O device → IOCS .