

# Syntax and Parsing

---

- ***Syntax***: the form of a program
- ***Semantics***: the meaning of a program
- Two parts to syntax analysis:
  - **lexical rules**: define legal characters and how they can be combined to form symbols  
("lexemes").
  - **syntax rules**: define how categories of lexemes  
("tokens") can be combined to form legal programs.

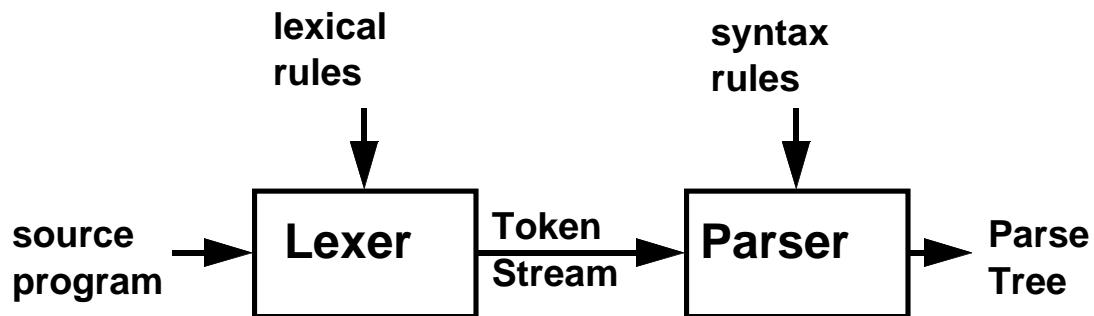
# Syntax and Parsing

---

- **More formal definition of Syntax:**
  - **Syntax defines the symbols of a language, and the spatial relationships between symbols in that language**
- **Semantics**
  - **Semantics describes the meanings to be ascribed to symbols, and structures of symbols, in a language.**

# Compiler Front End

---



## Syntax Analysis

- **We won't make a strict distinction, but will generally deal with syntax rules.**

# How to Describe the Syntax of a Language?

---

- **English description**
  - lengthy, tedious, ambiguous
- **Formal description**
  - recognizer: given a string, a recognizer for a language  $L$  tells whether or not the string is in  $L$  (ex: Compiler)
  - generator: a generator for  $L$  will produce an arbitrary string in  $L$  on demand. (ex: Grammar)
- **Recognition and generation are useful for different things, but are closely related.**
- **First, we'll talk about an important generation tool: BNF.**

# BNF - A Generative Tool

---

- **Backus-Naur Form (BNF) is a *metalanguage* for describing the syntax of programming languages.**
  - developed by John Backus and Peter Naur
  - first used to describe ALGOL60
- **A language description in BNF is called a *grammar*.**

Note: In formal linguistics a grammar is described as the syntax of a language together with an knowledge of the allowable constructs that provide “sensible” meanings. Obviously the syntax of a programming language does not contain “knowledge”, but we persist in this misuse of the term “grammar”.

# Grammars

---

- **A grammar is made up of productions, or rules, e.g.,**

`<sentence> --> <subject><verb><object>`

`<verb> --> see | love`

`<subject> --> I`

`<object> --> him | her`

- **Five components:**

**The rules using:**

**--> : "is defined as"**

**| : "or"**

**implied concatenation**

**terminals: see, love, I, him, her.**

**non-terminals: <sentence>, <verb>, <subject>, <object>**

**root symbol: <sentence>**

**the algorithm of generation or analysis**

Note: Except where specifically required, spaces are assumed to be inserted between “lexemes” in natural languages.

# Recursion

---

- **Need recursion to define strings of indefinite length:**

$\langle \text{ones} \rangle \rightarrow 1 \mid 1\langle \text{ones} \rangle$   
 $\Rightarrow 1, 11, 111, 1111, \dots$

$\langle \text{ablist} \rangle \rightarrow ab \mid a\langle \text{ablist} \rangle b$   
 $\Rightarrow ab$   
OR  $a\langle \text{ablist} \rangle b$   
 $\Rightarrow aabb$   
OR  $aa\langle \text{ablist} \rangle bb$   
 $\Rightarrow aaabbb$   
OR  $aaa\langle \text{ablist} \rangle bbb$   
 $\dots$

Describe the language resulting from this syntax!

Write a syntax that describes a language composed of strings of arbitrary numbers of “a”s followed by arbitrary numbers of “b”s

# Derivations

---

- The steps in generating a string using a grammar are called a *derivation*.

$$\langle \text{exp} \rangle \rightarrow \langle \text{id} \rangle \mid \langle \text{exp} \rangle + \langle \text{exp} \rangle \mid \langle \text{exp} \rangle * \langle \text{exp} \rangle \mid (\langle \text{exp} \rangle)$$
$$\langle \text{id} \rangle \rightarrow A \mid B \mid C$$

---

Note: “  $\rightarrow$  ”  
“  $\Rightarrow$  ”

$$\begin{aligned} \langle \text{exp} \rangle &\Rightarrow \langle \text{exp} \rangle + \langle \text{exp} \rangle \\ &\Rightarrow \langle \text{exp} \rangle * \langle \text{exp} \rangle + \langle \text{exp} \rangle \\ &\Rightarrow \langle \text{id} \rangle * \langle \text{exp} \rangle + \langle \text{exp} \rangle \\ &\Rightarrow A * \langle \text{exp} \rangle + \langle \text{exp} \rangle \\ &\Rightarrow A * \langle \text{id} \rangle + \langle \text{exp} \rangle \\ &\Rightarrow A * B + \langle \text{exp} \rangle \\ &\Rightarrow A * B + \langle \text{id} \rangle \\ &\Rightarrow A * B + C \end{aligned}$$

- Each step produces a *sentential form*. It is a string of terminals and nonterminals.
- A *sentence* is a sentential form containing only terminals.

# Language

---

- **We are now ready to define a LANGUAGE:**
  - A language is a set of strings of terminals.
- **Note**
  - A language may be finite or infinite. The empty set is a language.
- **So now define the LANGUAGE GENERATED BY A GRAMMAR:**
  - The language generated by a grammar is the set of sentences derivable using the grammar.

# Left-Most and Right-Most Derivations

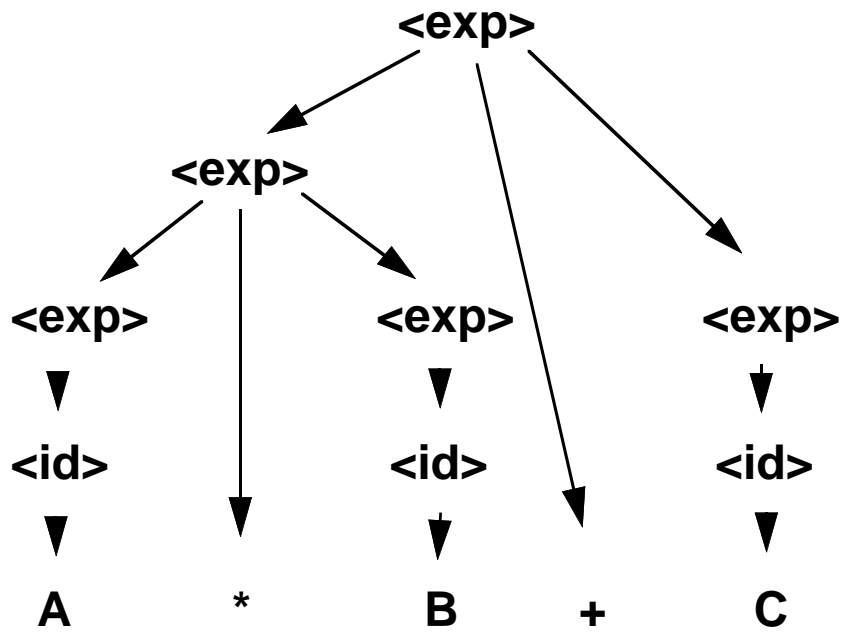
---

- **How to choose which non-terminal to replace next:**
  - left-most derivation: replace left-most NT first
  - right-most derivation: replace right-most NT first
- **Need not be either of these**
  - random replacement OK
    - can't affect **language** generated

# Parse Trees

---

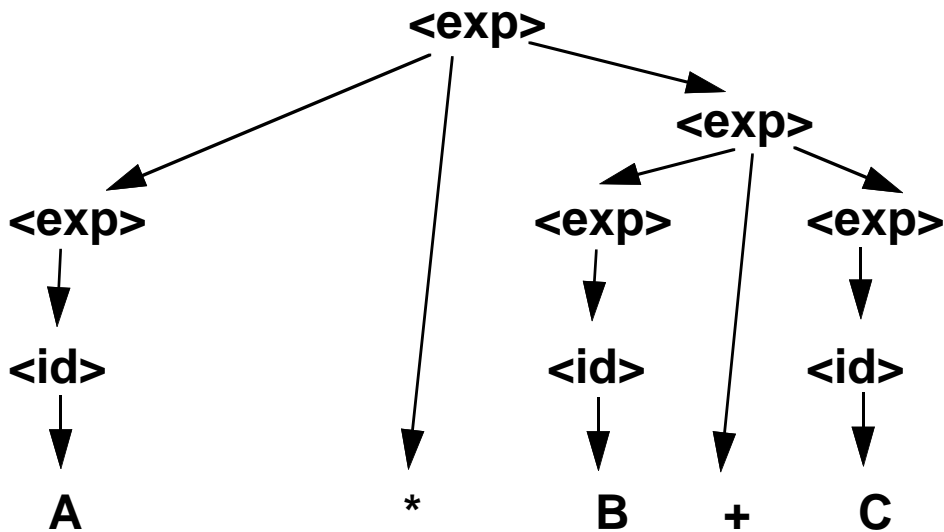
- Show the syntactic structure of sentences.



# Ambiguous Grammars

---

- A grammar is ambiguous if it generates the same sentence for which there are two or more parse trees.
- Another parse tree for  $A * B + C$ :



# Disambiguating the Grammar

---

- To disambiguate this grammar, change to:

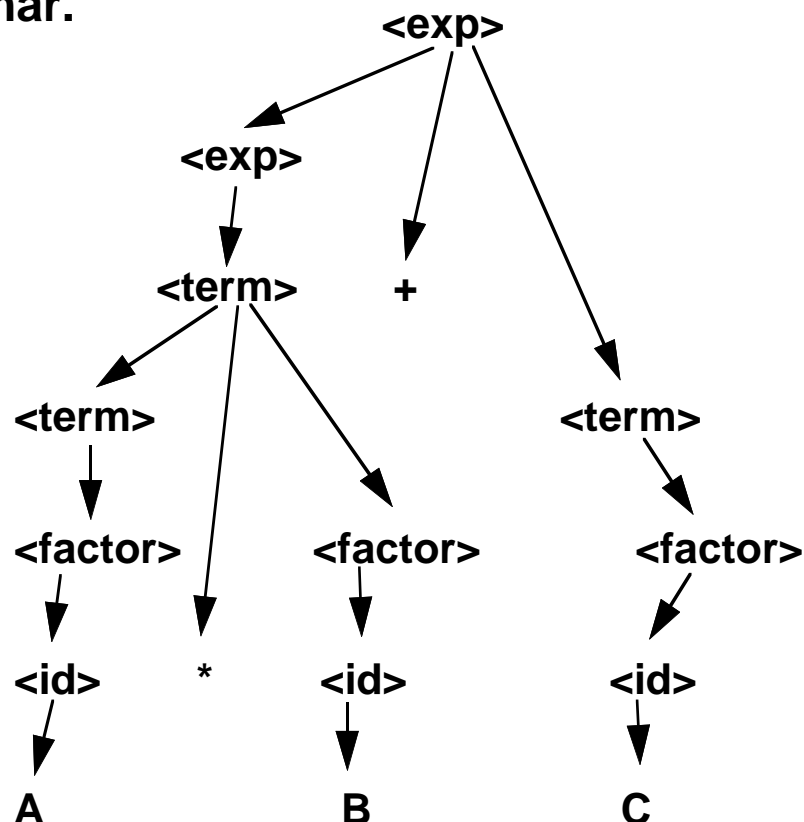
$\langle \text{exp} \rangle \rightarrow \langle \text{exp} \rangle + \langle \text{term} \rangle \mid \langle \text{term} \rangle$

$\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle * \langle \text{factor} \rangle \mid \langle \text{factor} \rangle$

$\langle \text{factor} \rangle \rightarrow (\langle \text{exp} \rangle) \mid \langle \text{id} \rangle$

$\langle \text{id} \rangle \rightarrow A \mid B \mid C$

- This forces  $*$  to have a higher precedence than  $+$ , even though it generates the same language as the first grammar.



# Chomsky Hierarchy of Syntax

---

- **0: Unrestricted grammar rules**

→  $\text{string} \rightarrow \text{string}$

- **1: Context Sensitive**

→  $ABC \rightarrow A\gamma C$

- **2: Context Free**

→  $A \rightarrow \text{string}$

- **3: Linear**

→  $A \rightarrow aB \mid b$

**where**

**N - set of non-terminals**

**T - set of terminals**

$A, B, C, \gamma \in N$

$a, b \in T$

$\text{string} \in (N \cup T)^*$

# Chomsky Hierarchy of Syntax

---

- **0: Unrestricted grammar**
  - $\text{string} \rightarrow \text{string}$  (arbitrary word processing change)
- **1: Context Sensitive**
  - $ABC \rightarrow A\gamma C$  (change B to  $\gamma$  in the context of A..C)
- **2: Context Free  $\rightarrow$  syntactic analyzer**
  - $A \rightarrow \text{string}$  (change A to string anywhere)
- **3: Linear  $\rightarrow$  lexical analyzer**
  - $A \rightarrow aB \mid b$  (there is a one-for-one replacement in N)

# Chomsky Hierarchy of Syntax

---

$$\mathbf{G}_0 \supset \mathbf{G}_1 \supset \mathbf{G}_2 \supset \mathbf{G}_3$$

**$G_3$  - linear, one-dimensional**

**$G_2$  - two-dimensional**

**$G_1$  - three-dimensional**

# Another Ambiguous Grammar

---

$\langle \text{stmt} \rangle \rightarrow \langle \text{assign} \rangle \mid \langle \text{if\_stmt} \rangle$

$\langle \text{assign} \rangle \rightarrow \langle \text{id} \rangle := \langle \text{exp} \rangle$

$\langle \text{if\_stmt} \rangle \rightarrow \text{IF } \langle \text{bool} \rangle \text{ THEN } \langle \text{stmt} \rangle \mid$

$\text{IF } \langle \text{bool} \rangle \text{ THEN } \langle \text{stmt} \rangle \text{ ELSE } \langle \text{stmt} \rangle$

- **Exercise:**

- **Prove that this is ambiguous.**
- **Write a grammar for the same language that is not ambiguous. (Use any published reference material)**

Note: This grammar (syntax) develops a language that demonstrates the “DANGLING ELSE” problem that occurred in ALGOL 60.

# Limitations of Context Free Grammars

---

- **Productions must always apply, regardless of context in which string appears.**

- **Can't handle some things:**

var x : integer;

y : boolean;

begin

x := 3;

OK

x := y;

types wrong

z := 5;

var z undeclared

- **Need "static semantics" . . .**

# Recognizers

---

- **How to generate a recognizer from a grammar?**
  - automatically (YACC)
  - by hand
- **There are many types of parsers:**
  - LL(0)
  - LL(k)
  - LR
  - LALR
  - recursive descent

# Extended BNF & Regular Expressions

---

## Bracket notations

- **[. . .] optional**

$\langle \text{if} \rangle \rightarrow \text{IF } \langle \text{bool} \rangle \text{ THEN } \langle \text{stmt} \rangle [\text{ELSE } \langle \text{stmt} \rangle]$

- **{. . .} zero or more times**

$\langle \text{ones} \rangle \rightarrow 1 \{ \langle \text{ones} \rangle \}$

- **(. . . | . . .) or (: : :) local choice**

$\langle \text{exp} \rangle \rightarrow \langle \text{id} \rangle | \langle \text{exp} \rangle (+ | * ) \langle \text{exp} \rangle$

## Regular Expressions

- $a^*$  equivalent to  $\{a\}$  - zero or more times
- $a^+$  equivalent to  $a\{a\}$  - one or more times

# Syntax and Parsing Summary

---

- **Recognizer vs. generator**
- **BNF**
  - **Four components**
  - **Recursion**
- **Derivation**
- **Ambiguous grammars**
- **Extended BNF**

# Semantics of Programming Languages

---

- **How to define the meaning of programs?**
- **Three approaches:**
  - **Operational**
  - **Axiomatic**
  - **Denotational**

# Operational Semantics

---

- Gives a program's meaning in terms of its implementation on a real or virtual machine
- Define two parts:
- machine
  - high level
  - low level
- translation from source code to "machine" code
- Let  $S$  be the *state*, having components mem, input, and output --- where mem consists of register values, memory values, etc.

$S \{ \text{instruction} \} S'$

$\text{Diff}(S, S') \Rightarrow \text{"meaning of instruction"}$

# Example

Pascal	Operational Semantics
<b>for i := x to y do</b> <b>begin</b> . . . <b>end</b>	<b>i := x</b> <b>loop: if i &gt; y goto out</b> . . . <b>i := i + 1</b> <b>goto loop</b> <b>out: ...</b>

- **Operational semantics could be much lower level, e.g.,**

```
    mov i,r1
    mov y,r2
    jmpifless(r2,r1,out)
    ...
out: ...
```

# Advantages and Disadvantages of Operational Semantics

---

- **Advantages:**
  - May be simple, intuitive for small examples
  - Useful for implementation
- **Disadvantages**
  - Very complex for large programs
  - Lacks mathematical rigor
- **Uses:**
  - Vienna Definition Language (VDL) used to define PL/I (Wegner 1972)
  - Compiler work

# Axiomatic Semantics

---

- Based on predicate calculus. Use assertions to show certain properties of programs.

$\{P\}$  statement  $\{Q\}$

- Compute precondition from postcondition:

$\{P\} \quad x := y + 1 \quad \{x > 0\}$

→ Possible Ps:

$y > 5$

$y = 37$

$y \geq 0$

etc.

**Weakest Precondition (WP)**



- *WP* → identifies *all possible* cases for which postcondition holds!

**WP  $\Rightarrow$  Least Restrictive**

# Finding the Weakest Precondition

---

- **Define function:**

wp: Stmt x Postcondition --> weakest precondition

stmt                  post condition  
    ↓                  ↓  
wp (x := e, P) =  $P_{x \rightarrow e}$  "substitute e for every x in P"

- **So:**

wp (x := y+1, x > 0)

=  $x > 0_{x \rightarrow y+1}$

=  $y+1 > 0$

=  $y \geq 0$

→ **basically, “undoing” the assignment and solving for y**

**y > -1 is WP, too**

**y > -1  $\wedge$  y  $\geq$  0 when y is an integer**

# Sequences of Statements

---

**$\{P\} \text{ S1; S2 } \{Q\}$**

- **Just apply wp twice**

$\text{wp } (x := y + 1; z := x + y, z > 5)$

$\text{wp } (z := x + y, P1)$   
 $= z > 5 \text{ }_{z \rightarrow x+y}$   
 $= x + y > 5$

What can we say about $x+y$ knowing that after the assignment is made $z > 5$ ?
---

$\text{wp } (x := y + 1, x + y > 5)$   
 $= x + y > 5 \text{ }_{x \rightarrow y+1}$   
 $= y + 1 + y > 5$   
 $= y > 2$

# Loops

---

- $\{P\}$  while B do S end  $\{Q\}$
- Need *loop invariant* I such that:
  - $P \implies I$
  - $\{I\} B \{I\}$
  - $\{I \ \& \ B\} S \{I\}$
  - $(I \ \& \ (\text{not } B)) \implies Q$
  - and the loop terminates

# Finding Loop Invariants

---

- **Work backwards through a few iterations and look for a pattern.**

**while  $y < x$  do  $y := y + 1$     $\{y = x\}$**

**$\text{wp } (y := y + 1, \{y = x\}) = \{y = x\}_{y \rightarrow y + 1}$**

**$= y = x - 1$**

**-- last time through**

**$\text{wp } (y := y + 1, \{y = x - 1\}) = \{y = x - 1\}_{y \rightarrow y + 1}$**

**$= y = x - 2$**

**-- next to last time**

**Now, by extension we get  $I = \{y < x\}$**

**When we factor in that the loop may not be executed even once (implying wp could be  $y = x$ ) then we get**

**$I = \{y \leq x\}$**

- **This also satisfies loop termination, so**

**$P = I = \{y \leq x\}$**

## Finding Loop Invariants (cont.)

---

- Does  $\{y \leq x\}$  satisfy conditions (1) - (5)?

(1)  $y = x \implies y \leq x$  ?

(2) if  $y \leq x$  and  $y \neq x$  is then evaluated, is it still the case that  $y \leq x$  ?

(3) if  $y \leq x$  and  $y \neq x$  are both true and then  $y := y+1$  is executed, then is  $y \leq x$  ?

(4) does  $y \leq x$  and  $y = x \implies y = x$  ?

(5) Can you argue convincingly that the program segment terminates?

# Finding Loop Invariants (a harder example)

---

- However, finding the WP is not always so easy!  
Consider:

**$\{P\}$  while  $y < x + 1$  do  $y := y + 1$   $\{y > 5\}$**

$$y > 5 \text{ }_{y \rightarrow y+1} \Rightarrow y > 4$$

$$y > 4 \text{ }_{y \rightarrow y+1} \Rightarrow y > 3$$

etc.

- really tells us *nothing* relative to  $x$  because  $x$  is not in  $Q \equiv \{y > 5\}$

- Try Using Boolean

$$I \ \& \ (\text{not } B) \Rightarrow Q$$

$$? \ \& \ y \geq x + 1 \Rightarrow y > 5$$

$$? \ \& \ y > x \Rightarrow y > 5$$

any  $x \geq 5$  satisfies implication

so . . . let  $I \equiv x \geq 5$

- Do the 4 Axioms hold?

# Advantages, Disadvantages, and Uses of Axiomatic Semantics

---

- **Advantages**

- Can be very abstract
- May be useful in proofs of correctness
- Solid theoretical foundations

- **Disadvantages**

- Predicate transformers are hard to define
- Hard to give complete meaning
- Does not suggest implementation

- **Uses of Axiomatic Semantics**

- Semantics of Pascal
- Reasoning about correctness

# HOMEWORK FOR AXIOMATIC SEMANTICS

---

Consider

$$\{P\} \ X := X * 3 \ \{X^2 = 36\}$$

Determine Weakest Precondition P

# Denotational Semantics

---

- **Define a function that maps a program (a syntactic object) to its meaning (a semantic object).**
- **Sort of like a high-level operational semantics, except**
  - machine is gone
  - language is  $\lambda$ -calculus
- **More abstract.**

# Example: Decimal Numbers

---

**Valuation function:  $V: \text{Number} \rightarrow \text{Integers}$**

↑  
syntax

↑  
meaning

- Syntax:**

$\langle \text{num} \rangle \rightarrow \langle \text{num} \rangle \langle \text{digit} \rangle \mid \langle \text{digit} \rangle$

$\langle \text{digit} \rangle \rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9$

- Semantics:**

→ **Let  $n \in \langle \text{num} \rangle$ ,  $d \in \langle \text{digit} \rangle$**

Applying the  
Function **V**  
to Syntax

$$\left\{ \begin{array}{l} V[\llbracket nd \rrbracket] = 10 * V[\llbracket n \rrbracket] + V[\llbracket d \rrbracket] \\ V[\llbracket 0 \rrbracket] = 0 \\ V[\llbracket 1 \rrbracket] = 1 \end{array} \right.$$

← integers as we know them

- Consider  $V[\llbracket 237 \rrbracket]$  :**

$$\begin{aligned} V[\llbracket 237 \rrbracket] &= 10 * V[\llbracket 23 \rrbracket] + V[\llbracket 7 \rrbracket] \\ &= 10 * (10 * V[\llbracket 2 \rrbracket] + V[\llbracket 3 \rrbracket]) + V[\llbracket 7 \rrbracket] \\ &= 10 * (10 * 2 + 3) + 7 \\ &= 10 * (20 + 3) + 7 \\ &= 10 * (23) + 7 \\ &= 230 + 7 \\ &= 237 \end{aligned}$$

# Expressions

---

- But for *real* programming languages we need more info:

E: Expression --> Integer

$E(\llbracket x \rrbracket) = ?$       where x is a variable

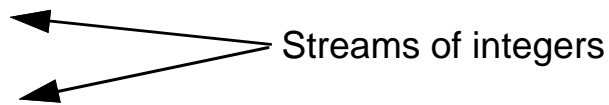
- Depends on the current *state*

→ **STATE** = <mem, input, output>

mem: Identifier --> Integer

input: Integer \*

output: Integer \*



- Now

E: Expression x STATE --> Integer

$E(\llbracket x \rrbracket, s) = \text{mem}(\llbracket x \rrbracket)$  where  $s = \langle \text{mem}, i, o \rangle$

$E(\llbracket e_1 + e_2 \rrbracket, s) = E(\llbracket e_1 \rrbracket, s) + E(\llbracket e_2 \rrbracket, s)$

# Statements

---

- **Expressions denote a value, but statements denote a state.**

ST: (Stmnt x STATE)  $\rightarrow$  STATE

ST (  $\llbracket x := e \rrbracket$ , s ) =  $\langle \text{mem}', i, o' \rangle$  where

s =  $\langle \text{mem}, i, o \rangle$

$\text{mem}' \llbracket x \rrbracket = E(\llbracket e \rrbracket, s)$

$\text{mem}' \llbracket y \rrbracket = \text{mem} \llbracket y \rrbracket$  for all  $y \neq x$

ST (  $\llbracket \text{write}(e) \rrbracket$ , s ) =  $\langle \text{mem}, i, o' \rangle$  where

s =  $\langle \text{mem}, i, o \rangle$

$o' = o \bullet ( E(\llbracket e \rrbracket, s) )$

# Sequences of Statements

---

- **Basic (sequential statement evaluation)**

$$\begin{aligned} \text{ST}(\llbracket \text{stmt}_1; \text{stmt}_2 \rrbracket, s) = \\ \text{ST}(\llbracket \text{stmt}_2 \rrbracket, s') \text{ where} \\ s' = \text{ST}(\llbracket \text{stmt}_1 \rrbracket, s) \end{aligned}$$

- **Parallel statement evaluation**

$$\begin{aligned} \text{ST}(\llbracket \text{stmt}_1; \text{stmt}_2; \text{stmt}_3 \rrbracket, s) = \{s_1, s_2, s_3\} \text{ where} \\ s_1 = \text{ST}(\llbracket \text{stmt}_1; \text{stmt}_2; \text{stmt}_3 \rrbracket, s) \\ s_2 = \text{ST}(\llbracket \text{stmt}_1; \text{stmt}_3; \text{stmt}_2 \rrbracket, s) \\ s_3 = \text{ST}(\llbracket \text{stmt}_3; \text{stmt}_1; \text{stmt}_2 \rrbracket, s) \end{aligned}$$

# Example

---

$$\begin{array}{l} \mathbf{P:} \quad \left\{ \begin{array}{l} \mathbf{x := 5;} \\ \mathbf{P':} \left\{ \begin{array}{l} \mathbf{y := x + 1;} \\ \mathbf{write(x * y);} \end{array} \right. \end{array} \right. \mathbf{P''} \end{array}$$

→ Initial state  $s = \langle \text{mem}, i, o \rangle$

$$\text{ST}(\llbracket P \rrbracket, s) = \text{ST}(\llbracket P' \rrbracket, (\text{ST}(\llbracket x := 5 \rrbracket, s)))$$

$\swarrow \quad \searrow$   
 $s'$

$s' = \langle \text{mem}', i', o' \rangle$  where

$$\text{mem}'(\llbracket x \rrbracket) = 5$$

$$\text{mem}'(\llbracket z \rrbracket) = \text{mem}(\llbracket z \rrbracket) \quad \text{for all } z \neq x$$

$$i' = i, o' = o$$

## Example (continued)

---

$$\rightarrow \text{ST}(\llbracket P' \rrbracket, s') = \text{ST}(\llbracket P' \rrbracket, \underbrace{\text{ST}(\llbracket y := x + 1 \rrbracket, s')}_{s''})$$

$s'' = \langle \text{mem}'', i'', o'' \rangle$  where

$$\text{mem}''(\llbracket y \rrbracket) = E(\llbracket x + 1 \rrbracket, s') = 6$$

$$\text{mem}''(\llbracket z \rrbracket) = \text{mem}'(\llbracket z \rrbracket) \quad \text{for all } z \neq y$$

$$i'' = i', o'' = o'$$

---

$$\rightarrow \text{ST}(\llbracket P'' \rrbracket, s'') = \text{ST}(\llbracket \text{write}(x * y) \rrbracket, s'') = s'''$$

$s''' = \langle \text{mem}''', i''', o''' \rangle$  where

$$\text{mem}''' = \text{mem}'', i''' = i''$$

$$o''' = o'' \cdot E(\llbracket x * y \rrbracket, s'') = o'' \cdot 30$$

---

$\rightarrow$  **So,**

$\text{ST}(\llbracket P \rrbracket, s) = \langle \text{mem}''', i''', o''' \rangle$  where

$$\text{mem}'''(\llbracket y \rrbracket) = 6$$

$$\text{mem}'''(\llbracket x \rrbracket) = 5$$

$$\text{mem}'''(\llbracket z \rrbracket) = \text{mem}(\llbracket z \rrbracket) \quad \text{for all } z \neq x, y$$

$$i''' = i$$

$$o''' = o \cdot 30$$

# Advantages, Disadvantages, and Uses of Denotational Semantics

---

- **Advantages (of denotational semantics)**
  - compact and precise
  - may help with implementation
  - solid mathematical foundations
- **Disadvantages**
  - Hard for programmer to use
- **Uses**
  - Semantics for Algol-60, Pascal, etc.
  - Compiler generation and optimization

# HOMWORK FOR DENOTATIONAL SEMANTICS

---

## Prefatory Consideration:

Prog. Langs. have conditional statements, e.g.

1) if b then stmt1, else stmt2

2) if b then exp1, else exp2

Assuming that conditionals only support expression evaluation and have no side effects, let's give meaning to 1) above:

$ST(\llbracket \text{if } b \text{ then stmt, else stmt2} \rrbracket, s) =$

if  $E(\llbracket b \rrbracket, s)$  then

$ST(\llbracket \text{stmt1} \rrbracket, s)$  else

$ST(\llbracket \text{stmt2} \rrbracket, s).$

**Note: use of previous defns T/F Assessment like in 'C' introduction of "IF THEN ELSE" in denotational language**

# UNDERSTAND/STUDY

---

1) ST ( $\llbracket \text{if } b \text{ then stmt1 else stmt2} \rrbracket, s$ )

definition and elaboration

2) Give denotational semantics for repeat until stmt

REPEAT stmt UNTIL b

You will need conditional statement like that specified above

**HINT:**

on RHS you might find recursive defn

# Summary

---

- **Each form of semantic description has its place:**

- **Operational**

- informal descriptions

- compiler work

- **Axiomatic**

- reasoning about particular properties

- proofs of correctness

- **Denotational**

- formal definitions

- provably correct implementations