

Another way to making an argument smaller for recursion:

$(\text{cdr } el)$ where el is a suitably transformed version of e .

Ex . ssort on Asn. 5

$(\text{ssort } L) \Rightarrow$ a list obtained by sorting the elements of L into ascending order.

list of real numbers e

$(\text{ssort } '(3 7 1 0 8 9 5)) \Rightarrow (0 1 3 5 7 8 9)$

In this case, let $el = (\underbrace{\text{Min-first } e}_{= e \text{ with the first occurrence of the least element moved to the first}})$.

least
 $= e \text{ with the first occurrence of the least element moved to the first}$

$el = (0 \underbrace{3 7 1 8 9 5}_{(\text{cdr } el)})$

Sometimes you should use different recursive strategies for different argument values.

Ex Merge-lists on Asn 5 (Asn = Assigned p, Assignment).

$(\text{merge-lists } L_1 \ L_2) \xrightarrow{\text{two lists of real members in ascending order}}$

a list of members in ascending order obtained by merging L_1 and L_2 .

$(\text{merge-lists } \underset{L_1}{'(2 7 9 12)} \ \underset{L_2}{' (0 1 8 10 15)}) \Rightarrow (0 1 2 7 8 9 10 12 15)$

Strategy 1: compute $(\text{merge-lists } L_1 \ L_2)$ from $(\text{merge-lists } (\text{cdr } L_1) \ L_2)$.

Strategy 2: compute $(\text{merge-lists } L_1 \ L_2)$ from $(\text{merge-lists } L_1 \ (\text{cdr } L_2))$.

Strategy 1 is appropriate of $(\text{car } L_1) < (\text{car } L_2)$

Strategy 2 is appropriate of $(\text{car } L_1) > (\text{car } L_2)$

Either strategy works if $(\text{car } L_1) = (\text{car } L_2)$

Another example unrepeated-elts (from Asn. 5)

$(\text{unrepeated-elts } '(C A D D A B B B A C)) \Rightarrow ((A A A))$

Strategy 1: compute $(\text{unrepeated } L)$ from $(\text{unrepeated elts } (\text{cdr } L))$

Strategy 2: $\sim \sim \sim \sim \sim$ (unrepeated $(\text{caddr } H)$)

CS316 3/17/2014 strategy ②

L = (CA DD ...) strategy 1

L = @ (AADD...) strategy 2

L = (AAA --) strategy 3

Programming paradigm.

Already discussed:

Imperative programming

object-oriented imperative programming

Functional programming

Another Large established paradigm in logic programming.

In logic programming we write collecting of facts facts and inference rules.

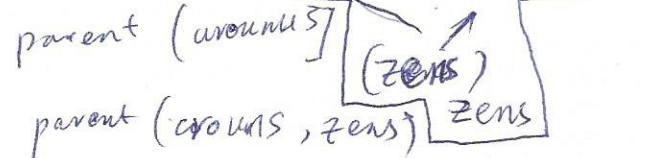
A logic program is used by writing queries which the logic programming system will find answers to by deduction from the facts and inference rules.

Prolog (developed in the early 1970s by colmerauer) is the best known logic programming language.

Simple Example:

Facts

parent (gaia, cronus)



parent (rhea, ~~zeus~~ zeus)

parent (cronus, zeus)

inference rule:

grand parent ()

parent (person1, person2)

i - parent (x, person2)

parent (person1, x).

The above is a very simple logic program (consisting of 4 fact clauses and 1 rule clause).

Question that use their program: ? - parent (gaia, X)

X = cronus; → try to find another answer

false = no more answers.

? - parent (X , zeus)

X = rhea;

X = cronos;

false,

? - grandparent (Y , zeus);

Y = gaia;

Y = uranus;

false;

Syntax of programming languages:

You will receive exercises soon - check email regularly (by Friday).

Syntax = form without regard to meaning

Semantics = meaning

Syntax rules have the property that the compiler can check whether or not they are violated - it is never necessary to execute the program to determine whether a syntax rule is violated.

Analogy with English.

The dog barked loudly.

A valid sentence. = syntactically valid has correct grammatical form
semantically valid ("meaningful")

The night barked loudly.

syntactically valid* — Has the same form/syntax as the above valid sentence
semantically invalid ("meaningless" — because nights cannot bark).

if it can be obtained from that sentence by replacing one name with another.

The dog loudly. → syntactically invalid — don't have correct grammatical form as it has no

For a program in a language such as C++ or Java, a syntactically valid program
wrongly speaking, a piece of text that can be obtained from some valid program by making
changes of the following kinds (zero or more times).

1. Remove a declaration (of a variable, function etc.)
2. Replace one identifier with another
3. Replace a constant in literal (e.g. "dog" or 237) with another (literal of the same kind).

Ex: A syntactically valid C program that is not a valid C++ program

```
int g()
{
    int x=7;
    return x[3]/((x-7)/(3-3));
}
```

This is syntactically valid because it can be obtained from the following valid program by making changes of the above kinds:

Exercise

Reminder

Do the reading assignment and exercises (on the syntax of expression) before Exam!

LispDouble Recursion ("divide-and-conquer")

A doubly recursive function computes its result from two different recursive calls (for at least some argument values).

No. of recursive calls grows exponentially with the depth of recursion. (But in many doubly recursive function the depth of recursion grows only logarithmically with argument size and so the no. of recursive calls does not grow exponentially with argument size).

Ex MSort from Asn 5

(msort $L \xrightarrow{\text{list of real numbers}} \text{ascending order.}$)

\Rightarrow a sorted list of the elements of L .

(msort '(3 7 1 6 5 8 2))

$\Rightarrow (1 2 3 4 5 6 7 8)$.

(msort '(2 7 3 0 4 5 8))

(msort '(2 3 4 8))

(2 3 4 8)

(msort '(7 0 5))

(0 5 7)

(merge-list '(2 3 4 8) '(0 5 7))

(0 2 3 4 5 7 8).

Another picture

(2 7 3 0 4 5 8).

split-list

(2 3 4 8)

msort \rightarrow (2 3 4 8)

(7 0 5)

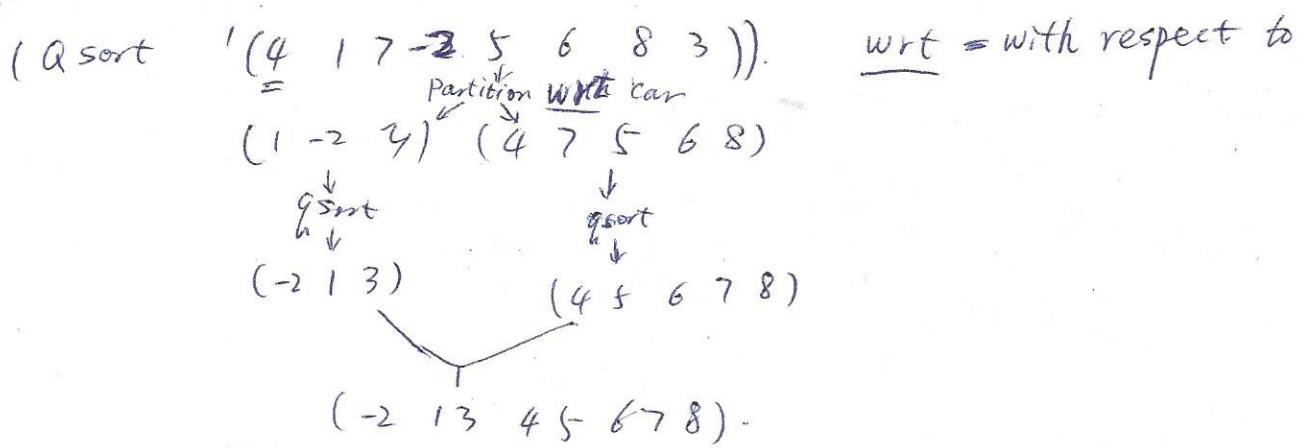
msort \rightarrow (0 5 7)

merge-lists

(0 2 3 4 5 7 8).

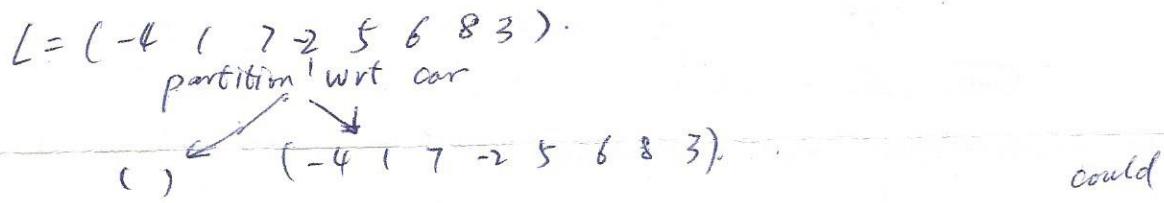
Another example

$(Qsort \ L) \Rightarrow$ same result as $(Msort \ L)$.



when does the above recursive strategy fail?

- ① L is NIL
- or ② First element of L is the first element is the smallest element of the list.



A special case of ② is when L is of length 1, but this ~~can't~~ be handled as a base case. A different recursive strategy is needed in case ②!

Functions that take functions as arguments

Ex: $\sum_{i=2}^4 i^2 = 2^2 + 3^2 + 4^2 = 29$

$(\text{defun} \ \text{sqr} \ (n) \ (\# \ n \ n))$

Ex: $(\text{sigma} \ #'\text{sqr} \ 2 \ 4) \Rightarrow 29$

How can we write function SIGMA that is like Σ ?

in Common Lisp [this does NOT apply to scheme] the value of a symbol as a variable (if any) is not related to the function definition of that symbol (if any)

SETF is one way to set the value

LET/LET+ is another way to set the value.

Most importantly when a function is called, the value of each formal parameter is set to the value of the corresponding actual argument.

DEFUN sets the function definition of a symbol. (DEFUN has no effect on the value of the symbol.)

```
(defun sqr (n) (* n n))
(setf sqr 5)
[> (sqr sqr)
  25]
```

#'f is an expression whose value is the definition function definition of f.
This explain why we wrote (sigma #'sqr 2 5) and NOT (sigma sqr 2 5).

~~Function~~ Funcall

(funcall g a₁ a₂ a₃ ... a_k) .

calls the function given by the value of g with a₁, ..., a_k as the actual arguments.

```
(setf g #'sqr)
```

```
(defun g (n) (* n 17))
```

```
>(g 3) | >(funcall g 3) ← g=sqr. | >(funcall #'g 3)
20   |   q                         | 20
```

In scheme, there is no difference between function definition and value. So FUNCALL and #' are not used in scheme.

Now we are ready to write the SIGMA function.

```
(defun sigma (f m n)
  (if(>m n))
    0
    (+ (funcall f m)
        (sigma f (+ m 1) n))))
```

$$\sum_{i=m}^n f(i) = f(m) + \dots + f(n)$$

$$= \begin{cases} \sum_{i=m+1}^n f(i) & \text{if } m \leq n \\ 0 & \text{if } m > n \end{cases}$$

MAPCAR (map in scheme)

(mapcar #'round '(7.1 6.8 3.2 3.1 5.9))
 $\Rightarrow (7 \ 7 \ 3 \ 3.6)$

(mapcar f L) \Rightarrow a list of same length L in which the i^{th} element is obtained by applying f to the i^{th} element of L.

(mapcar #'exp '(2 4 6 3 9)) $\Rightarrow (4 \ 16 \ 36 \ 81 \ 8)$

Let's write our own version of Mapcar.

(defun our-mapcar (f L))

(= e₁ e₂ ... e_k) our-mapcar f L) should $\Rightarrow (f(e_1), f(e_2), \dots, f(e_k))$

(defun our-mapcar (f L))

(if (endp L))

" cons (funcall f (car L))

(our-mapcar f (cdr L))))

The build-in ~~MAPCART~~ Mapcar (unlike our-MAPCAR) can also be used to map function of 2 or more arguments.

(mapcar #''+ '(7 6 3) '(2 4 5 8) '(0 1 6 1)) $\Rightarrow (3 \ 12 \ 17 \ 12)$

(mapcar #'cons '(A B C) '(a z) '(4 5) '(8)))

$\Rightarrow ((A \ a) (B \ 4) (C \ 5))$

~~today~~ Lambda Expression (also ~~available~~ available in Java, and any other res., Lambda expression evaluate to number nameless function).

Syntax is the same as for DEFUN but without the function name.

(lambda (n) (* n n))

((lambda (n) (* n n)) 5) $\Rightarrow 25$

(sigma (lambda (k) (* k k)) 2 4) $\Rightarrow 29$. $\sum_{k=2}^4 k^2$

Lambda Expression (contd.)

> ((lambda (x y z) (+ x (* y z))) 2 1 7)

≡ (optional)

> (mapcar (lambda (x) (+ x (* x x x)))
'(2 1 0 4))

(10 2 0 64)

> (mapcar (lambda (x y) (+ x (* y z)))

'(1 0 2 3)

'(0 1 1 2)) \Rightarrow (1 2 4 7)

Σ (sigma (lambda (i) (+ (* i i) (* i 3))) 0 3) $\sum_{i=0}^3 (i^2 + 3i)$.

$$\begin{array}{cccc} \Rightarrow 0^2 + 3*0 & + 1^2 + 3*1 & + 2^2 + 3*2 & + 3^2 + 3*3 \\ 0 & 4 & 10 & 18 \end{array}$$

$$\Rightarrow 32$$

Important A lambda expression can use variables defined outside that lambda expression relevant to FOO problem (Q14 on Asn 5)

Ex: Write a function SUM-POWERS such that

$$(\text{sum-powers } k \text{ m n}) = \sum_{i=m}^n i^k = m^k + (m+1)^k + \dots + n^k$$

(defun sum-powers (k m n)

(sigma (lambda (i)(expt i k)) m n))

Ex: INC-LIST-2: from assignment.

(inc-list-2 '(2 7 1 6 5) 3) \Rightarrow (5 10 4 9 8)

(defun inc-list-2 (L n)

(mapcar (lambda (i) (+ i n)) L)).

Apply Apply is like funcall but allows a list of arguments to be passed.
 *last argument must be a list.

(apply #'+ '(2 1 7 9))

= (funcall #'+ 2 1 7 9)

= (+ 2 1 7 9) = 19.

= (apply #'+ 2 '(1 7 9))

= (apply #'+ 2 1 '(7 9))

= (apply #'+ 2 1 7 '(9))

= (apply #'+ 2 1 7 9())

(apply f i₁ i₂ ... i_k) \leftarrow a list

calls f with the values of i₁, i₂ ... i_k and the elements of the values of L as arguments.

(funcall f i₁ i₂ ... i_k)

= (apply f (list i₁ i₂ ... i_k))

= (apply f i₁ (list i₂ ... i_k)) ... etc.

Ex. Sum from Asn-4

[(sum '(2 7 16)) \Rightarrow 16]

(defun sum (L)

 apply #'+ L).

(apply #'cons '((+ 2 3) (+ 4 A))) \Rightarrow ((+ 2 3) + 4 A)
 this is not F here.

Ex: (relevant to 16c).

(apply #'mapcar #'+ '((1 2 3)
 * (4 0 7)
 * (0 1 1)
 * (3 3 2)))

= (mapcar #'+ '((1 2 3) ' (4 0 7) ' (0 1 1) ' (3 3 2))) \Rightarrow (8 6 13)

SyntaxSyntax of Expression.

Expression can be written in many different notation / syntaxes.

Consider the expression : $f(g(h(1,2), f(3,4)), 5)$ (in java/c++ syntax)

This expression can be written in,

① infix notation as : $((1 \cdot h \cdot 2) \cdot g \cdot (3 \cdot f \cdot 4)) \cdot f \cdot 5$

② infix notation, with f and g having equal precedence. and belonging to a left associative precedence class, h having higher precedence than f and g .

$1 \cdot h \cdot 2 \cdot g \cdot (3 \cdot f \cdot 4) \cdot f \cdot 5.$ (if $f = +, g = -, h = \cdot$)

$1 \cdot 2 \cdot - \cdot (3 + 4) + 5$

③ Lisp notation $(f(g(h \cdot 1 \cdot 2) (f \cdot 3 \cdot 4)) \cdot 5)$

④ prefix notation = Lisp notation without parentheses. $[f \cdot g \cdot h \cdot 1 \cdot 2 \cdot (f \cdot 3 \cdot 4) \cdot 5]$

[This assumes you know how many arguments each operator/function takes]

⑤ "anti-lisp" notation — like Lisp but function names appear at the ends of lists (instead of ~~at~~ at the beginings of list).

$((1 \cdot 2 \cdot h) \cdot (3 \cdot 4 \cdot f) \cdot g) \cdot 5 \cdot f)$

⑥ Postfix notation = anti-lisp without parentheses.

$1 \cdot 2 \cdot h \cdot 3 \cdot 4 \cdot f \cdot g \cdot 5 \cdot f.$

Syntax of infix notation

Terminology operator = symbol that denotes a function ~~the~~

The arity of an operator is the no. of arguments taken by the function

Each argument maybe called an operand.

An operator of arity n is called an n -ary operator

binary = 2-ary

unary = 1-ary

ternary = 3-ary.

Syntactically valid infix expressions.

(S_{VIE}s) can be called defined as follows:

e is a S_{VIE} just if one of the following is true:

- ① e is a constant or a variable
- ② e is (e_1) where e_1 is an S_{VIE}.
- ③ e is $e_1 op e_2$ where each of e_1 and e_2 is an S_{VIE} and op is a binary operator
- ④ e is $op e_1$ where e_1 is an S_{VIE} and op is a prefix unary operator
- ⑤ e is $e_1 op$ where e_1 is an S_{VIE} and op is a postfix unary operator.

prefix unary operators
-y
-s
--x
*s

postfix unary
#++

This definition of S_{VIE}'s as a serious weakness. Some syntax decompositions of S_{VIE}s that are suggested by this definition are inconsistent with the desired semantics.

$$\text{Ex: } e = x + y * z - w$$

Can be decomposed using value ③ into

$$\begin{array}{ccc} e_1 * e_2 & | & \text{This is inconsistent with the usual semantics of} \\ // & | & \text{arithmetic expressions.} \\ x+y & z-w & \end{array}$$

After the exam, you will study sec 2.5 which gives another way to specify S_{VIE}s that is unambiguous specification that is consistent with given precedence and associativity rules.

One way to fix the problem is to add to ③, ④, and ⑤, the rule that op should be the operator that is applied last, and to give rules that determine which operator applied last.

$$2.(a) \quad a+b+c = * + a b c$$

$$b) a * (b + c) = + * b c a$$

$$c) a + b + c + d = + * a b + b d$$

$$d) a * b (b+c) * d = * * a + b c d$$

$$g) \left(b/2 + \sqrt{((b/2)*(b/2) - a*c)} \right) / a =$$

$$= \pm / - + 1/b^2 \text{ sign} \pm 1/b^2 / b^2 \neq a/c$$

$$2.2. \text{ a) } a * b + c = ab + c +$$

$$b) a*(b+c) = a \cdot b + a \cdot c.$$

$$c) a * b + c * d = ab + cd +.$$

$$d) a * (b + c) * d = b * c + a * d *$$

$$e) \frac{(b/z + \sqrt{((b/z)*(b/z) - a*c)})}{a}$$

$$= \cancel{b^2} b^2 / b^2 / b^2 + ac - 8\sqrt{t} + a /$$

2.3. a)  b)  c) 

d)

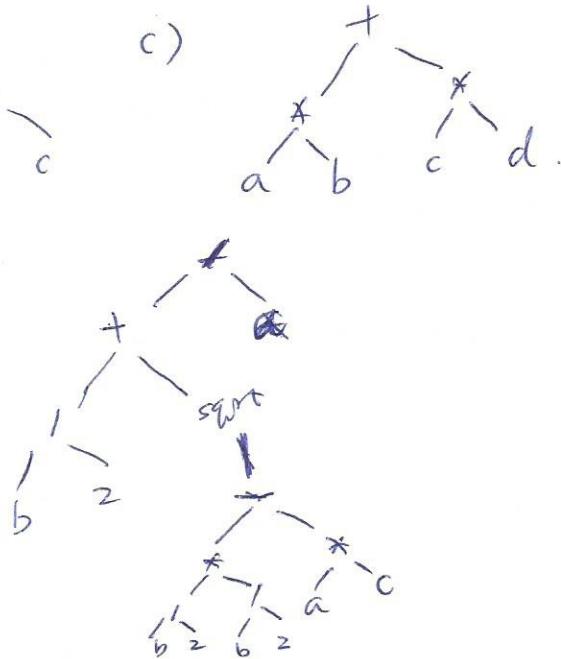
```

graph TD
    d((d*)) --- a((a))
    d --- c((c))
    a --- b((b))
    a --- c
  
```

$$\begin{array}{c}
 2+3 \\
 2 \quad 3 \\
 \text{a) } + \\
 \diagdown \quad \diagup \\
 2 \quad 3 \\
 b) \quad \overset{(2+3)}{+} \\
 \diagdown \quad \diagup \\
 2 \quad 3 \\
 c) \quad \overset{2+3+5}{+} \\
 \diagdown \quad \diagup \quad \diagup \\
 2 \quad 3 \quad 5 \\
 d) \quad \overset{(2+3)+5}{+} \\
 \diagdown \quad \diagup \quad \diagup \\
 2 \quad 3 \quad 5 \\
 e) \quad \overset{2+(3+5)}{+} \\
 \diagdown \quad \diagup \quad \diagup \\
 2 \quad 3 \quad 5
 \end{array}$$

$$2.8. \quad 77423+-$$

$$\begin{array}{c|c|c|c|c|c|c|c|c} 7 & 7 & 4 & 2 & 3 & - \\ \hline 7 & 7 & 49 & 49 & 49 & \\ & 7 & 49 & 49 & 49 & \\ & & 49 & 49 & 49 & \\ & & & 49 & 49 & \\ & & & & 49 & \\ & & & & & 49 & \\ & & & & & & 49 & \\ & & & & & & & 49 \end{array}$$



e is e_1 op e_2
 \downarrow SVIE SVIE

① $a+b+c$ $\xrightarrow{\substack{e_1 \\ \sim \\ e_2}} \text{bad decomposition} \times$
 $\xrightarrow{\substack{\text{op} \\ \text{op}}} a+b+c$
 $\xrightarrow{\substack{e_1 \text{ op} \\ e_2}} \text{good decomposition. } \checkmark$

$a+(b*c)/d$
 $\uparrow \quad \uparrow$
 $\text{op} \quad \text{op}$

decomposition $a+b+c$

How precedence and associativity Rules of an infix Notation are used to determine which operator is Applied Last.

Warning Precedence and associativity rules don't always uniquely determine which operator is applied first.

Ex: $y/z * --y$
 $\uparrow \quad \downarrow \quad \uparrow$
 applied last

Assuming $--$ belongs to a higher precedence class than $*$ and $/$ and that $*$ and $/$ belong to the same precedence class which is a left associative class.

The precedence and associativity rules imply $*$ is applied last, but do not say which operator is applied first.

Some languages may have additional rules that do determine which operator is applied first. Ex: Java has a left to right evaluation rule, which implies that $/$ is applied first in the above expression. [C++ , on the other hand, does not have this rule in C++ it is upto the compiler whether $--$ or $/$ is applied first in the above expr.

int y=4

System.out.println(y/z*--y) prints 6 in java

[would print 3 if $--$ were applied first].

Operators are partitioned into ranked precedence classes.

Each precedence class is specified as being left-associative or right-associative.

To determine which operator in an expression e is applied first:

① Let $e' = e$ without its ~~surrounding~~ parentheses if any.

Ex. $e = (x + (z * (y - 7/w)) * u)$

$e' = x + (z * (y - 7/w)) * u$.
 $\uparrow \quad \uparrow$

top level means not in parentheses.

- ② Find all the top-level operators in e' [top-level means "not surrounded by parentheses"]
 - ③ Among the operator is found at top ②, find the operators of lowest precedence.
 - ④ if just one operator is found at step ③, then that operator is the operator that is applied last. Otherwise, the operator that is applied last is the left most (right most) of the operators found at step ③ / according to whether their precedence class is right - (left-) associative.

	prefix unary	postfix unary	binary	associativity v.
Highest class 1	~			right
class 2	+ -		+ -	left
class 3		& ^ @		right
class 4.			# \$	left

which operator is applied last in left associativity

$$\text{Which operator is applied last in } \text{① } -x + z \# (xy^z) \& a @ z^y \stackrel{?}{=} x \# y - 1$$

$$\textcircled{2} \quad + x @ (y & \sim y \wedge y) + (a @ \sim y \wedge x) \& y \dashv.$$

right associativity

Postfix and Prefix Notations

prefix notation = Lisp notation without parentheses.

postfix notation* = "anti-lisp" notation without parentheses.

* also called reverse Polish notation or RPN.

In these notation you need to know the arity of each operator. You cannot use the same symbol to denote operators of different arities.

$- \infty y$ could mean $(-(-x)y)$ or $(-(-_2 x y))$
 $(-_2 (-x)y)$ $(-_1 (-_2 x y))$
 $\nearrow -_2 -, x y$ $(-_1 -_2 x y)$

We should use different symbols for binary and unary - to avoid this problem.

Differences between Postfix/Prefix Notations and infix Notation

- ① Operator may have any arity in postfix/prefix notation [they may only have arity 1 or 2 in infix notation]
- ② No parentheses in postfix/prefix notation
- ③ No precedence classes are used with postfix/prefix notation.
- ④ In postfix/prefix notation the same symbol should not be used to denote 2 operators of different arities.

Syntactically valid prefix/postfix expression can be defined as follows:

- ① A constant (s) or a variable (w) is a syntactically valid postfix/prefix expression.
- ② If op is an operator of arity k , then.
 - (a) $op \ e_1 \dots e_k$ is a syntactically valid prefix if each of the e 's is such an expr.
 - (b) $e_1 \dots e_k \ op$ is a syntactically valid postfix if each of the e 's is such an expr.

Semantics In case ②, the expr. is evaluated by evaluating each of the e 's and applying op to the values of the e 's

Postfix expression can be evaluated using a stack; ~~Reading the expr. from left to right.~~

- Reading the expr from left to right
- whenever you see a constant or variable, push its value
- whenever you see an operator of arity k , pop k values, apply the operator to those k values and push the result.

When the entire expr has been processed, its value will be the only thing on the stack.

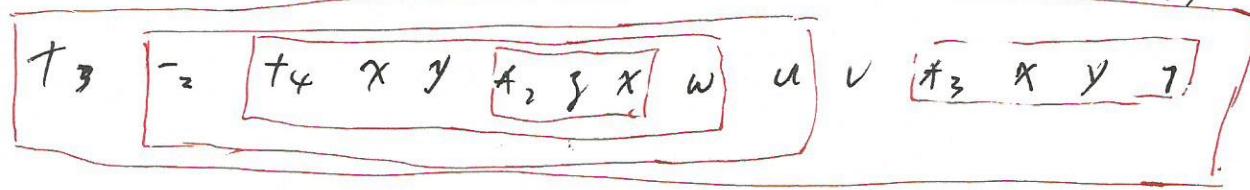
$\xrightarrow{x_{i^{\text{th}}}} \xleftarrow{f^{\text{th}}} \text{value to be popped is the } i^{\text{th}} - \text{last operand.}$

Prefix expression can be evaluated in the same way, except that the expression is read from right to left and when k operands are popped the i^{th} operand to be popped is the i^{th} operand of op (not the $i^{\text{th}} - \text{last}$)

You can use essentially the same algorithms to "put the parentheses back in a postfix/prefix expr. (to produce anti-lisp/lisp notation).

Lisp $(t_3 (-_2 (t_4 x y (+_2 z x) w) u) v (x_3 x y z))$

prefix.



TinyJ Assignment 1 email: Do the installation before next Monday.

* with assigned reading and exercises on syntax and method overriding in C++.

Syntax of Programming Languages.

The syntax of programming language is usually specified in terms of

- ① Units called tokens, such as while, for, IDENTIFIER (x, while, for) UNSIGNED-INT-LITERAL. (72. 47193) ; (>= - . (java/c++ tokens)).

Some tokens (e.g. IDENTIFIER and UNSIGNED-INT-LITERAL in Java/C++) have more than one instance.

Different instances of the same token are semantically different (i.e. have different meanings) but syntax specification written in BNF or some similar notation such as EBNF will not distinguish between them.

Note: It is also possible for the same instance of token to have more than one spelling,

Apple	73.0
APPLE	73.00

A token (like while or for in java/c++) that looks like an identifier but cannot be used as identifiers are called reserved words. [Some languages, such as

Fortran and Lisp, has no reserved words, in such languages words that have special

[`(defun nil (if (+ if 1))` is a legal lisp function definition] meaning (like IF in lisp) are called

keywords Some author refer to reversed words as keywords.

- ② Certain chosen language constructs (e.g. <while stmt>, <stmt> or <expression>). In a syntax specification of an entire language (rather just some part of a language) one of these constructs represents acceptable input for a compiler or interpreter - e.g. it might be called <program>.

The purpose of a syntax specification is to answer the following 2 questions:

- Q1: For each kind of token exactly which sequences of characters constitute a valid instance of that token.

Q₂ For each of the language constructs of ②, exactly which sequences of tokens constitute a syntactically valid instance of that language construct?

[Reminder] A syntactically valid ($\gamma = \text{UX}(Y-W)$) instance of a construct need not be a valid instance — it may contain errors such as ^{understand} ^{undeclared} ~~undeclared~~ identifier type ^{mismatch} or divided by zero.

Q₁ and Q₂ are usually answered using context-free grammar notation (BNF or some variant of this (such as EBNF))

Context-Free grammar notation is a relatively concise (and yet completely precise) notation for specifying a (possibly infinite) set of finite sequences of symbols. The symbols are called the Terminal of the grammar.

In instance answering Q₁ (for tokens that have many instances, like IDENTIFIER or Unsigned-float-literal) we can use a grammar whose ~~terminals~~ terminals are character. [Regular expression notation is another way to do this]

In answering Q₂, we can ~~use~~ use a grammar whose terminal are tokens. It is also possible to use a single grammar, whose ~~terminals~~ terminals are character, to answer both both Q₁ and Q₂.

Context-free grammar notation was intended by Chomsky, who invented a hierarchy of 4 types of grammar for specifying the syntax of English and other similar natural languages. Context-free grammars are type 2 in Chomsky's Hierarchy.

The idea of a context-free grammar was independently reinvented by Backus a few years later. Backus proposed the use of his grammar notation for specifying the syntax of Algol 60.

The Algol 60 Report (which specified the language Algol 60 in a way that was widely admired by computer scientists) adopted Backus's proposal.

Naur was the editor of the Algol 60 report. Backus's notation is now call BNF
(Backus Naur Form)

$\langle E \rangle ::= \langle E \rangle + \langle T \rangle$ (BNF)

$E \rightarrow E + T$. (not BNF)

Like many authors today, we will use BNF to mean some ~~part~~ commonly used notation for writing a context-free grammar.

④ But unlike Sethi:

cf Fig 2.6 on — P42.

and Fig 2.10 on P46.

Sethi would not consider Fig 2.6 to be BNF.
But we will consider it to be BNF.

Ex of a Grammar.

UNSIGNED-FLOAT-LITERAL tokens

can be specified in BNF as follows: (assuming they have the form: one or more digits on one or more digits . one or more digits) (Fig 2.3 on P36.)

① $ufpl ::= ip.f$ $ip ::= ip$ (integer part) $f ::= f$ (fraction) ← fraction.

② $ip ::= d$ $d ::= d / ip$ d. ("d" = "digit").

③ $ip ::= ip d$ ⑤

④ $f ::= d | df$

$d ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$ ⑥

undefined constant symbols

The 11 symbols 0, 1, ..., 9 are the terminals of this grammar.

The 4 symbols $ufpl$, ip , f , d are the nonterminals of this grammar.

Variables, each of which denotes of finite sequences of terminals.

There are 15 rules called productions ~~each~~ each of the form

$N ::= \alpha$ —
a single nonterminal

a sequence of zero or more terminals and/or nonterminals.

→ indicated, the starting nonterminal on the left is the nonterminal on the left side of the first production.

The nonterminal ufpl in the starting nonterminal

This nonterminal denotes the set of sequences that the grammar is intended to define — this set is called the language of / generated by the grammar. Unless otherwise

CS316 ① 4/7/2014.
(Final exam score) + $\frac{25}{40}$

replace the lower of your scores on exam 1 and 2, if it is higher.

parse Trees

- Last time it was mentioned that: Each nonterminal of a grammar represents a set of
grammar = context-free grammar = BNF specification. (in this course).
- finite sequences of terminals.

The set of sequences represented by the starting nonterminal is called the language of / generated by the grammar. [This is regarded as the set of sequences that is defined by the grammar]

① $utpl ::= ip \cdot f$

24. 1293.

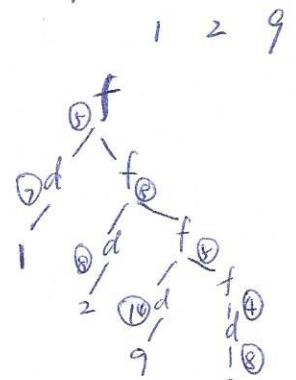
② $ip ::= \overset{\textcircled{1}}{d} / \overset{\textcircled{2}}{ip} \overset{\textcircled{3}}{d}$. "1" means "the left side is the same as in the previous production"
 $f ::= \overset{\textcircled{4}}{d} / df \overset{\textcircled{5}}{f}$
 $d ::= 0 / 1 / 2 / 3 / \dots / 9 \overset{\textcircled{6}}{f}$.

the set of sequence of terminals that is represented by a nonterminal N
is the set of all sequences of terminals that have parse trees rooted at N .

important

If the root of a parse tree is not specified, then it should be assumed
that the root is the starting nonterminal.

Here is a parse tree that proves:



A parse tree rooted at a nonterminal N is an ordered rooted tree with the following properties:

1. The root is N
2. each leaf is either a terminal or the symbol `<empty>`; `<empty>` leaf has no siblings.
3. each ~~internal~~ node (i.e. each internal node) is a nonterminal.

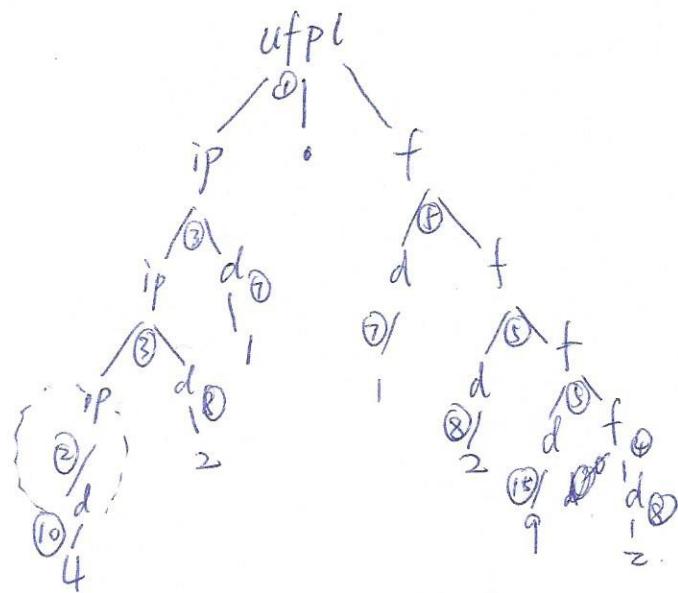
4. the sequence of children of each nonleaf node is the right side of a production whose left side is that node.

A parse tree whose sequence of non-empty leaves is t_1, t_2, \dots, t_k is called a parse tree of or a parse tree that generates t_1, t_2, \dots, t_k .

Here is a parse tree that parses

$$4 \cdot 2 \cdot 1 \cdot 1 \cdot 2 \cdot 9 \cdot 2 \in \text{ufpl}$$

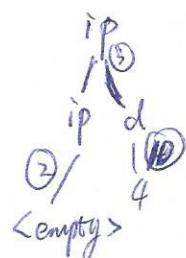
[i.e. $4 \cdot 2 \cdot 1 \cdot 1 \cdot 2 \cdot 9 \cdot 2 \in$ the language of the above grammar].

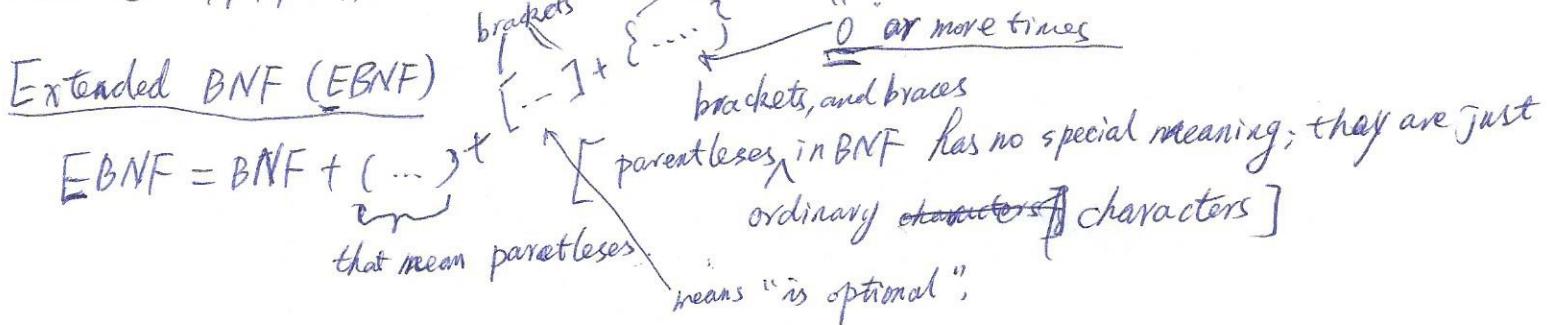


Suppose we change production ② to $\text{ip} ::= \langle \text{empty} \rangle$

Q1. Is the above tree still a correct parse tree? No

A correct parse tree of $4 \cdot 2 \cdot 1 \cdot 1 \cdot 2 \cdot 9 \cdot 2$ should have the following subtree instead of





Ex: $e ::= (+|-)t \quad (+|-)t.$ is a EBNF rule that is equivalent to the following 4 BNF productions

$$e ::= +t + t / -t + t / +t + t / -t - t$$

$e ::= [+|-]t [(+|-)t]$ is equivalent to 9 BNF productions: $[x] = \alpha | \langle \text{empty} \rangle$

$$e ::= t / t + t / t - t / +t / +t + t / +t - t / -t / -t + t / -t - t$$

$e ::= [+|-] + \{ (+|-)t \}$ is equivalent to infinitely many production including :

$$e ::= t + t + t - t - t - t + t.$$

$$e ::= -t - t + t + t + t + t - t - t$$

$$e ::= t.$$

EBNF specifications can always be rewritten as a finite set of BNF productions.

One way to do this is to replace each occurrence of (\dots) , $[\dots]$, and $\{ \dots \}$ with a new nonterminal, which you must define.

$$(\dots) \Rightarrow N \text{ where } N ::= \dots$$

$$[\dots] \Rightarrow N \text{ where } N ::= \dots \mid \langle \text{empty} \rangle$$

$$\{ \dots \} \Rightarrow N \text{ where } N ::= \langle \text{empty} \rangle \mid N \dots$$

$$\{ \alpha \mid \beta \mid \gamma \} \Rightarrow \langle \text{empty} \rangle \mid N(\alpha \mid \beta \mid \gamma) \Rightarrow N ::= \langle \text{Empty} \rangle \mid N\alpha \mid N\beta \mid N\gamma.$$

Work from ^{the} inside outwards

$$e ::= [+|-] + \{ (+|-)t \} \Rightarrow e ::= \underbrace{\text{opt-sign}}_{\substack{\text{option-sign} \\ \text{rest}}} + \underbrace{\text{rest}}_{\substack{\text{op} \\ \text{rest}}} \Rightarrow e ::= \text{opt-sign} + \text{rest}$$

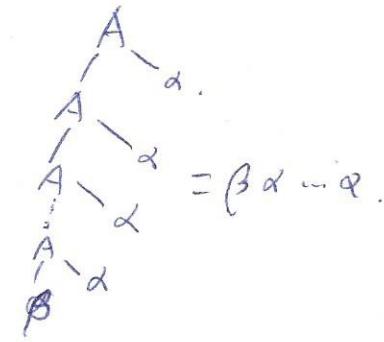
$\text{opt-sign} ::= + \mid - \mid \langle \text{empty} \rangle$

$\text{op} ::= + \mid -$

$\text{rest} ::= \langle \text{empty} \rangle \mid \text{rest op t.}$

Useful Equivalences

$$A ::= A\alpha/\beta \equiv A ::= \beta \{ \alpha \}$$



$$A ::= \alpha A / \beta \equiv A ::= \{\alpha\} \beta$$

Ex : A seq of one or more $\langle \text{stmt} \rangle$'s in which consecutive $\langle \text{stmt} \rangle$'s are separated by a semicolon.

EBNF $\langle \text{stmt-seq} \rangle ::= \langle \text{stmt} \rangle \{ ; \langle \text{stmt} \rangle \}$ (BNF: $A ::= A\alpha/\beta$)

$$\langle \text{stmt-seq} \rangle ::= \langle \text{stmt} \rangle ; \langle \text{stmt} \rangle \mid \langle \text{stmt} \rangle$$

seq of one or more d's

$$ip ::= d \{ d \} \Rightarrow ip ::= ip d / d.$$

~~$f ::= d \{ d \} \Rightarrow f ::= d d / f$~~

$$f ::= \{ d \} d \Rightarrow f ::= d f / d.$$

$$e ::= [+ | -] + \{ (+ | -) + \}$$

$$\Rightarrow e ::= e \overbrace{(+ | -)}^+ + | [+ | -] +$$

$$\Rightarrow e ::= e + t | e - t | + t | - t | t$$

euclid.cs.cuny.edu/~3161

Symbols.java

Each enum constant in symbols.java ~~represents~~ represents a
 → terminal / token.

→ or a non-terminal or a "fictitious token"

method vardecl ()

TJ. outputPrintSymbol (NTvarDecl)

inTreeDept () .

if (getCurrentToken () == (NT) {

nextToken ();

singleVarDecl ();

while (getCurrentToken () == (comma) {

nextToken ();

singleVarDecl ();

}

accept (semicolon);

}

else if (getCurrentToken () == Scanner) {

}

int X; Y = 3; W;

<VarDecl>

Reversed : int

<SingleVarDecl>

<SingleVarDecl>

<SingleVarDecl>

;

CS316 4/23/2014.

④

accept () → {
 accept current token
 nextToken().

nextToken()

There is a case for calling accept (...) even when nextToken() would be adequate:

accept (...) would do a redundant check in such cases, but the advantage is that changes to the code that are implemented later would be less likely to break the code.

TinyJ Virtual Machine (VM)

The memory of the TinyJ VM consists of 2 disjoint regions with separate address spaces:

① Code memory

② Data memory

Code memory is used to store instructions that can be fetched for execution

[For assignment 2, there will be no fetching of instructions for execution, your program will just compile the TinyJ source file into a sequence of TinyJ VM instructions]

Data memory is used to store:

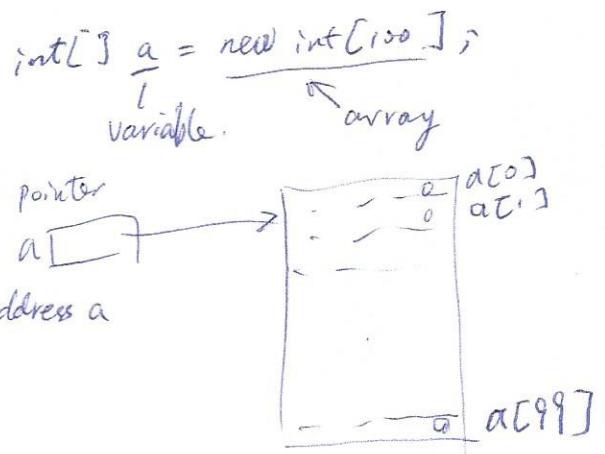
- Values of variables.
- contents of arrays.
- characters of string literals

Data memory address a is unrelated to code memory address a

Data memory consists of 3 ~~regions~~ regions:

a) A statically allocated regions (used to store values of static variables and characters of string literals)

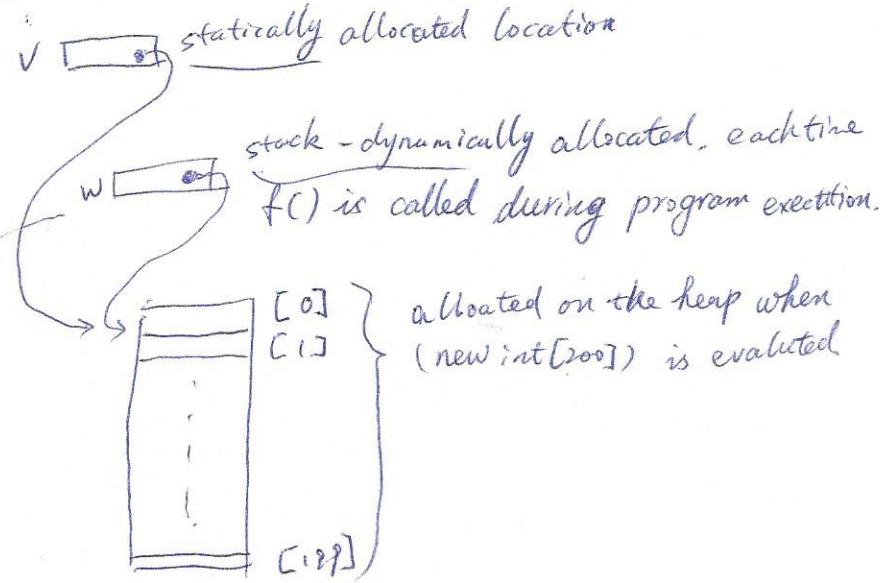
* In TinyJ, (unlike Java) strings are not objects, Arrays are the only objects in TinyJ



b) a stack-dynamically allocated region, used to store values of local variables of methods (including formal parameters of methods)

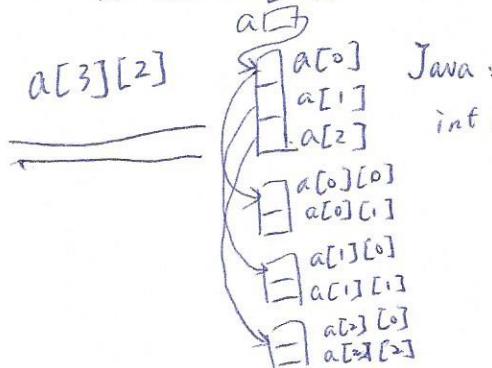
c) A heap-dynamically allocated region, used to store arrays (which are the objects in TinyT)

```
static int *v[ ];  
;  
void f () {  
    int w[ ];  
    w = new int[200];  
    &w = w;
```



class prog {

```
static int a[], b;
static int c;
static int d;
```

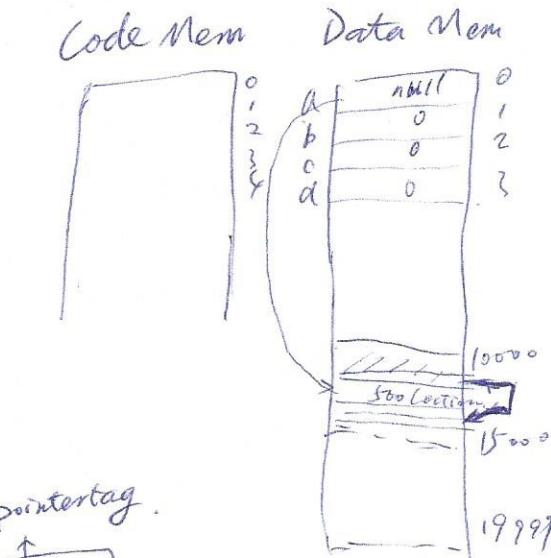
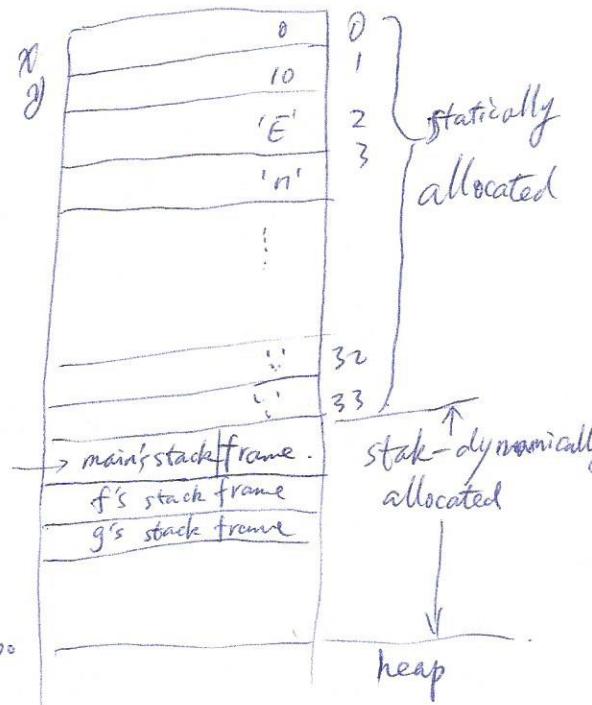
`a = new int[500]`

class E

Address values.

```
static int x, y = 10
```

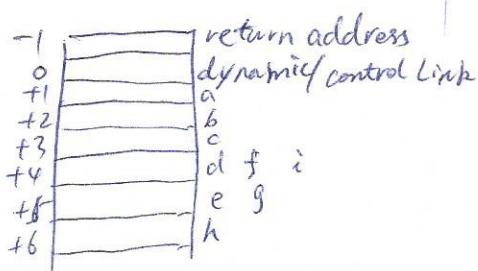
Address values.

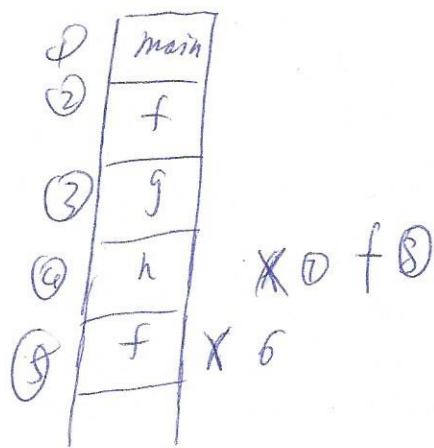
`System.out.println ("345678910...~");`example code
Data Mem

int func()

```
{ int +1, +2, +3  
  int a, b[], c;  
  ...  
  if (...)  
    int +4, +5  
    int d, e;  
  ...  
}  
else  
{ int +4, +5  
  int f, g;  
  int +6;  
}  
int +4 i;
```

main





Exprstack is

Exprstack is the TinyJ VM's Expression Evaluation Stack.

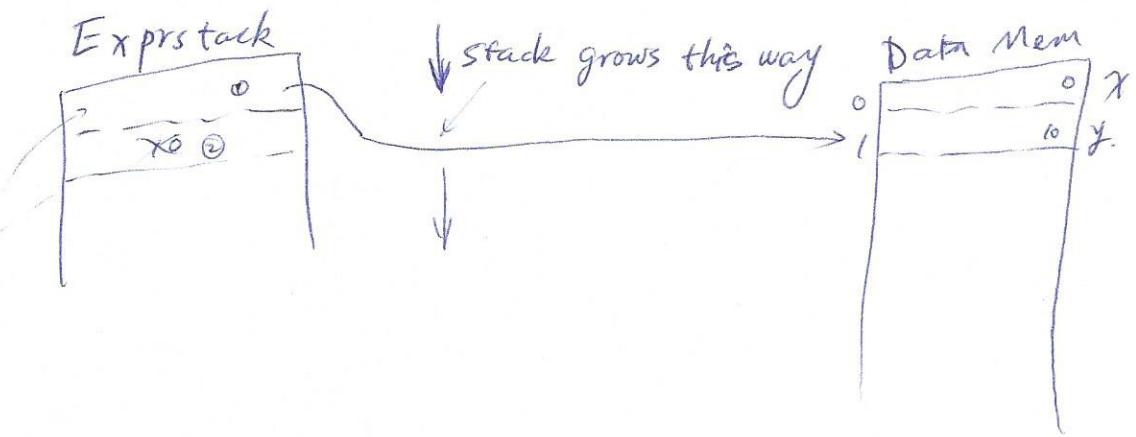
$$y = 10.$$

(y address is 1)

pushstatAddr 1

pushnum 10

savetoaddr



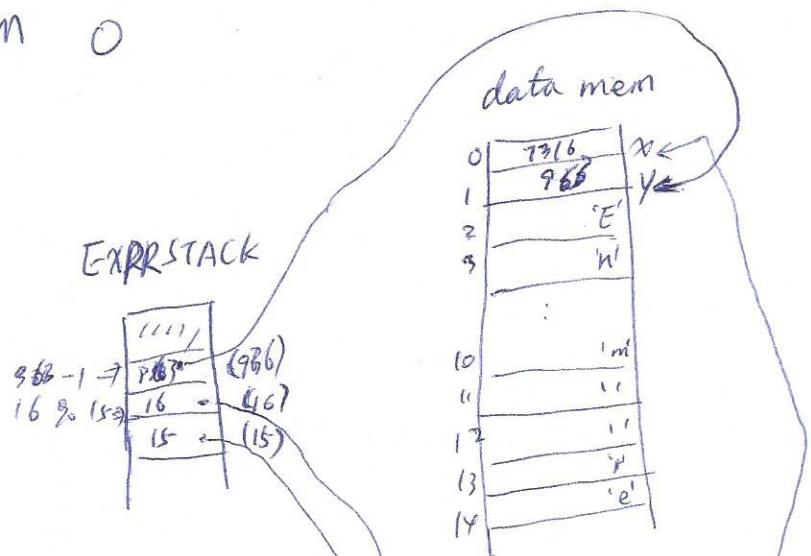
main() is translated to: INITSTKFRM 0

INITSTKFRM 0
WRITESTRING 2 12

PUSHSTATADDR 0

READINT

SAVETOADDR



evaluate expression: $y - a \% u$

② Translation of return $y - a \% u$; in method f()

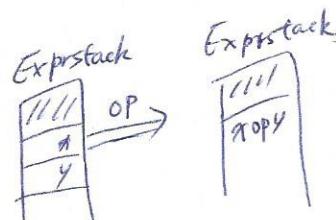
{ instruction that leave value of $y - a \% u$ on top of Exprstack.
return 3.

stack frame of each call of f() is as follows.

To leave value of $y - a \% u$ on top of Exprstack.

f1: | PUSHSTATADDR 1.
LOADFROMADDR
PUSHLOCADDR -4
LOADFROMADDR
PUSHLOCADDR +3
LOADFROMADDR
MOD
SUB

a	16	-4 ←
b	419	-3
c	3	-2
	return address	-1
	dynamic link	0
v	7619	+1
w	24	+2
u	15	+3 ←



OP Maybe: ADD, SUB, DIV, MUL
MOD,
AND, OR
EQ, NE, LE, LT, GE, GT

Translation of $f(21, 22, 23)$

PUSHNUM 21

PASSPARAM

PUSHNUM 22

PASSPARAM

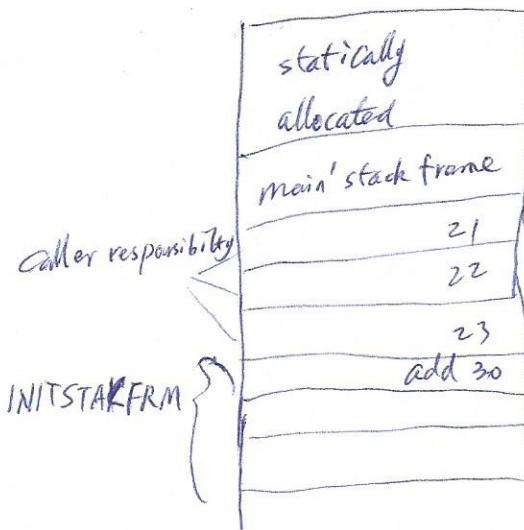
PUSHNUM 23

PASSPARAM

CALLSTATMETHOD 34.

in main()

Data MEM



Translation of

 $y + f(21, 22, 23)$:

PUSHSTATADDR 1

LOADFROMADDR

Code to call $f(21, 22, 23)$

ADD

Translation of

 $f(17, y, (x-y))$;

pushNUM 17

PASSPARAM

y } PUSHSTATADDR 1

LOADSTATADDR

PASSPARAM

x-y } PUSHSTATADDR 0

LOADSTATADDR

PUSHSTATADDR 1

LOADSTATADDR

SUB

PASSPARAM.

callSTATMETHOD 34

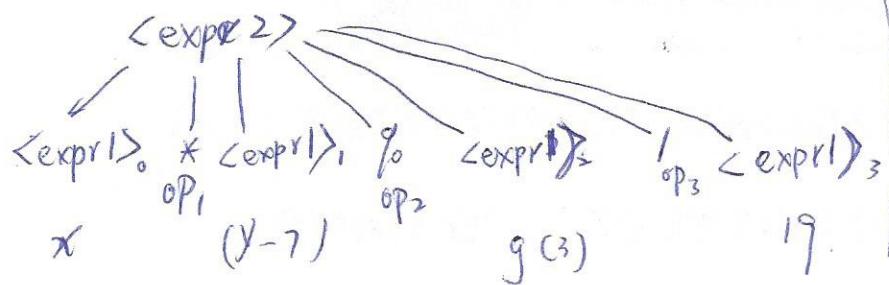
DISCARDVALUE

Notation:

For any node $\langle n \rangle$ in the TinyJ same file parse tree, let $\langle n \rangle$ codemean the VM instruction generated by the compiler for $\langle n \rangle$.Ex: if $\langle n \rangle$ is the $\langle \text{statement} \rangle$ node for $f(17, y, x-y)$ in the above program, then $\langle n \rangle$.code = instructions 8-20 in the generated code

Example of an <expr2>:

$$x * (y - 7) \% g(3) / 19$$



Variable
Assume x is static with address 3
Variable y is local with stack frame offset +4
 g is a method whose code starts at address 29.

In this ~~exp~~ example

<expr2>.code = <expr1>0.code ~~PUSHSTATADDR 3.~~
 LOADFROMADDR

<expr1>1.code ~~PUSHSTATADDR 4.~~
 LOADFROMADDR.
 PUSHNUM 7
 SUB

~~<expr1>i.instr~~
<expr1>2.code
 MUL
 PUSHNUM 3
 PASSPARAM.

CALLSTATMETHOD 29

op2.instr
MOD

<expr1>3.code.
 PUSHNUM 19

op3.instr
~~PUSHDIV~~

new ADDinstr() generates ADD.

new PUSHNUMinstr

new PUSHNUMinstr(23) generates PUSHNUM 23.

new WRITESTRINGinstr(4, 12) generates 4 12.

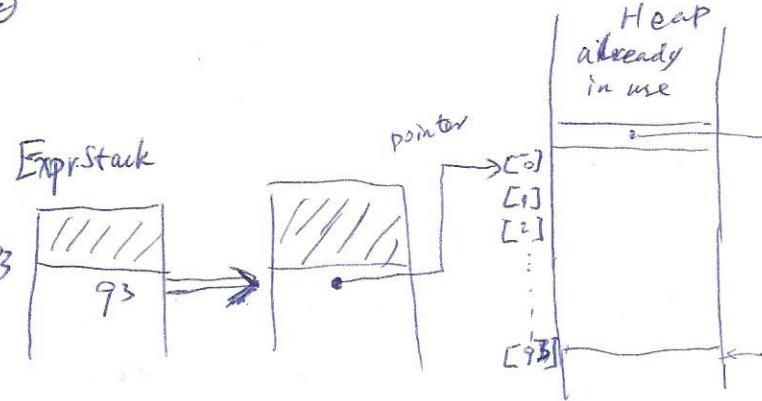
@ in <expr2>.while (ge

See Parse AND Translator.java. Txt. <expr2> changes.

new int[93]

Code generated:

`<expr3>.code = PUSHNUM 93
HEAPALLOC`

while ($x < 4$) {

The method whilestmt() should:

Instruction: `getNextCodeAddress()`

returns the code memory address of the next instruction to be off generated

1. Save `Instruction.getNextCodeAddress()` in an int variable a before calling `expr()`.
2. Generate the JumpOnfalse using something like:

`JUMPONFALSEinstr jFinstr = new JUMPONFALSEinstr(instruction.OPERAND
NOT_YETKNOWN)`

3. After calling `statement()`, generate

`JUMP a`

using the variable a of 1.

4. Fix up the `JUMPONFALSE` using

~~for~~ `jFinstr.fixUpOperand(Instruction.getNextCodeAddress())`

Assignment Or Invoc.:

$x = 19;$ `PUSHSTA ADDR 3` }
 `PUSHNUM 19` }
 `SAVE TO ADDR` }