

Ambiguity is Clearly a Bad Thing

Mathematics is an **unambiguous** language, given an expression:

$$x^{z+2} + 100 + y + z(2 + ab)$$

It has **only one value** (for a collection of variable values)

Ambiguity is Clearly a Bad Thing

Mathematics is an **unambiguous** language, given an expression:

$$x^{z+2} + 100 + y + z(2 + ab)$$

It has **only one value** (for a collection of variable values)

Programming languages may permit **more than one way** to express results or outcomes...

<pre>print("Hello world") SELECT name FROM students; x = x + 1;</pre>	<pre>Python SQL C/C++</pre>	<pre>sys.stdout.write("Hello world\n") SELECT s.name FROM students AS s; ++x;</pre>
---	-------------------------------------	---

There is still only one way to interpret the **intent of the programmer**.

Defining a Language

Here is a simple grammar that defines a language:

$$\begin{aligned} S &\rightarrow A \$ \mid x B x \$ \\ A &\rightarrow s B t \mid w \\ B &\rightarrow q s \mid s q \end{aligned}$$

By Convention...

1. the special symbol \$ means the **end of input**
2. UPPER case terms are **non-terminals**, they can appear on either side of the \rightarrow
3. terms *other than \$ and non-terminals* are called **terminals**, they can appear only on the right-hand side of \rightarrow
4. the vertical bar, |, is read as “OR”

Programs consist of *only* terminals.

Defining a Language

A simple grammar that defines a language:

$$\begin{aligned} S &\rightarrow A \$ \mid x B x \$ \\ A &\rightarrow s B t \mid w \\ B &\rightarrow q s \mid s q \end{aligned}$$

Is the single token w permitted by this grammar?

Defining a Language

A simple grammar that defines a language:

$$S \rightarrow A \$ \mid x B x \$$$

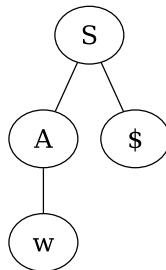
$$A \rightarrow s B t \mid w$$

$$B \rightarrow q s \mid s q$$

Is the single token w permitted by this grammar? **Yes**

$$S \rightarrow A \$$$

$$S \rightarrow w \$$$



← Parse Tree!

Defining a Language

A simple grammar that defines a language:

$$S \rightarrow A \$ \mid x B x \$$$

$$A \rightarrow s B t \mid w$$

$$B \rightarrow q s \mid s q$$

What about $x s q x$?

Defining a Language

A simple grammar that defines a language:

$$S \rightarrow A \$ \mid x B x \$$$

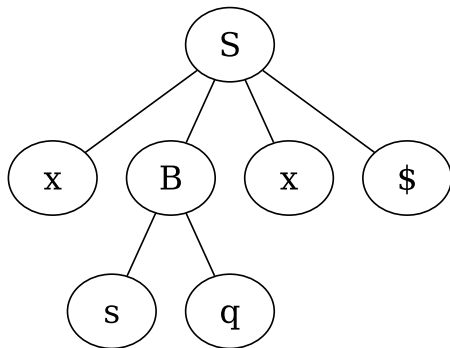
$$A \rightarrow s B t \mid w$$

$$B \rightarrow q s \mid s q$$

What about $x s q x$? **Yes**

$$S \rightarrow x B x \$$$

$$S \rightarrow x s q x \$$$



Defining a Language

A simple grammar that defines a language:

$$\begin{aligned} S &\rightarrow A \$ \mid x B x \$ \\ A &\rightarrow s B t \mid w \\ B &\rightarrow q s \mid s q \end{aligned}$$

Which of these are permitted?

$s s q t$

$w x q s x$

$x q s x$

Defining a Language

A simple grammar that defines a language:

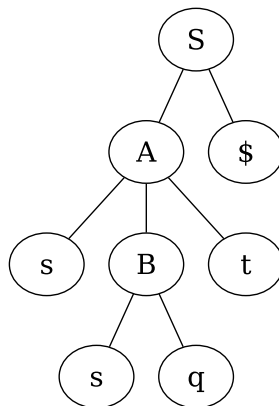
$$S \rightarrow A \$ \mid x B x \$$$

$$A \rightarrow s B t \mid w$$

$$B \rightarrow q s \mid s q$$

Which of these are permitted?

$s s q t$ **Yes**



Defining a Language

A simple grammar that defines a language:

$$\begin{aligned} S &\rightarrow A \$ \mid x B x \$ \\ A &\rightarrow s B t \mid w \\ B &\rightarrow q s \mid s q \end{aligned}$$

Which of these are permitted?

$$w x q s x \quad ?$$

Defining a Language

A simple grammar that defines a language:

$$\begin{aligned} S &\rightarrow A \$ \mid x B x \$ \\ A &\rightarrow s B t \mid w \\ B &\rightarrow q s \mid s q \end{aligned}$$

Which of these are permitted?

$w x q s x$ No

Syntax Error: did not expect x after w .

Defining a Language

$$\begin{aligned} S &\rightarrow A \$ \mid x B x \$ \\ A &\rightarrow s B t \mid w \\ B &\rightarrow q s \mid s q \end{aligned}$$

Which of these are permitted?

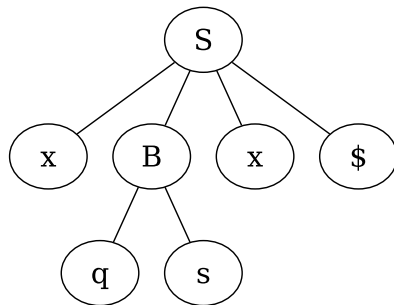
$$x \ q \ s \ x \quad ?$$

Defining a Language

$$\begin{aligned} S &\rightarrow A \$ \mid x B x \$ \\ A &\rightarrow s B t \mid w \\ B &\rightarrow q s \mid s q \end{aligned}$$

Which of these are permitted?

$x q s x$ **Yes**



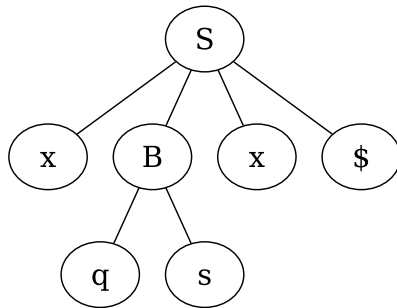
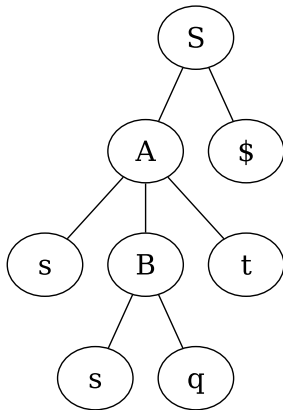
Defining a Language

- a. The parse tree **leaves** are special, what do they hold?
- b. What grammar parts are in the **non-leaf** nodes?

$$S \rightarrow A \$ \mid x B x \$$$

$$A \rightarrow s B t \mid w$$

$$B \rightarrow q s \mid s q$$



Defining a Language

w

s q s t

s s q t

x q s x

x s q x

$$S \rightarrow A \$ \mid x B x \$$$
$$A \rightarrow s B t \mid w$$
$$B \rightarrow q s \mid s q$$

Clearly, there are a limited number of **terminal sequences** permitted by this grammar.

It is **FINITE** — we certainly don't want programming languages with this property.

Recursive Language Definitions

$$\begin{array}{lll} S & \rightarrow & QLIST \$ \\ QLIST & \rightarrow & Q QLIST \\ & | & \lambda \\ Q & \rightarrow & a b c \\ & | & k \\ & | & s t u \end{array}$$

Here is a *recursive* grammar that permits an infinite collection of terminal sequences.

One Last Convention...

- the special symbol λ means an **empty sequence** of tokens — AKA “nothing”

Recursive Language Definitions

$$\begin{array}{ll} S & \rightarrow QLIST \$ \\ QLIST & \rightarrow Q QLIST \\ & | \lambda \\ Q & \rightarrow a b c \\ & | k \\ & | s t u \end{array}$$

Here is a *recursive* grammar that permits an infinite collection of terminal sequences.

What happens with multiple *Q*s?

k a b c k

Recursive Language Definitions

Program
k a b c k

Derivation

parsed by language



<i>S</i>	\rightarrow	<i>QLIST</i> \$
<i>QLIST</i>	\rightarrow	<i>Q</i> <i>QLIST</i>
		λ
<i>Q</i>	\rightarrow	<i>a b c</i>
		<i>k</i>
		<i>s t u</i>

$S \rightarrow QLIST \$$
 $S \rightarrow Q QLIST \$$
 $S \rightarrow k QLIST \$$
 $S \rightarrow k Q QLIST \$$
 $S \rightarrow k a b c QLIST \$$
 $S \rightarrow k a b c Q QLIST \$$
 $S \rightarrow k a b c k QLIST \$$
 $S \rightarrow k a b c k \lambda \$$

Recursive Language Definitions

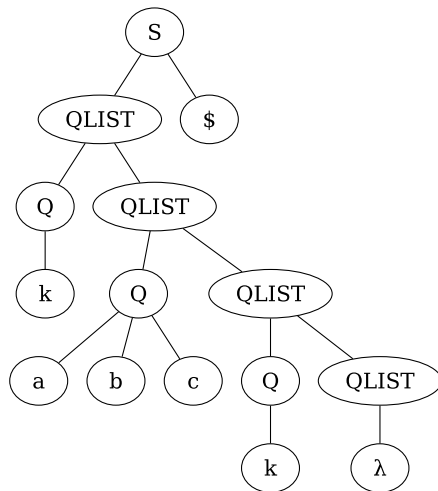
Program
k a b c k

parsed by language

S	\rightarrow	$QLIST\ \$$
$QLIST$	\rightarrow	$Q\ QLIST$
	$ $	λ
Q	\rightarrow	$a\ b\ c$
	$ $	k
	$ $	$s\ t\ u$



Parse Tree



A Simple Programming Language

PROGRAM → *SLIST* \$

SLIST → *S SLIST*
| λ

S → *var = EXPR*
| *if EXPR then (SLIST)*
| *if EXPR then (SLIST) else (SLIST)*
| *while EXPR do (SLIST)*
| *repeat (SLIST) while EXPR*
| *repeat (SLIST) until EXPR*

EXPR → *expr*
| *var*

A Simple Programming Language

PROGRAM → *SLIST* \$
SLIST → *S SLIST*
 | λ
S → *var = EXPR*
 | *if EXPR then (SLIST)*
 | *if EXPR then (SLIST) else (SLIST)*
 | *while EXPR do (SLIST)*
 | *repeat (SLIST) while EXPR*
 | *repeat (SLIST) until EXPR*
EXPR → *expr*
 | *var*

Which grammar rule means this is **not** a finite language?

Program Example A - Compilation Steps

- I. Lexical analysis detects **keywords**, variable names, expressions, and special symbols such as parenthesis and equal.
- II. The sequence of tokens is **parsed** using the grammar rules into a parse tree.
- III. The parse tree is simplified into a **sequence of assignments** with **comparisons and branches (jumps)**.
- IV. Assign memory locations for variables
- V. Assign registers and instructions for expressions
- VI. Generate comparison and jump instructions for control structures
- VII. ... and then we can generate machine instructions!

```
a = 3
c = (27)
b = a*10+1
repeat (
    b = b-a/2+1
    a = a+1
) until b<0
```

```

a = 3
c = (27)
b = a*10+1
repeat (
    b = b-a/2+1
    a = a+1
) until b<0

```

Lexical
Analysis
Describes
Tokens




```

{var,a} = {expr,3}
{var,c} = {expr,(27)}
{var,b} = {expr,a*10+1}
repeat (
    {var,b} = {expr,b-a/2+1}
    {var,a} = {expr,a+1}
) until {expr,b<0}

```

for
Grammar
Analysis



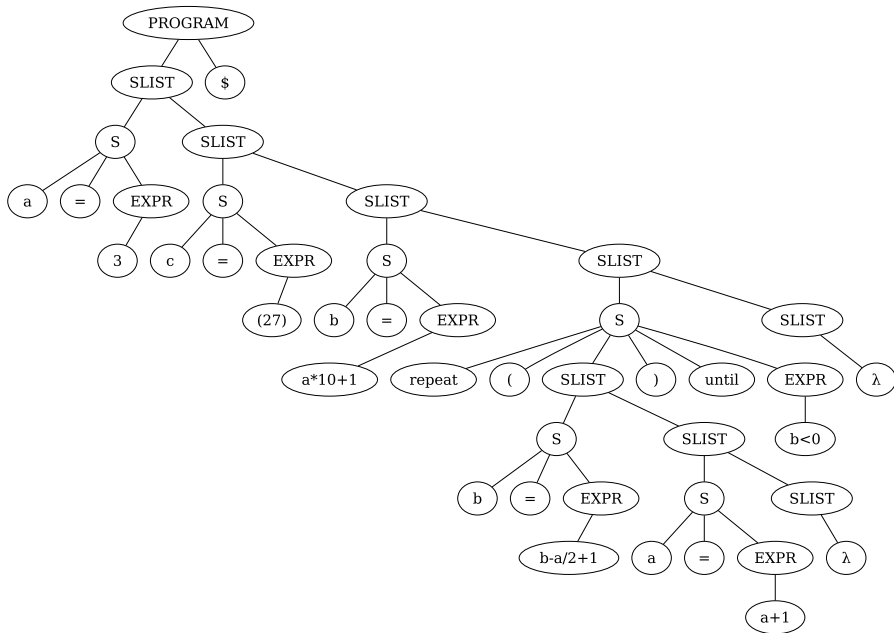
```

PROGRAM → SLIST $
SLIST  → S SLIST
        | λ
S      → var = EXPR
        | if EXPR then ( SLIST )
        | if EXPR then ( SLIST ) else ( SLIST )
        | while EXPR do ( SLIST )
        | repeat ( SLIST ) while EXPR
        | repeat ( SLIST ) until EXPR
EXPR   → expr
        | var

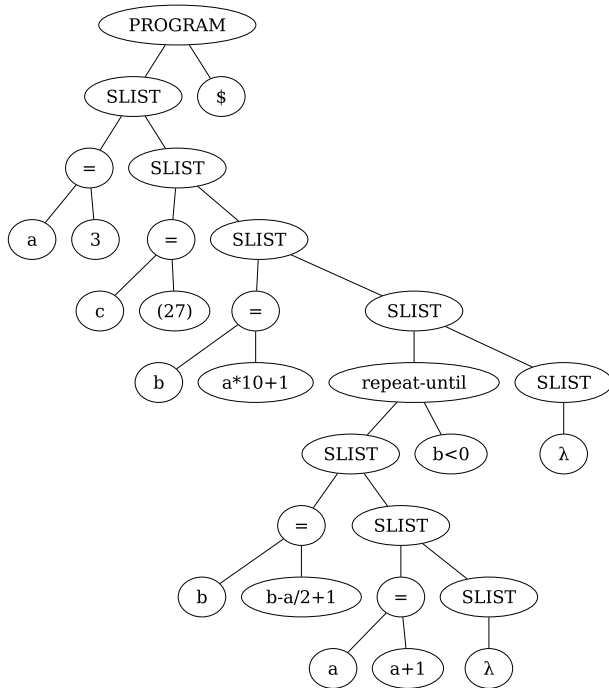
```

...that generates →

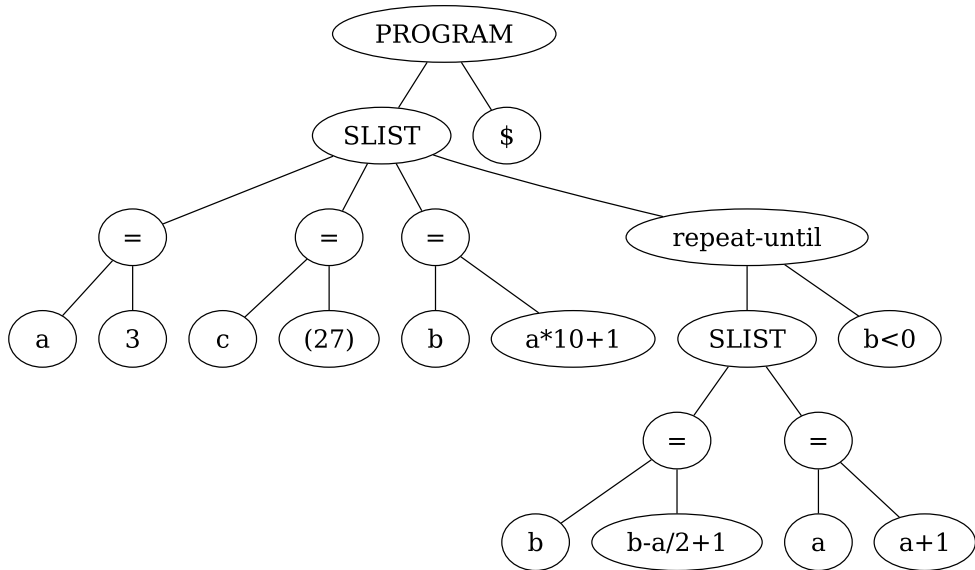
...a Parse Tree



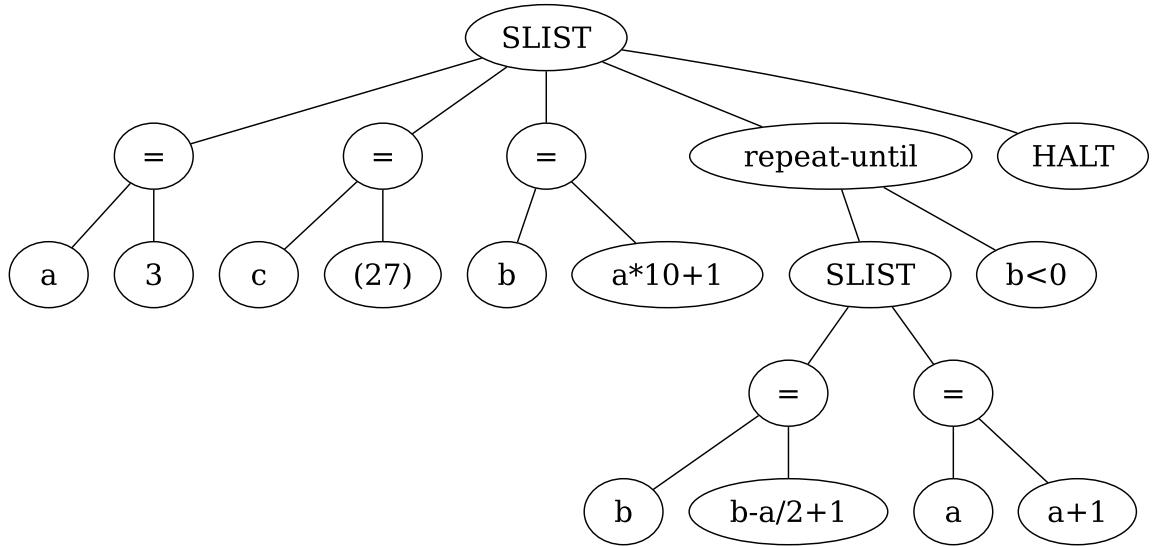
Simplify Statement S Nodes



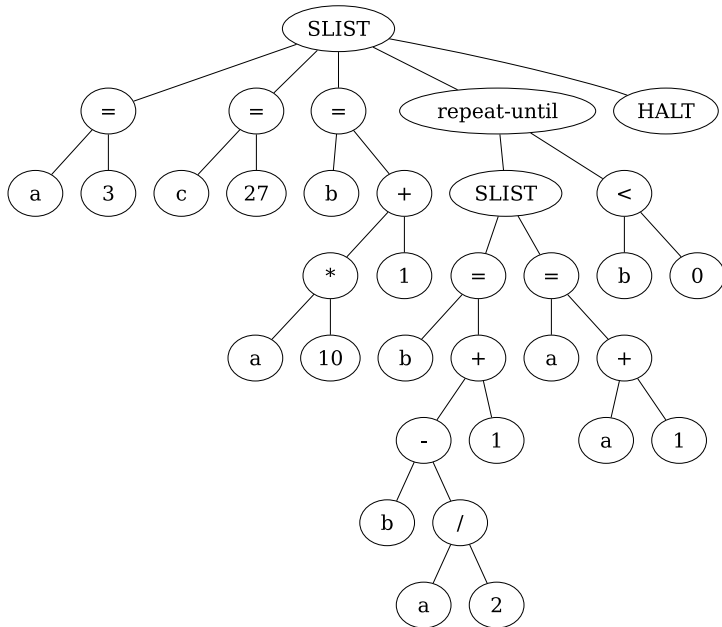
De-curse *SLIST*



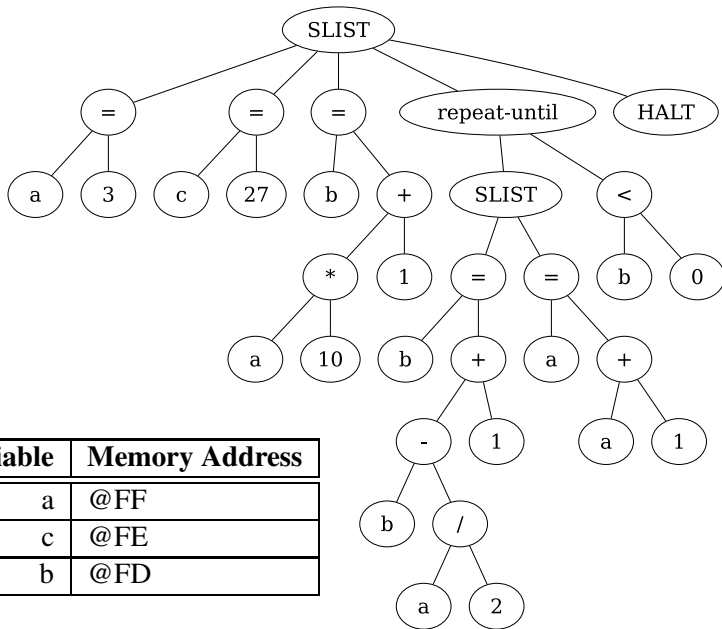
Programs are just an *SLIST* that HALTs



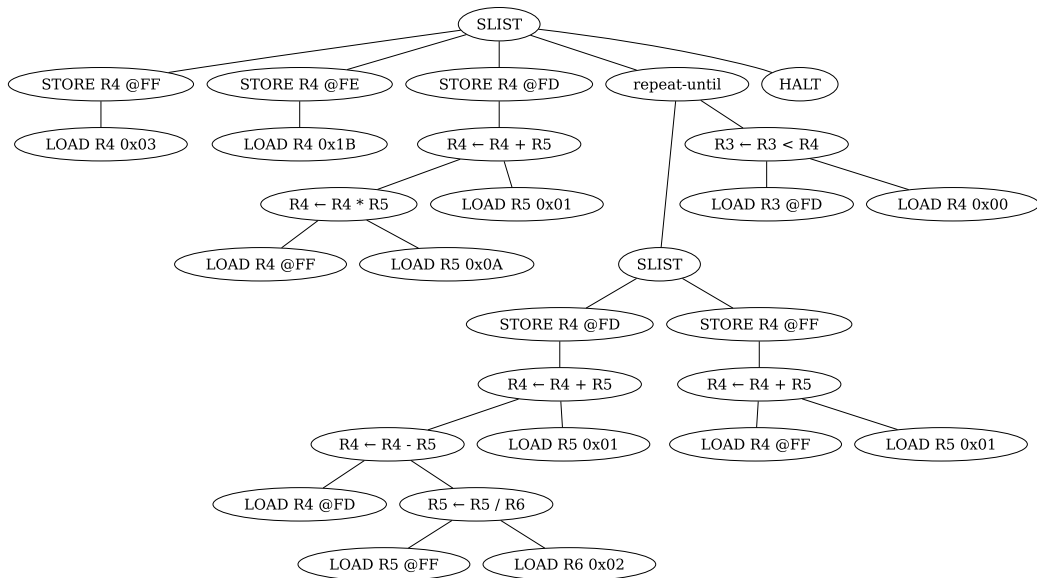
Generate RHS Expression Trees



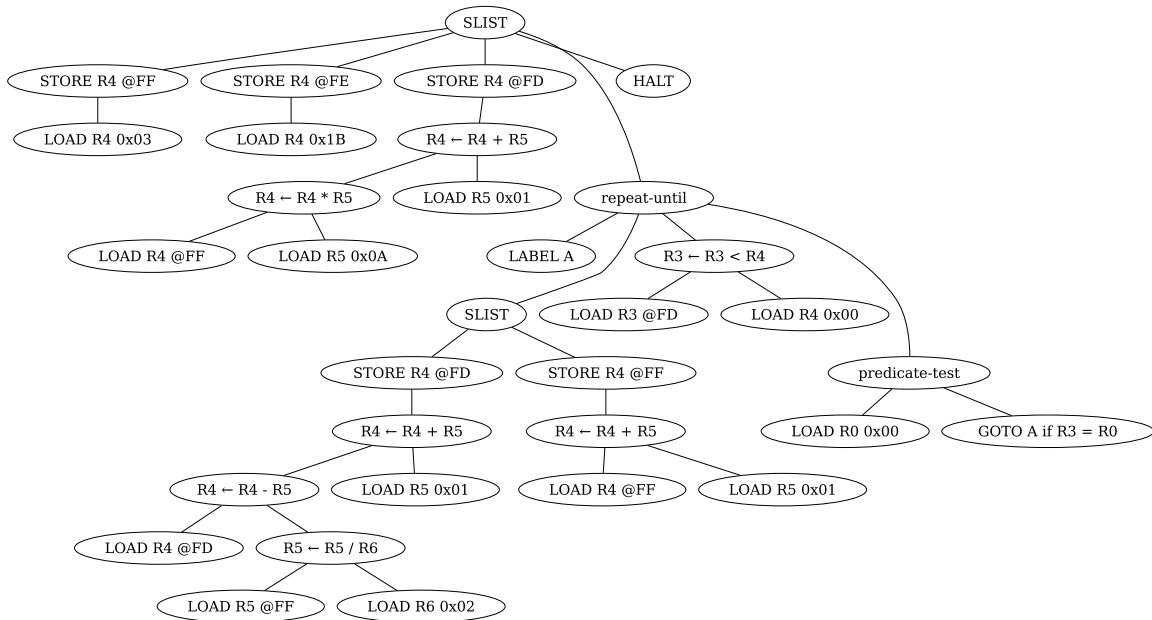
Assign Memory Locations for Variables



Assign Registers and Instructions for Expressions



Generate Instructions for Control Structures



(Pseudo) Assembly Code

```
LOAD R4 0x03
STORE R4 @FF
LOAD R4 0x1B
STORE R4 @FE
LOAD R4 @FF
LOAD R5 0x0A
R4 <- R4 * R5
LOAD R5 0x01
R4 <- R4 + R5
STORE R4 @FD
LABEL A
LOAD R4 @FD
LOAD R5 @FF
LOAD R6 0x02
R5 <- R5 / R6
```

```
R4 <- R4 - R5
LOAD R5 0x01
R4 <- R4 + R5
STORE R4 @FD
LOAD R4 @FF
LOAD R5 0x01
R4 <- R4 + R5
STORE R4 @FF
LOAD R3 @FD
LOAD R4 0x00
R3 <- R3 < R4
LOAD R0 0x00
GOTO A if R3 = R0
HALT
```

An **assembler** will take these instructions and generate the actual machine code, resolving **LABELs** and **GOTOs**, and perhaps performing some low-level optimizations such as removing redundant **LOADs** and **STOREs**.

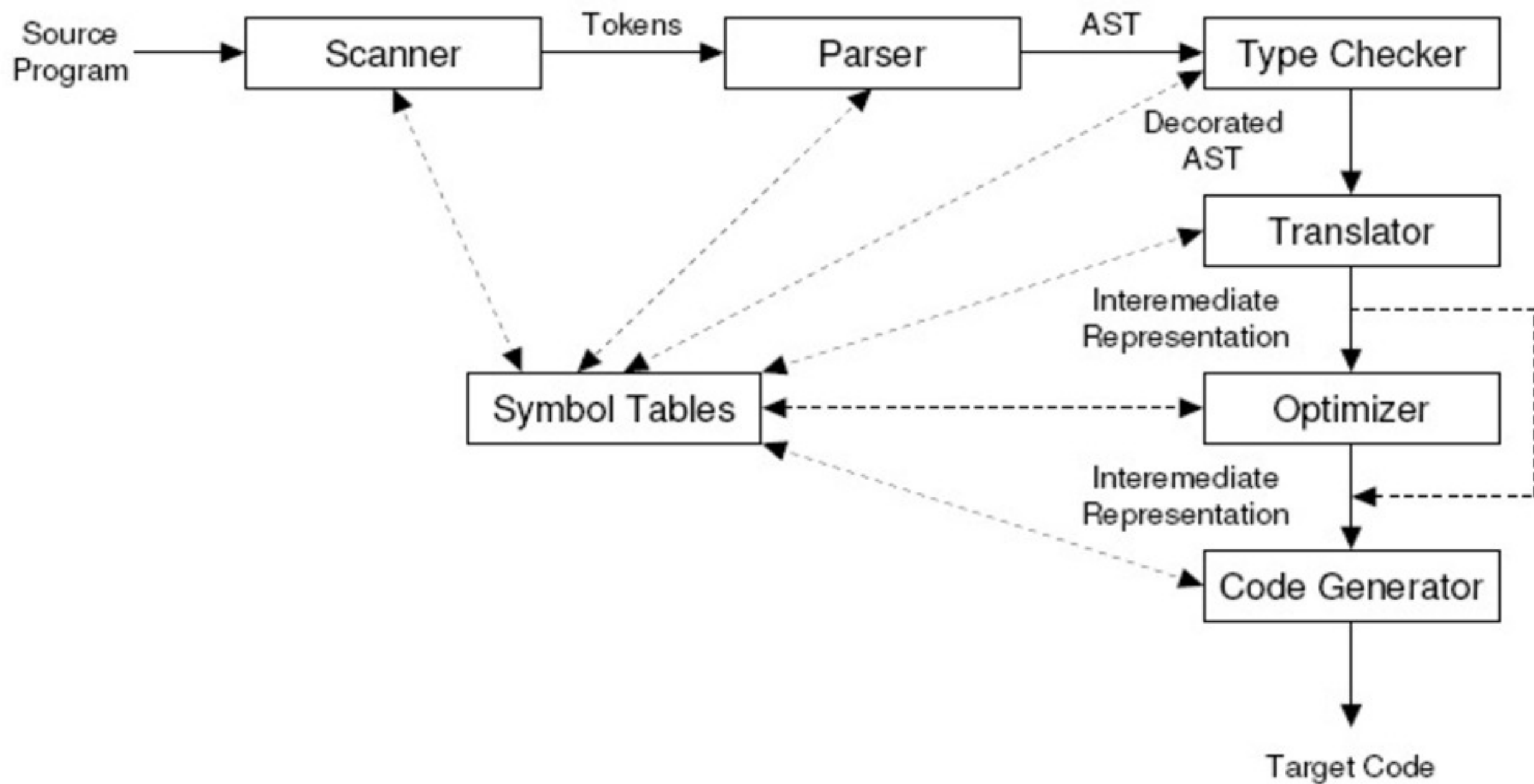


Figure 1.4: A syntax-directed compiler. AST denotes the Abstract Syntax Tree.

Compilers

We've looked at the essential steps taken to generate machine code from a high-level programming language.

Whether you use an **interpreted language** within a “virtual machine” such as Python, Ruby, Lisp, Java, ... or a language compiled “down to machine code” such as C, C++, or Fortran — you use a compiler of some sort, and that compiler pretty much follows all of these steps in order to (eventually) execute *your instructions* on a CPU.

The design of computer languages and the parsing algorithms associated with them is one of the classic and fundamental topics in Computer Science.

Program Example B - Compilation Steps

- I. Lexical analysis detects **keywords**, variable names, expressions, and special symbols such as parenthesis and equal.
- II. The sequence of tokens is **parsed** using the grammar rules into a parse tree.
- III. The parse tree is simplified into a **sequence of assignments with comparisons and branches (jumps)**.
- IV. Assign memory locations for variables
- V. Assign registers and instructions for expressions
- VI. Generate comparison and jump instructions for control structures
- VII. ... and then we can generate machine instructions!

```
x = (a+b)/23
if x<3 then (
    s = 1
) else (
    t = 2
)
```

```

x = (a+b)/23
if x<3 then (
    s = 1
) else (
    t = 2
)

```

Lexical
Analysis
Describes
Tokens

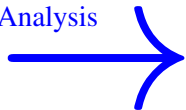


```

{var,x} = {expr,(a+b)/23}
if {expr,x<3} then (
    {var,s} = {expr,1}
) else (
    {var,t} = {expr,2}
)

```

for
Grammar
Analysis



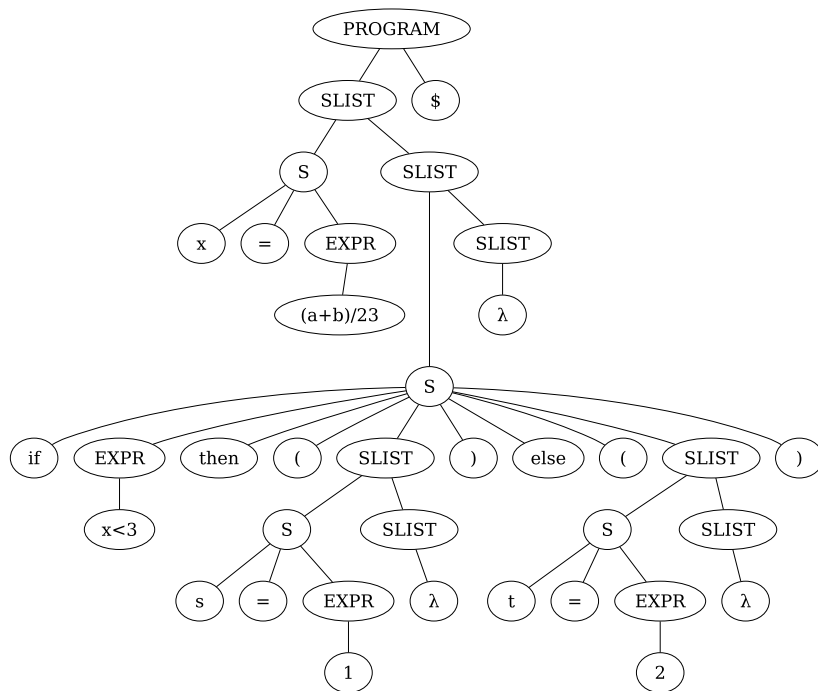
```

PROGRAM → SLIST $
SLIST → S SLIST
        | λ
S → var = EXPR
   | if EXPR then ( SLIST )
   | if EXPR then ( SLIST ) else ( SLIST )
   | while EXPR do ( SLIST )
   | repeat ( SLIST ) while EXPR
   | repeat ( SLIST ) until EXPR
EXPR → expr
      | var

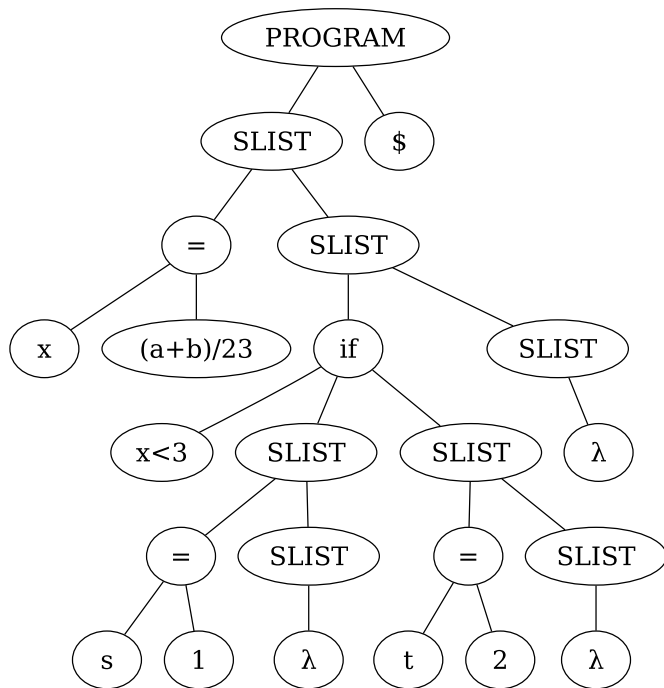
```

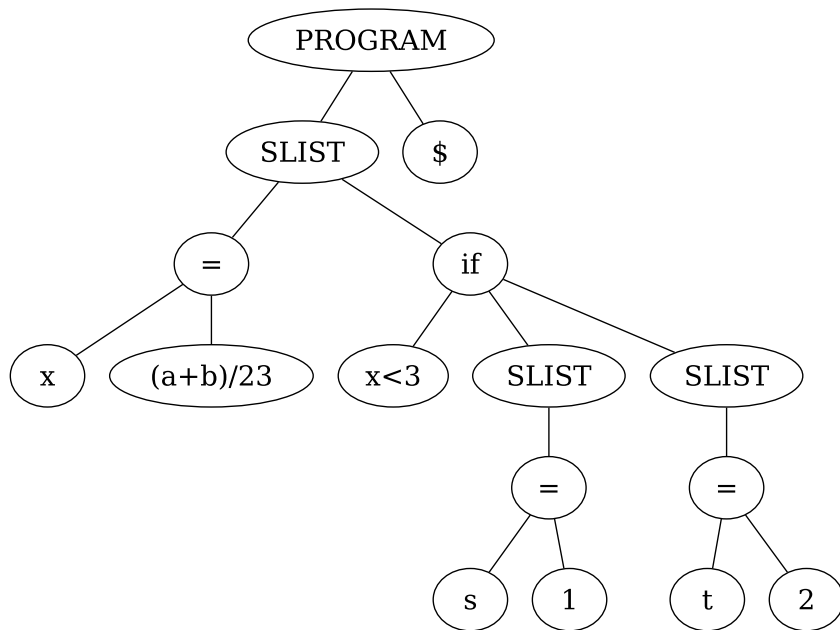
...that generates →

...a Parse Tree

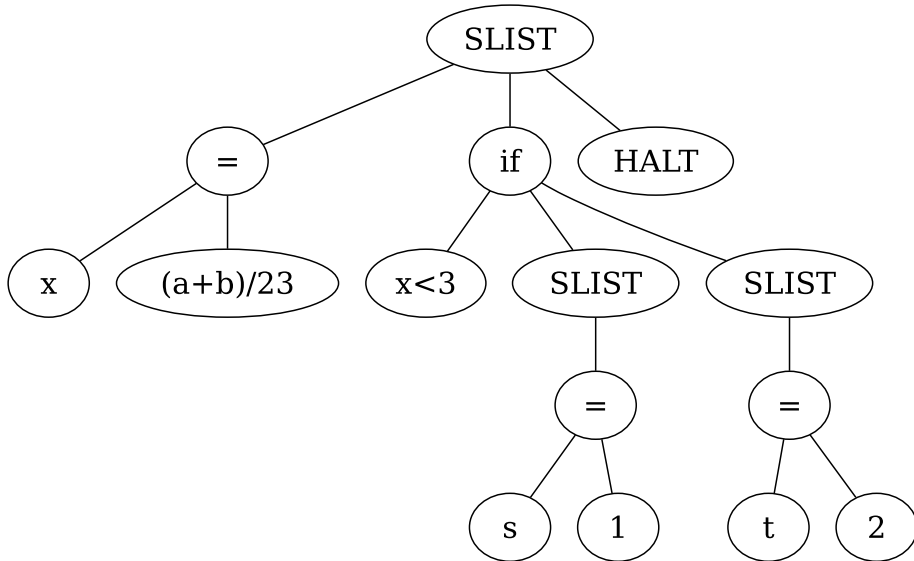


Simplify Statement S Nodes

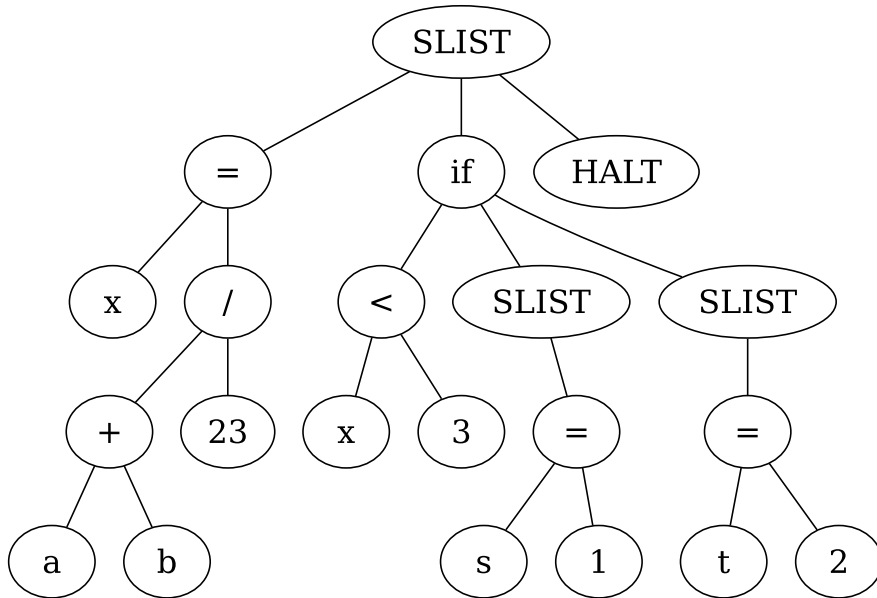




Programs are just an *SLIST* that HALTs

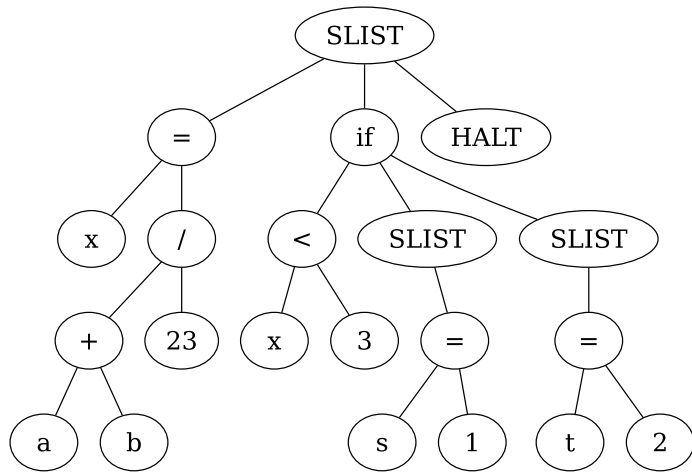


Generate RHS Expression Trees

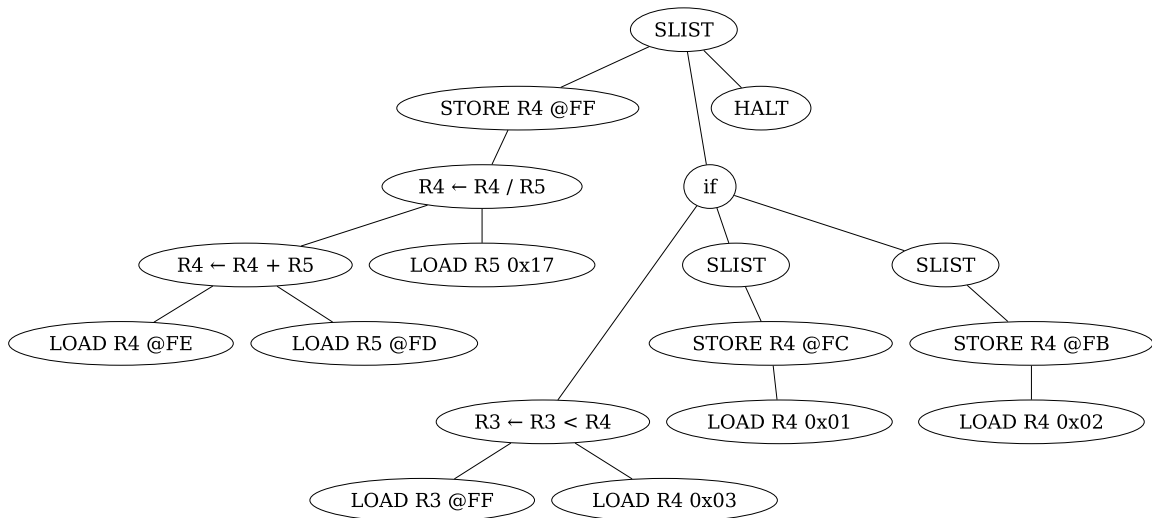


Assign Memory Locations for Variables

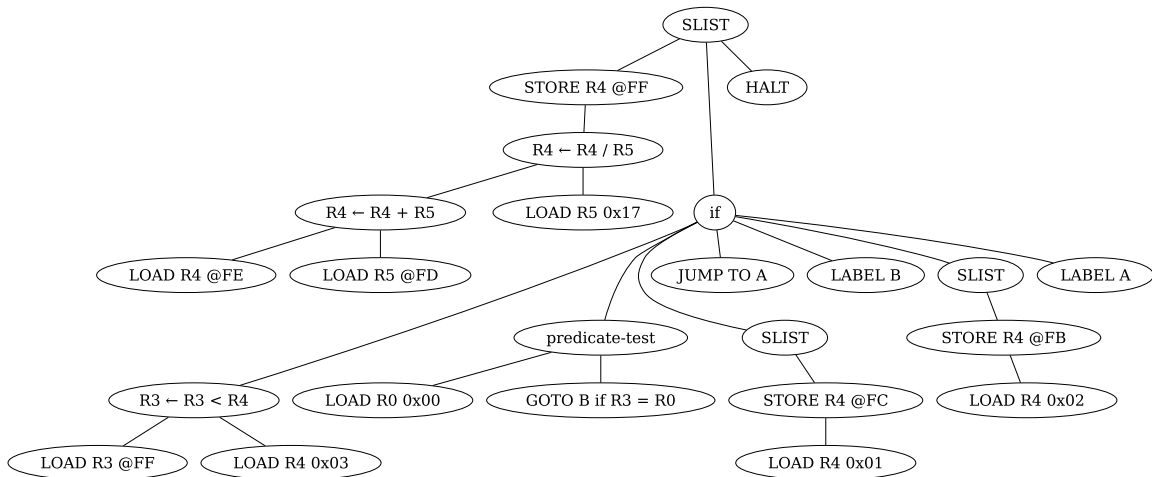
Variable	Memory Address
x	@FF
a	@FE
b	@FD
s	@FC
t	@FB



Assign Registers and Instructions for Expressions



Generate Instructions for Control Structures



(Pseudo) Assembly Code

```
LOAD R4 @FE
LOAD R5 @FD
R4 <- R4 + R5
LOAD R5 0x17
R4 <- R4 / R5
STORE R4 @FF
LOAD R3 @FF
LOAD R4 0x03
R3 <- R3 < R4
LOAD R0 0x00
GOTO B if R3 = R0
LOAD R4 0x01
STORE R4 @FC
JUMP TO A
LABEL B
```

```
LOAD R4 0x02
STORE R4 @FB
LABEL A
HALT
```

An **assembler** will take these instructions and generate the actual machine code, resolving **LABELs** and **GOTOs**, and perhaps performing some low-level optimizations such as removing redundant **LOADs** and **STOREs**.