

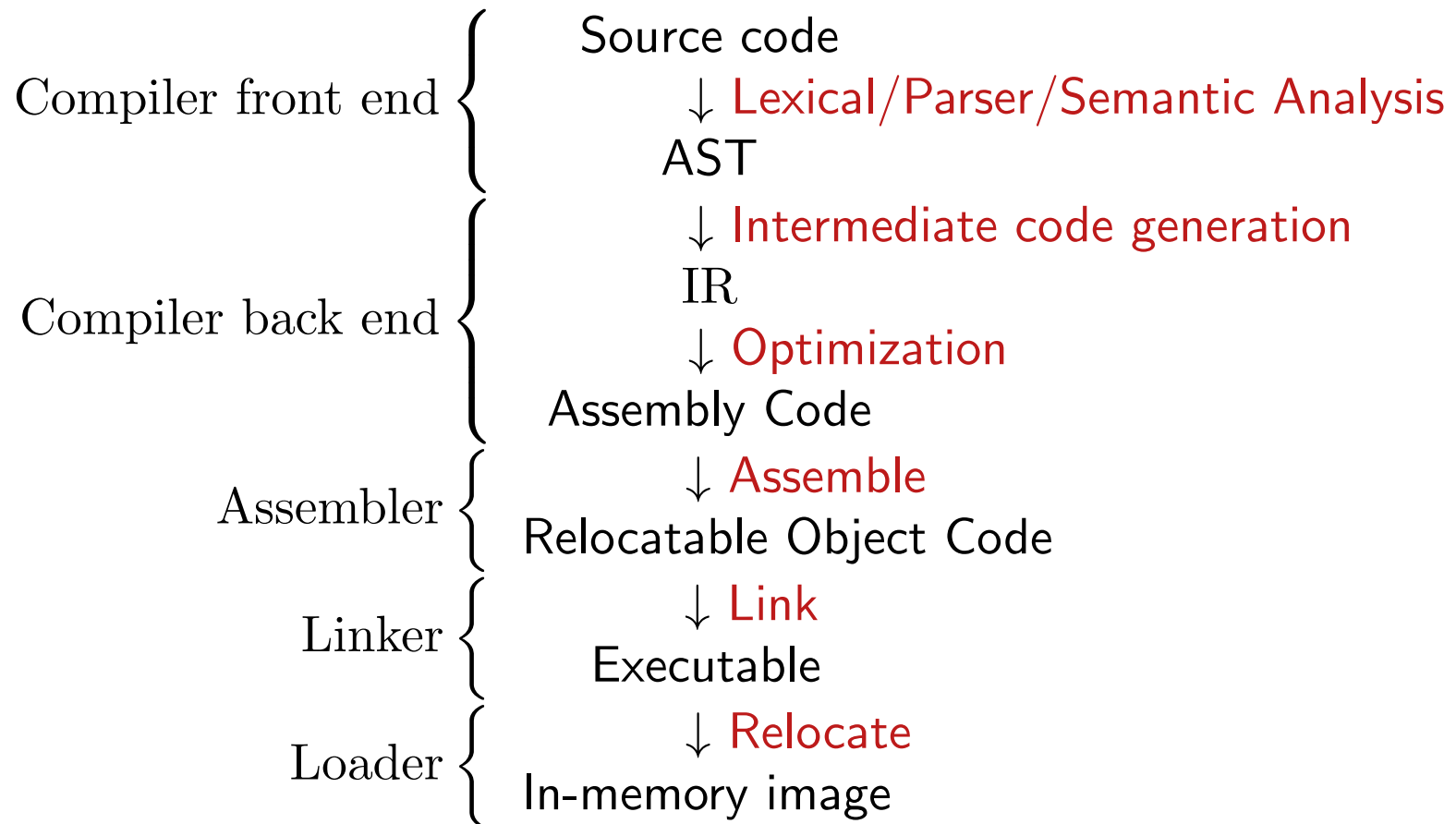
Intermediate code generation

Jordi Cortadella

Dept. CS, UPC

May 4, 2022

From source to machine code



Compiler Frontends and Backends

The front-end focuses on *analysis*:

- Lexical analysis
- Parsing
- Static semantic checking
- AST generation

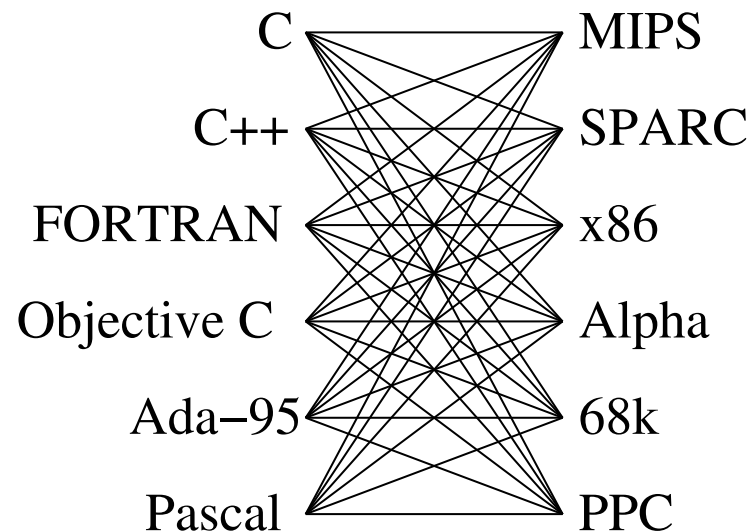
The back-end focuses on *synthesis*:

- Translation of the AST into intermediate code
- Optimization
- Assembly code generation

Portable Compilers

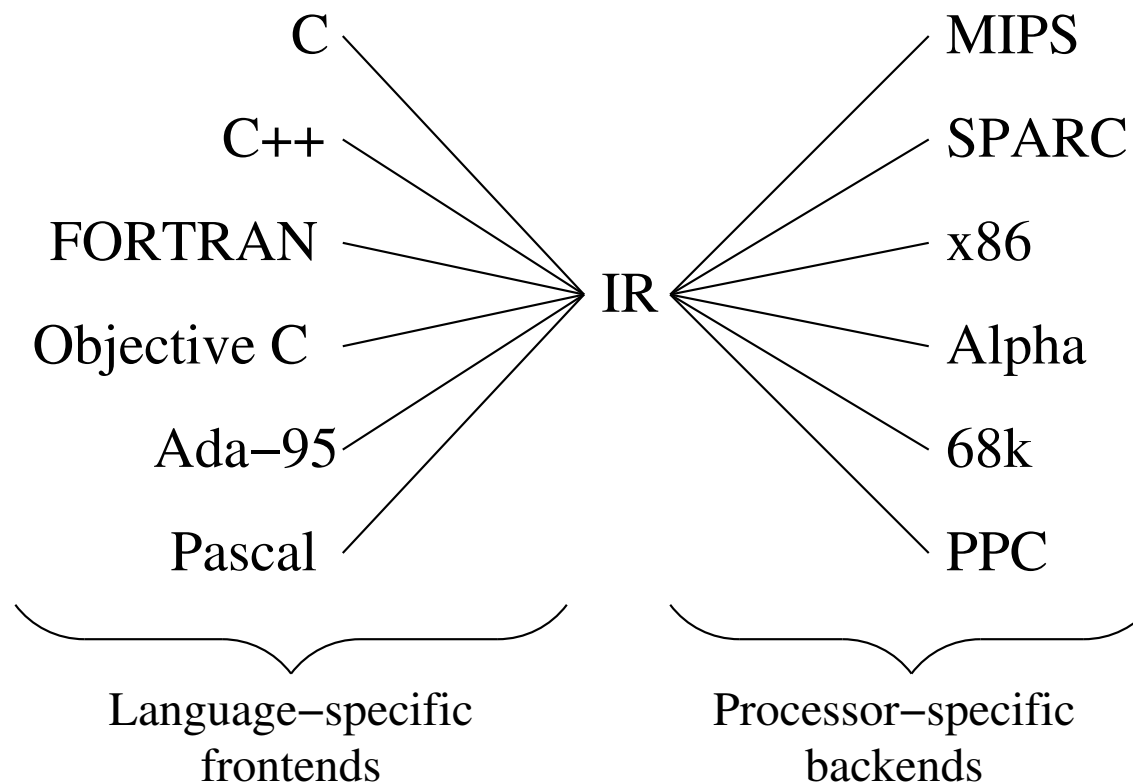
Building a compiler a large undertaking; most try to leverage it by making it portable.

Instead of



Portable Compilers

Use a common intermediate representation



Intermediate language

There are various intermediate languages used in compilers. They are languages of abstract machines used for analysis, optimization and machine language generation.

Some intermediate languages:

- Syntax trees
- Three-address codes
- Stack machines

Three-address code will be used in this course for code generation and optimization.

Three-address code

Most instructions have three addresses:

$$x = y \ op \ z$$

Only one operator allowed in the right side of an instruction, e.g., the instruction $a = x + y * z$ must be translated into this sequence of instructions:

$$\begin{aligned} t_1 &= y * z \\ a &= x + t_1 \end{aligned}$$

where t_1 is a compiler-generated temporary name.

Addresses and instructions

An **address** can be one of the following:

- A **name**. Source-program names appear as addresses. In an implementation, every source name is replaced by a pointer to its symbol-table entry.
- A **constant**.
- A **compiler-generated temporary**. For code optimization, it is convenient to use a different name each time a temporary is needed.

Three-address instructions:

- $x = y \text{ op } z$, where *op* is an arithmetic or logical operator and *x*, *y* and *z* are addresses.
- $x = \text{op } y$, where *op* is a unary operation (minus, logical negation, ...).
- $x = y$ (copy instructions).
- **goto** *L* (unconditional jump).
- **if** *x* **goto** *L* and **ifFalse** *x* **goto** *L* (conditional jump).
- **if** *x* *relop* *y* **goto** *L*, where *relop* is a relational operator (<, ==, >=, ...).

Addresses and instructions

To call the function $p(x_1, x_2, \dots, x_n)$:

```
param  x1  
param  x2  
...  
param  xn  
call   p, n
```

The integer n is not redundant because calls can be nested.
The call can return a result as follows: $t = \text{call } p, n$.

To return from a function:

```
return y
```

where y is the returned result (optional)

Addresses and instructions

Indexed copy instructions:

- $x = y[i]$ and $x[i] = y$,
where $x[i]$ represents the location i memory units beyond location x .

Address and pointer assignments:

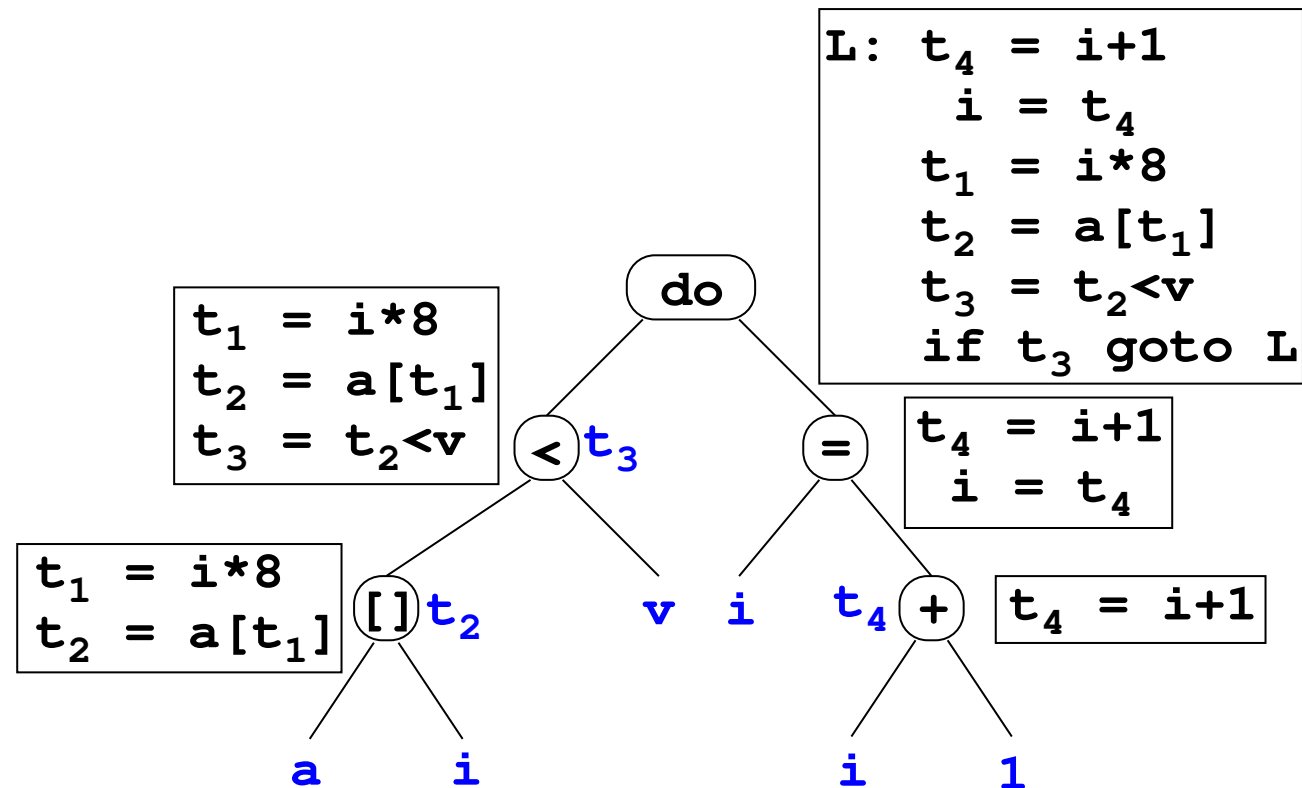
- $x = \&y$, sets x to point to the location of y .
- $x = *y$, sets x to the value of the location pointed to by y .
- $*x = y$, sets the value of the location pointed to by x to the value of y .

Code generation: example

The statement: `do i = i+1; while (a[i] < v);`

can be translated into:

```
L: t4 = i + 1  
   i = t4  
   t1 = i * 8  
   t2 = a[t1]  
   t3 = t2 < v  
   if t3 goto L
```



Declarations

Variables need to be allocated in memory regions.

- How much storage do we need for each type?
- What are the relative addresses of each component inside a type?

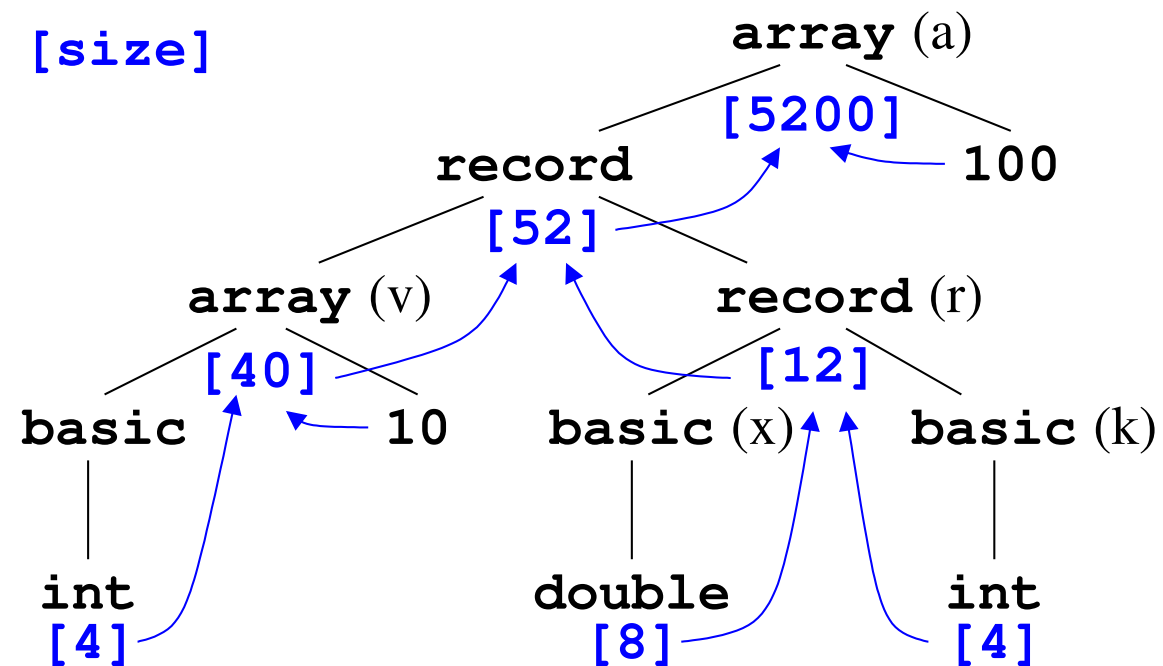
We will use this simple grammar for declarations:

$$\begin{aligned}\text{Decl} &\rightarrow \text{Type } \mathbf{id} ; \text{Decl} \mid \epsilon \\ \text{Type} &\rightarrow \text{Basic} \mid \text{Array} \mid \text{Record} \\ \text{Basic} &\rightarrow \mathbf{int} \mid \mathbf{char} \mid \mathbf{float} \mid \dots \\ \text{Array} &\rightarrow \text{Type} [\mathbf{num}] \\ \text{Record} &\rightarrow \mathbf{record} \{ \text{Decl} \}\end{aligned}$$

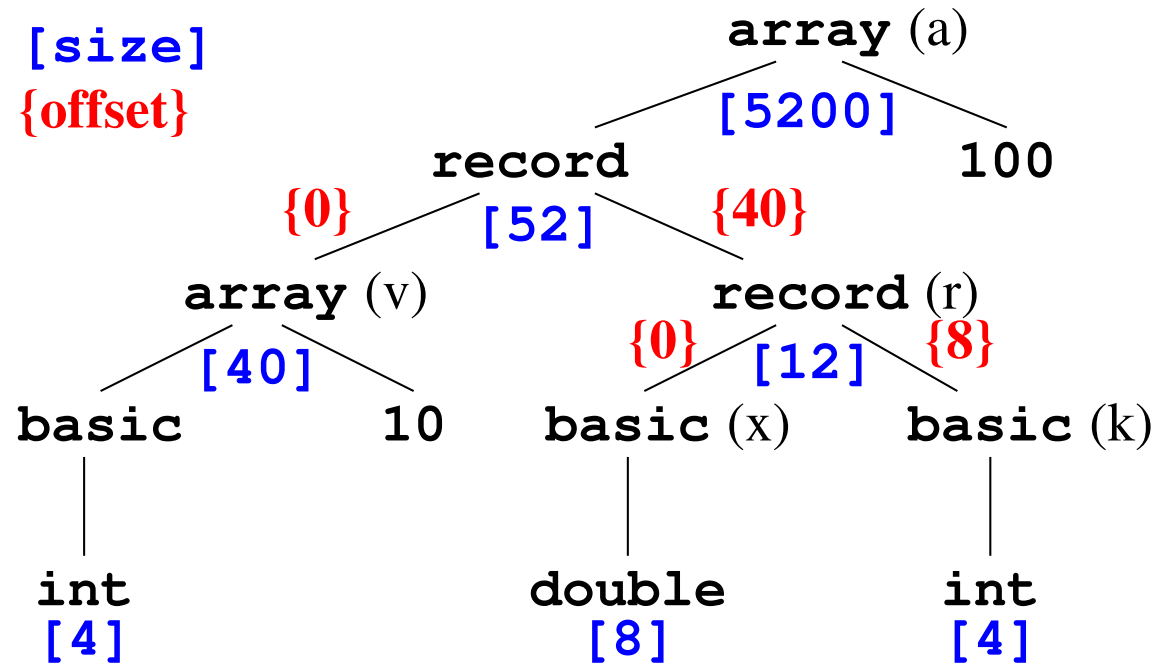
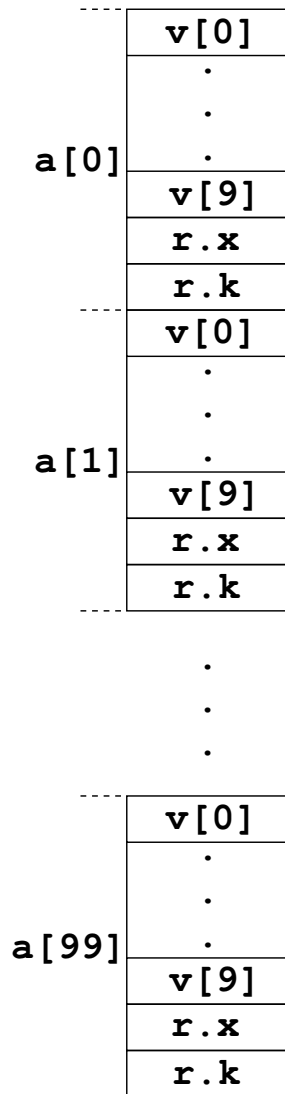
Declarations

What is the size of each type component?

```
record {  
  int [10] v;  
  record {  
    double x;  
    int k;  
  } r;  
} [100] a;
```



Declarations



Q: What is the memory location for `a[13].r.k` ?

Declarations

Let us assume an array $A[l \dots u]$ containing the elements $A[l], A[l + 1], \dots, A[u]$:

$$\text{Address}(A[i]) = \text{Base}(A) + (i - l) \times \text{SizeElem}(A)$$

For arrays in C, C++ and Java ($l = 0$),

$$\text{Address}(A[i]) = \text{Base}(A) + i \times \text{SizeElem}(A)$$

Example: let us consider the array `int A[100][25]`. What is the address of $A[i][j]$?

$$\text{Size}(\text{int}) = 4$$

$$\text{Size}(A[i]) = 25 \times \text{Size}(\text{int}) = 100$$

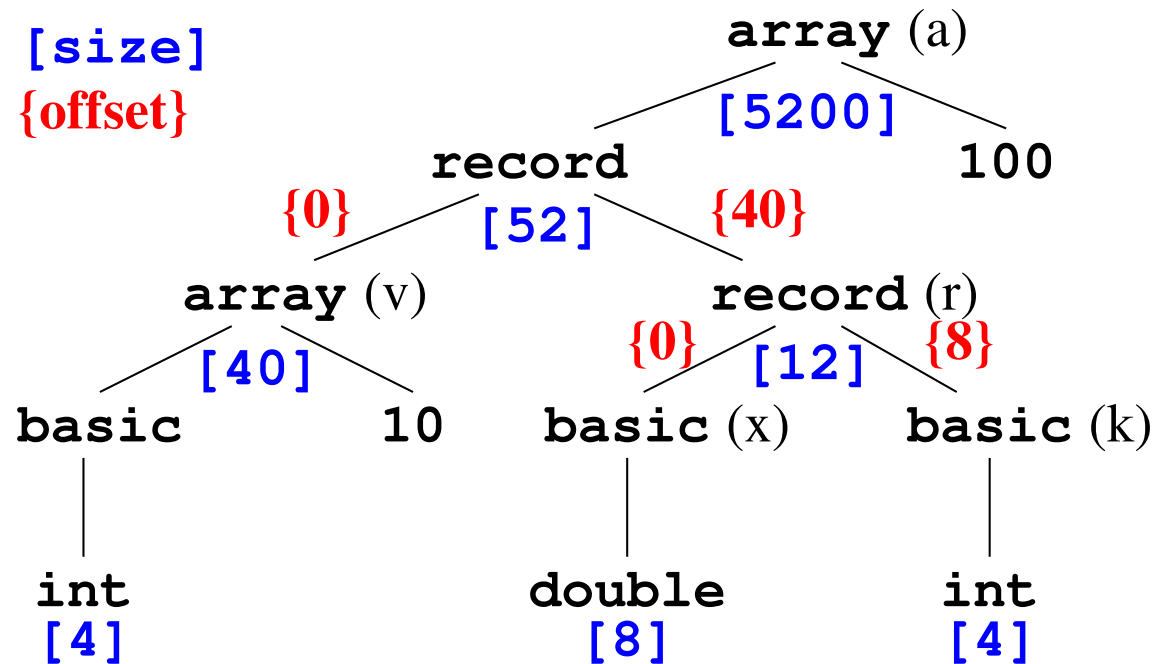
$$\text{Address}(A[i][j]) = \text{Address}(A[i]) + j \times 4 = \text{Base}(A) + i \times 100 + j \times 4$$

Records:

$$\text{Address}(A.z) = \text{Base}(A) + \text{Offset}(A, z)$$

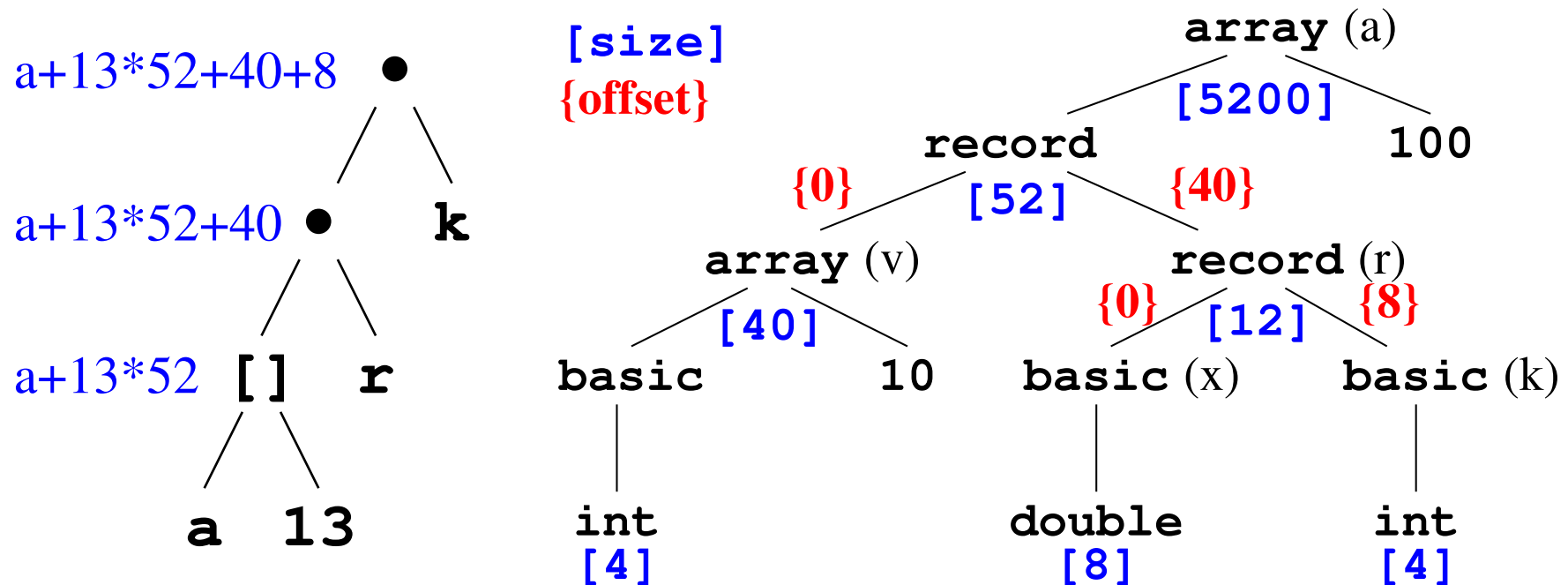
Declarations

```
record {  
  int [10] v;  
  record {  
    double x;  
    int k;  
  } r;  
} [100] a;
```



Q: What is the memory location for $a[13].r.k$?

Declarations

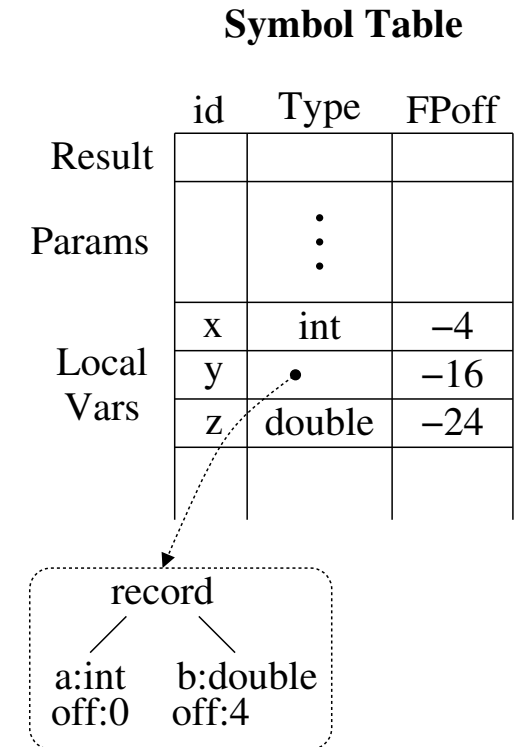
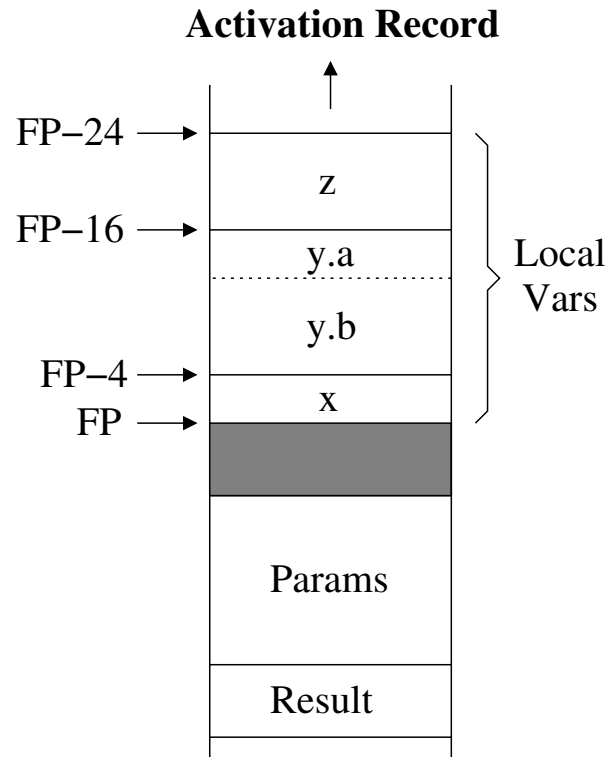


Q: What is the memory location for $a[13].r.k$?

$$\text{Address}(a[13].r.k) = \text{Base}(a) + 13 \times 52 + 40 + 8 = \text{Base}(a) + 724$$

Storage allocation for declarations

```
int f(int p, int q) {
    int x;
    record {
        int a;
        double b;
    } y;
    double z;
    ...
}
```



The symbol table contains information on how the variables are stored in the activation record (offsets with regard to the frame pointer).

Storage allocation for declarations

Ex.: declaration of local variables (grammar and semantic actions)

```
LocalVars  →  { FPoffset = 0; }  ListDecls
ListDecls  →  Decl ListDecls | ε
      Decl  →  Type id ;
               { FPoffset = FPoffset – Type.size;
                 SYMTAB.put (id.lexeme, Type.type, FPoffset); }
```

Exercise

Consider the semantic actions for storage allocation of:

- Parameters and results of a function.
- Nested scopes within a function.

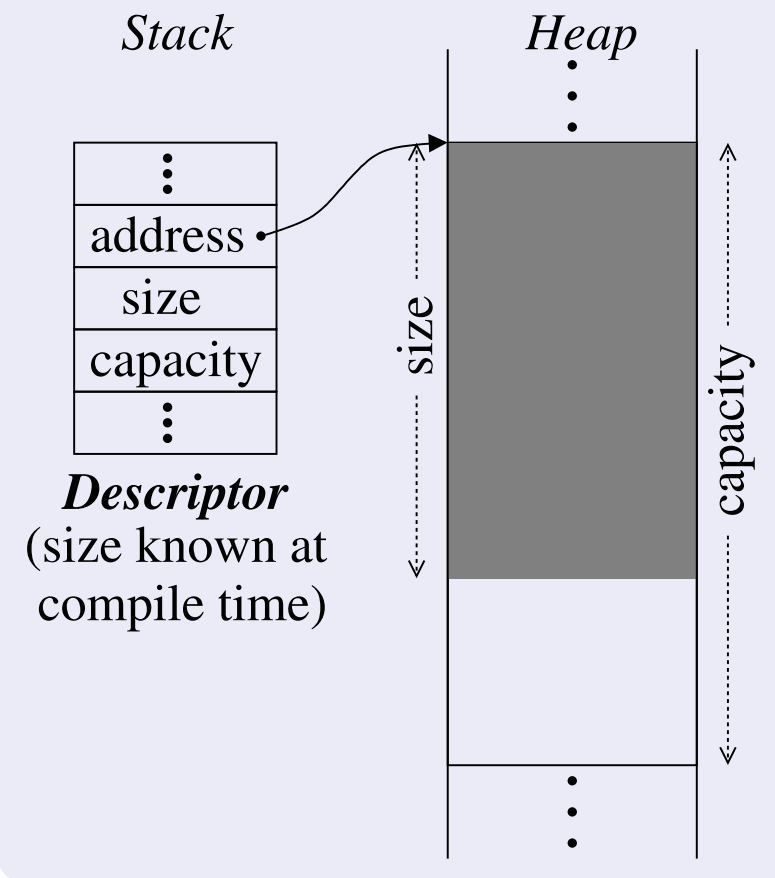
Dynamic storage allocation

Memory allocation for dynamic data structures is not easy. A combination of *static* and *dynamic* allocation is typically used:

- A fixed-sized descriptor of the data structure is allocated statically at compile time in the stack.
- The dynamic data storage is allocated and managed in the heap at runtime.

The heap is also used to allocate variables that outlive the call to the procedure that created them. These variables are usually managed by the program (e.g., *new/delete* statements in C++).

Ex.: STL vectors in C++

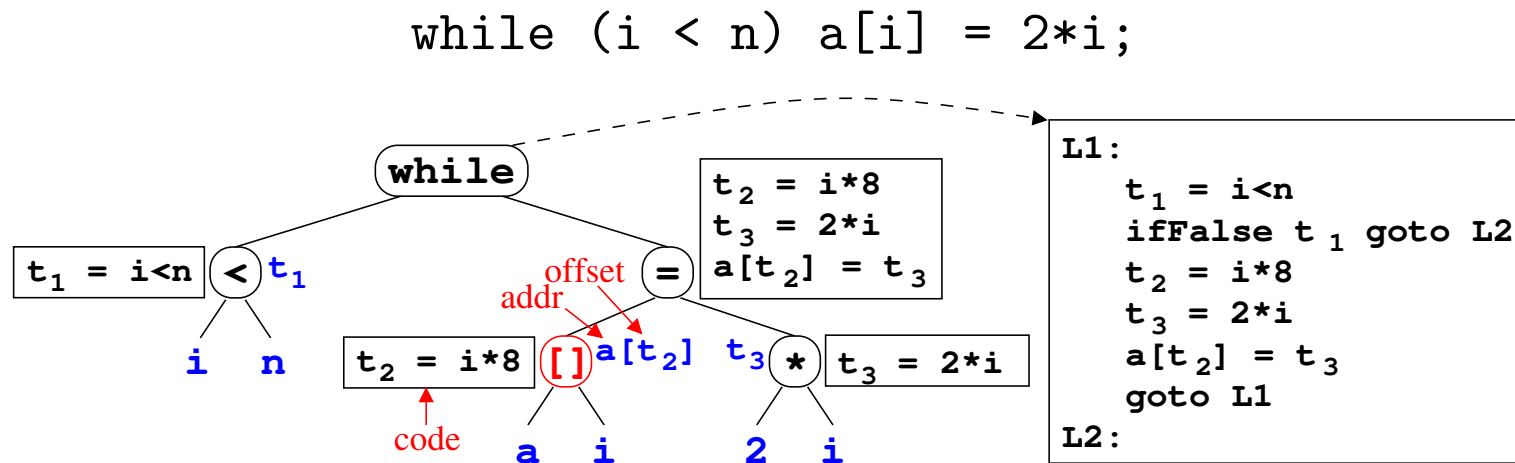


Attributes for code generation

Let us assume every node N of the syntax tree is annotated with three attributes:

- $N.addr$: denotes the address that will hold the value of N (e.g., an entry in the symbol table), or a constant.
- $N.offset$: denotes the offset applied to the address (for array/record access).
- $N.code$: a string holding the three-address code associated to N .

Let us assume that top is the current symbol table and $top.get(name)$ returns the symbol table entry associated to the string $name$.

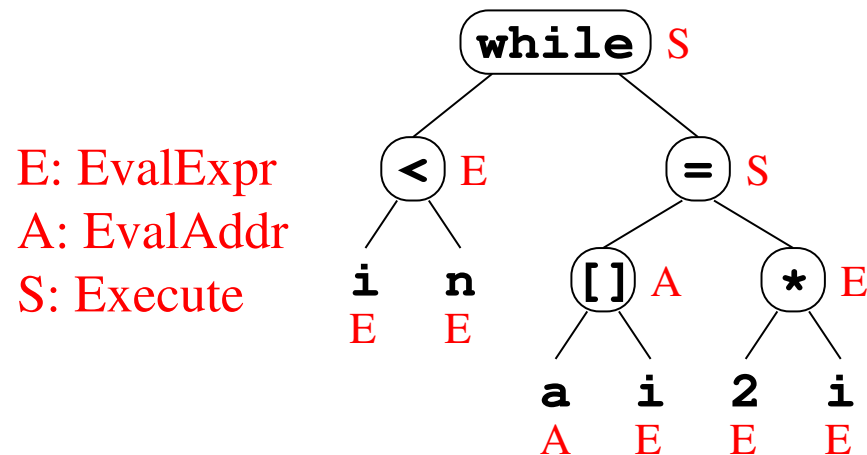


Functions for code generation

Three functions will be used for code generation that will define the previous attributes:

- **EvalExpr(E)** will generate code to evaluate an expression.
- **EvalAddr(E)** will generate code to calculate the address of an expression.
- **Execute(S)** will generate code to execute a statement.

```
while (i < n) a[i] = 2*i;
```



Code generation for expressions

`EvalExpr(E)` calculates the attributes $E.addr$ and $E.code$ of an expression. A grammar rule as parameter denotes the semantic contents of node E , e.g.,

`EvalExpr($E \rightarrow E_1 + E_2$)`

denotes the call to `EvalExpr(E)`, where E is a syntax-tree node representing the sum of two expressions.

A simple expression

`EvalExpr($E \rightarrow E_1 + E_2$)`

```
EvalExpr( $E_1$ ); EvalExpr( $E_2$ );  $E.addr$  = new Temp();
```

```
 $E.code$  =  $E_1.code$  ||  $E_2.code$  || " $E.addr$  =  $E_1.addr$  +  $E_2.addr$ ";
```

Let us assume that `type(E_1)=float` and `type(E_2)=int`. Then:

`EvalExpr($E \rightarrow E_1 + E_2$)`

```
EvalExpr( $E_1$ ); EvalExpr( $E_2$ );
```

```
 $E.addr$  = new Temp();  $t$  = new Temp();
```

```
 $E.code$  =  $E_1.code$  ||  $E_2.code$  ||
```

```
" $t$  = int2float  $E_2.addr$ " || " $E.addr$  =  $E_1.addr$  +  $t$ ";
```

Code generation for expressions

Pointers:

```
EvalExpr( $E \rightarrow *E_1$ )  
  EvalExpr( $E_1$ );  
   $E.addr = \text{new Temp}()$ ;  
   $E.code = E_1.code \parallel "E.addr = *E_1.addr";$ 
```

```
EvalExpr( $E \rightarrow \&E_1$ )  
  EvalExpr( $E_1$ );  
   $E.addr = \text{new Temp}()$ ;  
   $E.code = E_1.code \parallel "E.addr = \&E_1.addr";$ 
```


Code generation to calculate addresses

Identifiers:

EvalAddr($E \rightarrow id$)

```
E.addr = top.get(id.lexeme); E.code = ""; E.offset = 0;
```

Arrays:

EvalAddr($E \rightarrow E_1[E_2]$)

```
t = new Temp(); EvalAddr(E1); EvalExpr(E2);  
E.addr = E1.addr; E.offset = t;  
E.code = E1.code || E2.code ||  
    "t = E2.addr * SizeElem(E1)" ||  
    "t = t + E1.offset";
```

Records:

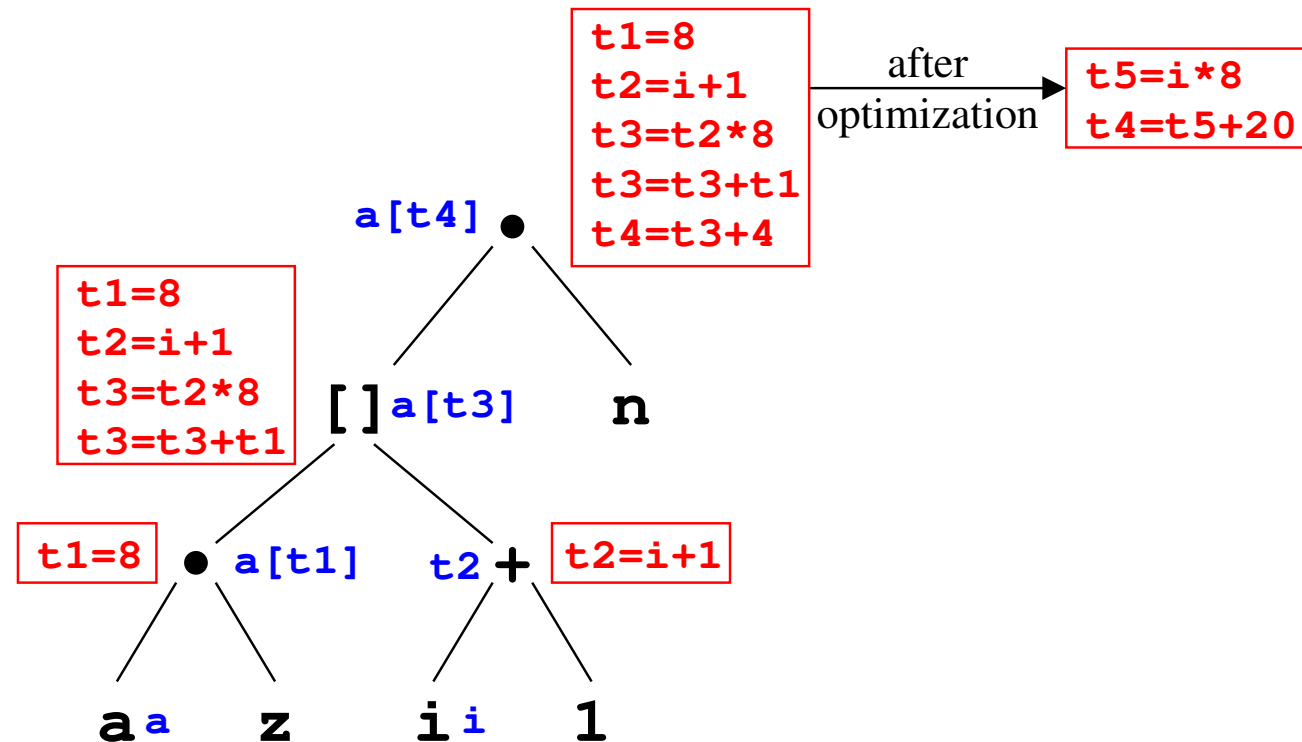
EvalAddr($E \rightarrow E_1.field$)

```
t = new Temp(); EvalAddr(E1); E.addr = E1.addr; E.offset = t;  
E.code = E1.code || "E.offset = E1.offset + Offset(E1, field);"
```

Code generation to calculate addresses: example

```
struct {  
  double x;  
  struct {  
    int m;  
    int n;  
  } z[20];  
} a;
```

Calculate address of:
`a.z[i+1].n`



For each tree node:

- `addr[offset]` in blue
- code in red boxes

Code generation for expressions

Identifiers:

```
EvalExpr( $E \rightarrow id$ )  $\equiv$  EvalAddr( $E$ );
```

Arrays and records:

```
EvalExpr( $E \rightarrow E_1[E_2] \mid E_1.field$ )  
  EvalAddr( $E$ ); t = new Temp();  
  E.code = E.code || "t = E.addr[E.offset]";  
  E.addr = t; E.offset = 0;
```

Code generation for statements: assignment

Assignment

```
Execute( $S \rightarrow E_1 := E_2$ )  
  EvalAddr( $E_1$ ); EvalExpr( $E_2$ );  
   $S.code = E_1.code \parallel E_2.code \parallel$   
     $"E_1.addr[E_1.offset] = E_2.addr";$ 
```

In case $E_1.offset = 0$, this reduces to

```
Execute( $S \rightarrow E_1 := E_2$ )  
  EvalAddr( $E_1$ ); EvalExpr( $E_2$ );  
   $S.code = E_1.code \parallel E_2.code \parallel "E_1.addr = E_2.addr";$ 
```

Execution order: are these two generated codes equivalent?

```
 $S.code = E_1.code \parallel E_2.code \parallel "E_1.addr = E_2.addr";$ 
```

```
 $S.code = E_2.code \parallel E_1.code \parallel "E_1.addr = E_2.addr";$ 
```

Code generation for statements: the order is important

Consider the following C++ code

```
int a[3] = {0, 0, 0};  
int i = 0;  
a[++i] = ++i;  
int n = ++i + i;
```

What is the contents of array `a` and the value of `n` after executing the code?

Code generation for statements: the order is important

Consider the following C++ code

```
int a[3] = {0, 0, 0};  
int i = 0;  
a[++i] = ++i;  
int n = ++i + i;
```

What is the contents of array `a` and the value of `n` after executing the code?

Suggestion: try `g++` and `clang++` with `--std=c++11` and `--std=c++17`.

Before C++17

```
a = [0, 0, ?]  
n = ?
```

After C++17

```
a = [0, 0, 1]  
n = ?
```

The semantics of any language should clearly specify the evaluation order of expressions and the cases in which the behavior is *undefined*. Ideally, no behavior should be undefined.

Code generation for statements: while

while B do S: Two possible schemes

```
L1: t = eval(B)
      ifFalse t goto L2
      S
      goto L1
L2:
```

Code generation for statements: while

while B do S: Two possible schemes

```
L1: t = eval(B)
      ifFalse t goto L2
      S
      goto L1
L2:
```

```
      goto L2
L1: S
L2: t = eval(B)
      if t goto L1
```


Code generation for statements: while

while B do S: Two possible schemes

```
L1: t = eval(B)
      ifFalse t goto L2
      S
      goto L1
L2:
```

```
      goto L2
L1: S
L2: t = eval(B)
      if t goto L1
```

Code generation for loops

Execute($I \rightarrow \text{while } B \text{ do } S$)

```
L1 = new Label(); L2 = new Label();
EvalExpr(B); Execute(S);
I.code = "L1:" || B.code ||
        "ifFalse B.addr goto L2" ||
        S.code || "goto L1" || "L2";
```

Code generation for statements: if-then

if B then S

```
    t = eval(B)
    ifFalse t goto L
    S
L:
```

Code generation for if-then

Execute($I \rightarrow$ if B then S)

```
L = new Label();
EvalExpr(B); Execute(S);
I.code = B.code || "ifFalse B.addr goto L" ||
        S.code || "L:";
```

Code generation for statements: if-then-else

if B then S_1 else S_2

```
t = eval(B)
ifFalse t goto L1
S1
goto L2
L1:
S2
L2:
```

Code generation for if-then-else

Execute($I \rightarrow$ if B then S_1 else S_2)

```
L1 = new Label(); L2 = new Label();
EvalExpr(B); Execute( $S_1$ ); Execute( $S_2$ );
I.code = B.code || "ifFalse B.addr goto L1" || S1.code
|| "goto L2" || L1.code || S2.code || "L2";
```

Code generation for a function call

$f(E_1, \dots, E_n)$

$t_1 = \text{EvalExpr/EvalAddr}(E_1)$

...

$t_n = \text{EvalExpr/EvalAddr}(E_n)$

param t_1

...

param t_n

$t = \text{call } f, n$

- The value or address of each expression will be evaluated depending on whether the parameter is passed by value or reference.
- Only **addressable** expressions can be passed by reference.
- No result is returned in case of a procedure call.
- The parameter count is required. Function calls can be nested, e.g.

$f(w, g(x, y), z).$

Evaluating Boolean expressions

Two schemes (semantically different):

- Eager evaluation: all operands are evaluated.
- Lazy (short-circuit) evaluation: the second operand is only evaluated if the value of the first one is not sufficient to determine the value of the expression.

E_1 and E_2

Eager evaluation

```
t1 = EvalExpr(E1)  
t2 = EvalExpr(E2)  
t3 = t1 and t2
```

E_1 or E_2

Eager evaluation

```
t1 = EvalExpr(E1)  
t2 = EvalExpr(E2)  
t3 = t1 or t2
```

Evaluating Boolean expressions

Two schemes (semantically different):

- Eager evaluation: all operands are evaluated.
- Lazy (short-circuit) evaluation: the second operand is only evaluated if the value of the first one is not sufficient to determine the value of the expression.

E_1 and E_2

Eager evaluation

```
t1 = EvalExpr(E1)  
t2 = EvalExpr(E2)  
t3 = t1 and t2
```

Lazy evaluation

```
t = EvalExpr(E1)  
if False t goto L  
t = EvalExpr(E2)  
L:
```

E_1 or E_2

Eager evaluation

```
t1 = EvalExpr(E1)  
t2 = EvalExpr(E2)  
t3 = t1 or t2
```

Lazy evaluation

```
t = EvalExpr(E1)  
if t goto L  
t = EvalExpr(E2)  
L:
```

Boolean conditions and jumps

Example: **while** not $(E_1 < E_2)$ and $(E_3 = E_4)$ **do** S;

Naïve code

While:

$t_1 = \text{EvalExpr}(E_1)$

$t_2 = \text{EvalExpr}(E_2)$

$t_5 = (t_1 < t_2)$

$t_6 = \text{not } t_5$

ifFalse t_6 goto EndCond

$t_3 = \text{EvalExpr}(E_3)$

$t_4 = \text{EvalExpr}(E_4)$

$t_6 = (t_3 = t_4)$

EndCond:

ifFalse t_6 goto Endwhile

S

goto While

EndWhile:

Boolean conditions and jumps

Example: **while** not $(E_1 < E_2)$ and $(E_3 = E_4)$ **do** S;

Naïve code

While:

```
t1 = EvalExpr(E1)
t2 = EvalExpr(E2)
t5 = (t1 < t2)
t6 = not t5
ifFalse t6 goto EndCond
t3 = EvalExpr(E3)
t4 = EvalExpr(E4)
t6 = (t3 = t4)
```

EndCond:

```
ifFalse t6 goto Endwhile
S
goto While
```

EndWhile:

Efficient code

While:

```
t1 = EvalExpr(E1)
t2 = EvalExpr(E2)
if t1 < t2 goto EndWhile
t3 = EvalExpr(E3)
t4 = EvalExpr(E4)
ifFalse t3 = t4 goto EndWhile
S
goto While
```

EndWhile:

How to generate efficient code for Boolean conditions combined with jumps?

Backpatching

A method to generate efficient code for Boolean expressions.

EvalBoolExpr(E , L_t , L_f)

Evaluates a Boolean expression generating code to jump to L_t or L_f in case the condition evaluates to true or false, respectively.

EvalBoolExpr($E \rightarrow E_1$ and E_2 , L_t , L_f)

(lazy evaluation)

```
Lf1 = new Label();  
EvalBoolExpr(E1, Lf1, Lf);  
EvalBoolExpr(E2, Lt, Lf);  
E.code = E1.code || "Lf1:" || E2.code;
```

EvalBoolExpr($E \rightarrow E_1$ or E_2 , L_t , L_f)

(lazy evaluation)

```
Lf1 = new Label();  
EvalBoolExpr(E1, Lt, Lf1);  
EvalBoolExpr(E2, Lt, Lf);  
E.code = E1.code || "Lf1:" || E2.code;
```

Backpatching

EvalBoolExpr($E \rightarrow \text{not } E_1, L_t, L_f$)

```
EvalBoolExpr( $E_1, L_f, L_t$ );  
E.code =  $E_1$ .code;
```

EvalBoolExpr($E \rightarrow E_1 \text{ relop } E_2, L_t, L_f$)

```
EvalExpr( $E_1$ );  
EvalExpr( $E_2$ );  
E.code =  $E_1$ .code ||  
          $E_2$ .code ||  
         "if  $E_1$ .addr relop  $E_2$ .addr goto  $L_t$ " ||  
         "goto  $L_f$ ";
```

Jump optimization

Some generated jumps are redundant and can be easily optimized, e.g., jumps to the next instruction, jumps to unconditional/conditional jumps, etc.
(see code optimization for more details)

Backpatching: revisiting **while** and **if-then-else**

Execute($I \rightarrow$ **while** B **do** S)

```
L1 = new Label();  
L2 = new Label();  
L3 = new Label();  
EvalBoolExpr(B, L2, L3);  
Execute(S);
```

```
I.code = "L1:" ||  
         B.code ||  
         "L2:" ||  
         S.code ||  
         "goto L1" ||  
         "L3";
```

Execute($I \rightarrow$ **if** B **then** S₁ **else** S₂)

```
L1 = new Label();  
L2 = new Label();  
L3 = new Label();  
EvalBoolExpr(B, L1, L2);  
Execute(S1); Execute(S2);
```

```
I.code =      B.code ||  
             "L1:" ||  
             S1.code ||  
             "goto L3" ||  
             "L2:" ||  
             S2.code ||  
             "L3";
```

Backpatching: example

while not (a < b) **and** (c = d) **do** S;

The code generator will call `EvalBoolExpr(B, L2, L3)`, where B is the condition of the loop, L2 is the label leading the loop body and L3 is the label at the end of the loop.

Before jump optimization

```
L1:
    if a < b goto L3
    goto L4
L4:
    if c != d goto L3
    goto L2
L2:
    S
    goto L1
L3:
```

After jump optimization

```
L1:
    if a < b goto L3
    if c != d goto L3
    S
    goto L1
L3:
```

Switch statements

- Implemented with a sequence of conditional jumps.
- It can be inefficient if the number of branches is large.

```
switch (E) {  
  case  $v_1$ :  $S_1$   
  case  $v_2$ :  $S_2$   
     $\vdots$   
  case  $v_{n-1}$ :  $S_{n-1}$   
  default:  $S_n$   
}
```

\Rightarrow

```
t = EvalExpr(E);  
if (t ==  $v_1$ )  $S_1$ ;  
else if (t ==  $v_2$ )  $S_2$ ;  
   $\vdots$   
else if (t ==  $v_{n-1}$ )  $S_{n-1}$ ;  
else  $S_n$ ;
```

Switch statements with computed jumps

Table of computed jumps

JumpTable:

v_1 :	L_1
	\dots
	L_n
	\dots
v_2 :	L_2
	\dots
	L_n
	\dots
v_3 :	L_3
	\vdots
	\vdots
	\vdots
v_{n-1} :	L_{n-1}

Code (assuming $v_1 < v_2 < \dots < v_{n-1}$)

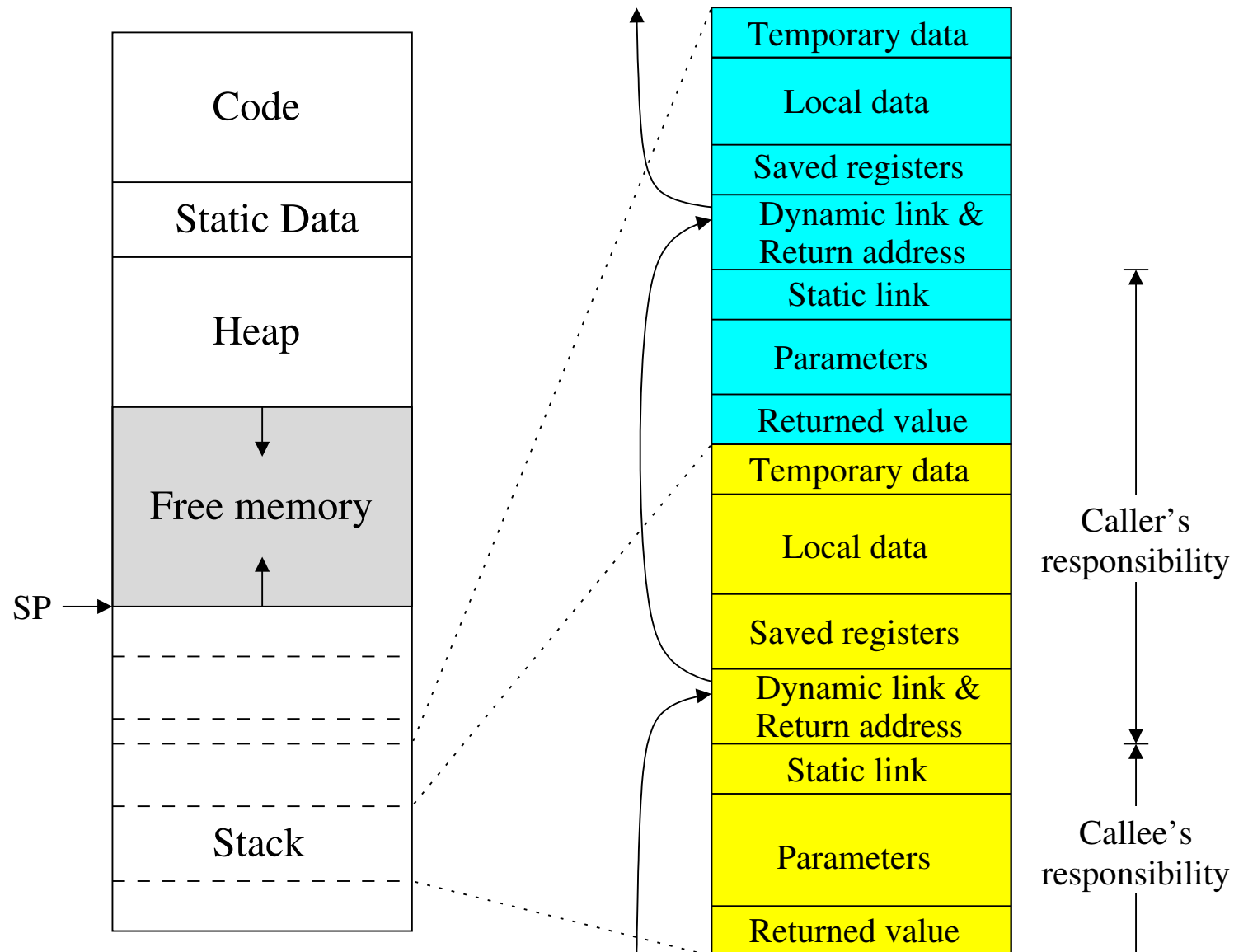
```
t = EvalExpr(E);  
if (t <  $v_1$  or t >  $v_{n-1}$ ) goto  $L_n$ ;  
goto JumpTable[t- $v_1$ ];
```

```
 $L_1$ :  $S_1$ ; goto endsw;  
 $L_2$ :  $S_2$ ; goto endsw;  
...  
 $L_{n-1}$ :  $S_{n-1}$ ; goto endsw;  
 $L_n$ :  $S_n$ ;
```

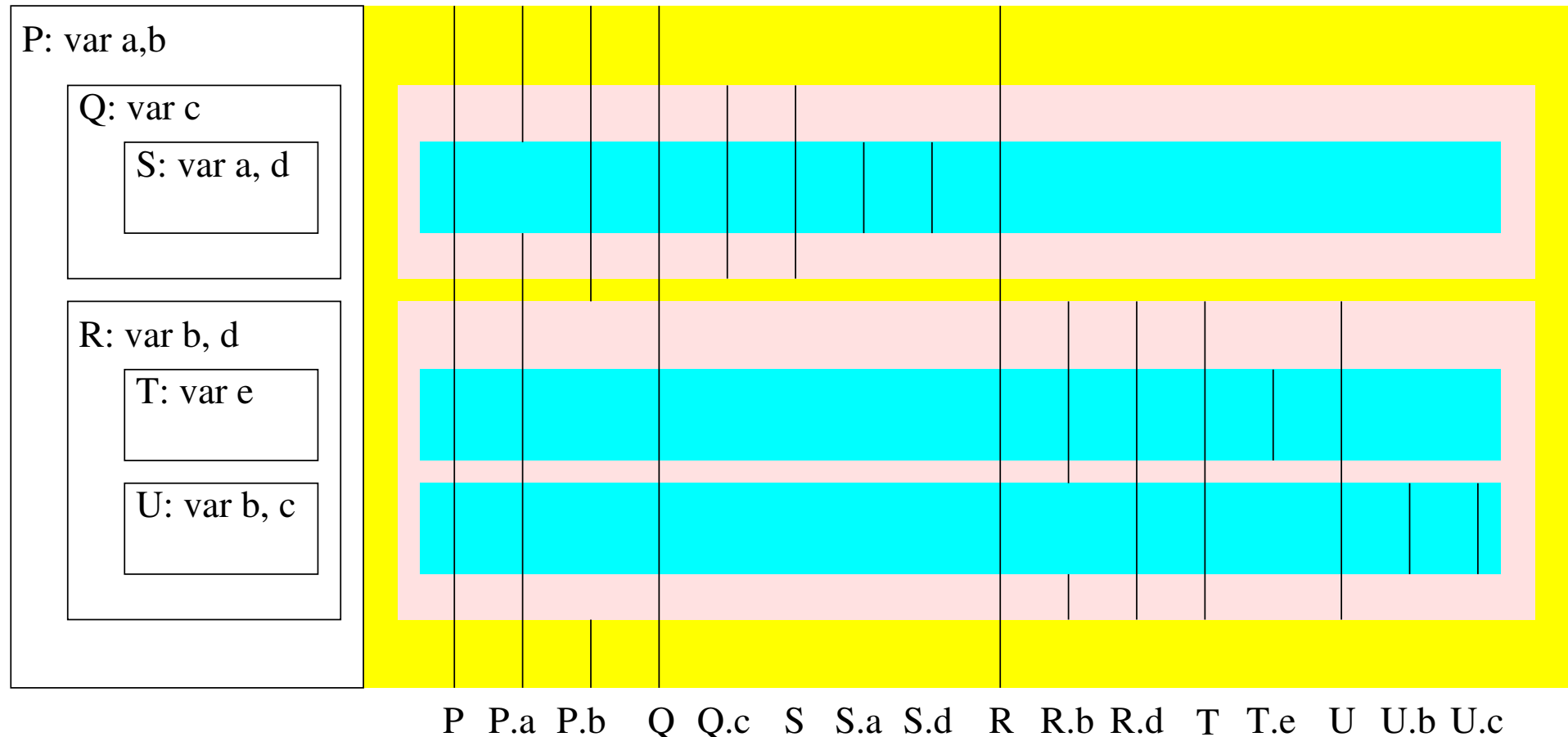
```
endsw:
```

Hybrid approaches are also possible, e.g., using hash tables or several tables with min/max buckets.

Storage organization

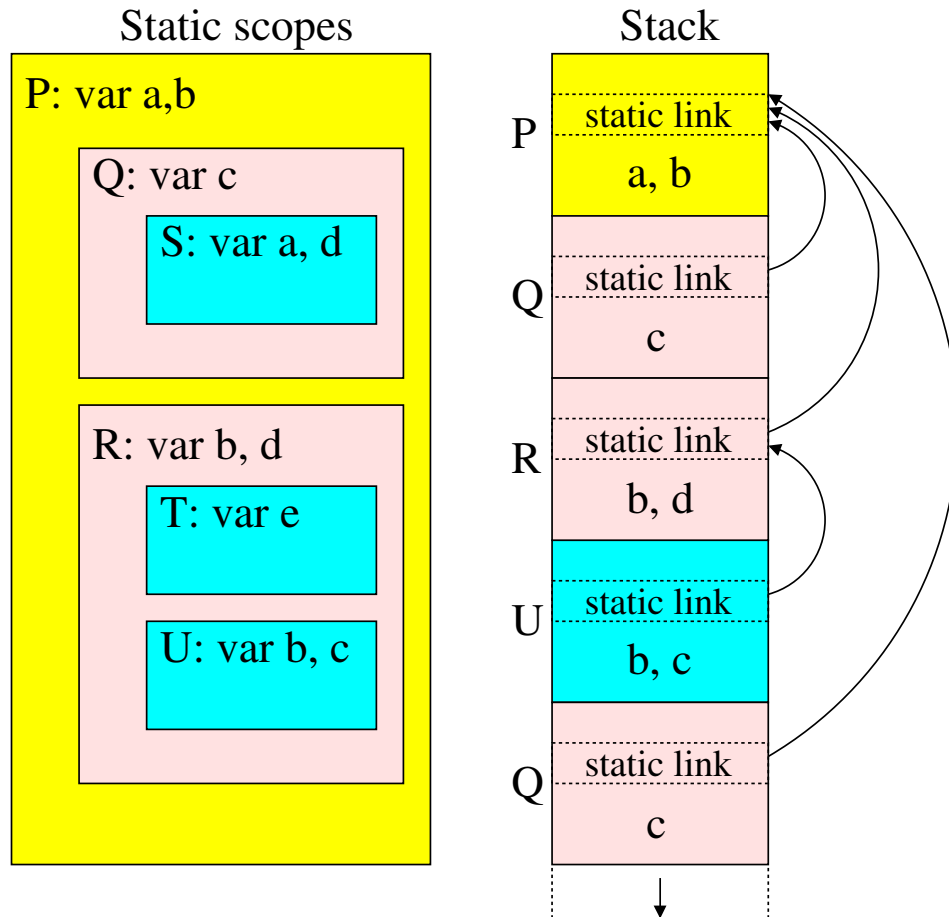


Static scopes



- How to read/write variables stored at different activation records?
- How to manage the static links?

Static scopes



- Each activation record has a static link pointing at the AR of the static ancestor in the source code.
- At every function call, the caller passes the static link as an *implicit* parameter.
- Every access to a non-local variable requires to traverse the chain of static links until the AR storing the variable is reached.

Static scopes

Function P calls $Q(x_1, \dots, x_n)$.

We know that $\text{depth}(Q) - \text{depth}(P) \leq 1$

Code

```
ts1 = &static_link  
ts1 = *ts1  
  ⋮  
ts1 = *ts1 }  $\text{depth}(P) - \text{depth}(Q) + 1$   
param ts1  
param t1  
  ⋮  
param tn  
call  $Q, n + 1$ 
```

Static scopes

Function P calls $Q(x_1, \dots, x_n)$.

We know that $\text{depth}(Q) - \text{depth}(P) \leq 1$

Code

```
ts1 = &static_link  
ts1 = *ts1  
  ⋮  
ts1 = *ts1 } depth(P) - depth(Q) + 1  
param ts1  
param t1  
  ⋮  
param tn  
call Q, n + 1
```

Case $\text{depth}(Q) = \text{depth}(P) + 1$

```
ts1 = &static_link  
param ts1  
  ⋮
```

Case $\text{depth}(Q) < \text{depth}(P) + 1$

```
ts1 = static_link  
ts1 = *ts1  
  ⋮  
ts1 = *ts1 } depth(P) - depth(Q)  
param ts1  
  ⋮
```

Static scopes

Access to variable v declared in P while executing function Q .

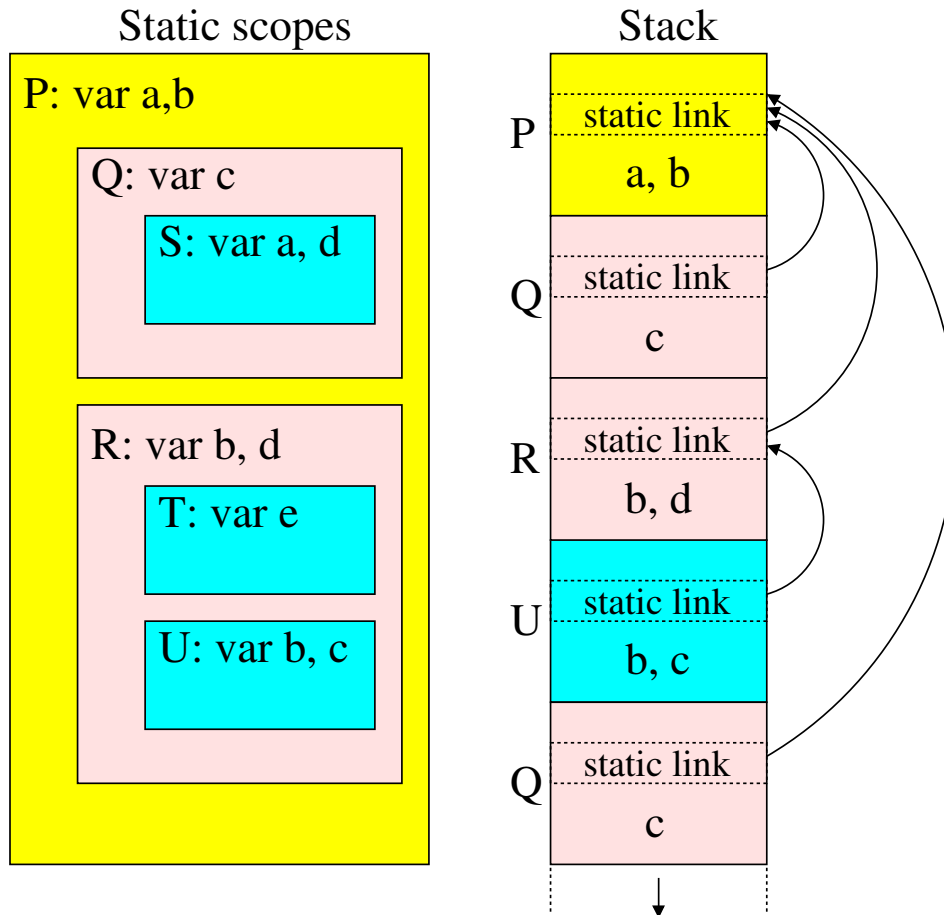
We know that $\text{depth}(Q) - \text{depth}(P) \leq 1$

Code (assume non-local access: $\text{depth}(Q) > \text{depth}(P)$)

```
ts1 = static_link  
ts1 = *ts1  
  ⋮  
ts1 = *ts1 } depth(Q) - depth(P) - 1  
ts1 = ts1 + offset(P :: v)
```

t_{s1} contains the address of variable v

Static scopes: example



Code inside U: $d = a + b$

```
t1 = static_link
t1 = *t1
t1 = t1 + offset(P::a)
t1 = *t1
t2 = t1 + b
t3 = static_link
t3 = t3 + offset(R::d)
*t3 = t2
```

Code inside U: $Q(b)$

```
t1 = static_link
t1 = *t1
param t1
param b
call Q, 2
```