

Reinhard Wilhelm + Helmut Seidl

Abstract Machines

Saarbrücken + Trier

Summer 2002

0 Introduction

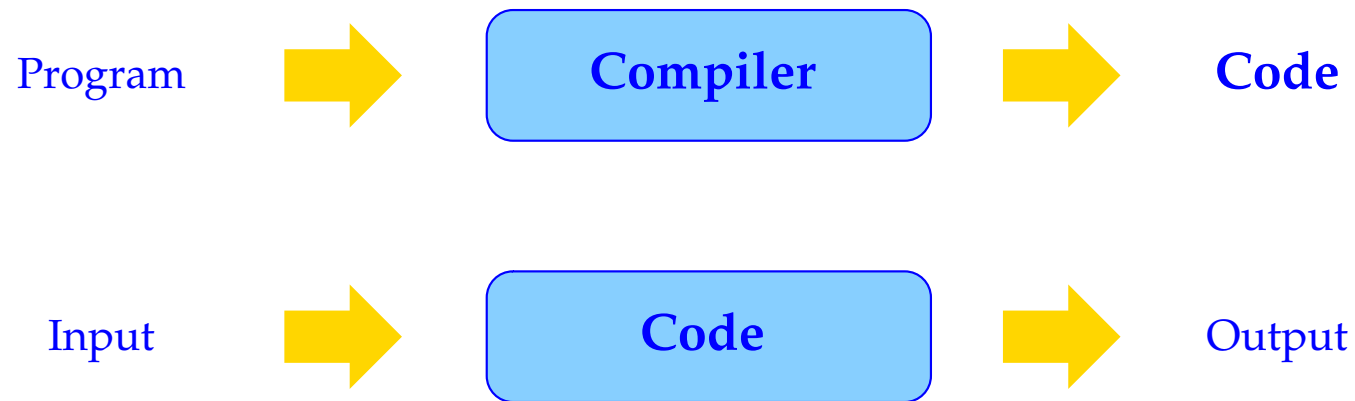
Principle of Interpretation:



Advantage: No precomputation on the program text \implies no/short startup-time

Disadvantages: Program parts are repeatedly analyzed during execution + less efficient access to program variables \implies slower execution speed

Principle of Compilation:



Two Phases (at two different Times):

- Translation of the source program into a machine program (at **compile time**);
- Execution of the machine program on input data (at **run time**).

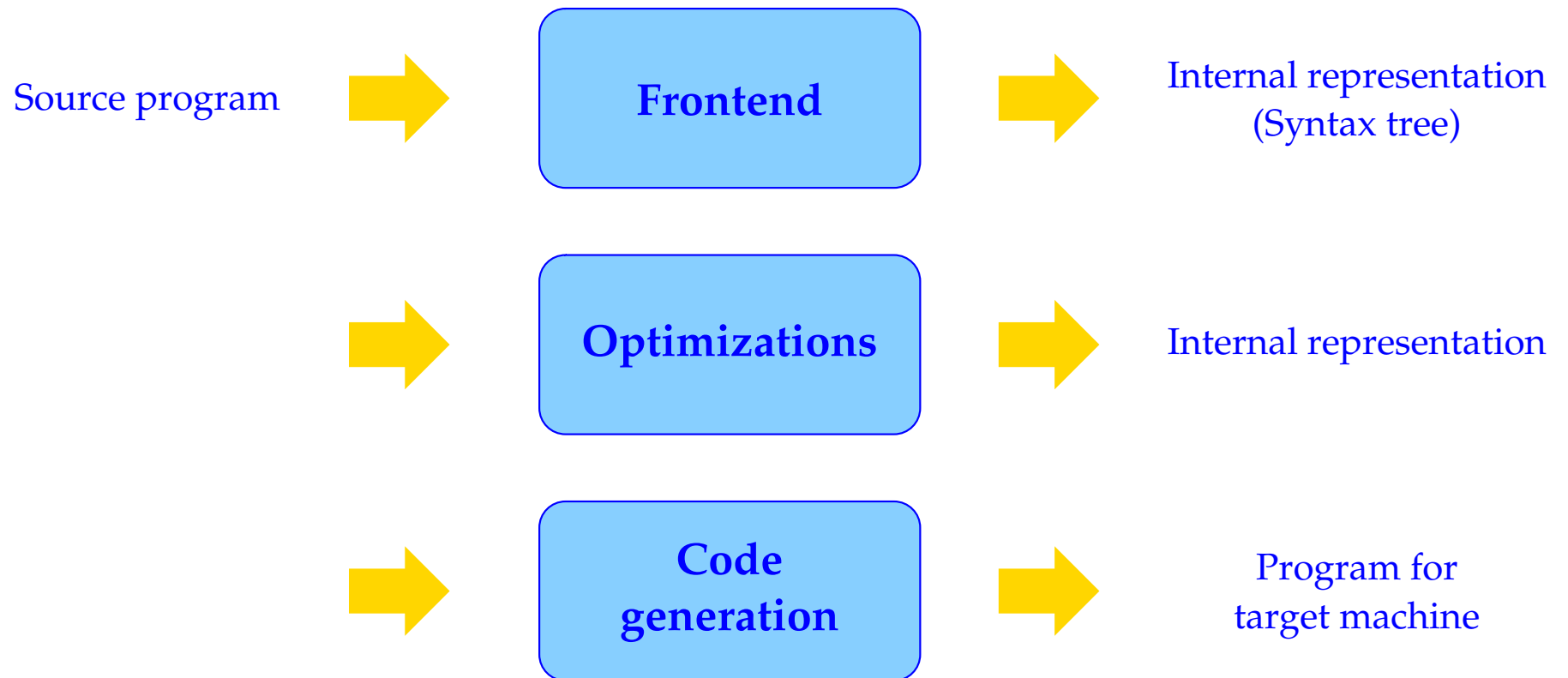
Preprocessing of the source program provides for

- efficient access to the values of program variables at run time
- global program transformations to increase execution speed.

Disadvantage: Compilation takes time

Advantage: Program execution is sped up \implies compilation pays off in long running or often run programs

Structure of a compiler:

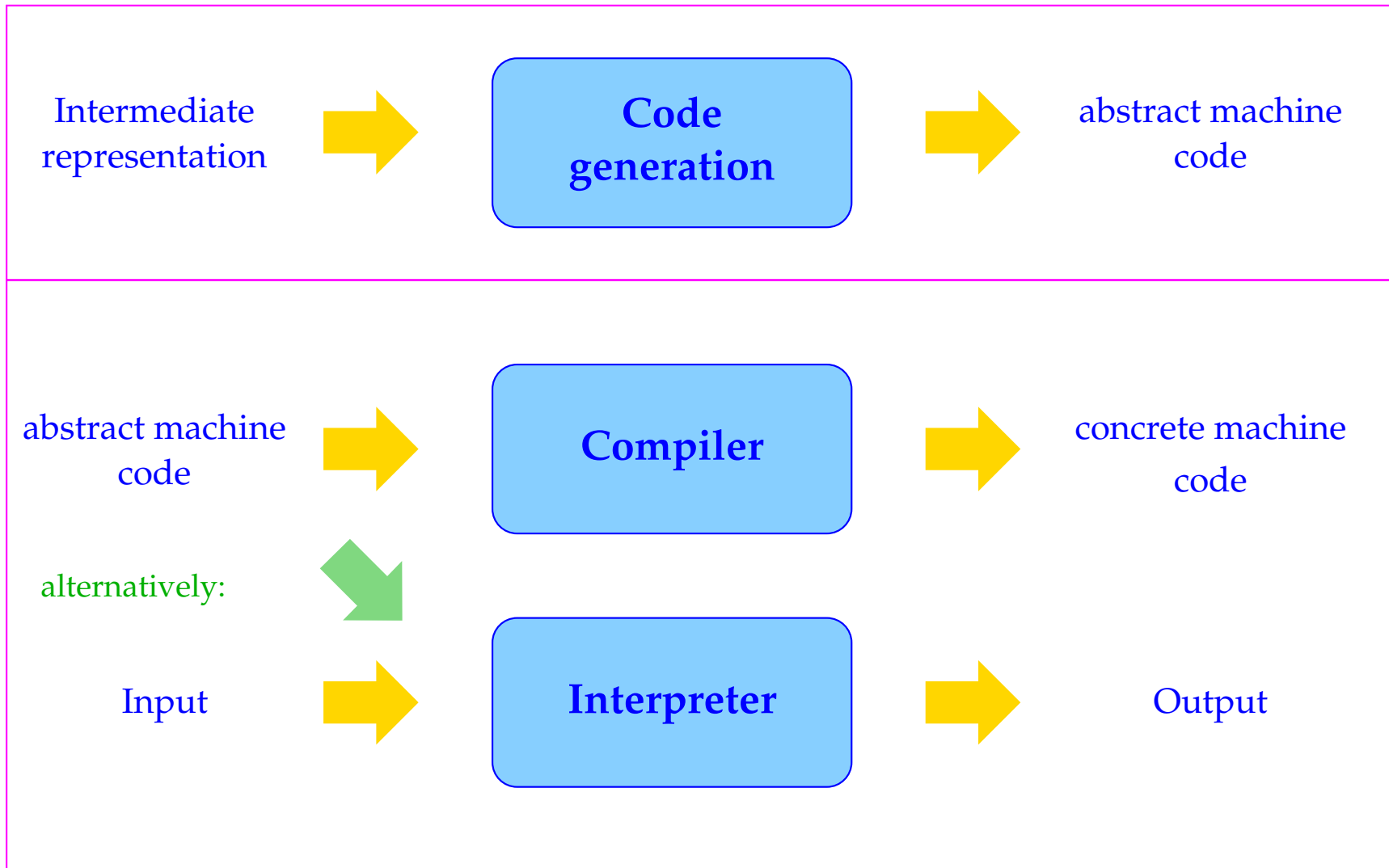


Subtasks in code generation:

Goal is a good exploitation of the hardware resources:

1. **Instruction Selection:** Selection of efficient, semantically equivalent instruction sequences;
2. **Register-allocation:** Best use of the available processor registers
3. **Instruction Scheduling:** Reordering of the instruction stream to exploit intra-processor parallelism

For several reasons, e.g. modularization of code generation and portability, code generation may be split into **two phases**:



Abstract machine

- idealized architecture,
- simple code generation,
- easily implemented on real hardware.

Advantages:

- Porting the compiler to a new target architecture is simpler,
- Modularization makes the compiler easier to modify,
- Translation of program constructs is separated from the exploitation of architectural features.

Abstract machines for some programming languages:

Pascal	→	P-machine	
Smalltalk	→	Bytecode	
Prolog	→	WAM	(“Warren Abstract Machine”)
SML, Haskell	→	STGM	
Java	→	JVM	

We will consider the following languages and abstract machines:

C	→	CMa	//	<i>imperative</i>
PuF	→	MaMa	//	<i>functional</i>
Proll	→	WiM	//	<i>logic based</i>
multi-threaded C	→	threaded CMa	//	<i>concurrent</i>

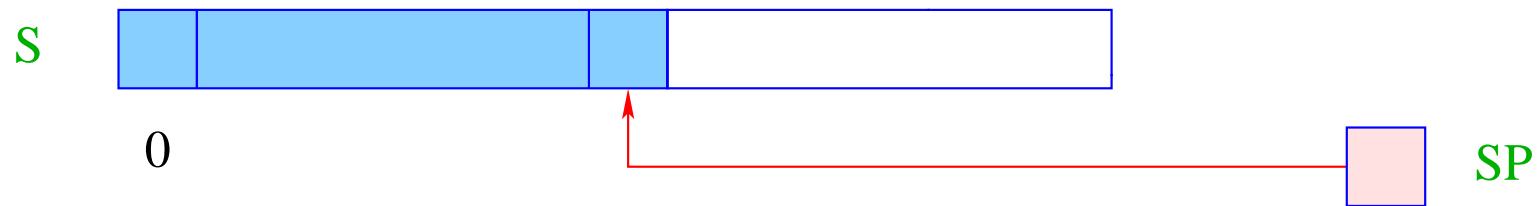
The Translation of C

1 The Architecture of the CMa

- Each abstract machine provides a set of **instructions**
- Instructions are executed on the abstract hardware
- This abstract hardware can be viewed as a set of data structures, which the instructions access
- ... and which are managed by the **run-time system**

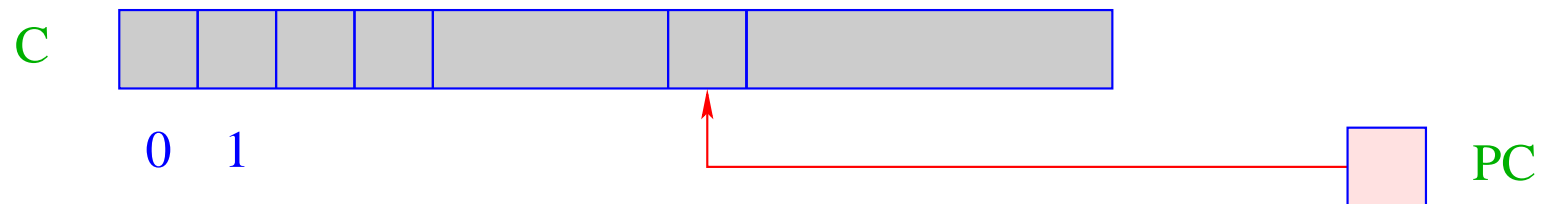
For the CMa we need:

The Data Store:



- **S** is the (data) store, onto which new cells are allocated in a LIFO discipline
⇒ **Stack**.
- **SP** ($\hat{=}$ **Stack Pointer**) is a register, which contains the address of the topmost allocated cell,
Simplification: All types of data fit into one cell of **S**.

The Code/Instruction Store:



- **C** is the Code store, which contains the program.
Each cell of field **C** can store exactly one abstract instruction.
- **PC** ($\hat{=}$ Program Counter) is a register, which contains the address of the instruction to be executed *next*.
- Initially, **PC** contains the address 0.
 \Rightarrow **C**[0] contains the instruction to be executed first.

Execution of Programs:

- The machine loads the instruction in $C[PC]$ into a **Instruction-Register IR** and executes it
- **PC** is incremented by 1 before the execution of the instruction

```
while (true) {  
    IR = C[PC]; PC++;  
    execute (IR);  
}
```

- The execution of the instruction may overwrite the **PC** (jumps).
- The **Main Cycle** of the machine will be halted by executing the instruction **halt** , which returns control to the environment, e.g. the operating system
- More instructions will be introduced **by demand**

2 Simple expressions and assignments

Problem: evaluate the expression $(1 + 7) * 3$!

This means: generate an instruction sequence, which

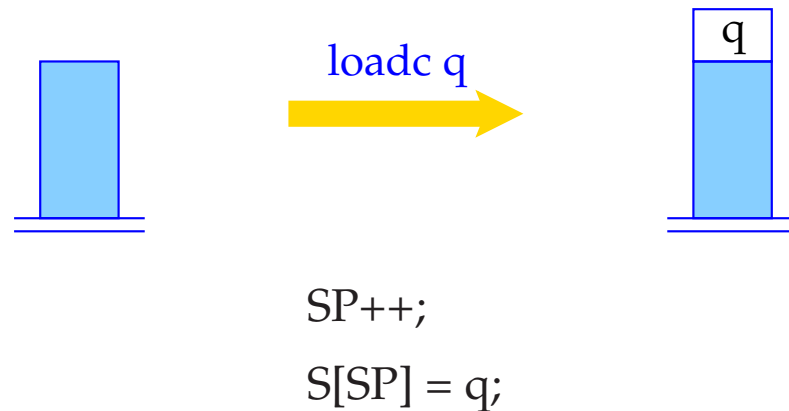
- determines the value of the expression and
- pushes it on top of the stack...

Idea:

- first compute the values of the subexpressions,
- save these values on top of the stack,
- then apply the operator.

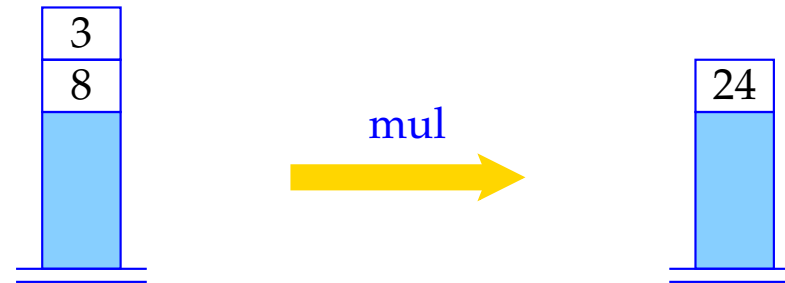
The general principle:

- instructions expect their arguments on top of the stack,
- execution of an instruction consumes its operands,
- results, if any, are stored on top of the stack.



Instruction `loadc q` needs no operand on top of the stack, pushes the constant `q` onto the stack.

Note: the content of register `SP` is only implicitly represented, namely through the height of the stack.



SP--;

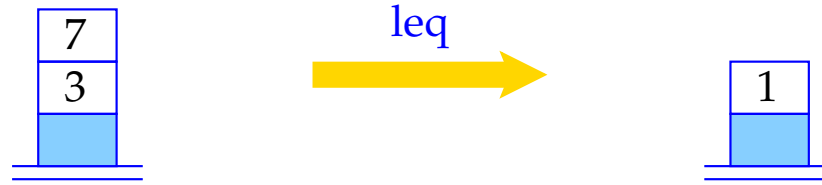
$S[SP] = S[SP] * S[SP+1];$

mul expects two operands on top of the stack, consumes both, and pushes their product onto the stack.

... the other binary arithmetic and logical instructions, **add**, **sub**, **div**, **mod**, **and**, **or** and **xor**, work analogously, as do the comparison instructions **eq**, **neq**, **le**, **leq**, **gr** and **geq**.

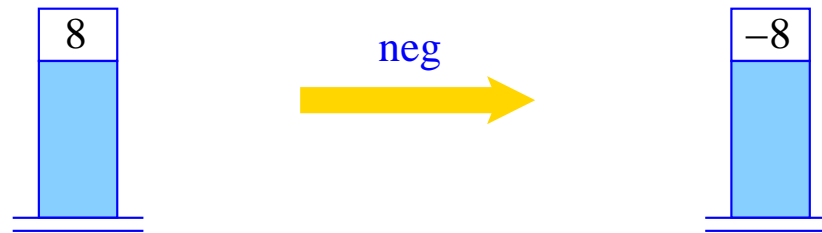
Example:

The operator `leq`



Remark: 0 represents *false*, all other integers *true*.

Unary operators `neg` and `not` consume one operand and produce one result.



$$S[SP] = -S[SP];$$

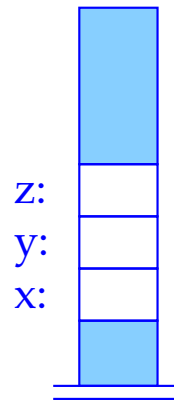
Example: Code for $1 + 7$:

loadc 1 loadc 7 add

Execution of this code sequence:



Variables are associated with cells in S :



Code generation will be described by some **Translation Functions**, $code$, $code_L$, and $code_R$.

Arguments: A program construct and a function ρ . ρ delivers for each variable x the relative address of x . ρ is called **Address Environment**.

Variables can be used in two different ways:

Example: $x = y + 1$

We are interested in the **value** of y , but in the **address** of x .

The syntactic position determines, whether the **L-value** or the **R-value** of a variable is required.

L-value of x = **address** of x

R-value of x = **content** of x

$\text{code}_R e \rho$	produces code to compute the R-value of e in the address environment ρ
$\text{code}_L e \rho$	analogously for the L-value

Note:

Not every expression has an L-value (Ex.: $x + 1$).

We define:

$$\begin{aligned} \text{code}_R (e_1 + e_2) \rho &= \text{code}_R e_1 \rho \\ &\quad \text{code}_R e_2 \rho \\ &\quad \text{add} \end{aligned}$$

... analogously for the other binary operators

$$\begin{aligned} \text{code}_R (-e) \rho &= \text{code}_R e \rho \\ &\quad \text{neg} \end{aligned}$$

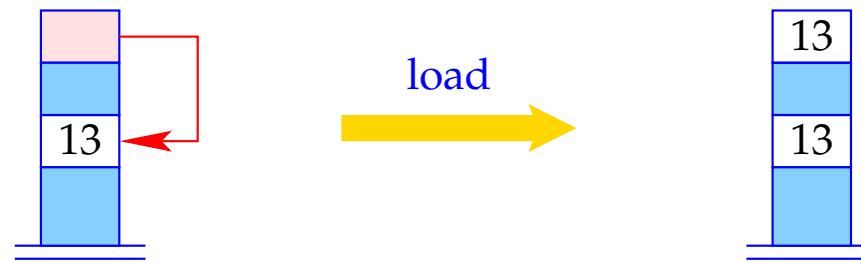
... analogously for the other unary operators

$$\begin{aligned} \text{code}_R q \rho &= \text{loadc } q \\ \text{code}_L x \rho &= \text{loadc } (\rho x) \\ &\quad \dots \end{aligned}$$

$$\text{code}_R \ x \ \rho = \text{code}_L \ x \ \rho$$

load

The instruction `load` loads the contents of the cell, whose address is on top of the stack.



$S[SP] = S[S[SP]];$

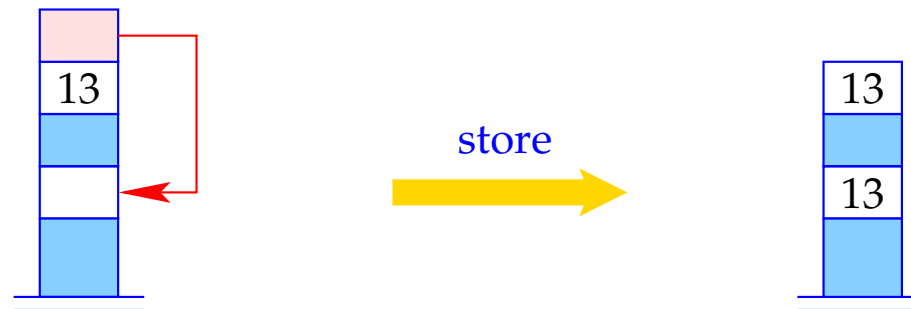
$$\text{code}_R (x = e) \rho = \text{code}_R e \rho$$

$$\text{code}_L x \rho$$

store

store writes the contents of the second topmost stack cell into the cell, whose address is on top of the stack, and leaves the written value on top of the stack.

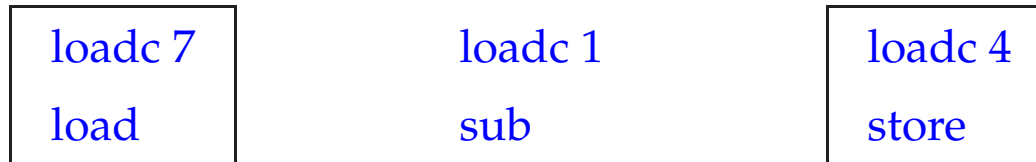
Note: this is different from the corresponding store-instruction of the P-machine in Wilhelm/Maurer!



$$S[S[SP]] = S[SP-1];$$

$$SP--;$$

Example: Code for $e \equiv x = y - 1$ with $\rho = \{x \mapsto 4, y \mapsto 7\}$.
 $\text{code}_R e \rho$ produces:



Improvements:

Introduction of special instructions for frequently used instruction sequences,
e.g.,

loada q	=	loadc q load
storea q	=	loadc q store

3 Statements and Statement Sequences

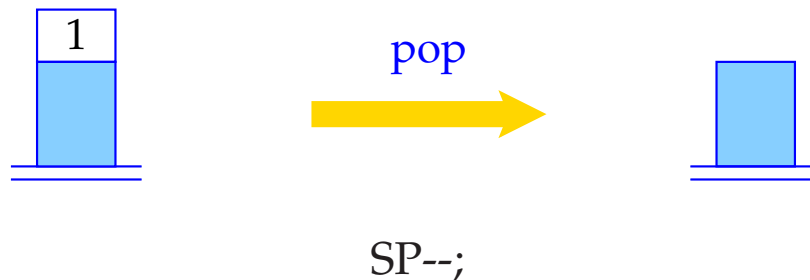
Is e an expression, then $e;$ is a statement.

Statements do not deliver a value. The contents of the **SP** before and after the execution of the generated code must therefore be the same.

$$\text{code } e; \rho = \text{code}_R e \rho$$

pop

The instruction **pop** eliminates the top element of the stack.

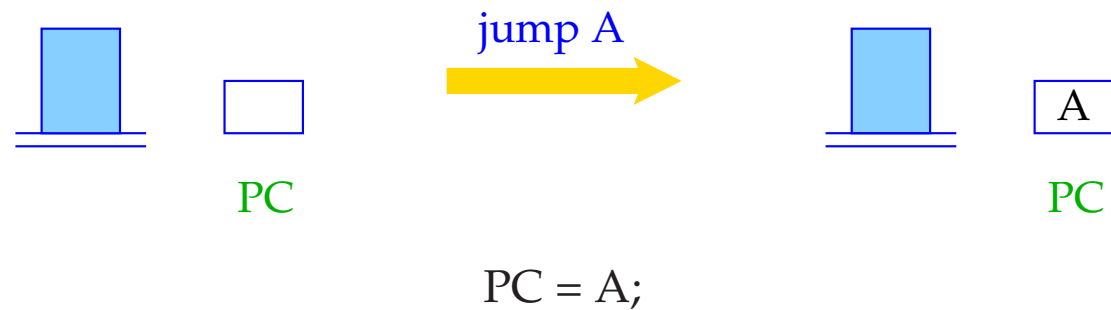


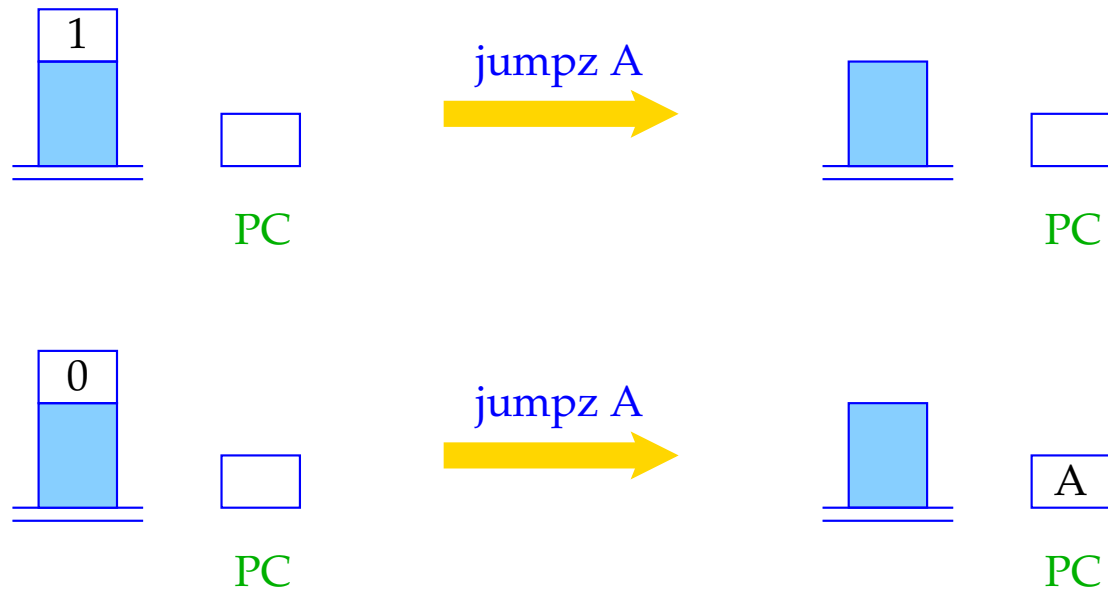
The code for a statement sequence is the concatenation of the code for the statements of the sequence:

$$\begin{aligned} \text{code } (s \text{ } ss) \rho &= \text{code } s \rho \\ &\quad \text{code } ss \rho \\ \text{code } \varepsilon \rho &= // \text{ empty sequence of instructions} \end{aligned}$$

4 Conditional and Iterative Statements

We need jumps to deviate from the serial execution of consecutive statements:





```
if (S[SP] == 0) PC = A;  
SP--;
```

For ease of comprehension, we use **symbolic jump targets**. They will later be replaced by absolute addresses.

Instead of absolute code addresses, one could generate **relative** addresses, i.e., relative to the actual **PC**.

Advantages:

- **smaller addresses** suffice most of the time;
- the code becomes **relocatable**, i.e., can be moved around in memory.

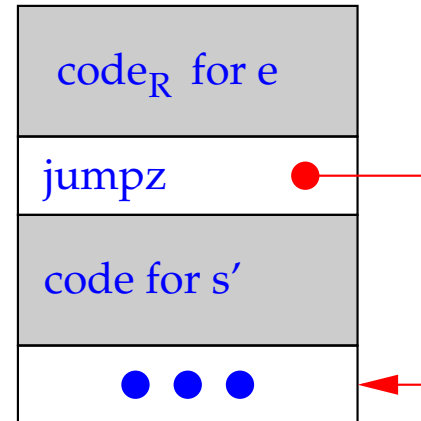
4.1 One-sided Conditional Statement

Let us first regard $s \equiv \mathbf{if} (e) s'$.

Idea:

- Put code for the evaluation of e and s' consecutively in the code store,
- Insert a conditional jump ([jump on zero](#)) in between.

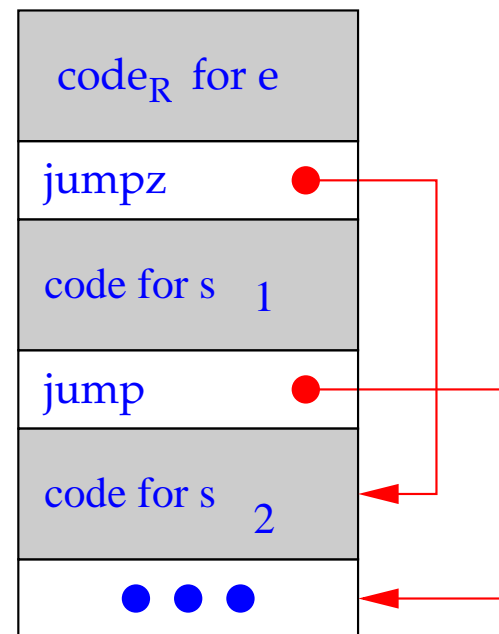
`code` $s \rho$ = `code`_R $e \rho$
`jumpz` A
`code` $s' \rho$
 A : ...



4.2 Two-sided Conditional Statement

Let us now regard $s \equiv \mathbf{if} (e) s_1 \mathbf{else} s_2$. The same strategy yields:

$\mathbf{code} s \rho = \mathbf{code}_R e \rho$
 $\mathbf{jumpz} A$
 $\mathbf{code} s_1 \rho$
 $\mathbf{jump} B$
A : $\mathbf{code} s_2 \rho$
B : ...



Example:

Be $\rho = \{x \mapsto 4, y \mapsto 7\}$ and

$s \equiv$ **if** $(x > y)$ (i)
 $x = x - y;$ (ii)
 else $y = y - x;$ (iii)

code s ρ produces:

loada 4
loada 7
ge
jumpz A

(i)

loada 4
loada 7
sub
storea 4
pop
jump B

(ii)

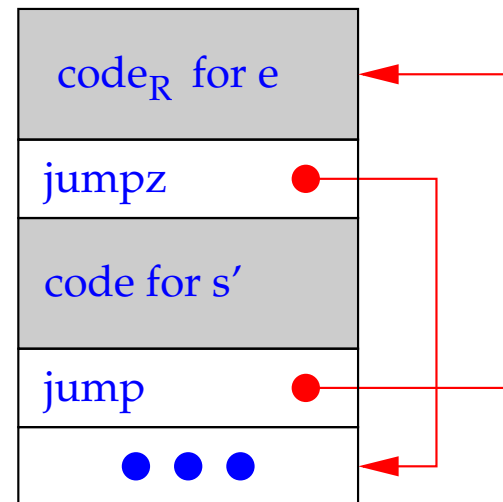
A: loada 7
loada 4
sub
storea 7
pop
B: ...

(iii)

4.3 while-Loops

Let us regard the loop $s \equiv \mathbf{while} (e) s'$. We generate:

```
code  $s \rho$  =  
A : codeR  $e \rho$   
    jumpz B  
    code  $s' \rho$   
    jump A  
B : ...
```



Example:

Be $\rho = \{a \mapsto 7, b \mapsto 8, c \mapsto 9\}$ and s the statement:

while $(a > 0) \{c = c + 1; a = a - b;\}$

code $s \rho$ produces the sequence:

A:	loada 7	loada 9	loada 7	B: ...
	loadc 0	loadc 1	loada 8	
	ge	add	sub	
	jumpz B	storea 9	storea 7	
		pop	pop	
			jump A	

4.4 for-Loops

The **for**-loop $s \equiv \mathbf{for} (e_1; e_2; e_3) s'$ is equivalent to the statement sequence $e_1; \mathbf{while} (e_2) \{s' e_3; \}$ – provided that s' contains no **continue**-statement.

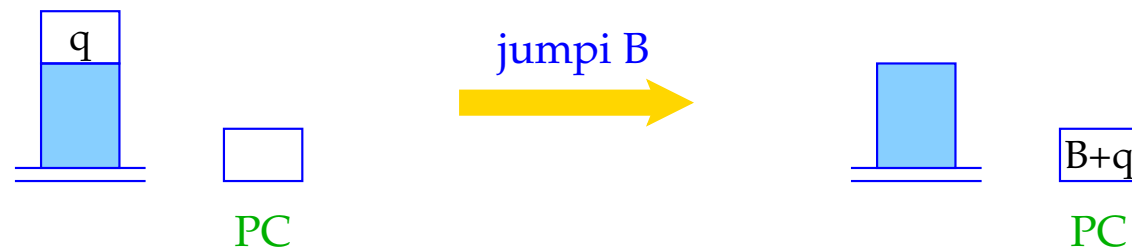
We therefore translate:

```
code s ρ = codeR e1
           pop
           A : codeR e2 ρ
               jumpz B
               code s' ρ
               codeR e3 ρ
               pop
               jump A
           B : ...
```

4.5 The switch-Statement

Idea:

- Multi-target branching in **constant time!**
- Use a **jump table**, which contains at its i -th position the jump to the beginning of the i -th alternative.
- Realized by **indexed jumps**.



$PC = B + S[SP];$
 $SP--;$

Simplification:

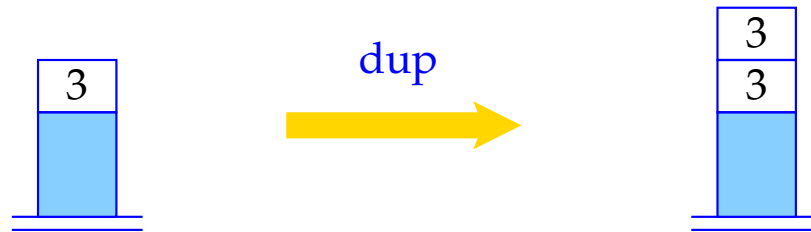
We only regard **switch**-statements of the following form:

$$s \equiv \text{switch } (e) \{$$
$$\quad \text{case } 0: \quad ss_0 \text{ break;}$$
$$\quad \text{case } 1: \quad ss_1 \text{ break;}$$
$$\quad \quad \quad \vdots$$
$$\quad \text{case } k - 1: \quad ss_{k-1} \text{ break;}$$
$$\quad \text{default: } \quad ss_k$$
$$\quad \quad \quad \}$$

s is then translated into the instruction sequence:

<code>check 0 k B</code>	=	<code>dup</code>	<code>dup</code>	<code>jumpi B</code>
		<code>loadc 0</code>	<code>loadc k</code>	<code>A: pop</code>
		<code>geq</code>	<code>le</code>	<code>loadc k</code>
		<code>jumpz A</code>	<code>jumpz A</code>	<code>jumpi B</code>

- The R-value of e is still needed for indexing after the comparison. It is therefore copied before the comparison.
- This is done by the instruction `dup`.
- The R-value of e is replaced by k before the indexed jump is executed if it is less than 0 or greater than k .



$S[SP+1] = S[SP];$

$SP++;$

Note:

- The jump table could be placed directly after the code for the Macro `check`. This would save a few unconditional jumps. However, it may require to search the `switch`-statement twice.
- If the table starts with u instead of 0, we have to decrease the R-value of e by u before using it as an index.
- If all potential values of e are `definitely` in the interval $[0, k]$, the macro `check` is not needed.

5 Storage Allocation for Variables

Goal:

Associate **statically**, i.e. at compile time, with each variable x a fixed (relative) address ρx

Assumptions:

- variables of basic types, e.g. **int**, ... occupy one storage cell.
- variables are allocated in the store in the order, in which they are declared, starting at address 1.

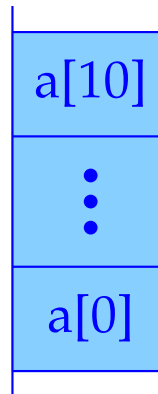
Consequently, we obtain for the declaration $d \equiv t_1 x_1; \dots t_k x_k$ (t_i basic type) the address environment ρ such that

$$\rho x_i = i, \quad i = 1, \dots, k$$

5.1 Arrays

Example: `int [11] a;`

The array a consists of 11 components and therefore needs 11 cells.
 ρa is the address of the component $a[0]$.



We need a function `sizeof` (notation: $|\cdot|$), computing the space requirement of a type:

$$|t| = \begin{cases} 1 & \text{if } t \text{ basic} \\ k \cdot |t'| & \text{if } t \equiv t'[k] \end{cases}$$

Accordingly, we obtain for the declaration $d \equiv t_1 x_1; \dots t_k x_k$

$$\begin{aligned} \rho x_1 &= 1 \\ \rho x_i &= \rho x_{i-1} + |t_{i-1}| && \text{for } i > 1 \end{aligned}$$

Since $|\cdot|$ can be computed at compile time, also ρ can be computed at compile time.

Task:

Extend `codeL` and `codeR` to expressions with accesses to array components.

Be `t[c] a;` the declaration of an array `a`.

To determine the start address of a component `a[i]`, we compute $\rho a + |t| * (R\text{-value of } i)$.

In consequence:

```
codeL a[e] ρ = loadc (ρ a)
               codeR e ρ
               loadc |t|
               mul
               add
```

... or more general:

$$\text{code}_L e_1[e_2] \rho = \begin{array}{l} \text{code}_R e_1 \rho \\ \text{code}_R e_2 \rho \\ \text{loadc } |t| \\ \text{mul} \\ \text{add} \end{array}$$

Remark:

- In **C**, an array is a **pointer**. A declared array a is a **pointer-constant**, whose R-value is the start address of the array.
- Formally, we define for an array e : $\text{code}_R e \rho = \text{code}_L e \rho$
- In **C**, the following are equivalent (as L-values):

$$2[a] \quad a[2] \quad a + 2$$

Normalization: Array names and expressions evaluating to arrays occur in front of index brackets, index expressions inside the index brackets.

5.2 Structures

In **Modula** and **Pascal**, structures are called **Records**.

Simplification:

Names of structure components are not used elsewhere.

Alternatively, one could manage a separate environment ρ_{st} for each structure type st .

Be `struct { int a ; int b ; } x ;` part of a declaration list.

- x has as relative address the address of the first cell allocated for the structure.
- The components have addresses **relative** to the start address of the structure. In the example, these are $a \mapsto 0, b \mapsto 1$.

Let $t \equiv \mathbf{struct} \{t_1 c_1; \dots t_k c_k\}$. We have

$$|t| = \sum_{i=1}^k |t_i|$$

$$\rho c_1 = 0 \quad \text{and}$$

$$\rho c_i = \rho c_{i-1} + |t_{i-1}| \quad \text{for } i > 1$$

We thus obtain:

$$\mathbf{code}_L(e.c) \rho = \mathbf{code}_L e \rho$$

$$\mathbf{loadc}(\rho c)$$

$$\mathbf{add}$$

Example:

Be `struct { int a; int b; } x;` such that $\rho = \{x \mapsto 13, a \mapsto 0, b \mapsto 1\}$.

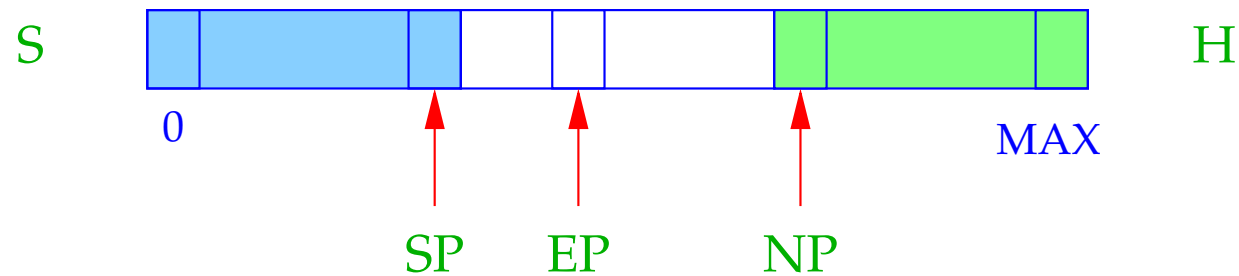
This yields:

```
codeL (x.b) ρ = loadc 13  
                loadc 1  
                add
```

6 Pointer and Dynamic Storage Management

Pointer allow the access to anonymous, dynamically generated objects, whose life time is not subject to the LIFO-principle.

⇒ We need another potentially unbounded storage area H – the Heap.



NP $\hat{=}$ **New Pointer**; points to the lowest occupied heap cell.

EP $\hat{=}$ **Extreme Pointer**; points to the uppermost cell, to which SP can point (during execution of the actual function).

Idea:

- Stack and Heap grow toward each other in S, but must not collide. ([Stack Overflow](#)).
- A collision may be caused by an increment of **SP** or a decrement of **NP**.
- **EP** saves us the check for collision at the stack operations.
- The checks at heap allocations are still necessary.

What can we do with pointers (pointer values)?

- **set** a pointer to a storage cell,
- **dereference** a pointer, access the value in a storage cell pointed to by a pointer.

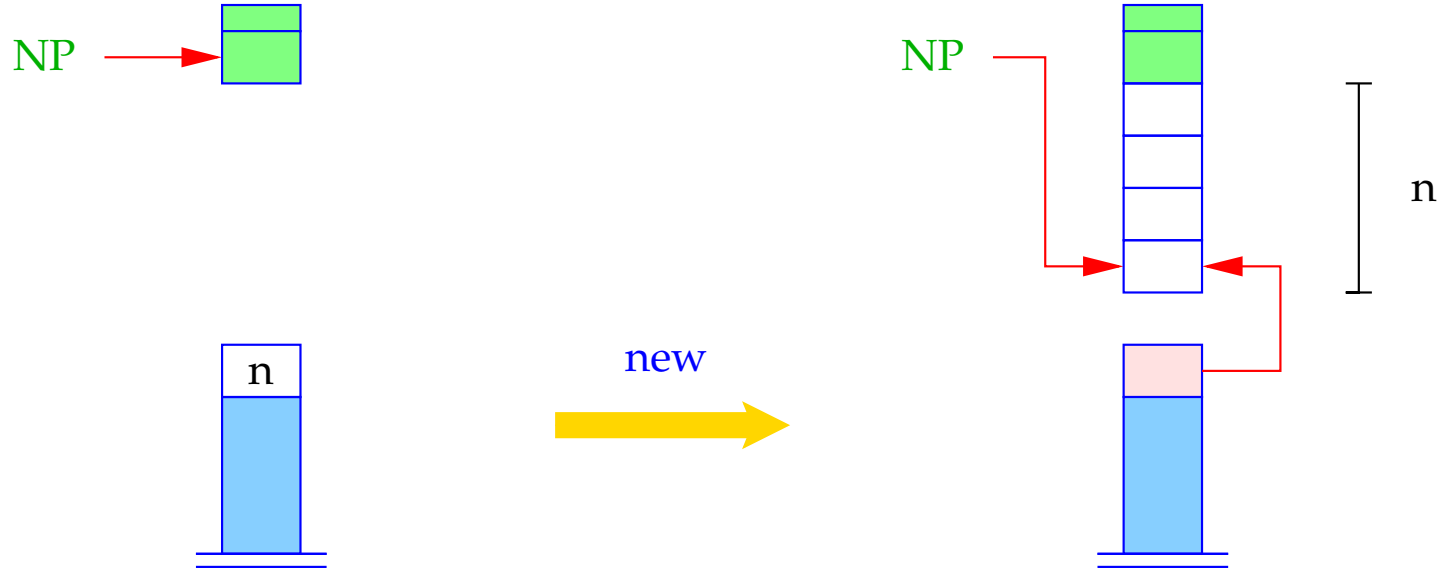
There are two ways to set a pointer:

- (1) A call **malloc**(e) reserves a heap area of the size of the value of e and returns a pointer to this area:

$$\text{code}_R \text{ malloc}(e) \rho = \text{code}_R e \rho \text{ new}$$

- (2) The application of the address operator **&** to a variable returns a **pointer** to this variable, i.e. its address ($\hat{=}$ **L-value**). Therefore:

$$\text{code}_R (\&e) \rho = \text{code}_L e \rho$$



```

if (NP - S[SP] ≤ EP)
    S[SP] = NULL;
else {
    NP = NP - S[SP];
    S[SP] = NP;
}

```

- NULL is a special pointer constant, identified with the integer constant 0.
- In the case of a collision of stack and heap the NULL-pointer is returned.

Dereferencing of Pointers:

The application of the operator `*` to the expression e returns the **contents** of the storage cell, whose address is the R-value of e :

$$\text{code}_L (*e) \rho = \text{code}_R e \rho$$

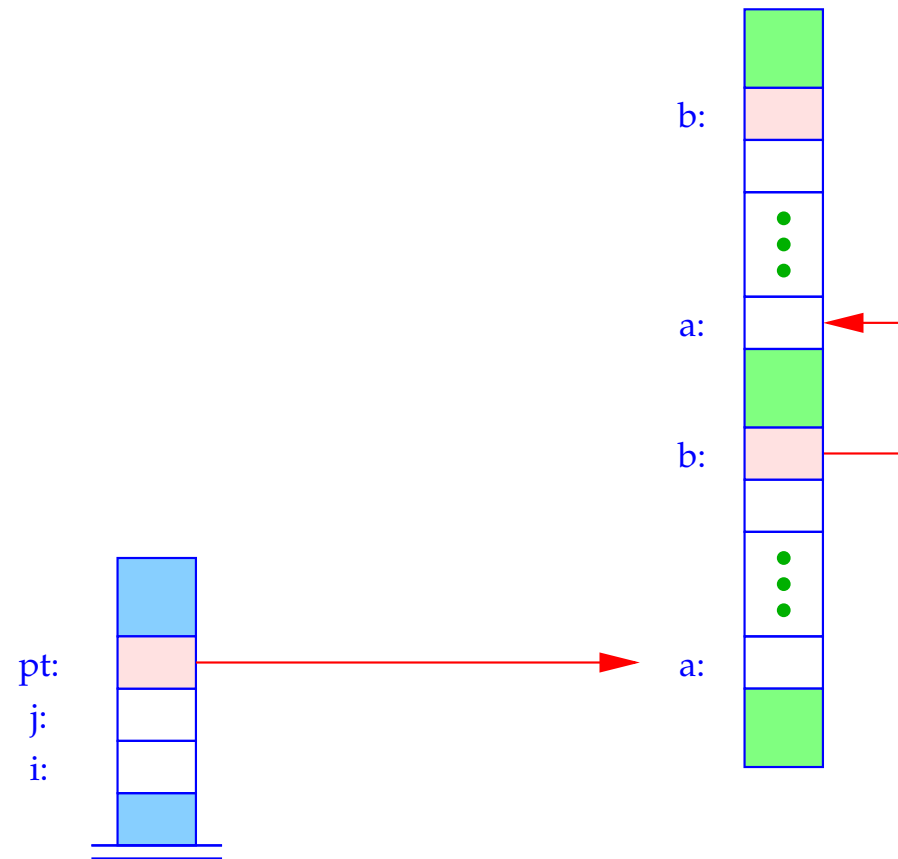
Example: Given the declarations

```
struct  $t$  { int  $a[7]$ ; struct  $t$   $*b$ ; };  
int  $i, j$ ;  
struct  $t$   $*pt$ ;
```

and the expression $((pt \rightarrow b) \rightarrow a)[i + 1]$

Because of $e \rightarrow a \equiv (*e).a$ holds:

$$\begin{aligned} \text{code}_L (e \rightarrow a) \rho &= \text{code}_R e \rho \\ &\quad \text{loadc} (\rho a) \\ &\quad \text{add} \end{aligned}$$



Be $\rho = \{i \mapsto 1, j \mapsto 2, pt \mapsto 3, a \mapsto 0, b \mapsto 7\}$. Then:

$$\begin{aligned}
 & \text{code}_L((pt \rightarrow b) \rightarrow a)[i + 1] \rho \\
 = & \text{code}_R((pt \rightarrow b) \rightarrow a) \rho & = & \text{code}_R((pt \rightarrow b) \rightarrow a) \rho \\
 & \text{code}_R(i + 1) \rho & & \text{loada 1} \\
 & \text{loadc 1} & & \text{loadc 1} \\
 & \text{mul} & & \text{add} \\
 & \text{add} & & \text{loadc 1} \\
 & & & \text{mul} \\
 & & & \text{add}
 \end{aligned}$$

For arrays, their R-value equals their L-value. Therefore:

$$\text{code}_R((pt \rightarrow b) \rightarrow a) \rho = \text{code}_R(pt \rightarrow b) \rho = \begin{array}{l} \text{loada } 3 \\ \text{loadc } 7 \\ \text{add} \\ \text{load} \\ \text{loadc } 0 \\ \text{add} \end{array}$$

In total, we obtain the instruction sequence:

loada 3	load	loada 1	loadc 1
loadc 7	loadc 0	loadc 1	mul
add	add	add	add

7 Conclusion

We tabulate the cases of the translation of expressions:

$$\begin{aligned} \text{code}_L (e_1[e_2]) \rho &= \text{code}_R e_1 \rho \\ &\quad \text{code}_R e_2 \rho \\ &\quad \text{loadc } |t| \\ &\quad \text{mul} \\ &\quad \text{add} \qquad \text{if } e_1 \text{ has type } t^* \text{ or } t[] \end{aligned}$$

$$\begin{aligned} \text{code}_L (e.a) \rho &= \text{code}_L e \rho \\ &\quad \text{loadc } (\rho a) \\ &\quad \text{add} \end{aligned}$$

$$\text{code}_L (*e) \rho = \text{code}_R e \rho$$

$$\text{code}_L x \rho = \text{loadc} (\rho x)$$

$$\text{code}_R (\&e) \rho = \text{code}_L e \rho$$

$$\text{code}_R e \rho = \text{code}_L e \rho \quad \text{if } e \text{ is an array}$$

$$\begin{aligned} \text{code}_R (e_1 \square e_2) \rho &= \text{code}_R e_1 \rho \\ &\quad \text{code}_R e_2 \rho \\ &\quad \text{op} \end{aligned}$$

op instruction for operator ' \square '

$\text{code}_R q \rho = \text{load } q \quad q \text{ constant}$

$\text{code}_R (e_1 = e_2) \rho = \text{code}_R e_2 \rho$
 $\text{code}_L e_1 \rho$
 store

$\text{code}_R e \rho = \text{code}_L e \rho$
 $\text{load} \quad \text{otherwise}$

Example: `int a[10], *b;` with $\rho = \{a \mapsto 7, b \mapsto 17\}$.

For the statement: `*a = 5;` we obtain:

$$\begin{aligned} \text{code}_L(*a) \rho &= \text{code}_R a \rho = \text{code}_L a \rho = \text{loadc } 7 \\ \text{code}(*a = 5;) \rho &= \text{loadc } 5 \\ &\quad \text{loadc } 7 \\ &\quad \text{store} \\ &\quad \text{pop} \end{aligned}$$

As an exercise translate:

$$s_1 \equiv b = (&a) + 2; \quad \text{and} \quad s_2 \equiv *(b + 3) = 5;$$


```

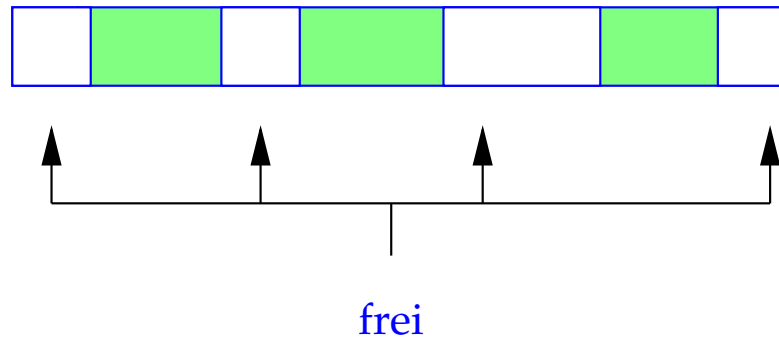
code (s1s2) ρ =   loadc 7           loadc 5
                   loadc 2           loadc 17
                   loadc 10 // size of int[10]  load
                   mul // scaling             loadc 3
                   add                       loadc 1 // size of int
                   loadc 17                  mul // scaling
                   store                      add
                   pop // end of s1         store
                                           pop // end of s2

```

8 Freeing Occupied Storage

Problems:

- The freed storage area is still referenced by other pointers ([dangling references](#)).
- After several deallocations, the storage could look like this ([fragmentation](#)):



Potential Solutions:

- Trust the programmer. Manage freed storage in a particular data structure (free list) \implies `malloc` or `free` may become expensive.
- Do nothing, i.e.:

`code free (e); ρ` = `codeR e ρ`
`pop`

\implies simple and (in general) efficient.

- Use an `automatic`, potentially “conservative” `Garbage-Collection`, which occasionally collects `certainly` inaccessible heap space.

9 Functions

The definition of a function consists of

- a **name**, by which it can be called,
- a specification of the **formal parameters**;
- maybe a **result type**;
- a **statement part**, the **body**.

For **C** holds:

$$\text{code}_R f \rho = \text{loadc } _f = \text{starting address of the code for } f$$

⇒⇒ Function names must also be managed in the address environment!

Example:

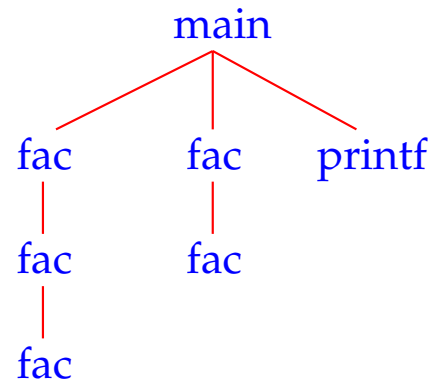
```
int fac (int x) {  
    if (x ≤ 0) return 1;  
    else return x * fac(x - 1);  
}
```

```
main () {  
    int n;  
    n = fac(2) + fac(1);  
    printf ("%d", n);  
}
```

At any time during the execution, several **instances** of one function may exist, i.e., may have started, but not finished execution.

An instance is created by a call to the function.

The recursion tree in the example:



We conclude:

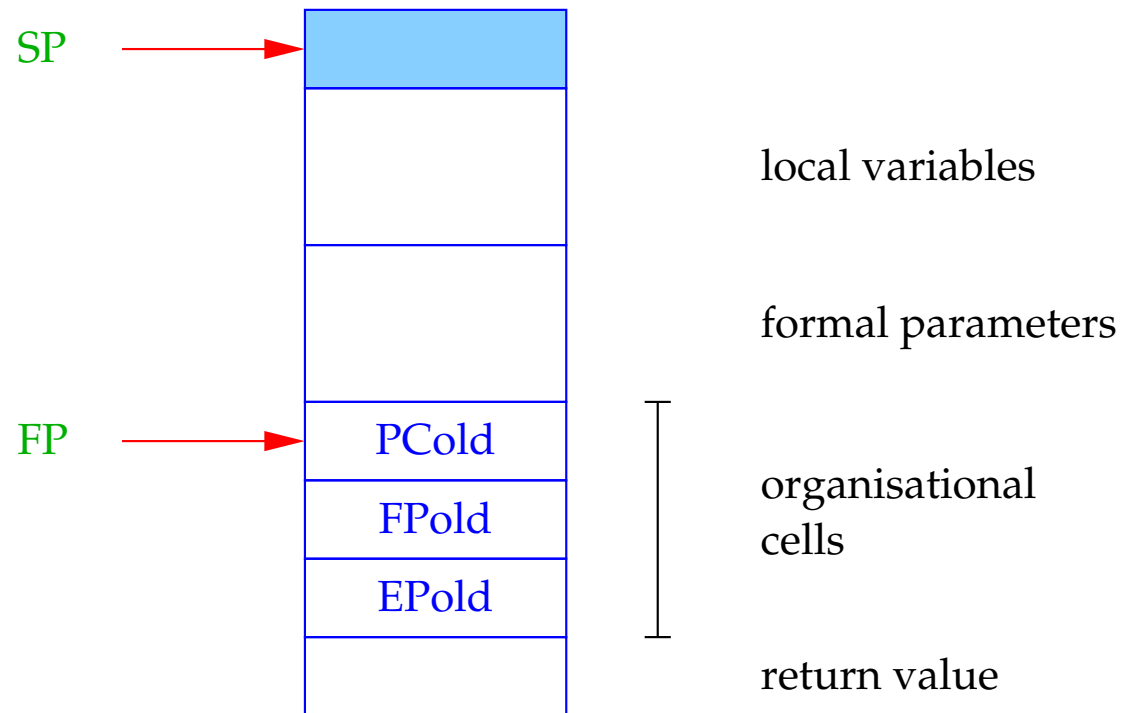
The **formal parameters** and **local variables** of the different **instances** of the same function must be kept separate.

Idea:

Allocate a special storage area for each instance of a function.

In sequential programming languages these storage areas can be managed on a stack. They are therefore called **Stack Frames**.

9.1 Storage Organization for Functions



FP $\hat{=}$ **Frame Pointer**; points to the last **organizational cell** and is used to address the formal parameters and the local variables.

The caller must be able to continue execution in its frame after the return from a function. Therefore, at a function call the following values have to be saved into [organizational cells](#):

- the [FP](#)
- the [continuation address](#) after the call and
- the actual [EP](#).

[Simplification](#): The return value fits into one storage cell.

[Translation tasks for functions](#):

- Generate code for the body!
- Generate code for calls!

9.2 Computing the Address Environment

We have to distinguish two different kinds of variables:

1. **globals**, which are defined externally to the functions;
2. **locals**/automatic (including formal parameters), which are defined internally to the functions.



The address environment ρ associates pairs $(tag, a) \in \{G, L\} \times \mathbb{N}_0$ with their names.

Note:

- There exist more refined notions of visibility of (the defining occurrences of) variables, namely **nested blocks**.
- The translation of different program parts in general uses different address environments!

Example (1):

```
0  int i;
    struct list {
        int info;
        struct list * next;
    } * l;

1  int ith (struct list * x, int i) {
    if (i ≤ 1) return x →info;
    else return ith (x →next, i - 1);
}
```

```
2  main () {
    int k;
    scanf ("%d", &i);
    scanlist (&l);
    printf ("\n\t%d\n", ith (l,i));
}
```

address	environment	at	0
ρ_0 :	i	\mapsto	$(G, 1)$
	l	\mapsto	$(G, 2)$
	ith	\mapsto	$(G, _ith)$
	$main$	\mapsto	$(G, _main)$
			...

Example (2):

```
0  int i;
    struct list {
        int info;
        struct list * next;
    } * l;

1  int ith (struct list * x, int i) {
    if (i ≤ 1) return x → info;
    else return ith (x → next, i - 1);
}
```

```
2  main () {
    int k;
    scanf ("%d", &i);
    scanlist (&l);
    printf ("\n\t%d\n", ith (l,i));
}
```

1	inside	of	ith:
$\rho_1 :$	i	\mapsto	$(L, 2)$
	x	\mapsto	$(L, 1)$
	l	\mapsto	$(G, 2)$
	ith	\mapsto	(G, ith)
	main	\mapsto	(G, main)
			...

Example (3):

```
0  int i;
    struct list {
        int info;
        struct list * next;
    } * l;

1  int ith (struct list * x, int i) {
    if (i ≤ 1) return x →info;
    else return ith (x →next, i - 1);
}
```

```
2  main () {
    int k;
    scanf ("%d", &i);
    scanlist (&l);
    printf ("\n\t%d\n", ith (l,i));
}
```

2	inside	of	main:
$\rho_2 :$	i	\mapsto	$(G, 1)$
	l	\mapsto	$(G, 2)$
	k	\mapsto	$(L, 1)$
	ith	\mapsto	$(G, _ith)$
	$main$	\mapsto	$(G, _main)$
			...

9.3 Calling/Entering and Leaving Functions

Be f the actual function, the **Caller**, and let f call the function g , the **Callee**.

The code for a function call has to be distributed among the Caller and the Callee:

The distribution depends on **who** has **which** information.

Actions upon **calling/entering** g :

1. Saving **FP, EP** } **mark**
 2. Computing the actual parameters
 3. Determining the start address of g
 4. Setting the new **FP**
 5. Saving **PC** and } **call**
jump to the beginning of g
 6. Setting the new **EP** } **enter**
 7. Allocating the local variables } **alloc**
- } available in f
- } available in g

Actions upon **leaving** g :

1. Restoring the registers **FP, EP, SP**
 2. Returning to the code of f , i.e. restoring the **PC**
- } **return**

Altogether we generate for a call:

$$\begin{aligned} \text{code}_R g(e_1, \dots, e_n) \rho &= \text{mark} \\ &\quad \text{code}_R e_1 \rho \\ &\quad \dots \\ &\quad \text{code}_R e_m \rho \\ &\quad \text{code}_R g \rho \\ &\quad \text{call } n \end{aligned}$$

where n = space for the actual parameters

Note:

- Expressions occurring as actual parameters will be evaluated to their **R-value** \implies **Call-by-Value**-parameter passing.
- Function g can also be an **expression**, whose **R-value** is the start address of the function to be called ...

- Function names are regarded as **constant pointers** to functions, similarly to declared arrays. The R-value of such a pointer is the start address of the function.

- For a variable `int (*)() g;`, the two calls

`(*g)()` und `g()`

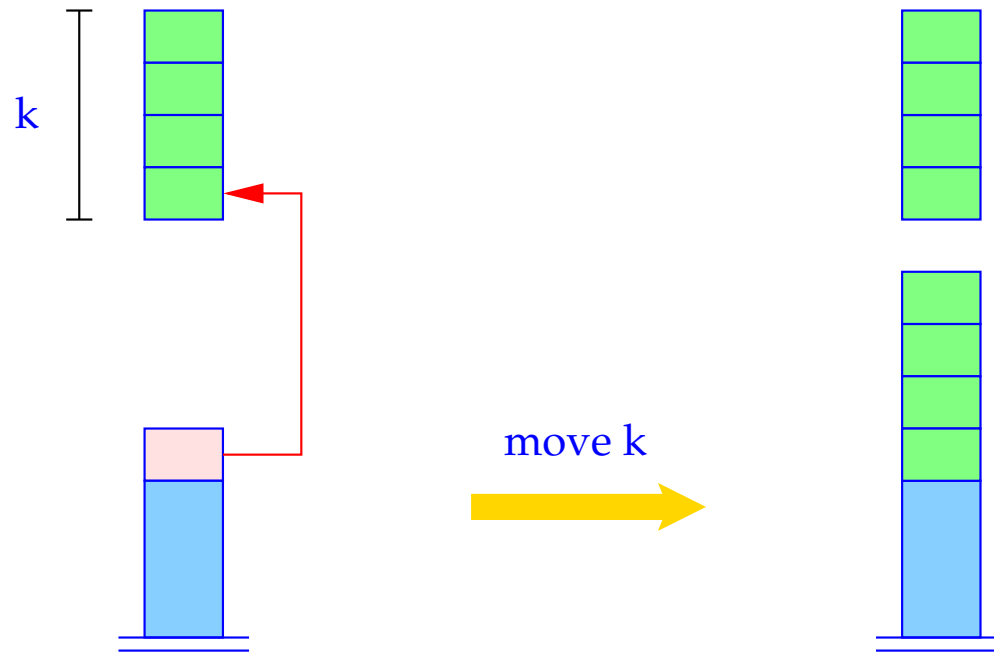
are equivalent :-)

Normalization: Dereferencing of a function pointer is ignored.

- Structures are copied when they are passed as parameters.

In consequence:

<code>code_R f ρ</code>	<code>= loadc (ρ f)</code>	<code>f</code> a function name
<code>code_R (*e) ρ</code>	<code>= code_R e ρ</code>	<code>e</code> a function pointer
<code>code_R e ρ</code>	<code>= code_L e ρ</code>	
	<code>move k</code>	<code>e</code> a structure of size <code>k</code>

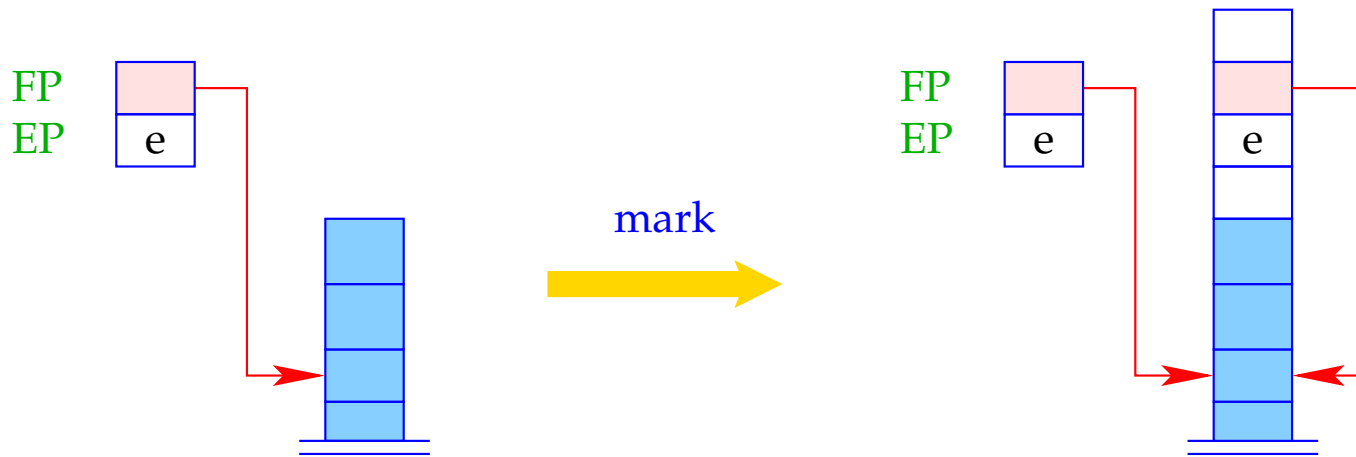


```

for (i = k-1; i ≥ 0; i--)
    S[SP+i] = S[S[SP]+i];
SP = SP+k-1;

```

The instruction `mark` allocates space for the return value and for the organizational cells and saves the FP and EP.

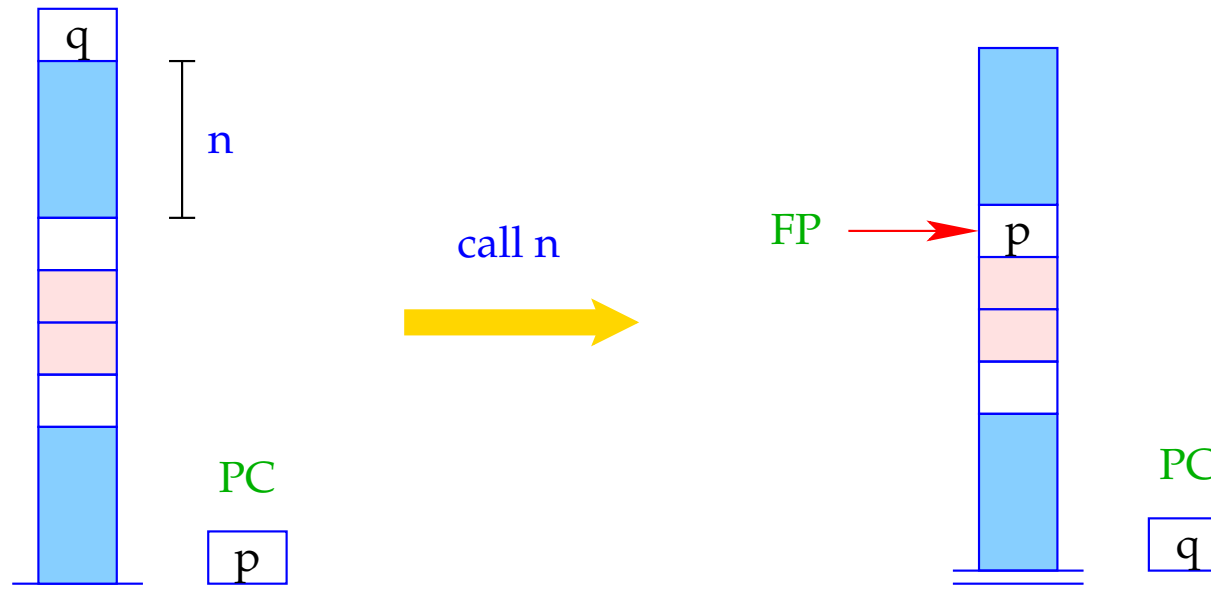


$$S[SP+2] = EP;$$

$$S[SP+3] = FP;$$

$$SP = SP + 4;$$

The instruction `call n` saves the continuation address and assigns `FP`, `SP`, and `PC` their new values.



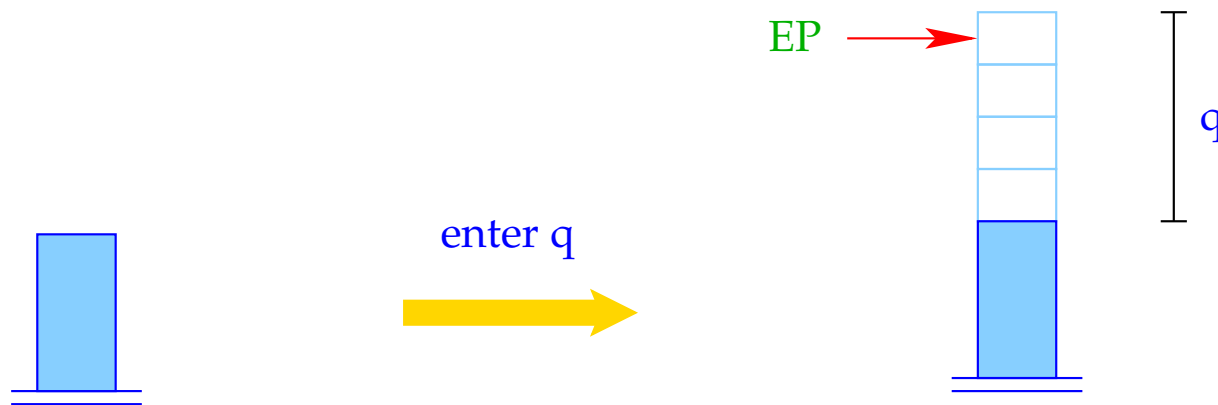
$FP = SP - n - 1;$
 $S[FP] = PC;$
 $PC = S[SP];$
 $SP--;$

Correspondingly, we translate a function definition:

```
code t f (specs){V_defs ss} ρ =  
    _f:  enter q      // Setting the EP  
        alloc k      // Allocating the local variables  
code ss ρf  
return // leaving the function
```

where t = return type of f with $|t| \leq 1$
 q = $maxS + k$ where
 $maxS$ = maximal depth of the local stack
 k = space for the local variables
 ρ_f = address environment for f
// takes care of *specs*, *V_defs* and ρ

The instruction `enter q` sets `EP` to its new value. Program execution is terminated if not enough space is available.

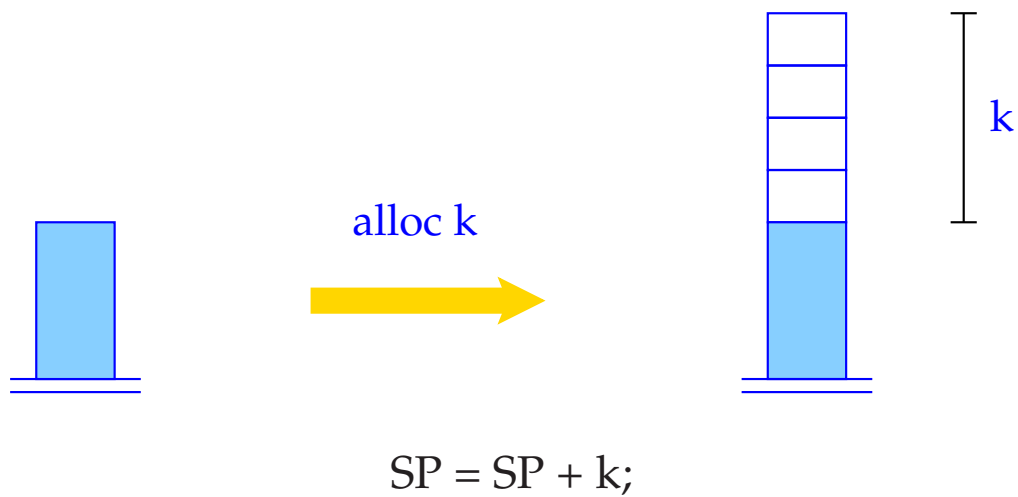


```
EP = SP + q;
```

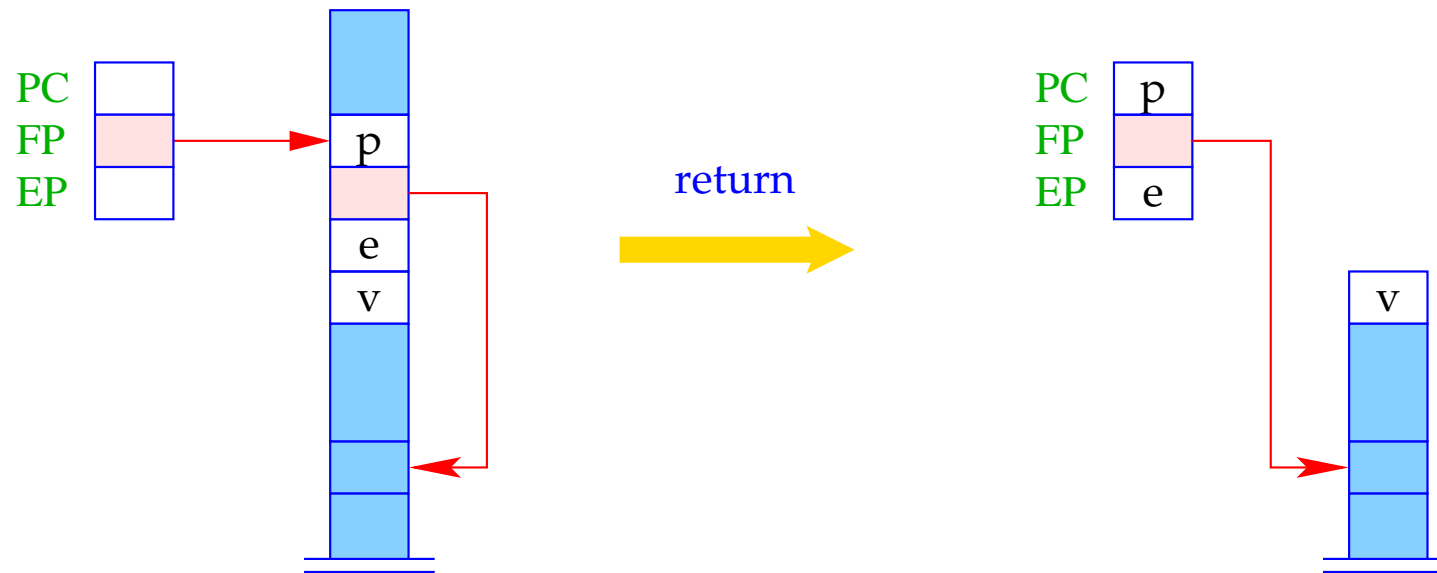
```
if (EP ≥ NP)
```

```
    Error ("Stack Overflow");
```

The instruction `alloc k` reserves stack space for the local variables.



The instruction `return` pops the actual stack frame, i.e., it restores the registers `PC`, `EP`, `SP`, and `FP` and leaves the return value on top of the stack.



$PC = S[FP]; EP = S[FP-2];$

if ($EP \geq NP$) Error ("Stack Overflow");

$SP = FP-3; FP = S[SP+2];$

9.4 Access to Variables and Formal Parameters, and Return of Values

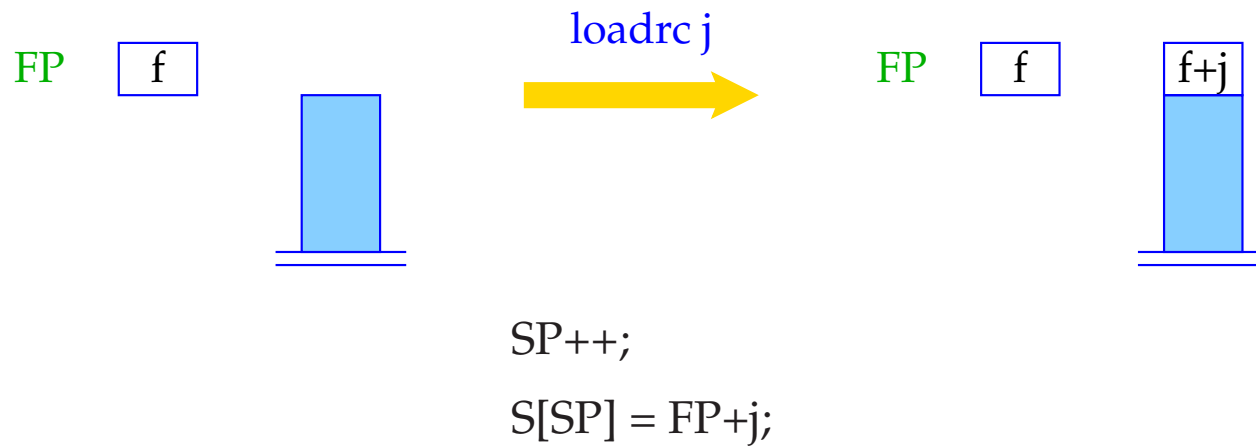
Local variables and formal parameters are addressed relative to the current FP.

We therefore modify `codeL` for the case of variable names.

For $\rho x = (tag, j)$ we define

$$\text{code}_L x \rho = \begin{cases} \text{loadc } j & tag = G \\ \text{loadrc } j & tag = L \end{cases}$$

The instruction `loadrc j` computes the sum of `FP` and `j`.



As an optimization one introduces the instructions `loadr j` and `storer j` .
This is analogous to `loada j` and `storea j`.

$$\text{loadr } j = \begin{array}{l} \text{loadrc } j \\ \text{load} \end{array}$$
$$\text{storer } j = \begin{array}{l} \text{loadrc } j \\ \text{store} \end{array}$$

The code for `return e;` corresponds to an assignment to a variable with relative address `-3`.

$$\text{code return } e; \rho = \begin{array}{l} \text{code}_R e \rho \\ \text{storer } -3 \\ \text{return} \end{array}$$

Example: For the function

```
int fac (int x) {  
    if (x ≤ 0) return 1;  
    else return x * fac (x - 1);  
}
```

we generate:

<code>_fac:</code>	<code>enter q</code>	<code>loadc 1</code>	<code>A:</code>	<code>loadr 1</code>	<code>mul</code>
	<code>alloc 0</code>	<code>storer -3</code>		<code>mark</code>	<code>storer -3</code>
	<code>loadr 1</code>	<code>return</code>		<code>loadr 1</code>	<code>return</code>
	<code>loadc 0</code>	<code>jump B</code>		<code>loadc 1</code>	<code>B:</code> <code>return</code>
	<code>leq</code>			<code>sub</code>	
	<code>jumpz A</code>			<code>loadc _fac</code>	
				<code>call 1</code>	

where $\rho_{\text{fac}} : x \mapsto (L, 1)$ and $q = 1 + 4 + 2 = 7$.

10 Translation of Whole Programs

The state before program execution starts:

$$SP = -1 \quad FP = EP = 0 \quad PC = 0 \quad NP = \text{MAX}$$

Be $p \equiv V_defs \ F_def_1 \dots F_def_n$, a program, where F_def_i defines a function f_i , of which one is named `main`.

The code for the program p consists of:

- Code for the function definitions F_def_i ;
- Code for allocating the global variables;
- Code for the call of `main()`;
- the instruction `halt`.

We thus define:

```
code  $p \ \emptyset$    =   enter (k + 6)
                   alloc (k + 1)
                   mark
                   loadc _main
                   call 0
                   pop
                   halt
                   _f1: code F_def1  $\rho$ 
                       ⋮
                   _fn: code F_defn  $\rho$ 
```

where \emptyset $\hat{=}$ empty address environment;
 ρ $\hat{=}$ global address environment;
 k $\hat{=}$ space for global variables
 $_main \in \{_f_1, \dots, _f_n\}$

The Translation of Functional Programming Languages

11 The language PuF

We only regard a mini-language PuF (“Pure Functions”).

We do not treat, as yet:

- Side effects;
- Data structures.

A Program is an expression e of the form:

$$\begin{aligned} e ::= & b \mid x \mid (\square_1 e) \mid (e_1 \square_2 e_2) \\ & \mid (\mathbf{if} \ e_0 \ \mathbf{then} \ e_1 \ \mathbf{else} \ e_2) \\ & \mid (e' \ e_0 \ \dots \ e_{k-1}) \\ & \mid (\mathbf{fn} \ x_0, \dots, x_{k-1} \Rightarrow e) \\ & \mid (\mathbf{let} \ x_1 = e_1; \dots; x_n = e_n \ \mathbf{in} \ e_0) \\ & \mid (\mathbf{letrec} \ x_1 = e_1; \dots; x_n = e_n \ \mathbf{in} \ e_0) \end{aligned}$$

An expression is therefore

- a basic value, a variable, the application of an operator, or
- a function-[application](#), a function-[abstraction](#), or
- a **let**-expression, i.e. an expression with [locally defined variables](#), or
- a **letrec**-expression, i.e. an expression with [simultaneously defined local variables](#).

For simplicity, we only allow [int](#) and [bool](#) as basic types.

Example:

The following well-known function computes the factorial of a natural number:

```
letrec fac    =    fn x  $\Rightarrow$  if x  $\leq$  1 then 1  
                    else x · fac (x - 1)  
  
    in fac 7
```

As usual, we only use the minimal amount of parentheses.

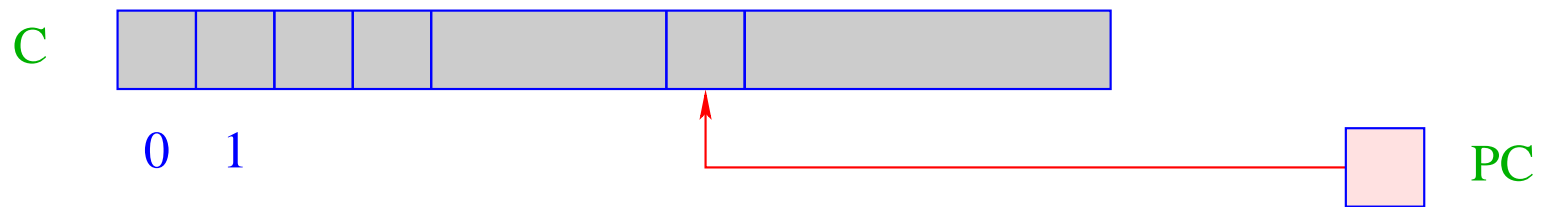
There are two **Semantics**:

CBV: Arguments are evaluated **before** they are passed to the function (as in SML);

CBN: Arguments are passed unevaluated; they are only evaluated when their value is needed (as in Haskell).

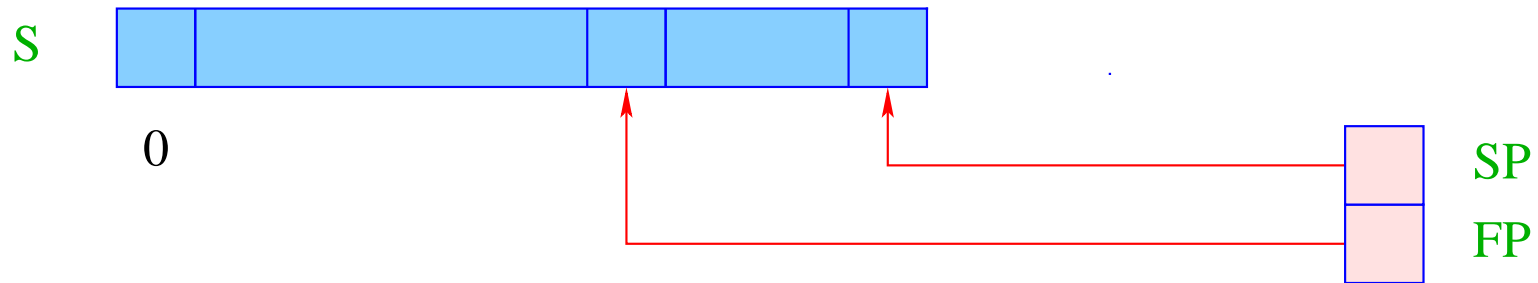
12 Architecture of the MaMa:

We know already the following components:



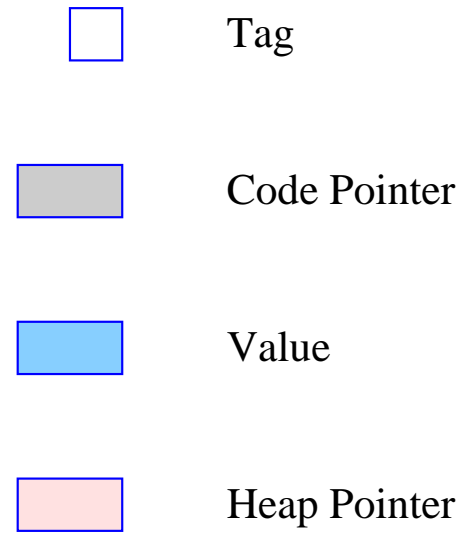
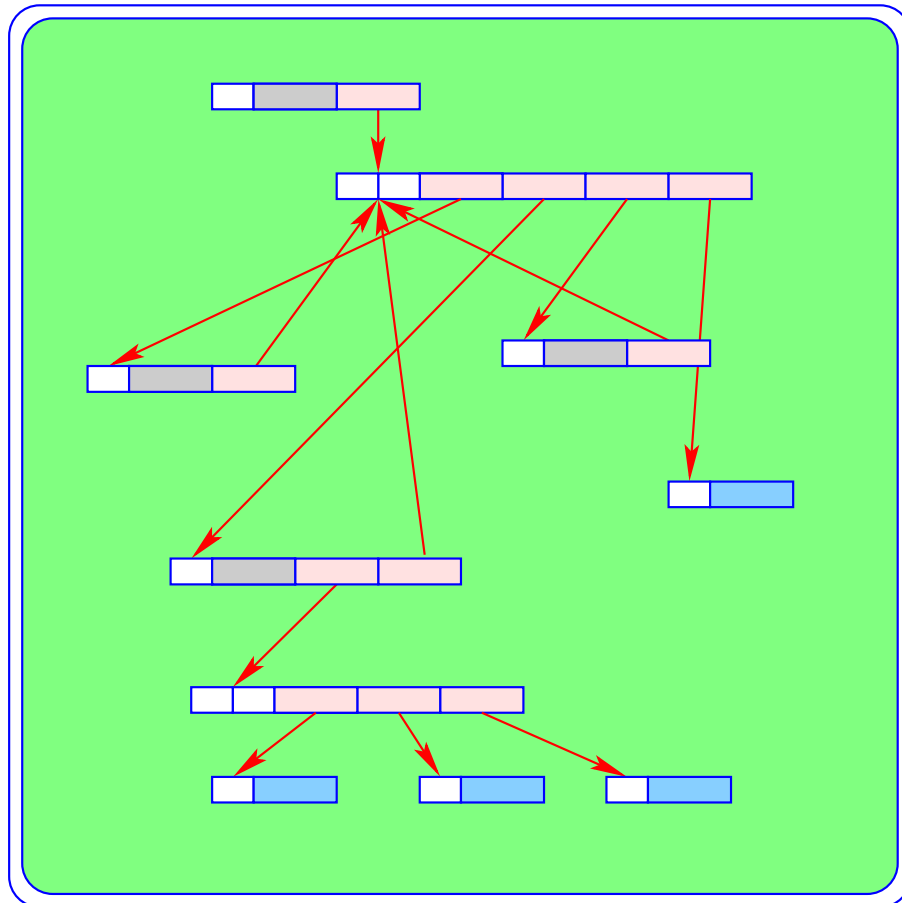
C = Code-store – contains the MaMa-program;
each cell contains one instruction;

PC = Program Counter – points to the instruction to be executed next;

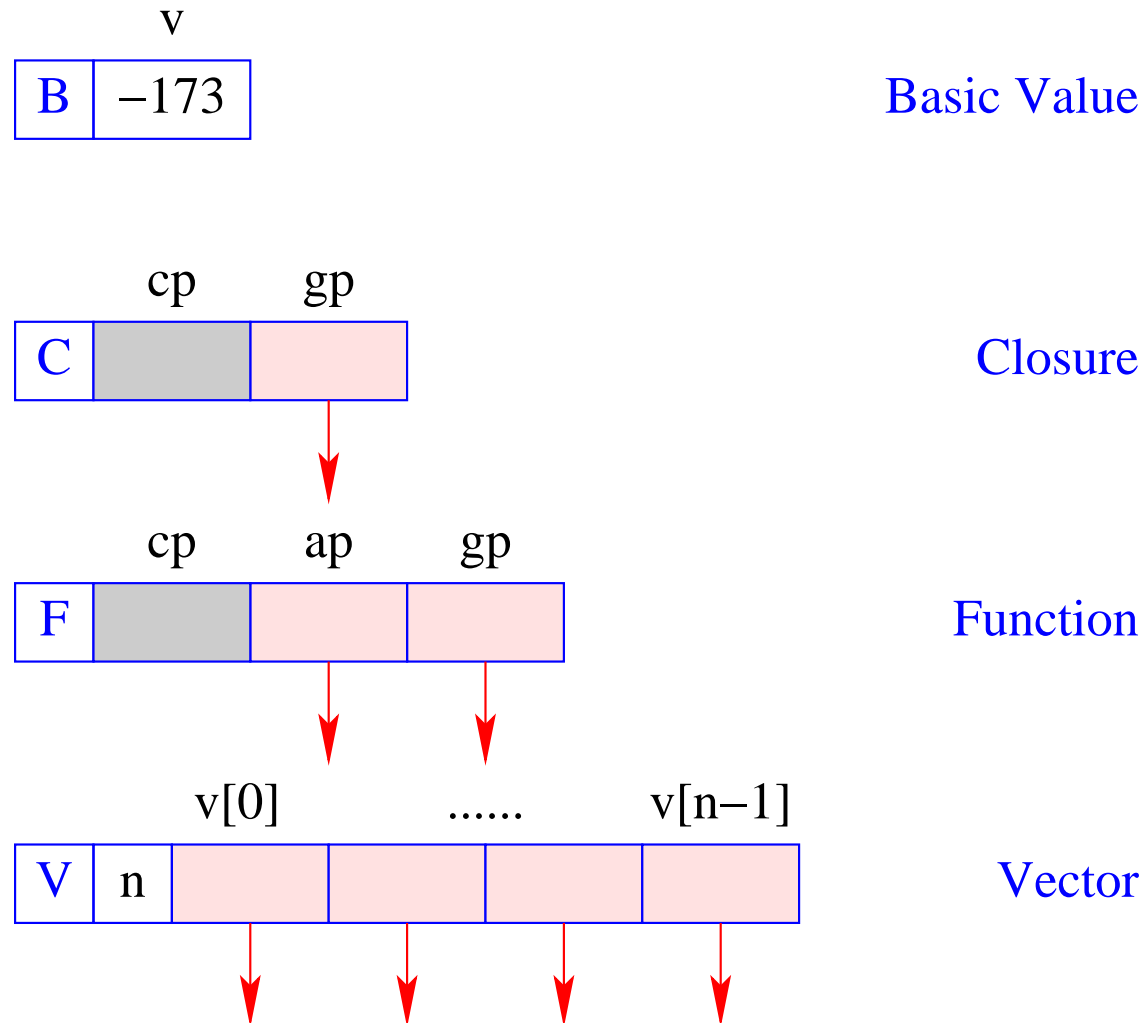


- S** = Runtime-Stack – each cell can hold a basic value or an address;
- SP** = Stack-Pointer – points to the topmost occupied cell;
as in the **CMa** implicitly represented;
- FP** = Frame-Pointer – points to the actual stack frame.

We also need a heap **H**:



... it can be thought of as an **abstract data type**, being capable of holding data objects of the following form:



The instruction `new (tag, args)` creates a corresponding object (B, C, F, V) in **H** and returns a reference to it.

We distinguish three different kinds of code for an expression e :

- `codeV e` — (generates code that) computes the **V**alue of e , stores it in the heap and returns a reference to it on top of the stack (the normal case);
- `codeB e` — computes the value of e , and returns it on the top of the stack (only for **B**asic types);
- `codeC e` — does **not** evaluate e , but stores a **C**losure of e in the heap and returns a reference to the closure on top of the stack.

We start with the code schemata for the first two kinds:

13 Simple expressions

Expressions consisting only of constants, operator applications, and conditionals are translated like expressions in imperative languages:

$$\begin{aligned} \text{code}_B b \rho \text{sd} &= \text{loadc } b \\ \text{code}_B (\square_1 e) \rho \text{sd} &= \text{code}_B e \rho \text{sd} \\ &\quad \text{op}_1 \\ \text{code}_B (e_1 \square_2 e_2) \rho \text{sd} &= \text{code}_B e_1 \rho \text{sd} \\ &\quad \text{code}_B e_2 \rho (\text{sd} + 1) \\ &\quad \text{op}_2 \end{aligned}$$

$\text{code}_B(\text{if } e_0 \text{ then } e_1 \text{ else } e_2) \rho \text{ sd} =$

- $\text{code}_B e_0 \rho \text{ sd}$
- jumpz A
- $\text{code}_B e_1 \rho \text{ sd}$
- jump B
- A: $\text{code}_B e_2 \rho \text{ sd}$
- B: ...

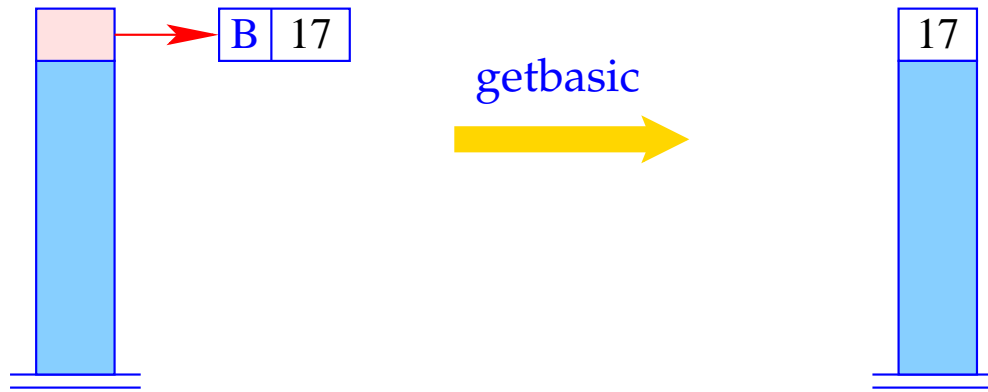
Note:

- ρ denotes the actual **address environment**, in which the expression is translated. Address environments have the form:

$$\rho : Vars \rightarrow \{L, G\} \times \mathbb{Z}$$

- The extra argument **sd**, the **stack difference**, *simulates* the movement of the **SP** when instruction execution modifies the stack. It is needed later to address variables.
- The instructions **op₁** and **op₂** implement the operators \square_1 and \square_2 , in the same way as the operators **neg** and **add** implement negation resp. addition in the **CMa**.
- For all other expressions, we first compute the value in the heap and then dereference the returned pointer:

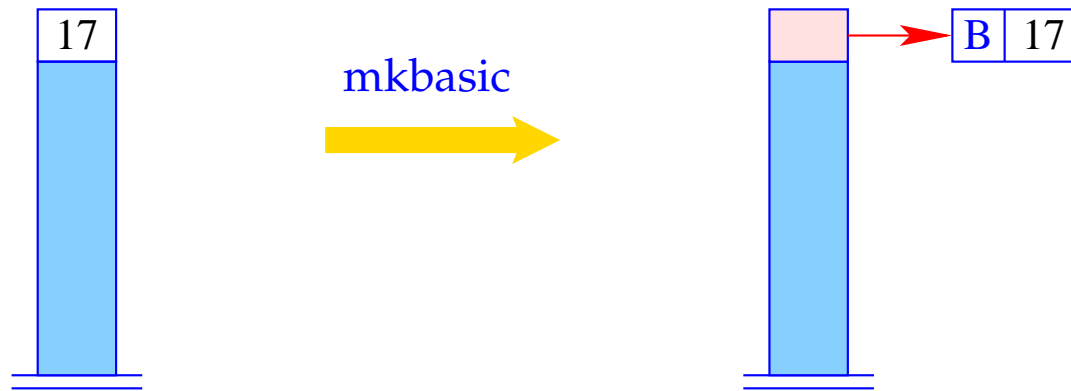
$$\text{code}_B e \rho \text{sd} = \text{code}_V e \rho \text{sd} \\ \text{getbasic}$$



```
if (H[S[SP]] != (B,_))  
    Error "not basic!";  
else  
    S[SP] = H[S[SP]].v;
```

For code_V and simple expressions, we define analogously:

$$\begin{aligned}
 \text{code}_V b \rho \text{sd} &= \text{loadc } b; \text{mkbasic} \\
 \text{code}_V (\square_1 e) \rho \text{sd} &= \text{code}_B e \rho \text{sd} \\
 &\quad \text{op}_1; \text{mkbasic} \\
 \text{code}_V (e_1 \square_2 e_2) \rho \text{sd} &= \text{code}_B e_1 \rho \text{sd} \\
 &\quad \text{code}_B e_2 \rho (\text{sd} + 1) \\
 &\quad \text{op}_2; \text{mkbasic} \\
 \text{code}_V (\text{if } e_0 \text{ then } e_1 \text{ else } e_2) \rho \text{sd} &= \text{code}_B e_0 \rho \text{sd} \\
 &\quad \text{jumpz } A \\
 &\quad \text{code}_V e_1 \rho \text{sd} \\
 &\quad \text{jump } B \\
 &\quad A: \text{code}_V e_2 \rho \text{sd} \\
 &\quad B: \dots
 \end{aligned}$$



`S[SP] = new (B,S[SP]);`

14 Accessing Variables

We must distinguish between **local** and **global** variables.

Example: Regard the function f :

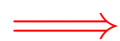
```
let c = 5
    f = fn a => let b = a * a
                in b + c
in f c
```

The function f uses the **global** variable c and the **local** variables a (as formal parameter) and b (introduced by the inner **let**).

The binding of a global variable is determined, when the function is **constructed** (**static scoping!**), and later only looked up.

Accessing Global Variables

- The bindings of global variables of an expression or a function are kept in a vector in the heap (**Global Vector**).
- They are addressed consecutively starting with 0.
- When an F-object or a C-object are constructed, the Global Vector for the function or the expression is determined and a reference to it is stored in the `gp`-component of the object.
- During the evaluation of an expression, the (**new**) register **GP** (**Global Pointer**) points to the actual Global Vector.
- In contrast, local variables should be administered on the stack ...



General form of the address environment:

$$\rho : Vars \rightarrow \{L, G\} \times \mathbb{Z}$$

Accessing Local Variables

Local variables are administered on the stack, in **stack frames**.

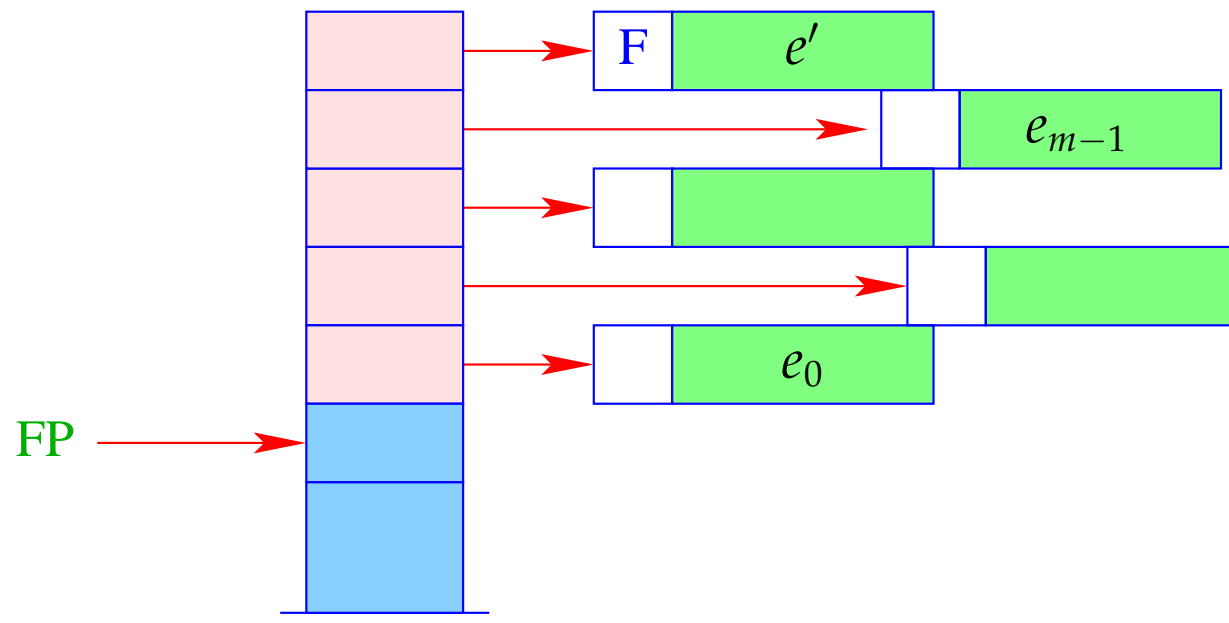
Let $e \equiv e' e_0 \dots e_{m-1}$ be the application of a function e' to arguments e_0, \dots, e_{m-1} .

Warning:

The arity of e' does not need to be m :-)

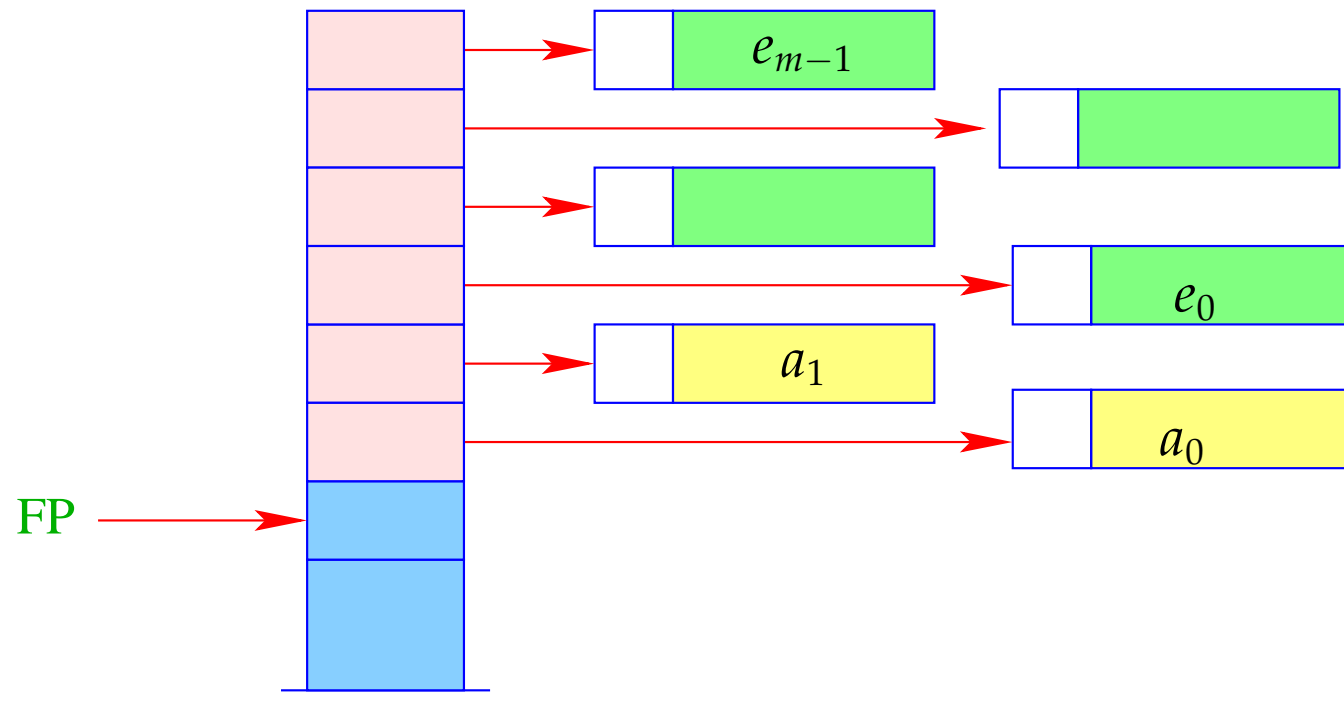
- PuF functions have **curried** types, $f : t_1 \rightarrow t_2 \rightarrow \dots \rightarrow t_n \rightarrow t$
- f may therefore receive less than n arguments (**under supply**);
- f may also receive more than n arguments, if t is a **functional type** (**over supply**).

Possible stack organisations:

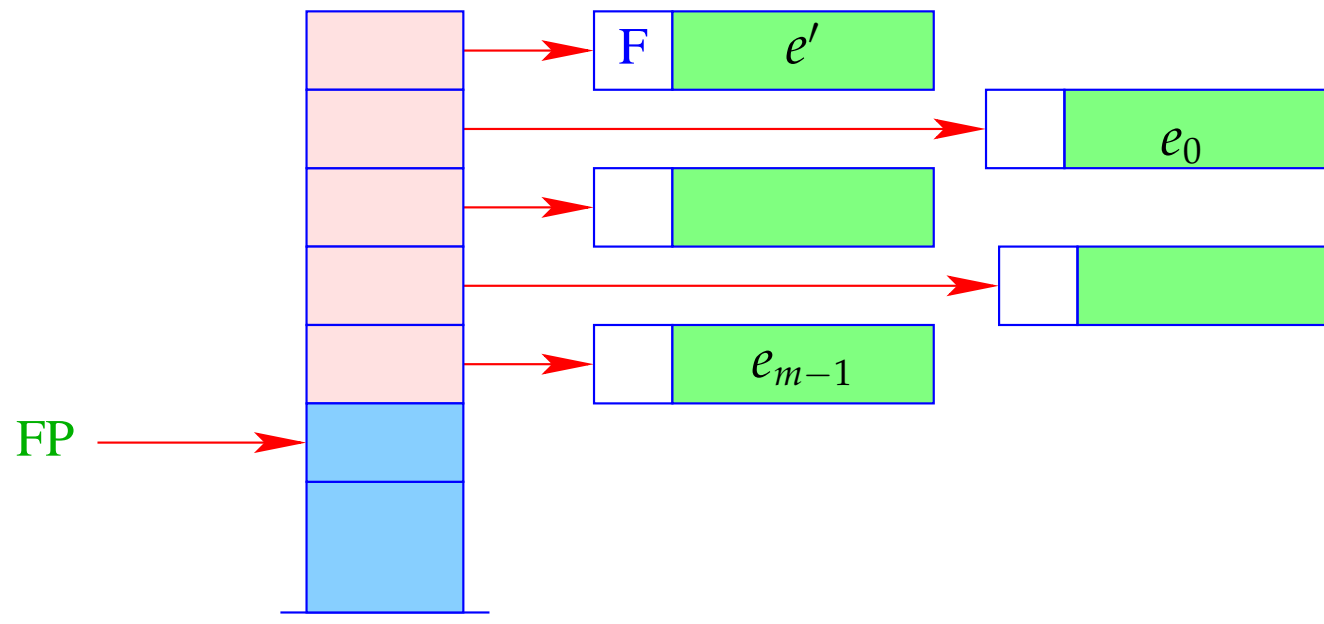


- + Addressing of the arguments can be done relative to **FP**
- The local variables of e' cannot be addressed relative to **FP**.
- If e' is an n -ary function with $n < m$, i.e., we have an over-supplied function application, the remaining $m - n$ arguments will have to be shifted.

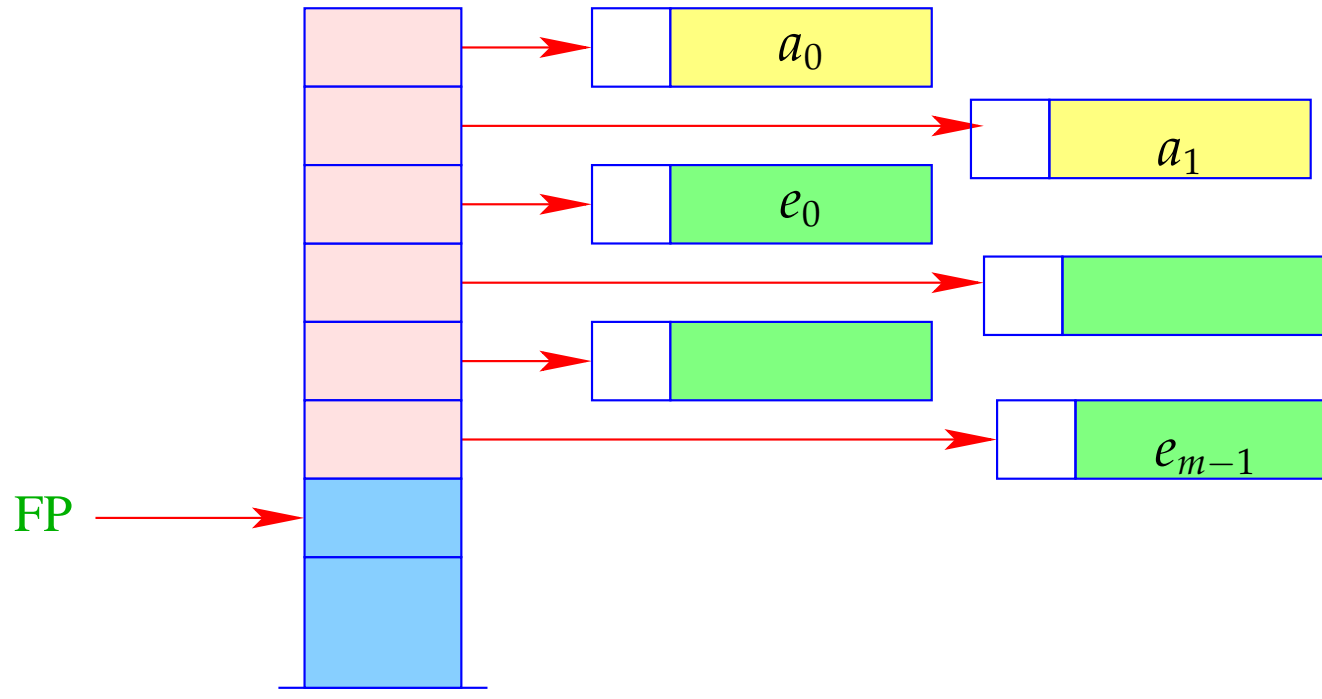
- If e' evaluates to a function, which has already been partially applied to the parameters a_0, \dots, a_{k-1} , these have to be sneaked in underneath e_0 :



Alternative:



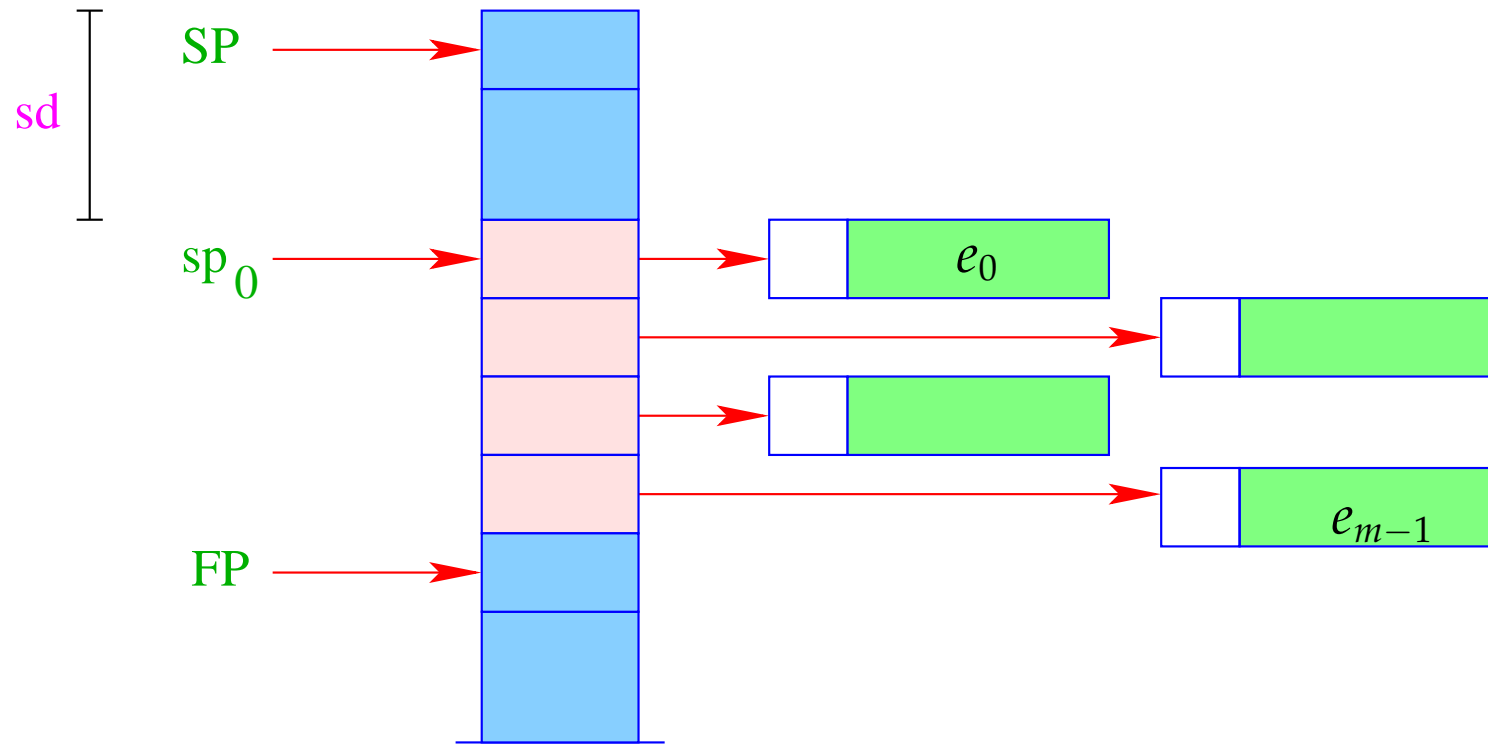
- + The further arguments a_0, \dots, a_{k-1} and the local variables can be allocated above the arguments.



- Addressing of arguments and local variables relative to **FP** is no more possible. (Remember: m is unknown when the function definition is translated.)

Way out:

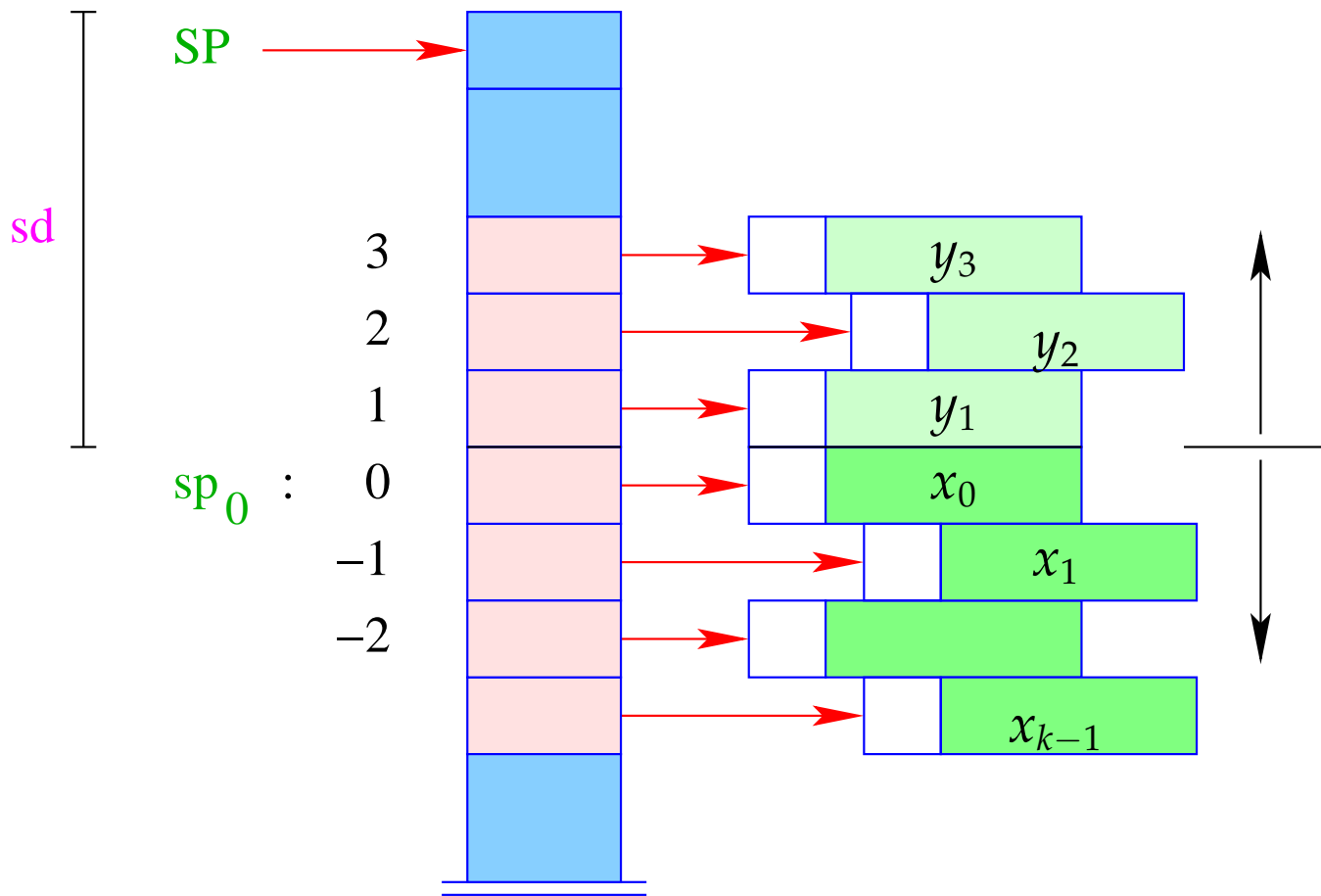
- We address both, arguments and local variables, relative to the stack pointer
SP !!!
- However, the stack pointer changes during program execution...



- The difference between the **current** value of **SP** and its value sp_0 at the entry of the function body is called the stack distance, **sd**.
- Fortunately, this stack distance can be determined at compile time for each program point, by **simulating the movement** of the **SP**.
- The formal parameters x_0, x_1, x_2, \dots successively receive the **non-positive** relative addresses $0, -1, -2, \dots$, i.e., $\rho x_i = (L, -i)$.
- The **absolute** address of the i -th formal parameter consequently is

$$sp_0 - i = (\mathbf{SP} - \mathbf{sd}) - i$$

- The local **let**-variables y_1, y_2, y_3, \dots will be successively pushed onto the stack:



- The y_i have **positive** relative addresses $1, 2, 3, \dots$, that is: $\rho y_i = (L, i)$.
- The absolute address of y_i is then $sp_0 + i = (SP - sd) + i$

With **CBN**, we generate for the access to a variable:

$$\text{code}_V x \rho \text{sd} = \text{getvar } x \rho \text{sd} \\ \text{eval}$$

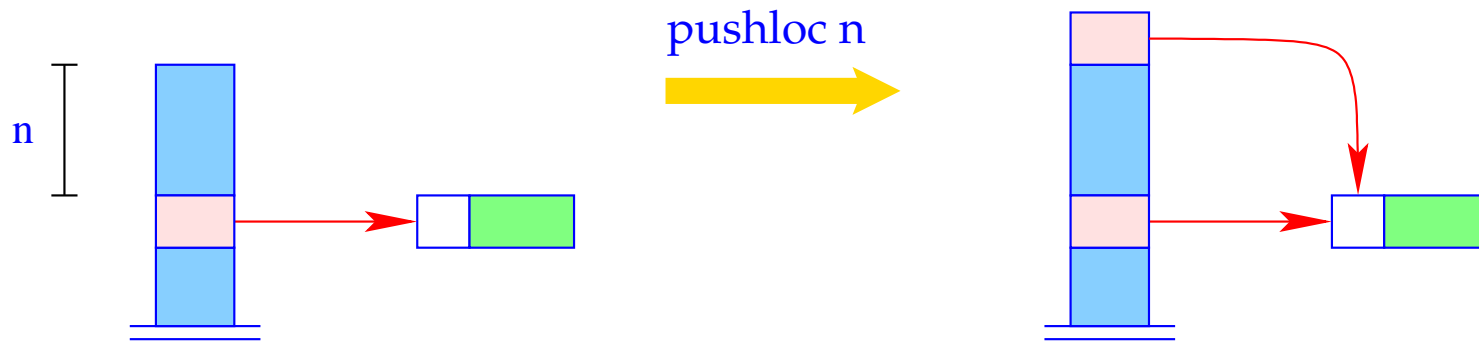
The instruction **eval** checks, whether the value has already been computed or whether its evaluation has to yet to be done (\implies will be treated later :-)

With **CBV**, we can just delete **eval** from the above code schema.

The (compile-time) macro **getvar** is defined by:

$$\text{getvar } x \rho \text{sd} = \text{let } (t, i) = \rho x \text{ in} \\ \text{case } t \text{ of} \\ \quad L \Rightarrow \text{pushloc } (\text{sd} - i) \\ \quad G \Rightarrow \text{pushglob } i \\ \text{end}$$

The access to local variables:



$S[SP+1] = S[SP - n]; SP++;$

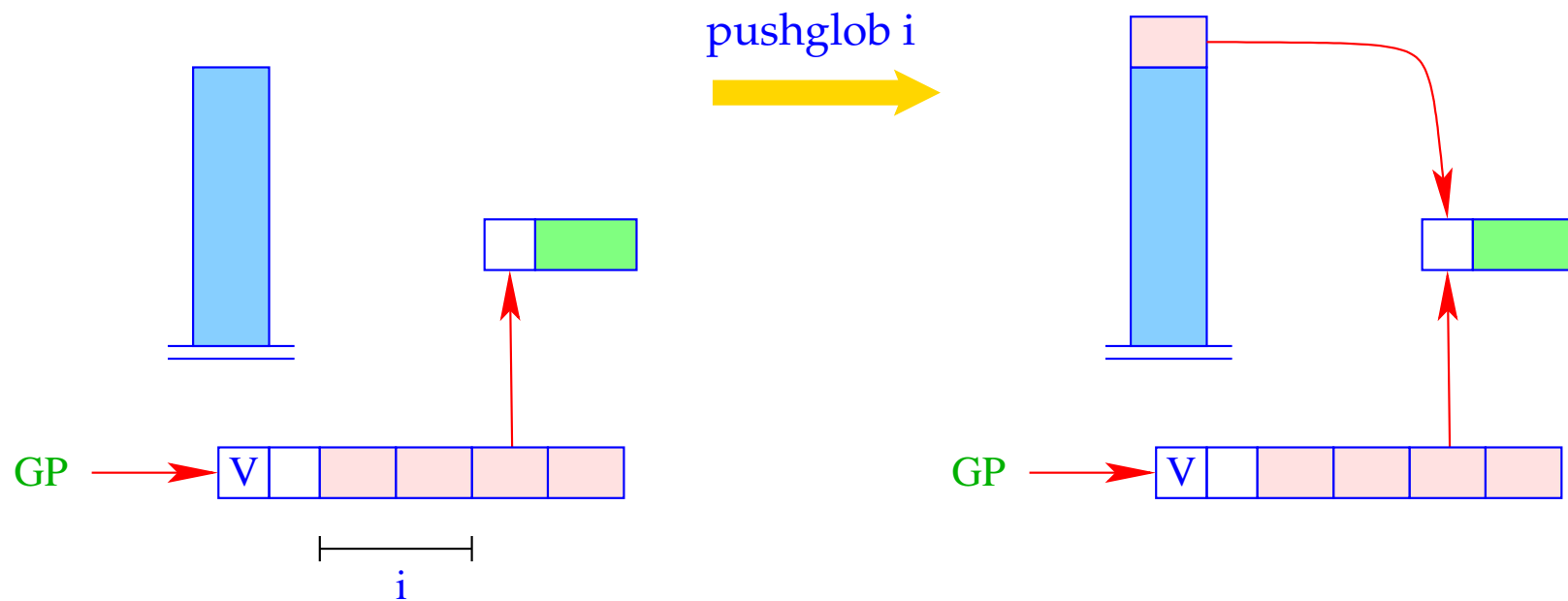
Correctness argument:

Let sp and sd be the values of the stack pointer resp. stack distance **before** the execution of the instruction. The value of the local variable with address i is loaded from $S[a]$ with

$$a = sp - (sd - i) = (sp - sd) + i = sp_0 + i$$

... exactly as it should be :-)

The access to global variables is much simpler:



$SP = SP + 1;$
 $S[SP] = GP \rightarrow v[i];$

Example:

Regard $e \equiv (b + c)$ for $\rho = \{b \mapsto (L, 1), c \mapsto (G, 0)\}$ and $sd = 1$.

With **CBN**, we obtain:

<code>code_v e ρ 1</code>	=	<code>getvar b ρ 1</code>	=	1	<code>pushloc 0</code>
		<code>eval</code>		2	<code>eval</code>
		<code>getbasic</code>		2	<code>getbasic</code>
		<code>getvar c ρ 2</code>		2	<code>pushglob 0</code>
		<code>eval</code>		3	<code>eval</code>
		<code>getbasic</code>		3	<code>getbasic</code>
		<code>add</code>		3	<code>add</code>
		<code>mkbasic</code>		2	<code>mkbasic</code>

15 let-Expressions

As a warm-up let us first consider the treatment of local variables :-)

Let $e \equiv \mathbf{let} \ y_1 = e_1; \dots; y_n = e_n \ \mathbf{in} \ e_0$ be a **let**-expression.

The translation of e must deliver an instruction sequence that

- allocates local variables y_1, \dots, y_n ;
- in the case of
 - CBV**: evaluates e_1, \dots, e_n and binds the y_i to their values;
 - CBN**: constructs closures for the e_1, \dots, e_n and binds the y_i to them;
- evaluates the expression e_0 and returns its value.

Here, we consider the **non-recursive** case only, i.e. where y_j only depends on y_1, \dots, y_{j-1} . We obtain for **CBN**:

```

codeV e ρ sd = codeC e1 ρ sd
                codeC e2 ρ1 (sd + 1)
                ...
                codeC en ρn-1 (sd + n - 1)
                codeV e0 ρn (sd + n)
                slide n // deallocates local variables

```

where $\rho_j = \rho \oplus \{y_i \mapsto (L, \text{sd} + i) \mid i = 1, \dots, j\}$.

In the case of **CBV**, we use `codeV` for the expressions e_1, \dots, e_n .

Warning!

All the e_i must be associated with the same binding for the global variables!

Example:

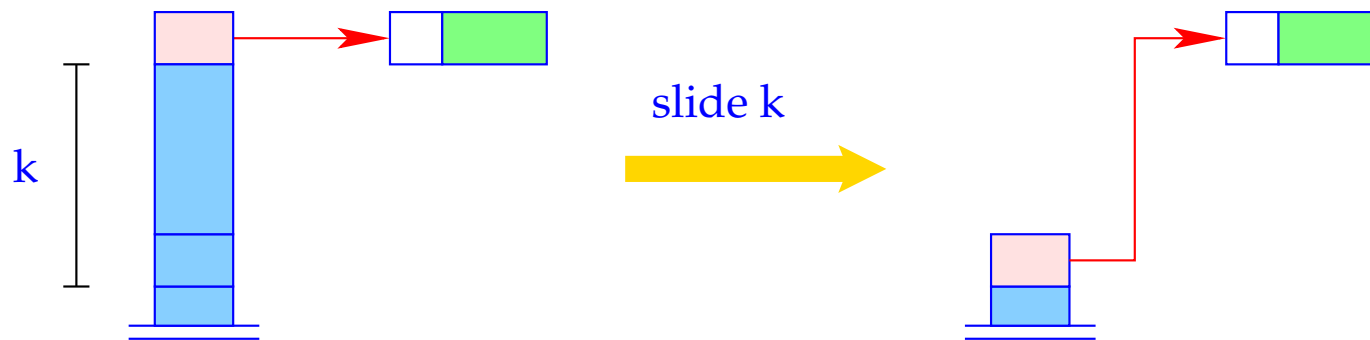
Consider the expression

$$e \equiv \mathbf{let} \ a = 19; b = a * a \ \mathbf{in} \ a + b$$

for $\rho = \emptyset$ and $\mathbf{sd} = 0$. We obtain (for **CBV**):

0	loadc 19	3	getbasic	3	pushloc 1
1	mkbasic	3	mul	4	getbasic
1	pushloc 0	2	mkbasic	4	add
2	getbasic	2	pushloc 1	3	mkbasic
2	pushloc 1	3	getbasic	3	slide 2

The instruction `slide k` deallocates again the space for the locals:



$S[SP-k] = S[SP];$
 $SP = SP - k;$

16 Function Definitions

The definition of a function f requires code that allocates a **functional value** for f in the heap. This happens in the following steps:

- Creation of a Global Vector with the binding of the free variables;
- Creation of an (initially empty) argument vector;
- Creation of an F-Object, containing references to these vectors and the start address of the code for the body;

Separately, code for the body has to be generated.

Thus:

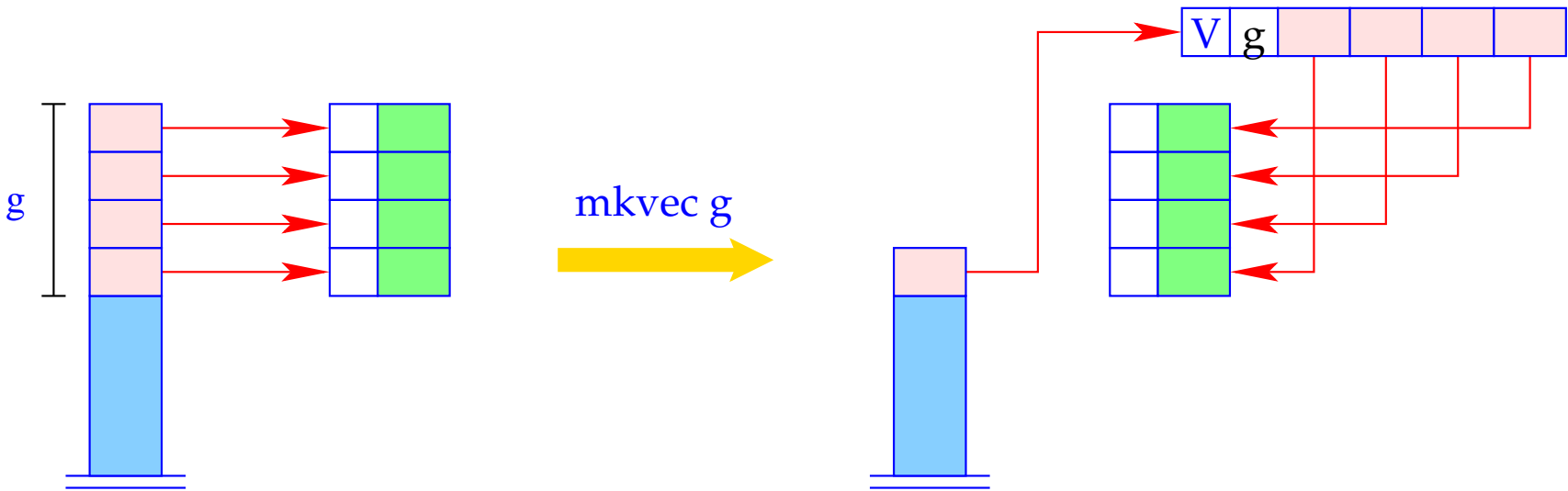
$$\text{code}_V(\mathbf{fn } x_0, \dots, x_{k-1} \Rightarrow e) \rho \text{ sd} =$$

```

getvar z0 ρ sd
getvar z1 ρ (sd + 1)
...
getvar zg-1 ρ (sd + g - 1)
mkvec g
mkfunval A
jump B
A : targ k
    codeV e ρ' 0
    return k
B : ...

```

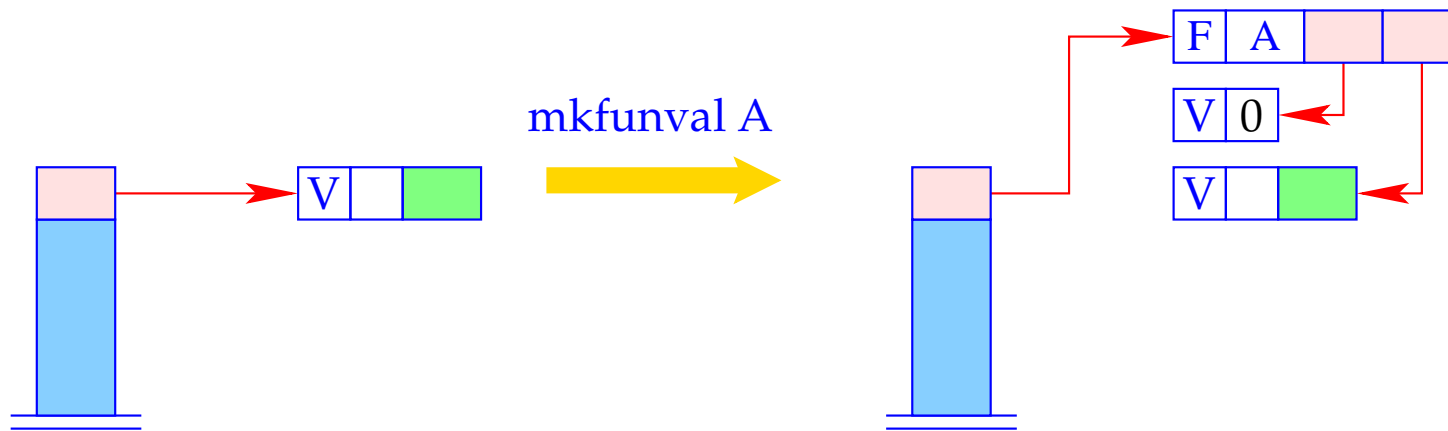
where $\{z_0, \dots, z_{g-1}\} = \text{free}(\mathbf{fn } x_0, \dots, x_{k-1} \Rightarrow e)$
and $\rho' = \{x_i \mapsto (L, -i) \mid i = 0, \dots, k-1\} \cup \{z_j \mapsto (G, j) \mid j = 0, \dots, g-1\}$



```

h = new (V, n);
SP = SP - g + 1;
for (i=0; i<g; i++)
    h->v[i] = S[SP + i];
S[SP] = h;

```



```

a = new (V,0);
S[SP] = new (F, A, a, S[SP]);

```

Example:

Regard $f \equiv \mathbf{fn} \ b \Rightarrow a + b$ for $\rho = \{a \mapsto (L, 1)\}$ and $\mathbf{sd} = 1$.

$\mathbf{code}_V f \ \rho \ 1$ produces:

1	pushloc 0	0	pushglob 0	2	getbasic
2	mkvec 1	1	eval	2	add
2	mkfunval A	1	getbasic	1	mkbasic
2	jump B	1	pushloc 1	1	return 1
0	A: targ 1	2	eval	2	B: ...

The secrets around $\mathbf{targ} \ k$ and $\mathbf{return} \ k$ will be revealed later :-)

17 Function Application

Function applications correspond to function calls in **C**.

The necessary actions for the evaluation of $e' e_0 \dots e_{m-1}$ are:

- Allocation of a stack frame;
- Transfer of the actual parameters, i.e. with:
 - CBV**: Evaluation of the actual parameters;
 - CBN**: Allocation of closures for the actual parameters;
- Evaluation of the expression e' to an F-object;
- Application of the function.

Thus for **CBN**:

```

codeV (e' e0 ... em-1) ρ sd = mark A // Allocation of the frame
                                codeC em-1 ρ (sd + 3)
                                codeC em-2 ρ (sd + 4)
                                ...
                                codeC e0 ρ (sd + m + 2)
                                codeV e' ρ (sd + m + 3) // Evaluation of e'
                                apply // corresponds to call
                                A: ...

```

To implement **CBV**, we use `codeV` instead of `codeC` for the arguments e_i .

Example: For $(f\ 42)$, $\rho = \{f \mapsto (L, 2)\}$ and $sd = 2$, we obtain with **CBV**:

```

2 mark A           6 mkbasic           7 apply
5 loadc 42        6 pushloc 4         3 A: ...

```

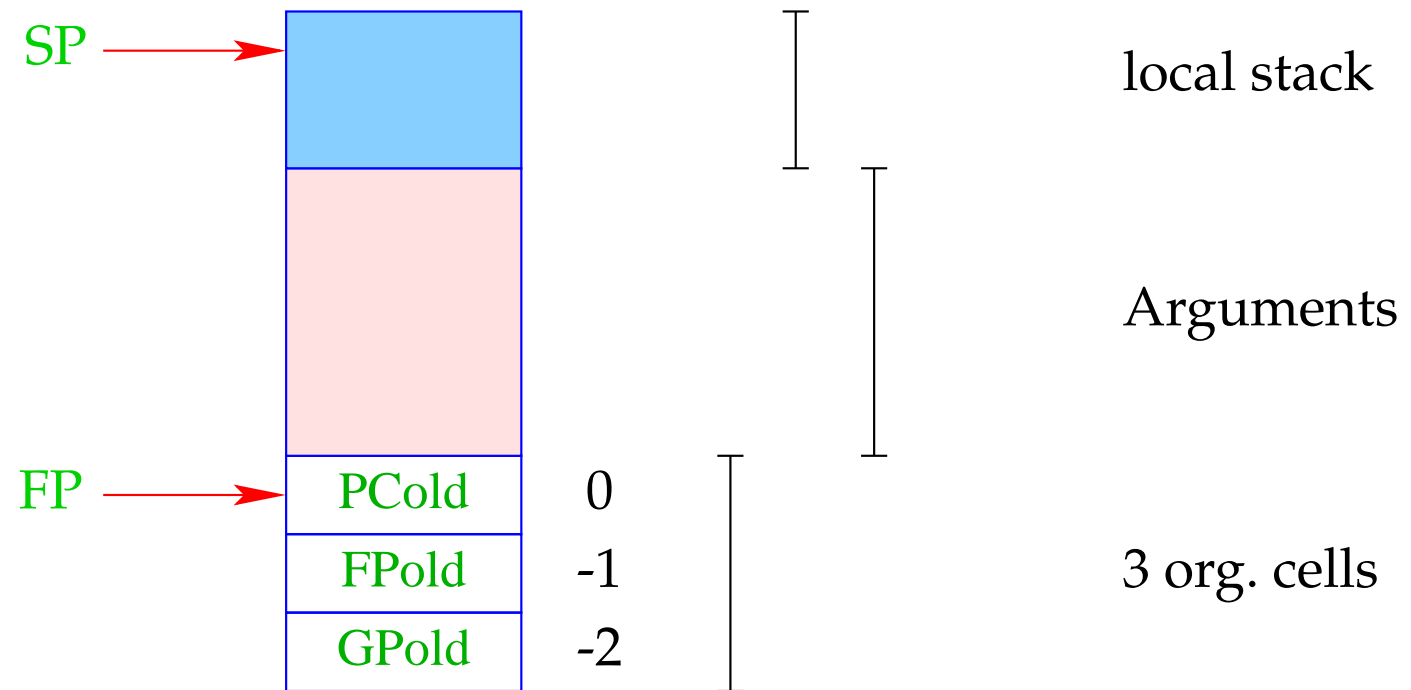
A Slightly Larger Example:

let $a = 17$; $f = \mathbf{fn}$ $b \Rightarrow a + b$ **in** f 42

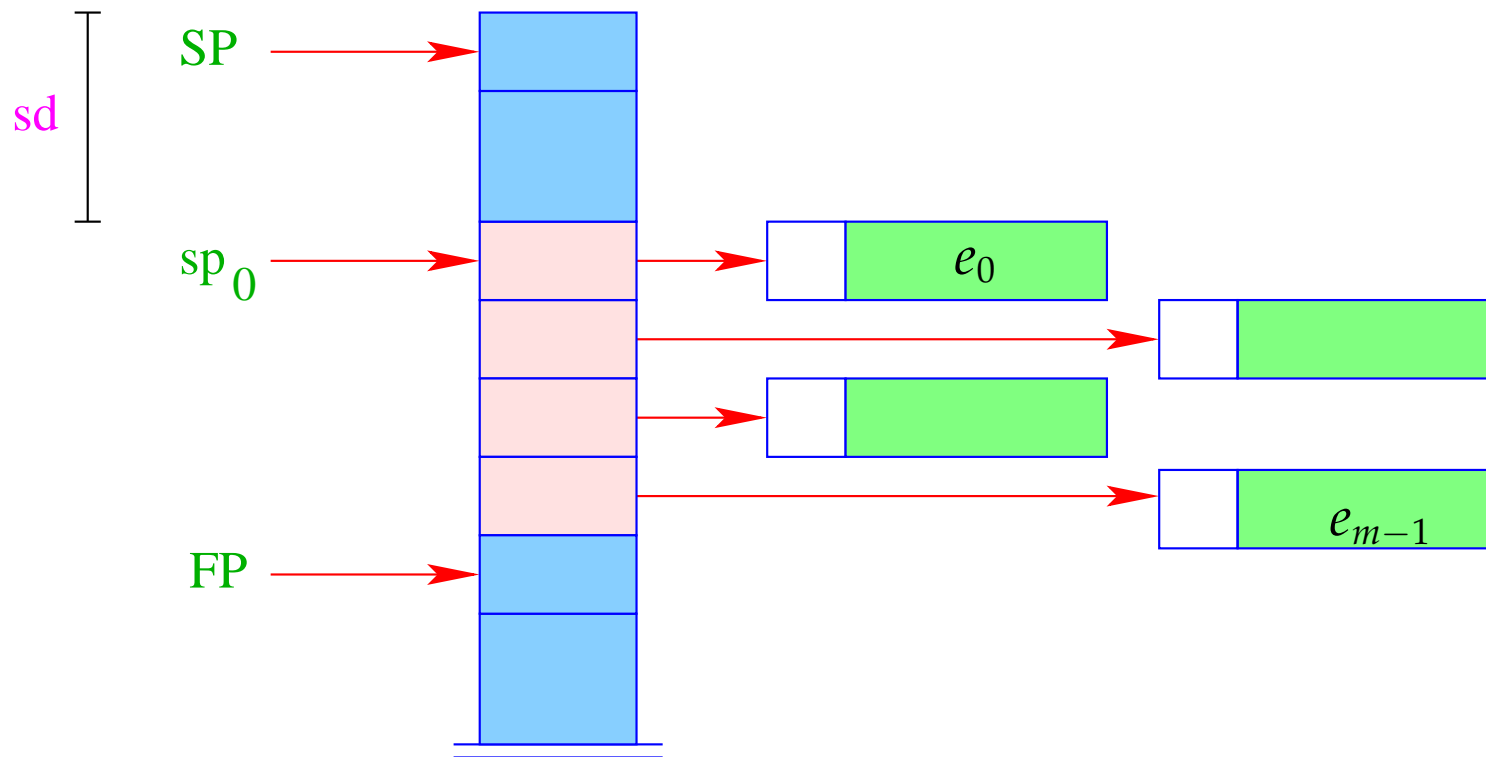
For **CBV** and $sd = 0$ we obtain:

0	loadc 17	2		jump B	2		getbasic	5		loadc 42
1	mkbasic	0	A:	targ 1	2		add	5		mkbasic
1	pushloc 0	0		pushglob 0	1		mkbasic	6		pushloc 4
2	mkvec 1	1		getbasic	1		return 1	7		apply
2	mkfunval A	1		pushloc 1	2	B:	mark C	3	C:	slide 2

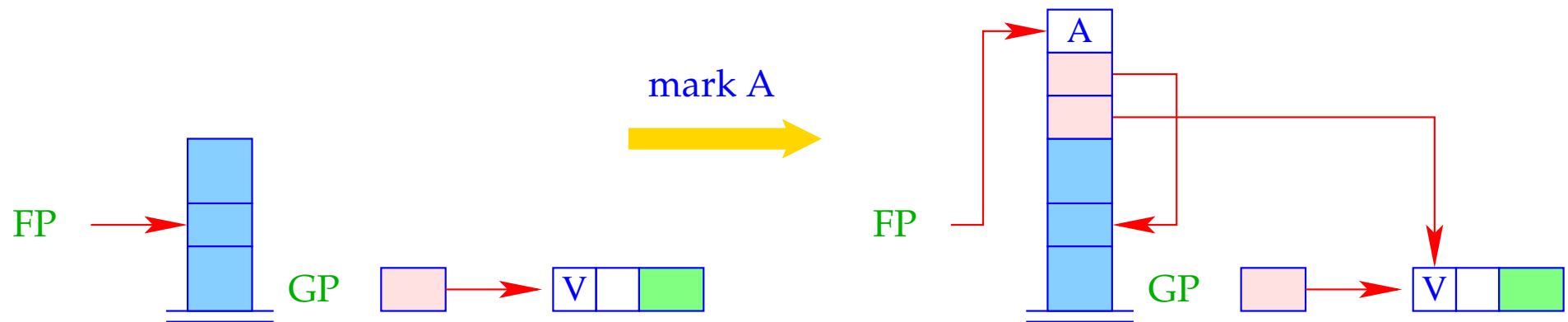
For the implementation of the new instruction, we must fix the organization of a stack frame:



Remember: Addressing of arguments and local variables

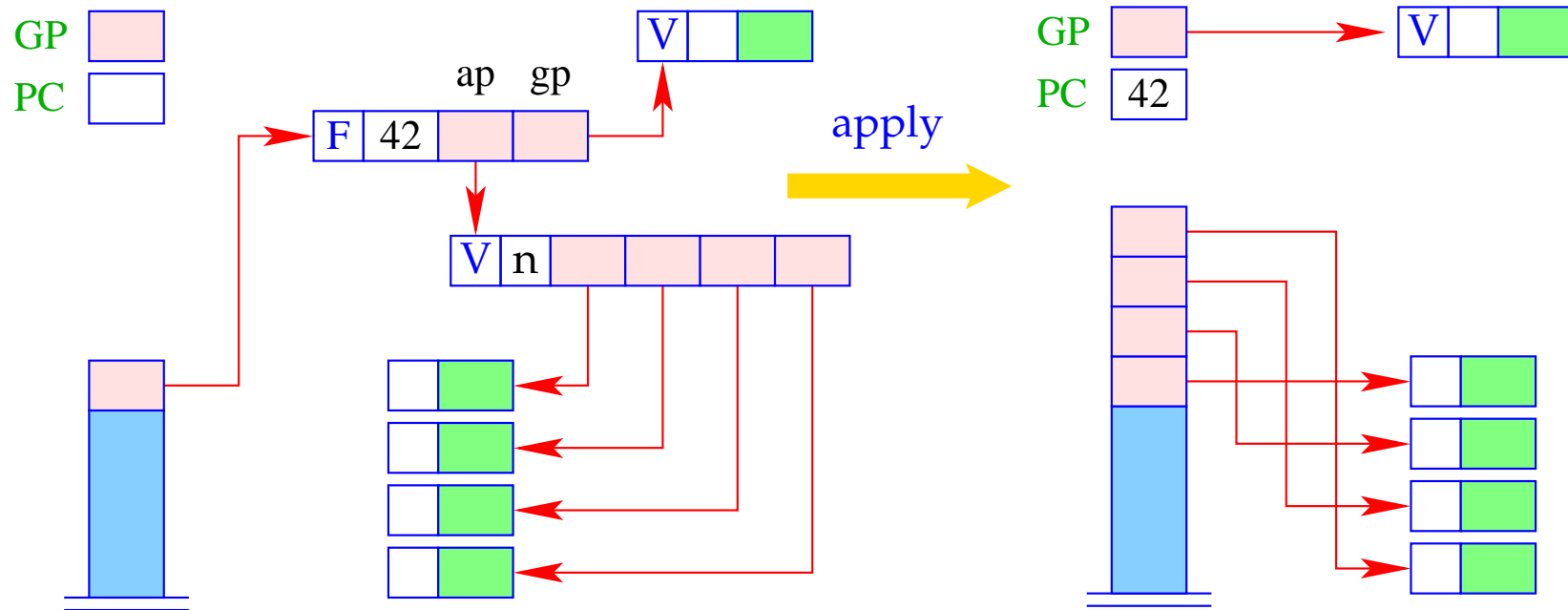


Different from the **CMa**, the instruction **mark A** already saves the return address:



$S[SP+1] = GP;$
 $S[SP+2] = FP;$
 $S[SP+3] = A;$
 $FP = SP = SP + 3;$

The instruction `apply` unpacks the F-object, a reference to which (hopefully) resides on top of the stack, and continues execution at the address given there:



```

h = S[SP];
if (H[h] != (F,_,_))
    Error "no fun";
else {

```

```

    GP = h->gp; PC = h->cp;
    for (i=0; i < h->ap->n; i++)
        S[SP+i] = h->ap->v[i];
    SP = SP + h->ap->n - 1;
}

```

Warning:

- The last element of the argument vector is the last to be put onto the stack. This must be the **first** argument reference.
- This should be kept in mind, when we treat the packing of arguments of an under-supplied function application into an F-object !!!

18 Over- and Undersupply of Arguments

The first instruction to be executed when entering a function body, i.e., after an `apply` is `targ k`.

This instruction checks whether there are enough arguments to evaluate the body.

Only if this is the case, the execution of the code for the body is started.

Otherwise, i.e. in the case of `under-supply`, a new F-object is returned.

The test for number of arguments uses: $SP - FP$

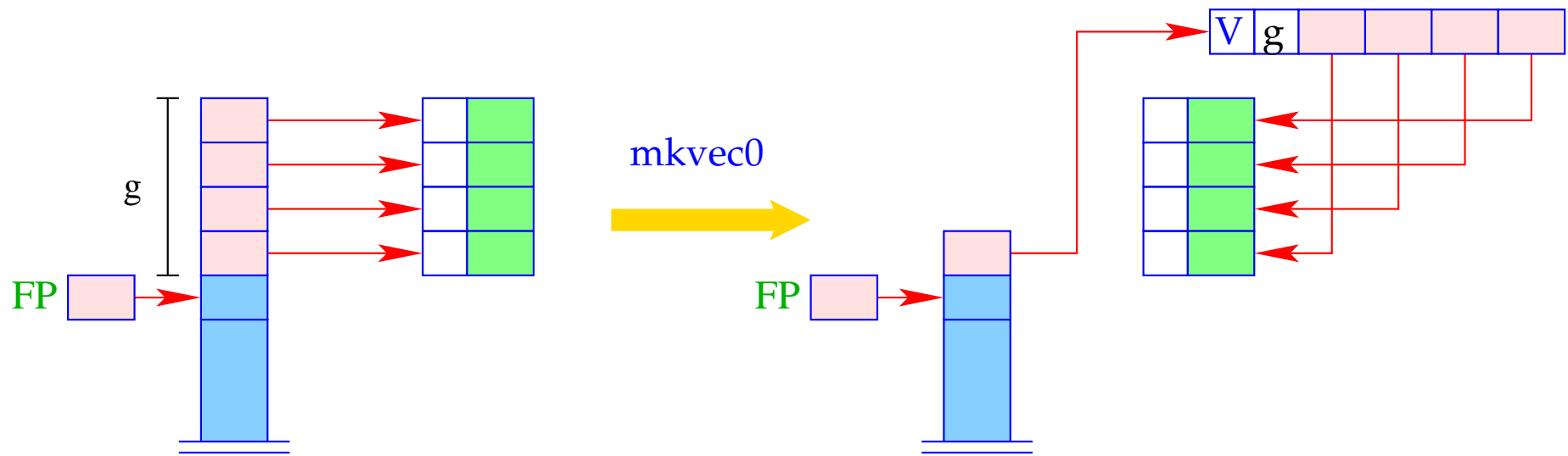
`targ k` is a complex instruction.

We decompose its execution in the case of `under-supply` into several steps:

```
targ k = if (SP - FP < k) {  
    mkvec0;           // creating the argumentvector  
    wrap;             // wrapping into an F - object  
    popenv;          // popping the stack frame  
}
```

The combination of these steps into one instruction is a kind of optimization :-)

The instruction `mkvec0` takes all references from the stack above `FP` and stores them into a vector:

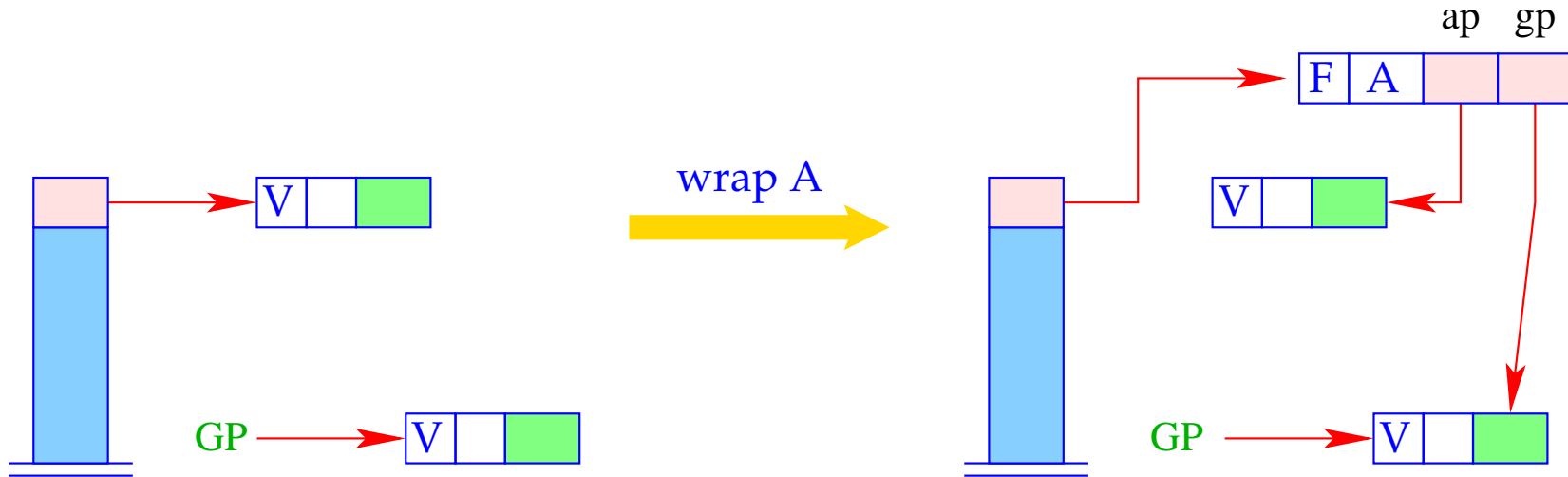


```

g = SP-FP; h = new (V, g);
SP = FP+1;
for (i=0; i<g; i++)
    h->v[i] = S[SP + i];
S[SP] = h;

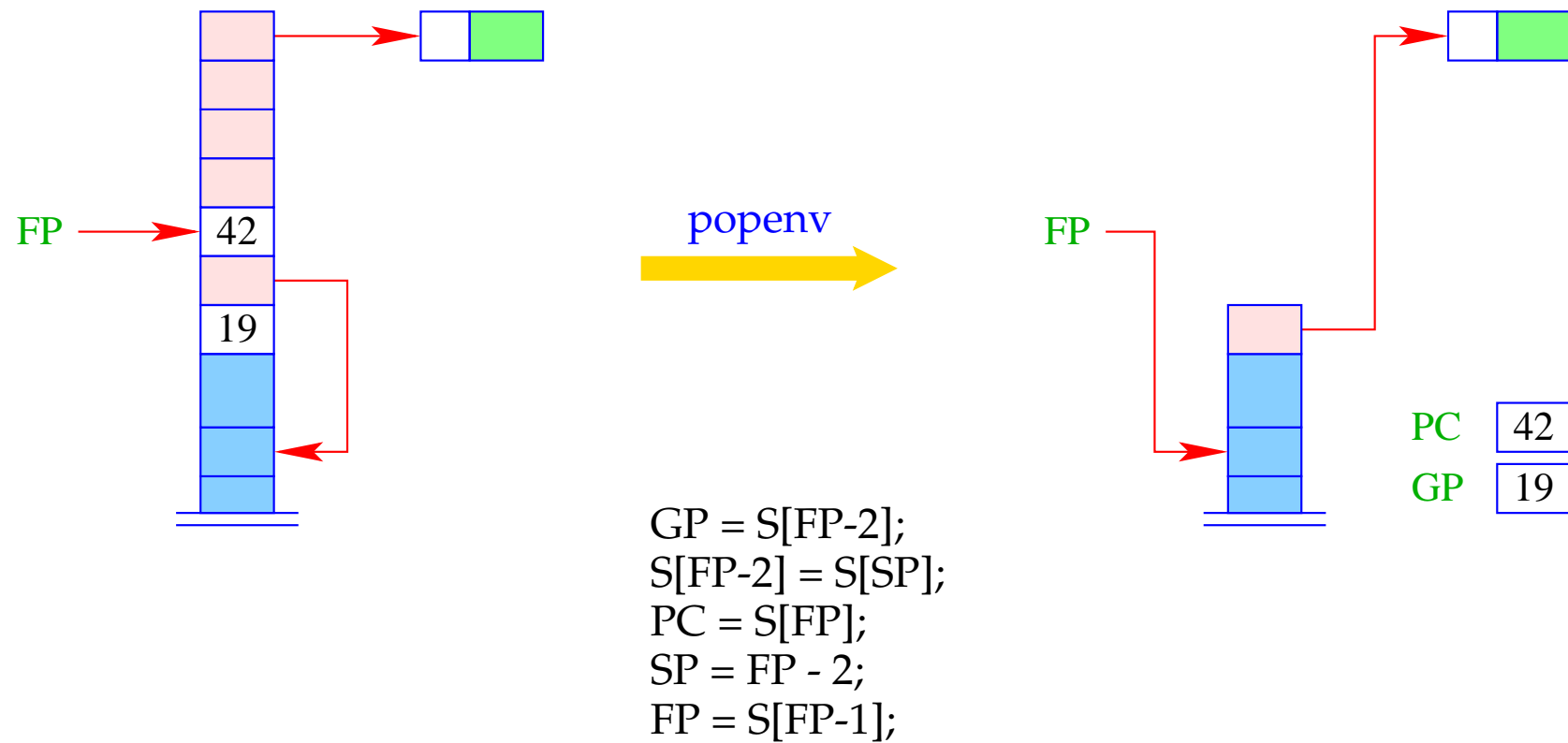
```

The instruction `wrap A` wraps the argument vector together with the global vector into an F-object:

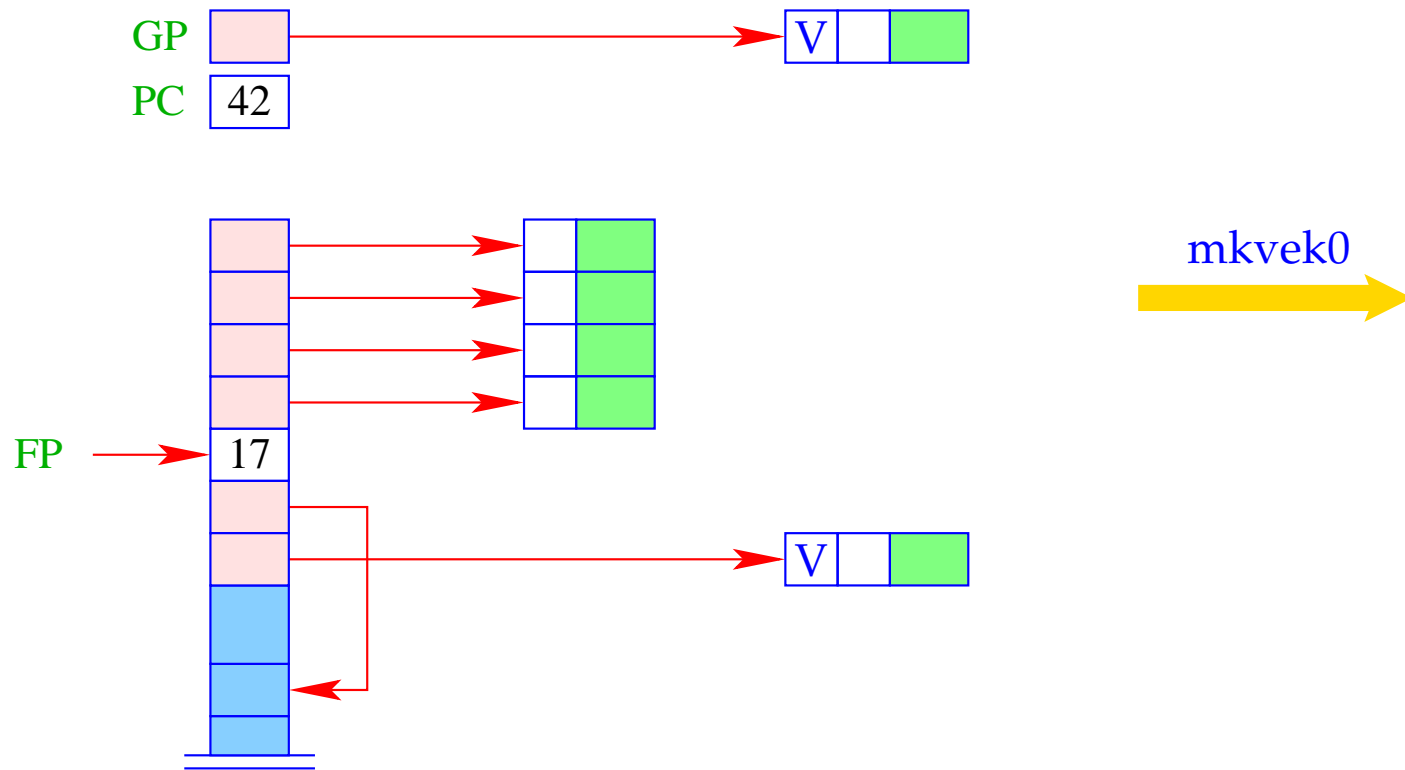


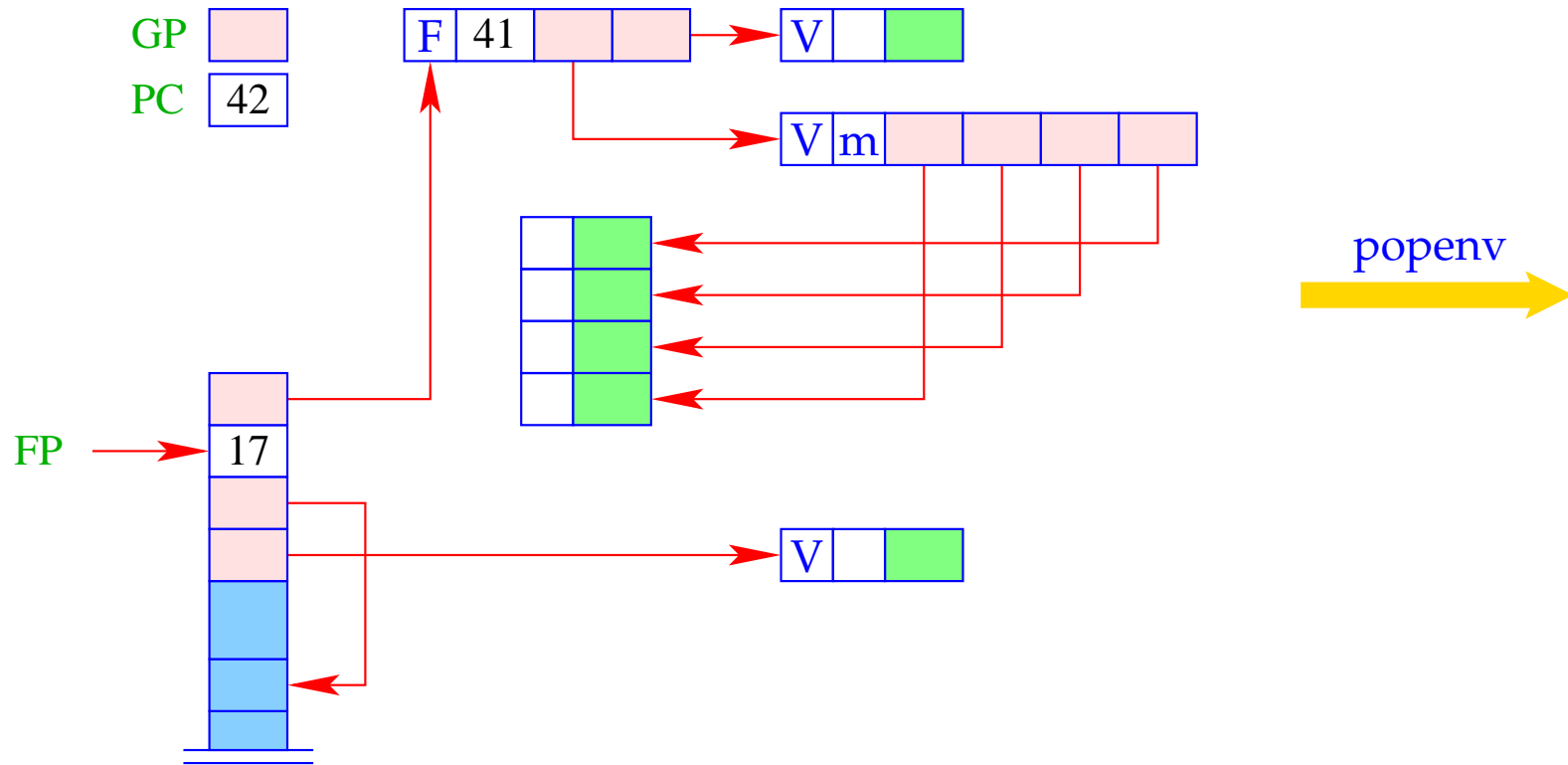
$S[SP] = \text{new } (F, A, S[SP], GP);$

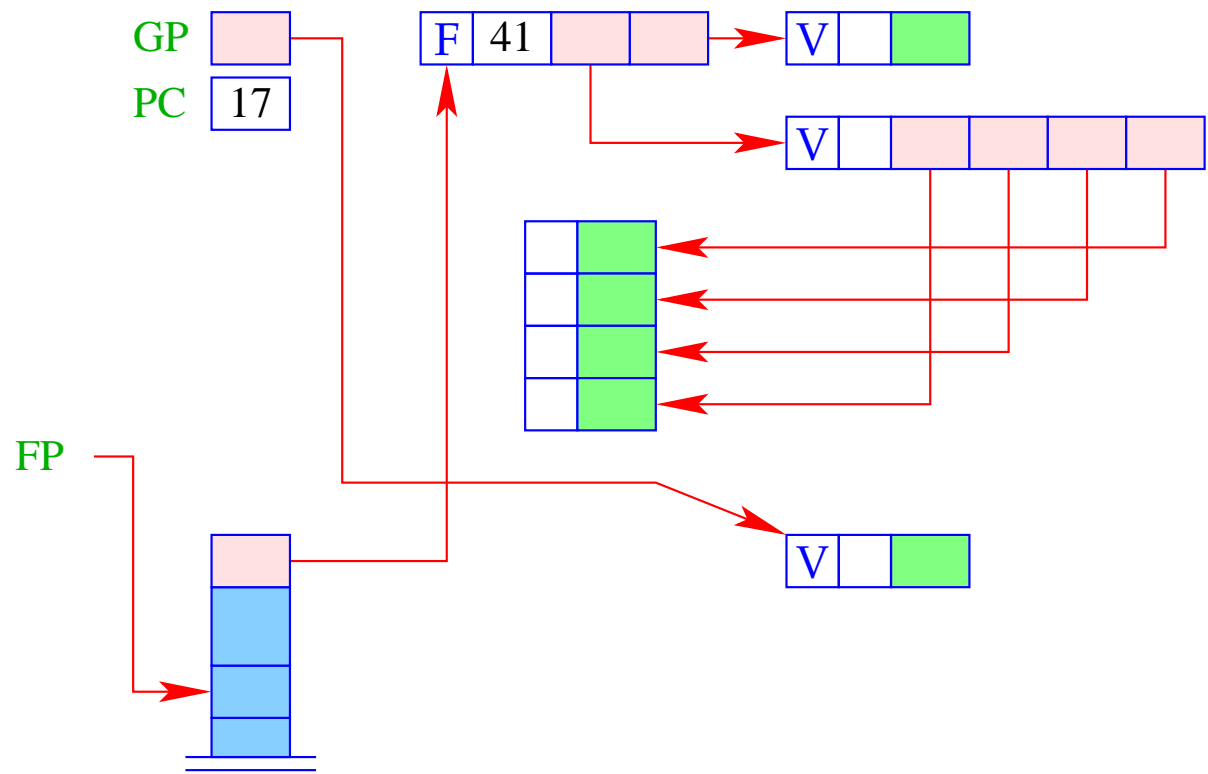
The instruction `popenv` finally releases the stack frame:

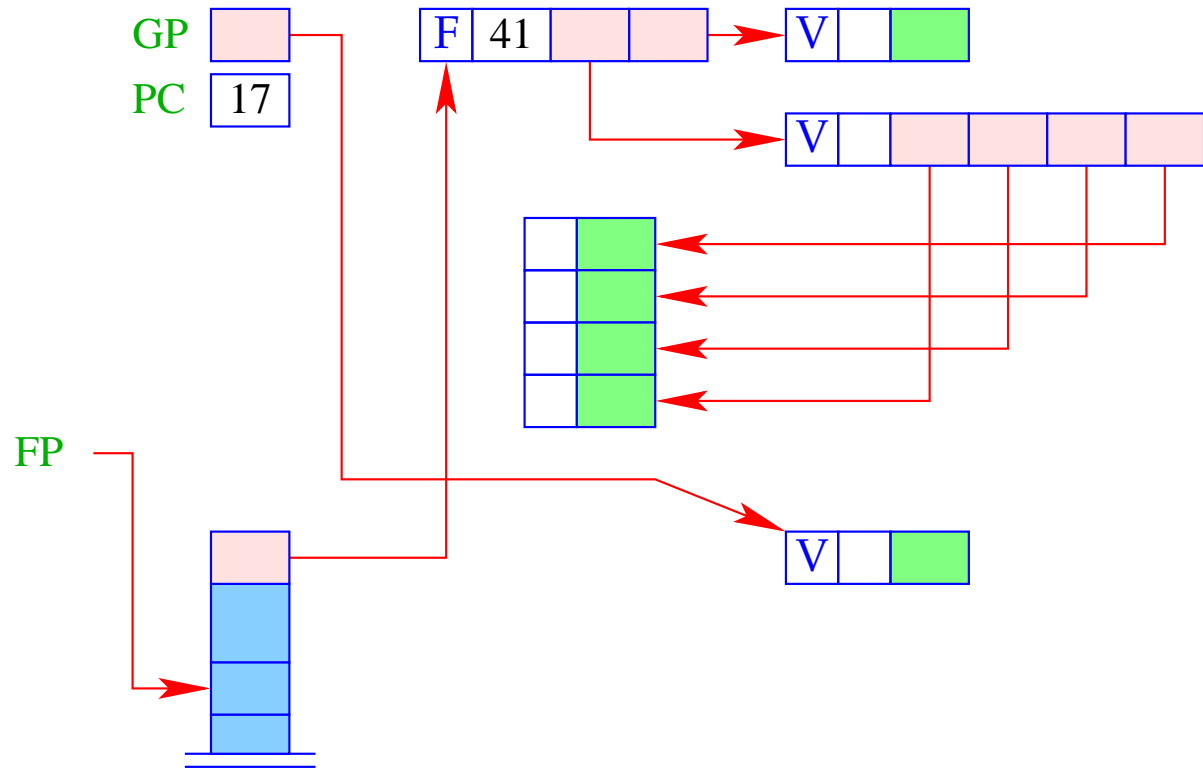


Thus, we obtain for `targ k` in the case of under supply:









- The stack frame can be released **after the execution of the body** if exactly the right number of arguments was available.
- If there is an **oversupply** of arguments, the body must evaluate to a function, which consumes the rest of the arguments ...
- The check for this is done by **return k**:

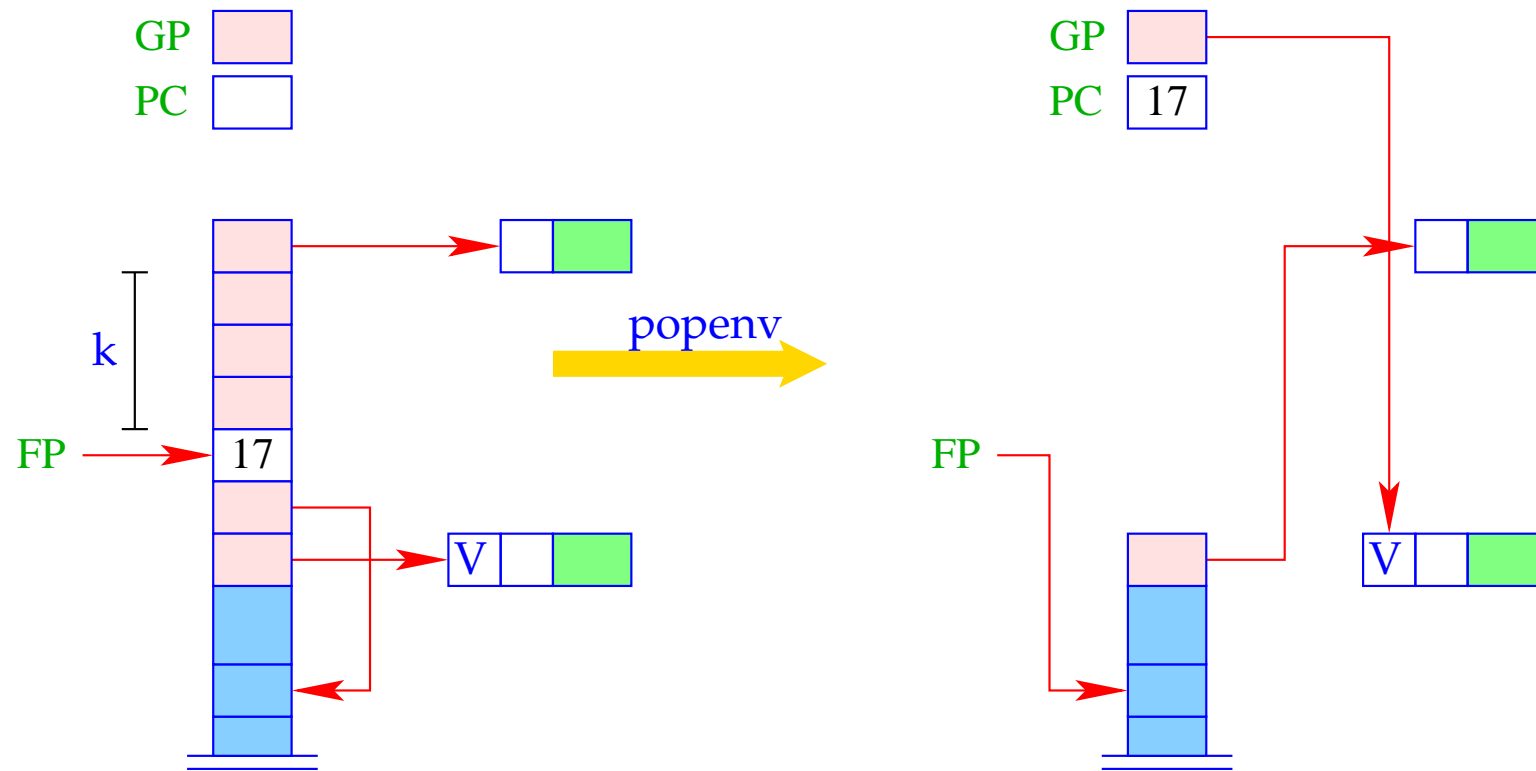
```

return k = if (SP - FP = k + 1)
    popenv;           // Done
else {               // There are more arguments
    slide k;
    apply;           // another application
}

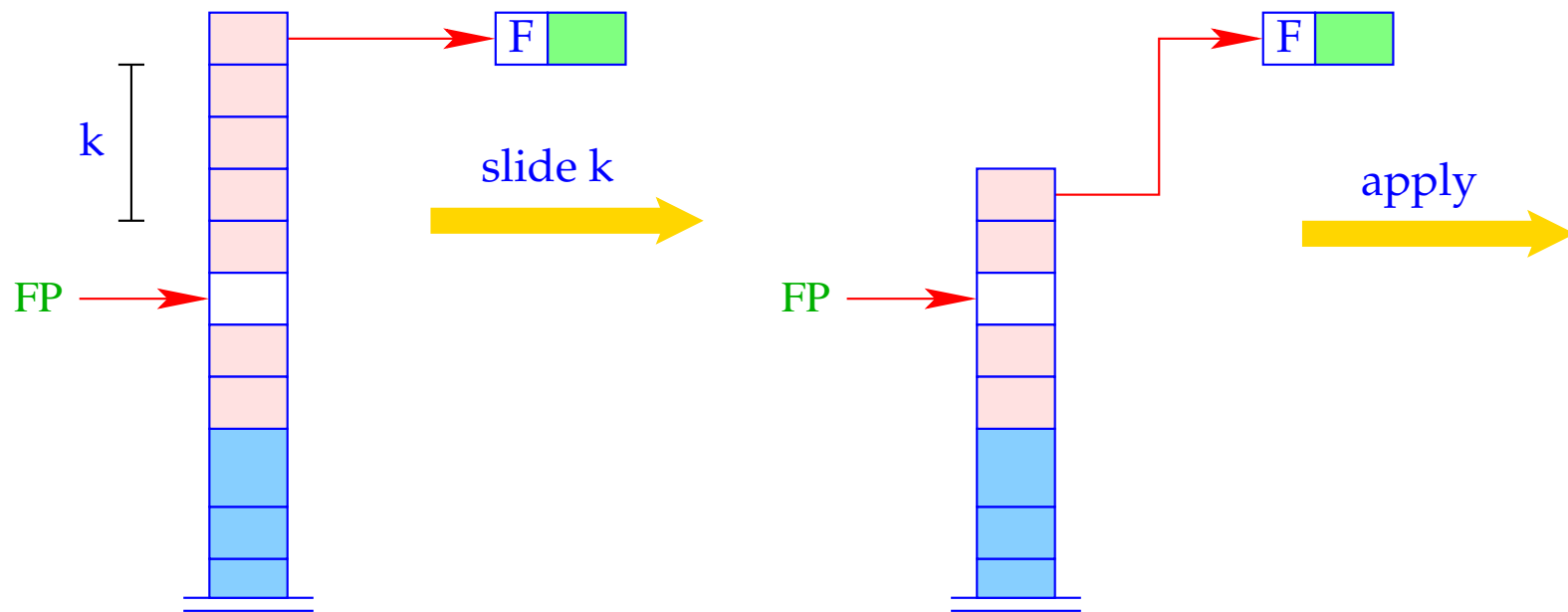
```

The execution of **return k** results in:

Case: Done



Case: Over-supply



19 letrec-Expressions

Consider the expression $e \equiv \mathbf{letrec} \ y_1 = e_1; \dots; y_n = e_n \ \mathbf{in} \ e_0$.

The translation of e must deliver an instruction sequence that

- allocates local variables y_1, \dots, y_n ;
- in the case of
 - CBV**: evaluates e_1, \dots, e_n and binds the y_i to their values;
 - CBN**: constructs closures for the e_1, \dots, e_n and binds the y_i to them;
- evaluates the expression e_0 and returns its value.

Warning:

In a **letrec**-expression, the definitions can use variables that will be allocated only **later!** \implies **Dummy**-values are put onto the stack before processing the definition.

For **CBN**, we obtain:

```
codeV e ρ sd = alloc n           // allocates local variables
                codeC e1 ρ' (sd + n)
                rewrite n
                ...
                codeC en ρ' (sd + n)
                rewrite 1
                codeV e0 ρ' (sd + n)
                slide n           // deallocates local variables
```

where $\rho' = \rho \oplus \{y_i \mapsto (L, \text{sd} + i) \mid i = 1, \dots, n\}$.

In the case of **CBV**, we also use `codeV` for the expressions e_1, \dots, e_n .

Warning:

Recursive definitions of basic values are **undefined** with **CBV!!!**

Example:

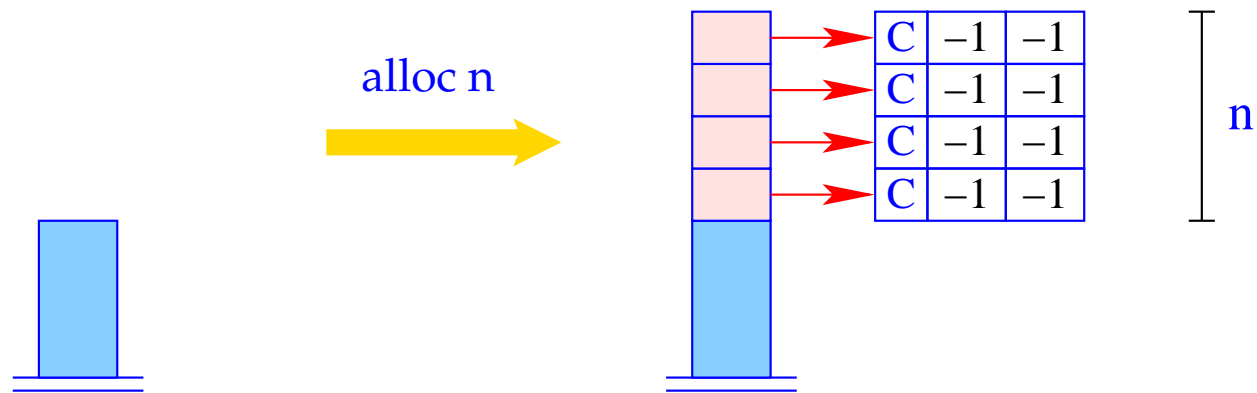
Consider the expression

$$e \equiv \mathbf{letrec} \ f = \mathbf{fn}x, y \Rightarrow \mathbf{if}y \leq 1 \ \mathbf{then} \ x \ \mathbf{else} \ f(x * y)(y - 1) \ \mathbf{in} \ f1$$

for $\rho = \emptyset$ and $\mathbf{sd} = 0$. We obtain (for **CBV**):

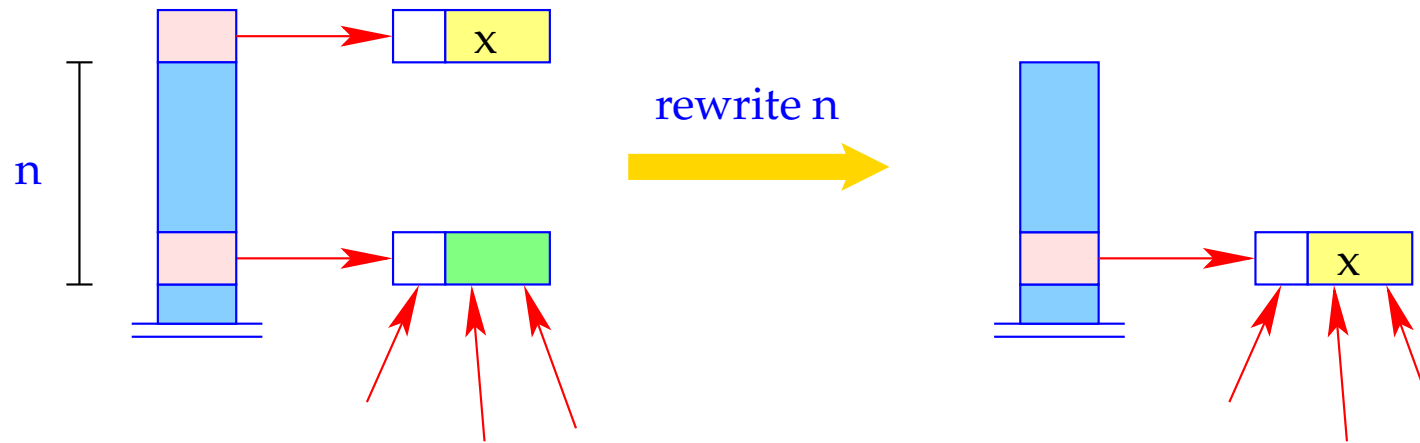
0	alloc 1	0	A:	targ 2	4	loadc 1
1	pushloc 0	0		...	5	mkbasic
2	mkvec 1	1		return 2	5	pushloc 4
2	mkfunval A	2	B:	rewrite 1	6	apply
2	jump B	1		mark C	2	C: slide 1

The instruction `alloc n` reserves n cells on the stack and initialises them with n dummy nodes:



```
for (i=1; i<=n; i++)  
    S[SP+i] = new (C,-1,-1);  
SP = SP + n;
```

The instruction `rewrite n` overwrites the contents of the heap cell pointed to by the reference at $S[SP-n]$:



$H[S[SP-n]] = H[S[SP]];$
 $SP = SP - 1;$

- The **reference** $S[SP - n]$ remains unchanged!
- Only its **contents** is changed!

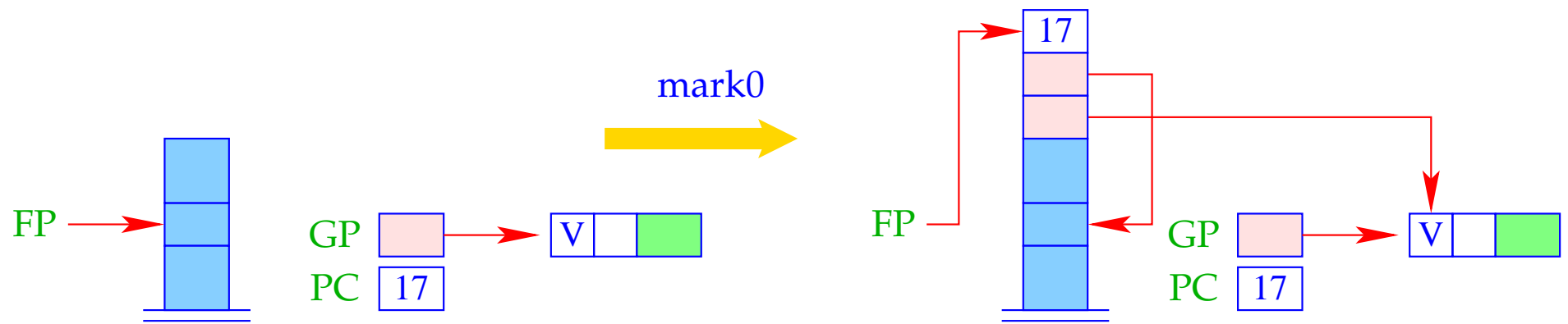
20 Closures and their Evaluation

- Closures are needed only for the implementation of CBN.
- Before the value of a variable is accessed (with CBN), this value **must** be available.
- Otherwise, a stack frame must be created to determine this value.
- This task is performed by the instruction **eval**.

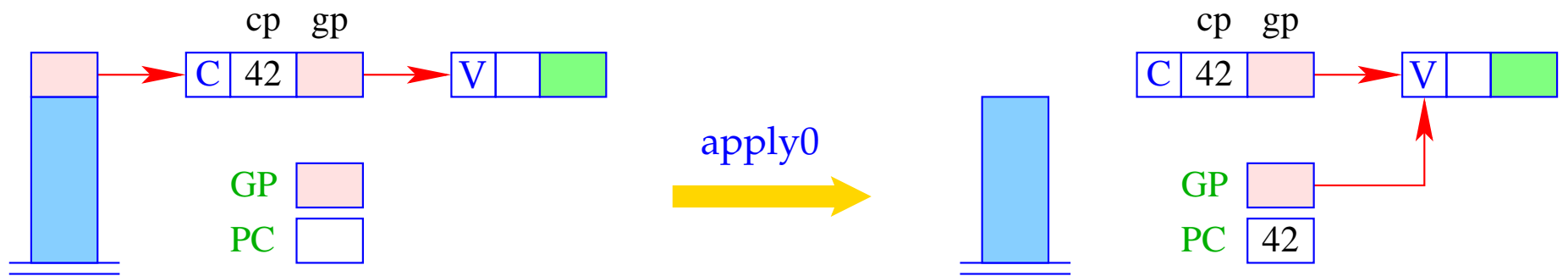
`eval` can be decomposed into small actions:

```
eval = if (H[S[SP]] ≡ (C, _, _)) {  
    mark0;           // allocation of the stack frame  
    pushloc 3;      // copying of the reference  
    apply0;         // corresponds to apply  
}
```

- A closure can be understood as a parameterless function. Thus, there is no need for an `ap`-component.
- Evaluation of the closure thus means evaluation of an application of this function to 0 arguments.
- In contrast to `mark A`, `mark0` dumps the current `PC`.
- The difference between `apply` and `apply0` is that no argument vector is put on the stack.



$S[SP+1] = GP;$
 $S[SP+2] = FP;$
 $S[SP+3] = PC;$
 $FP = SP = SP + 3;$

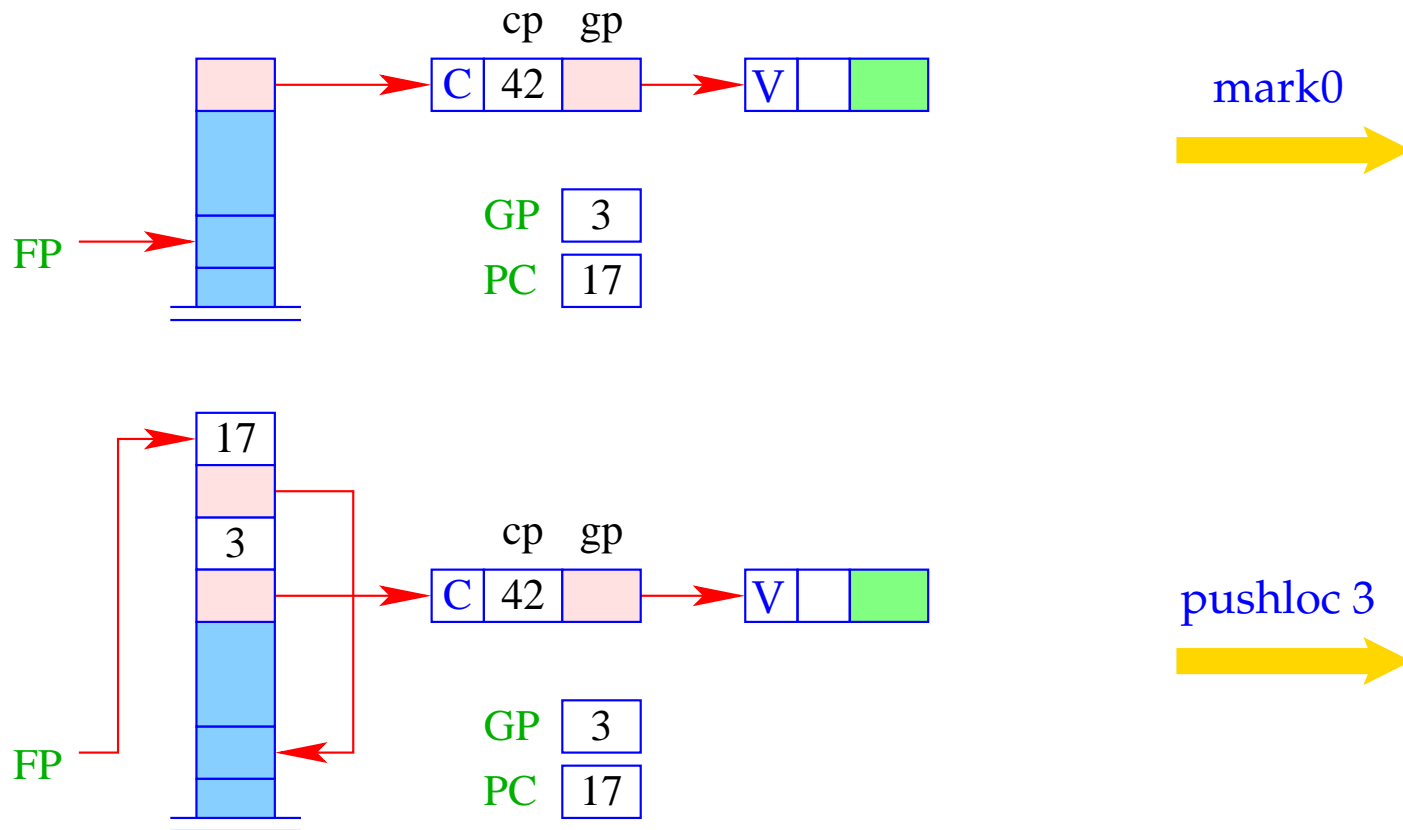


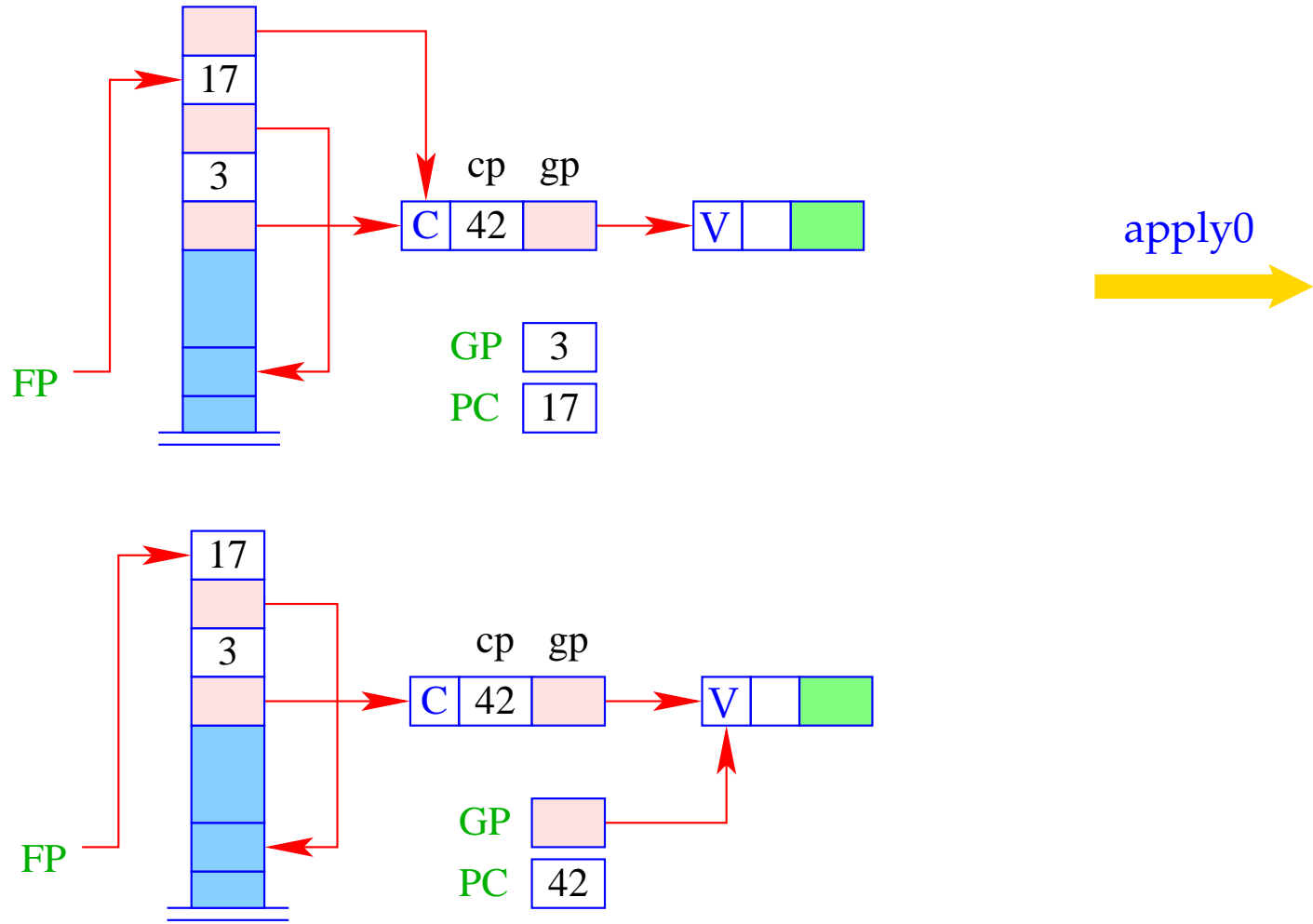
```

h = S[SP]; SP--;
GP = h->gp; PC = h->cp;

```

We thus obtain for the instruction `eval`:





The **construction** of a closure for an expression e consists of:

- Packing the bindings for the free variables into a vector;
- Creation of a C-object, which contains a reference to this vector and to the code for the evaluation of e :

```

codeC e ρ sd =      getvar z0 ρ sd
                     getvar z1 ρ (sd + 1)
                     ...
                     getvar zg-1 ρ (sd + g - 1)
                     mkvec g
                     mkclos A
                     jump B
A : codeV e ρ' 0
    update
B : ...

```

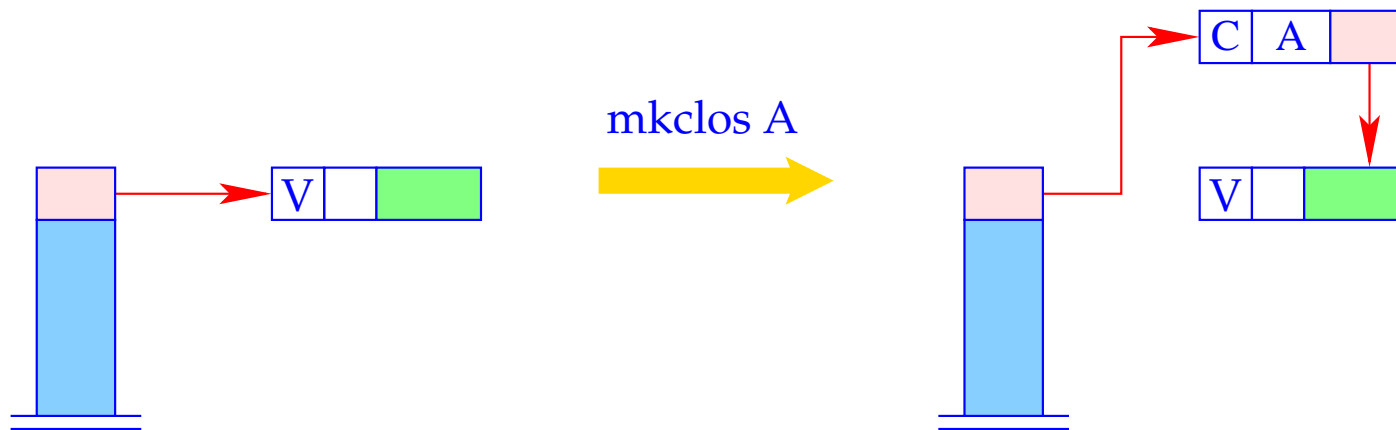
where $\{z_0, \dots, z_{g-1}\} = \text{free}(e)$ and $\rho' = \{z_i \mapsto (G, i) \mid i = 0, \dots, g - 1\}$.

Example:

Consider $e \equiv a * a$ with $\rho = \{a \mapsto (L, 0)\}$ and $sd = 1$. We obtain:

1	pushloc 1	0	A:	pushglob 0	2	getbasic
2	mkvec 1	1		eval	2	mul
2	mkcloc A	1		getbasic	1	mkbasic
2	jump B	1		pushglob 0	1	update
		2		eval	2	B: ...

- The instruction `mkclos A` is analogous to the instruction `mkfunval A`.
- It generates a C-object, where the included code pointer is `A`.

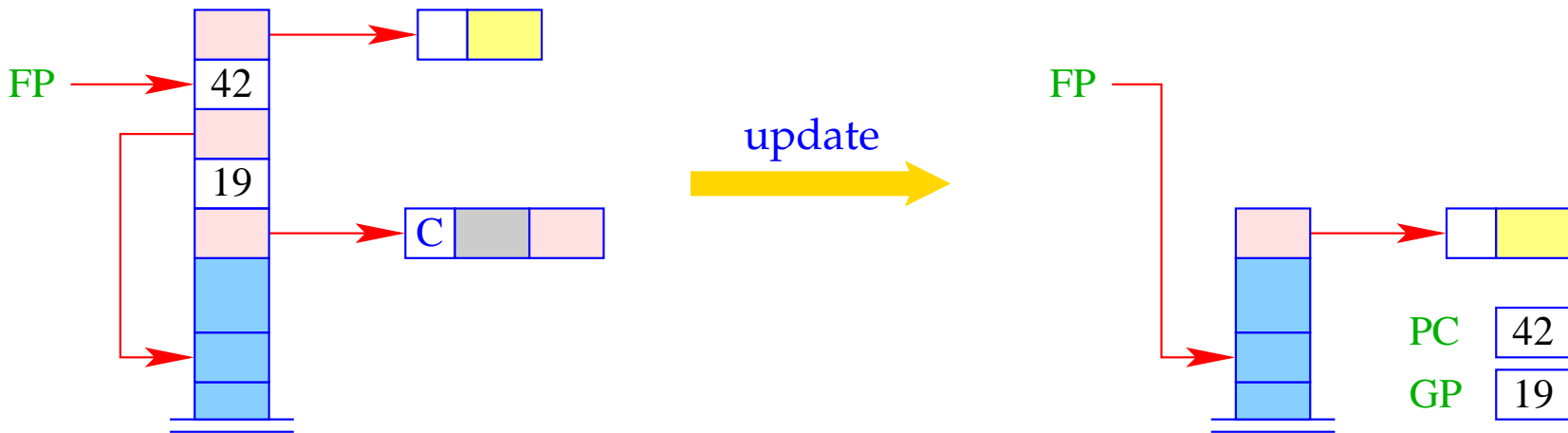


`S[SP] = new (C, A, S[SP]);`

In fact, the instruction `update` is the combination of the two actions:

`popenv`
`rewrite 1`

It overwrites the closure with the computed value.



21 Optimizations I: Global Variables

Observation:

- Functional programs construct many F- and C-objects.
- This requires the inclusion of (the bindings of) all global variables.
Recall, e.g., the construction of a closure for an expression e ...

```

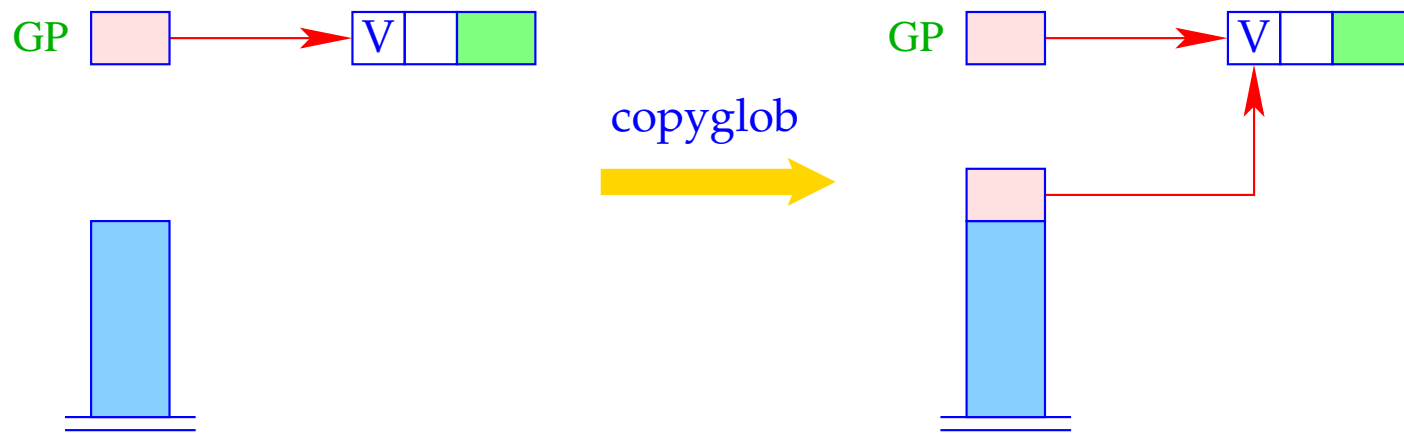
codeC e ρ sd =   getvar z0 ρ sd
                   getvar z1 ρ (sd + 1)
                   ...
                   getvar zg-1 ρ (sd + g - 1)
                   mkvec g
                   mkclos A
                   jump B
A : codeV e ρ' 0
    update
B : ...

```

where $\{z_0, \dots, z_{g-1}\} = \text{free}(e)$ and $\rho' = \{z_i \mapsto (G, i) \mid i = 0, \dots, g-1\}$.

Idea:

- **Reuse** Global Vectors, i.e. share Global Vectors!
- Profitable in the translation of **let**-expressions or function applications: Build one Global Vector for the union of the free-variable sets of all let-definitions resp. all arguments.
- Allocate (references to) global vectors with multiple uses in the stack frame like local variables!
- Support the access to the current **GP** by an instruction **copyglob** :



```
SP++;  
S[SP] = GP;
```

- The optimization will cause Global Vectors to contain **more** components than just references to the free the variables that occur in one expression ...

Disadvantage: Superfluous components in Global Vectors prevent the deallocation of already useless heap objects \implies **Space Leaks :-)**

Potential Remedy: Deletion of references at the end of their life time.

22 Optimizations II: Closures

In some cases, the construction of closures can be avoided, namely for

- Basic values,
- Variables,
- Functions.

Basic Values:

The construction of a closure for the value is at least as expensive as the construction of the B-object itself!

Therefore:

$$\text{code}_C b \rho sd = \text{code}_V b \rho sd = \begin{array}{l} \text{loadc b} \\ \text{mkbasic} \end{array}$$

This replaces:

mkvec 0		jump B	mkbasic	B:	...
mkclos A	A:	loadc b	update		

Variables:

Variables are either bound to values or to C-objects. Constructing another closure is therefore superfluous. Therefore:

$$\text{code}_C x \rho \text{sd} = \text{getvar } x \rho \text{sd}$$

This replaces:

<code>getvar</code> $x \rho \text{sd}$	<code>mkclos</code> A	A:	<code>pushglob</code> 0	<code>update</code>
<code>mkvec</code> 1	<code>jump</code> B		<code>eval</code>	B: ...

Example:

$e \equiv \text{letrec } a = b; b = 7 \text{ in } a.$ `codeV` $e \emptyset 0$ produces:

0	<code>alloc</code> 2	3	<code>rewrite</code> 2	3	<code>mkbasic</code>	2	<code>pushloc</code> 1
2	<code>pushloc</code> 0	2	<code>loadc</code> 7	3	<code>rewrite</code> 1	3	<code>eval</code>
						3	<code>slide</code> 2

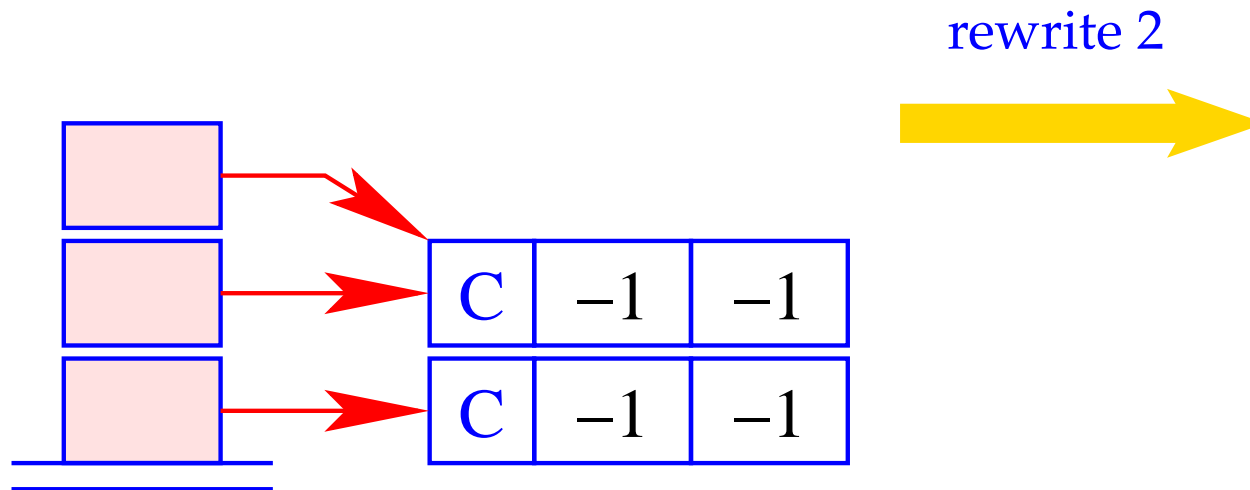
The execution of this instruction sequence should deliver the basic value 7 ...

0	alloc 2	3	rewrite 2	3	mkbasic	2	pushloc 1
2	pushloc 0	2	loadc 7	3	rewrite 1	3	eval
						3	slide 2

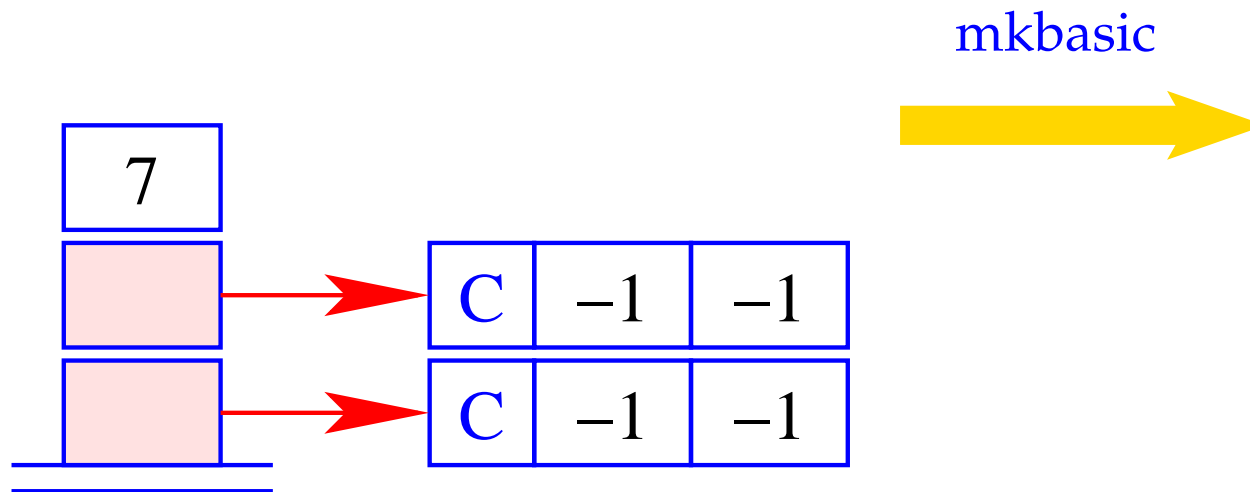
alloc 2



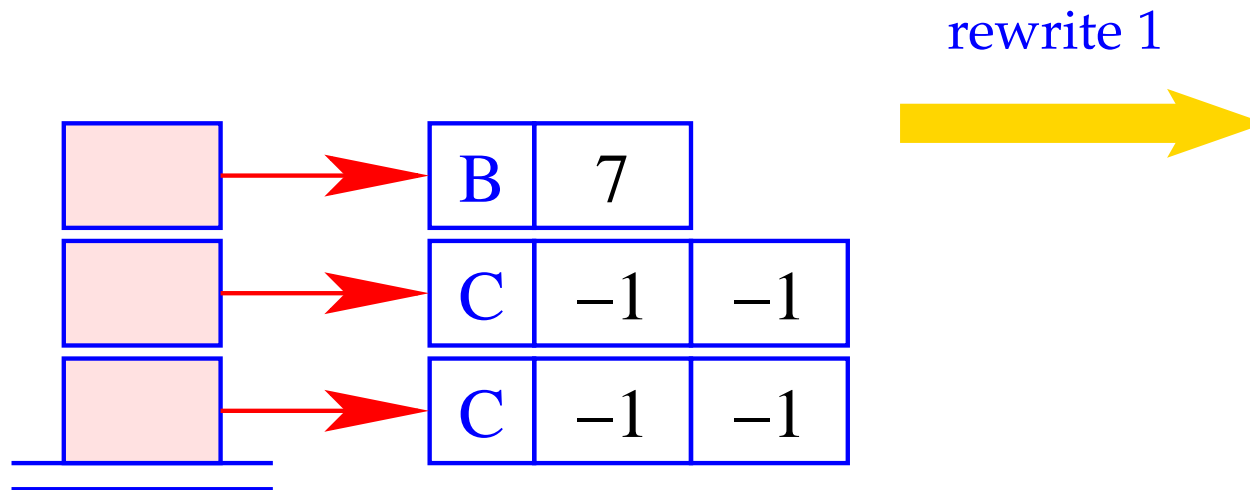
0	alloc 2	3	rewrite 2	3	mkbasic	2	pushloc 1
2	pushloc 0	2	loadc 7	3	rewrite 1	3	eval
						3	slide 2



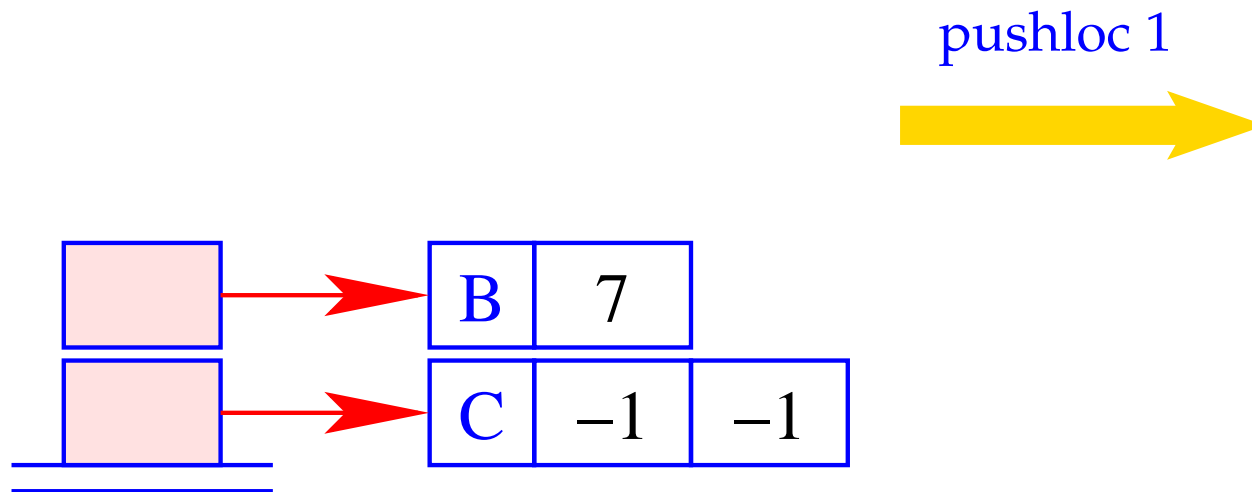
0	alloc 2	3	rewrite 2	3	mkbasic	2	pushloc 1
2	pushloc 0	2	loadc 7	3	rewrite 1	3	eval
						3	slide 2



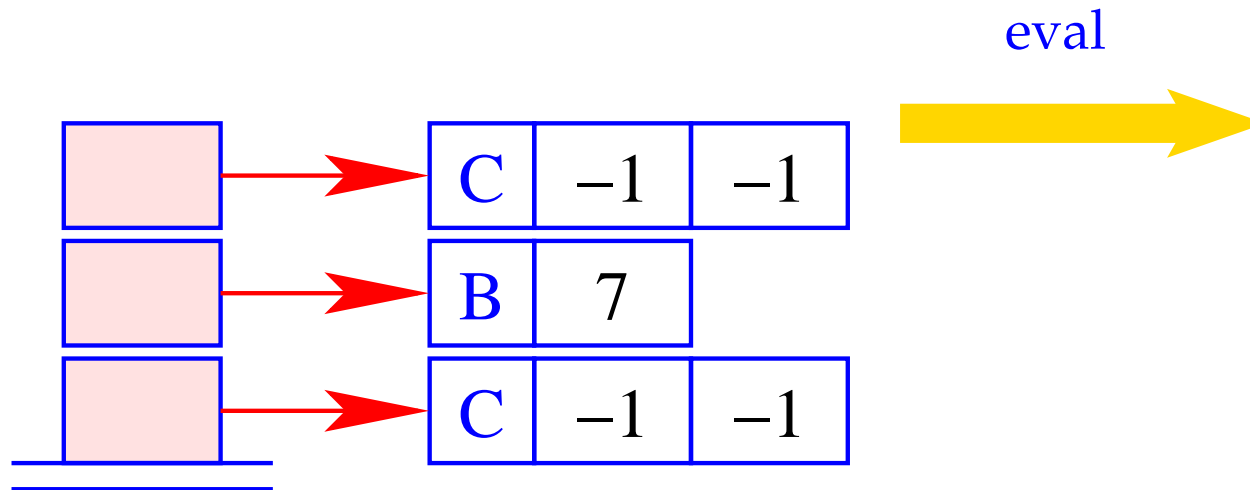
0	alloc 2	3	rewrite 2	3	mkbasic	2	pushloc 1
2	pushloc 0	2	loadc 7	3	rewrite 1	3	eval
						3	slide 2



0	alloc 2	3	rewrite 2	3	mkbasic	2	pushloc 1
2	pushloc 0	2	loadc 7	3	rewrite 1	3	eval
						3	slide 2



0	alloc 2	3	rewrite 2	3	mkbasic	2	pushloc 1
2	pushloc 0	2	loadc 7	3	rewrite 1	3	eval
						3	slide 2



0	alloc 2	3	rewrite 2	3	mkbasic	2	pushloc 1
2	pushloc 0	2	loadc 7	3	rewrite 1	3	eval
						3	slide 2

Segmentation Fault !!

Apparently, this optimization was not quite **correct** :-)

The Problem:

Binding of variable y to variable x **before** x 's dummy node is replaced!!



The Solution:

cyclic definitions: reject sequences of definitions like

let $a = b; \dots b = a$ **in** \dots

acyclic definitions: order the definitions $y = x$ such that the dummy node for the right side of x is already overwritten.

Functions:

Functions are values, which are not evaluated further. Instead of generating code that constructs a closure for an F-object, we generate code that constructs the F-object directly.

Therefore:

$$\text{code}_C (\text{fn } x_0, \dots, x_{k-1} \Rightarrow e) \rho \text{sd} = \text{code}_V (\text{fn } x_0, \dots, x_{k-1} \Rightarrow e) \rho \text{sd}$$

23 The Translation of a Program Expression

Execution of a program e starts with

$$PC = 0 \quad SP = FP = GP = -1$$

The expression e must not contain **free variables**.

The value of e should be determined and then a **halt** instruction should be executed.

$$\text{code } e = \text{code}_V e \ \emptyset \ 0 \\ \text{halt}$$

Remarks:

- The code schemata as defined so far produce **Spaghetti code**.
- Reason: Code for function bodies and closures placed directly behind the instructions **mkfunval** resp. **mkclos** with a jump over this code.
- Alternative: Place this code somewhere else, e.g. **following** the **halt**-instruction:

Advantage: Elimination of the direct jumps following **mkfunval** and **mkclos**.

Disadvantage: The code schemata are more complex as they would have to accumulate the code pieces in a **Code-Dump**.



Solution:

Disentangle the Spaghetti code in a subsequent optimization phase :-)

Example: **let** $a = 17$; $f = \mathbf{fn}$ $b \Rightarrow a + b$ **in** f 42

Disentanglement of the jumps produces:

0	loadc 17	2	mark B	3	B:	slide 2	1	pushloc 1
1	mkbasic	5	loadc 42	1		halt	2	eval
1	pushloc 0	6	mkbasic	0	A:	targ 1	2	getbasic
2	mkvec 1	6	pushloc 4	0		pushglob 0	2	add
2	mkfunval A	7	eval	1		eval	1	mkbasic
		7	apply	1		getbasic	1	return 1

24 Structured Data

In the following, we extend our functional programming language by some datatypes.

24.1 Tuples

Constructors: $(., \dots, .)$, k -ary with $k \geq 0$;

Destructors: $\#j$ for $j \in \mathbb{N}_0$ (Projections)

We extend the syntax of expressions correspondingly:

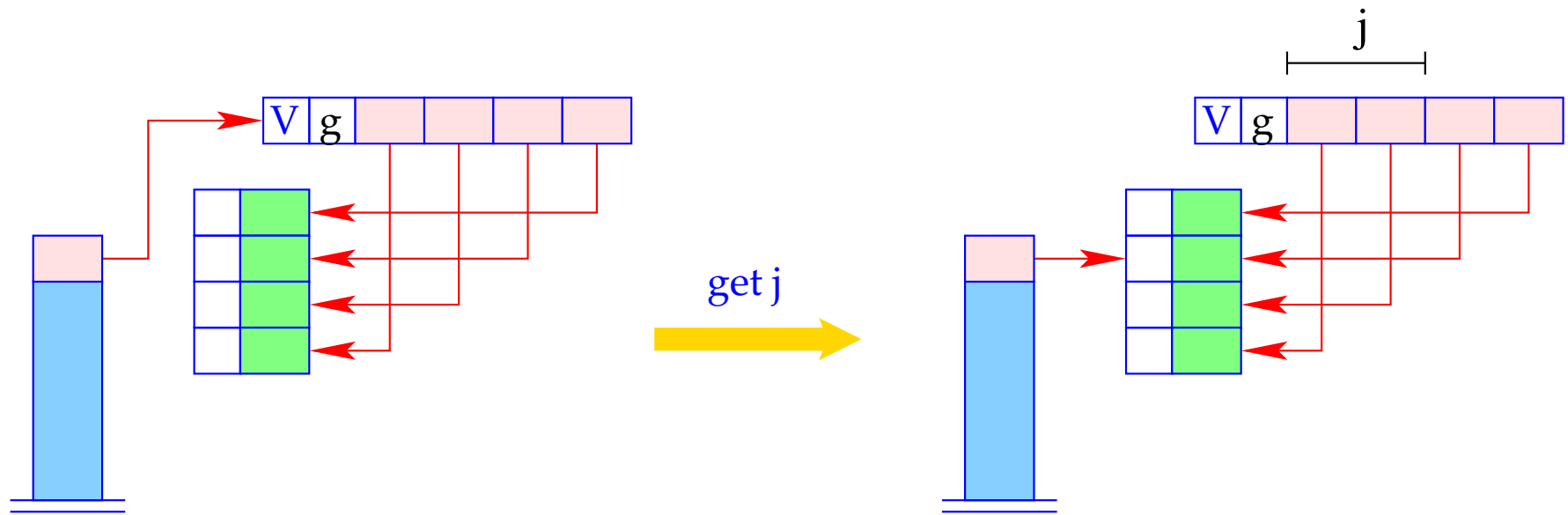
$$e ::= \dots \mid (e_0, \dots, e_{k-1}) \mid \#j e \\ \mid \mathbf{let} (x_0, \dots, x_{k-1}) = e_1 \mathbf{in} e_0$$

- In order to **construct** a tuple, we collect sequence of references on the stack. Then we construct a vector of these references in the heap using **mkvec**
- For returning **components** we use an indexed access into the tuple.

$$\begin{aligned}
 \text{code}_V (e_0, \dots, e_{k-1}) \rho \text{sd} &= \text{code}_C e_0 \rho \text{sd} \\
 &\quad \text{code}_C e_1 \rho (\text{sd} + 1) \\
 &\quad \dots \\
 &\quad \text{code}_C e_{k-1} \rho (\text{sd} + k - 1) \\
 &\quad \text{mkvec } k
 \end{aligned}$$

$$\begin{aligned}
 \text{code}_V (\#j e) \rho \text{sd} &= \text{code}_V e \rho \text{sd} \\
 &\quad \text{get } j
 \end{aligned}$$

In the case of **CBV**, we directly compute the values of the e_i .



```

if (S[SP] == (V,g,v))
  S[SP] = v[j];
else Error "Vector expected!";

```

Inversion: Accessing all components of a tuple simultaneously:

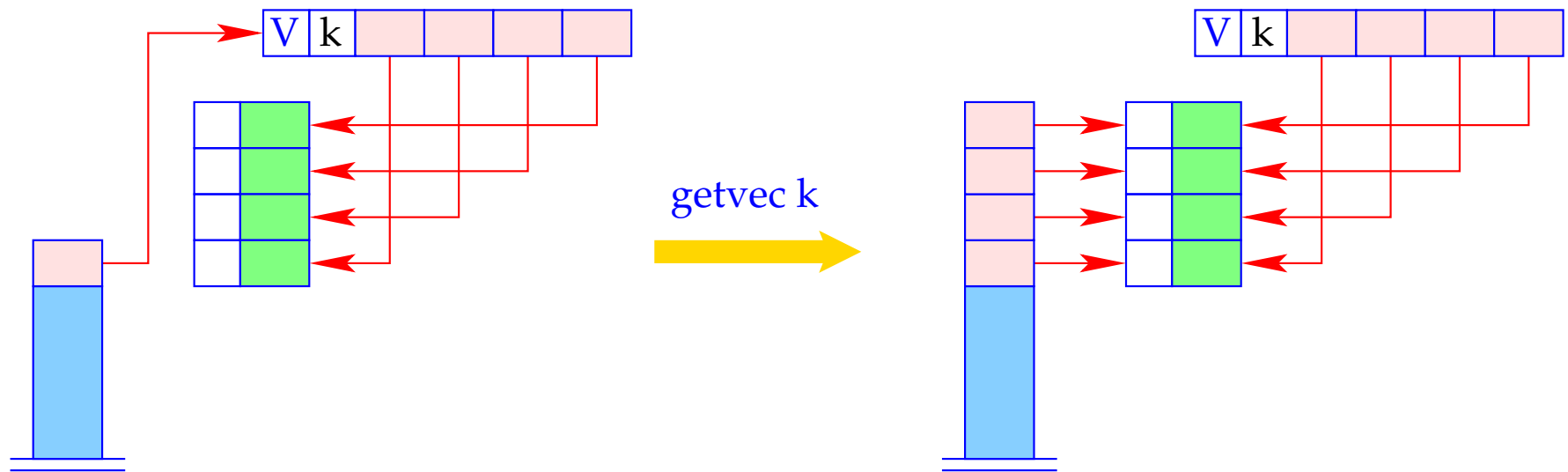
$$e \equiv \mathbf{let} (y_0, \dots, y_{k-1}) = e_1 \mathbf{in} e_0$$

This is translated as follows:

$$\begin{aligned} \mathbf{code}_V e \rho \mathbf{sd} &= \mathbf{code}_V e_1 \rho \mathbf{sd} \\ &\quad \mathbf{getvec} \ k \\ &\quad \mathbf{code}_V e_0 \rho' (\mathbf{sd} + k) \\ &\quad \mathbf{slide} \ k \end{aligned}$$

where $\rho' = \rho \oplus \{y_i \mapsto (L, \mathbf{sd} + i) \mid i = 0, \dots, k - 1\}$.

The instruction $\mathbf{getvec} \ k$ pushes the components of a vector of length k onto the stack:



```

if (S[SP] == (V,k,v)) {
    SP--;
    for(i=0; i<k; i++) {
        SP++; S[SP] = v[i];
    }
} else Error "Vector expected!";

```

24.2 Lists

Lists are constructed by the **constructors**:

`[]` “Nil”, the empty list;

`“:”` “Cons”, right-associative, takes an element and a list.

Access to list components is possible by **case-expressions** ...

Example: The append function `app`:

$$\begin{aligned} \text{app} &= \text{fn } l, y \Rightarrow \text{case } l \text{ of} \\ &\quad [] \quad \rightarrow \quad y \\ &\quad h : t \quad \rightarrow \quad h : (\text{app } t \ y) \end{aligned}$$

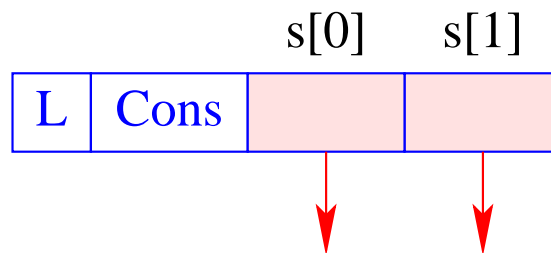
accordingly, we extend the syntax of expressions:

$$e ::= \dots \mid [] \mid (e_1 : e_2) \\ \mid (\mathbf{case} \ e_0 \ \mathbf{of} \ [] \rightarrow e_1; \ h : t \rightarrow e_2)$$

Additionally, we need new heap objects:



empty list



non-empty list

24.3 Building Lists

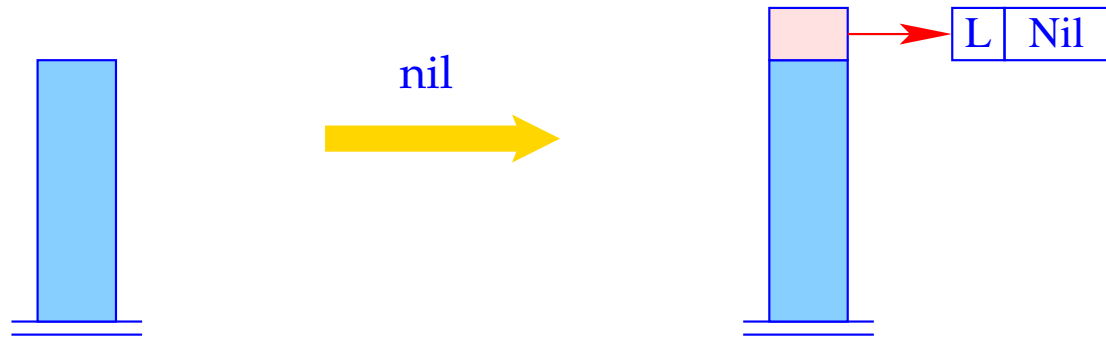
The new instructions `nil` and `cons` are introduced for building list nodes.

We translate for **CBN**:

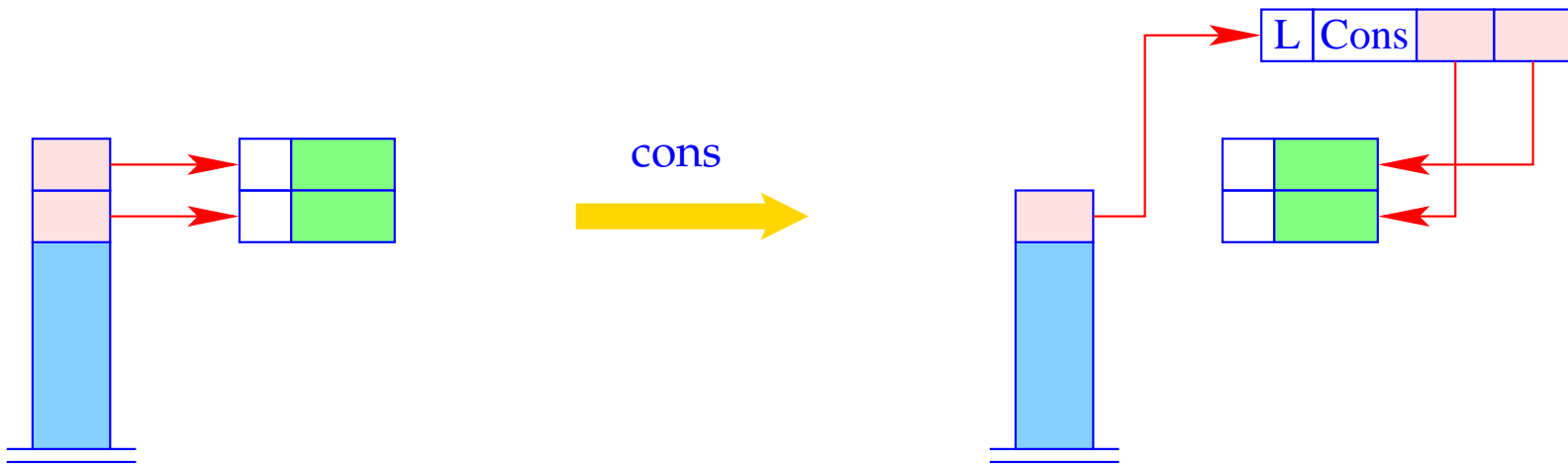
$$\begin{aligned} \text{code}_V [] \rho \text{sd} &= \text{nil} \\ \text{code}_V (e_1 : e_2) \rho \text{sd} &= \text{code}_C e_1 \rho \text{sd} \\ &\quad \text{code}_C e_2 \rho (\text{sd} + 1) \\ &\quad \text{cons} \end{aligned}$$

Note:

- With **CBN**: Closures are constructed for the arguments of ":";
- With **CBV**: Arguments of ":" are evaluated :-)



$S[SP] = SP++; S[SP] = \text{new}(L, \text{Nil});$



$S[SP-1] = \text{new } (L, \text{Cons}, S[SP-1], S[SP]);$
 $SP--;$

24.4 Pattern Matching

Consider the expression $e \equiv \mathbf{case} \ e_0 \ \mathbf{of} \ [] \rightarrow e_1; \ h : t \rightarrow e_2$.

Evaluation of e requires:

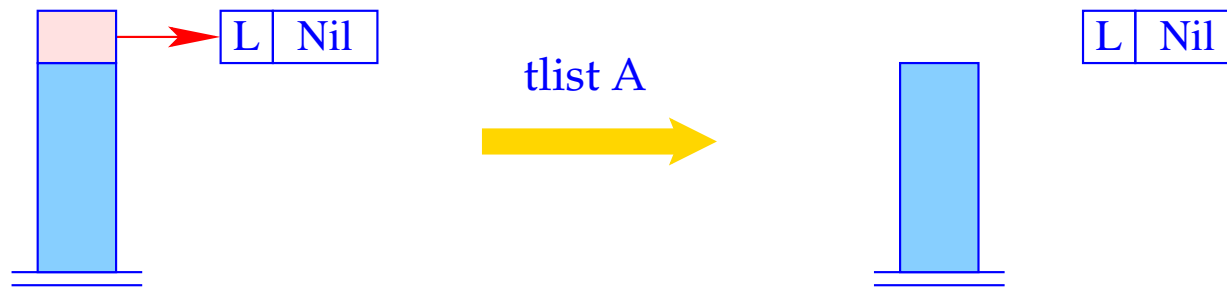
- evaluation of e_0 ;
- check, whether resulting value v is an L-object;
- if v is the empty list, evaluation of e_1 ...
- otherwise storing the two references of v on the stack and evaluation of e_2 .
This corresponds to **binding** h and t to the two components of v .

In consequence, we obtain (for CBN as for CBV):

$$\begin{aligned} \text{code}_V e \rho \text{sd} &= && \text{code}_V e_0 \rho \text{sd} \\ &&& \text{tlist A} \\ &&& \text{code}_V e_1 \rho \text{sd} \\ &&& \text{jump B} \\ &&& \text{A : } \text{code}_V e_2 \rho' (\text{sd} + 2) \\ &&& \text{slide 2} \\ &&& \text{B : } \dots \end{aligned}$$

where $\rho' = \rho \oplus \{h \mapsto (L, \text{sd} + 1), t \mapsto (L, \text{sd} + 2)\}$.

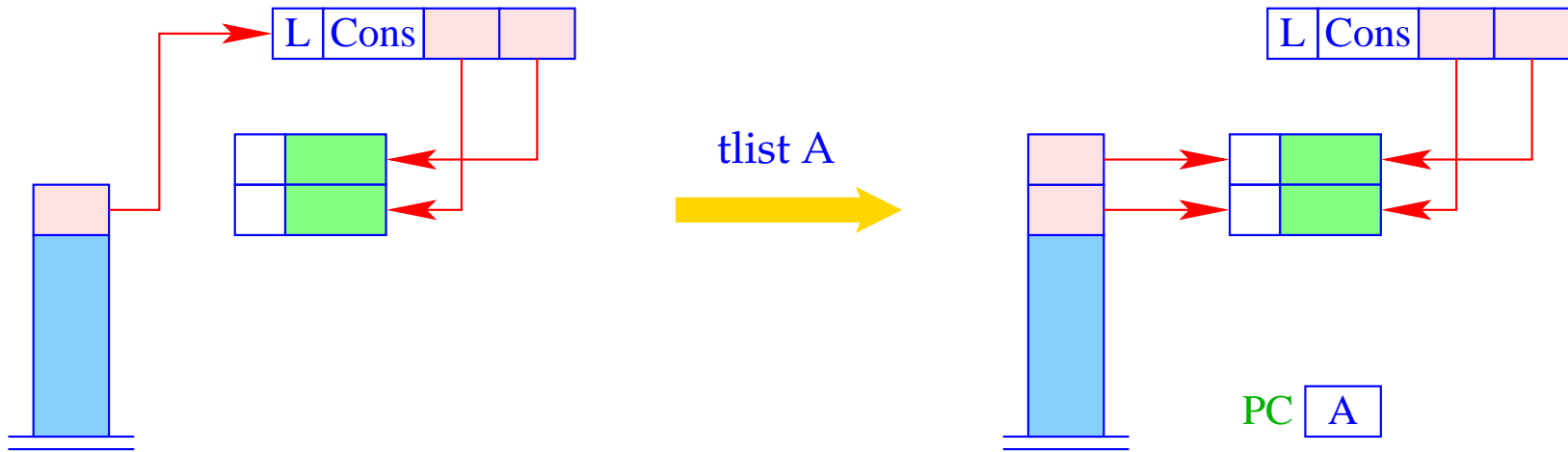
The new instruction `tlist A` does the necessary checks and (in the case of Cons) allocates two new local variables:



```

h = S[SP];
if (H[h] != (L,...))
    Error "no list!";
if (H[h] == (_,Nil)) SP- -;
...

```



```

... else {
  S[SP+1] = S[SP]→s[1];
  S[SP] = S[SP]→s[0];
  SP++; PC = A;
}

```

Example: The (disentangled) body of the function `app` with
`app` \mapsto $(G, 0)$:

0	targ 2	3	pushglob 0	0	C:	mark D
0	pushloc 0	4	pushloc 2	3		pushglob 2
1	eval	5	pushloc 6	4		pushglob 1
1	tlist A	6	mkvec 3	5		pushglob 0
0	pushloc 1	4	mkclos C	6		eval
1	eval	4	cons	6		apply
1	jump B	0	slide 2	1	D:	update
2	A: pushloc 1	3	B: return 2			

Note:

Datatypes with more than two constructors need a generalization of the `tlist` instruction, corresponding to a `switch`-instruction :-)

24.5 Closures of Tuples and Lists

The general schema for `codeC` can be optimized for tuples and lists:

$$\begin{aligned} \text{code}_C (e_0, \dots, e_{k-1}) \rho \text{sd} &= \text{code}_V (e_0, \dots, e_{k-1}) \rho \text{sd} = \text{code}_C e_0 \rho \text{sd} \\ & \quad \text{code}_C e_1 \rho (\text{sd} + 1) \\ & \quad \dots \\ & \quad \text{code}_C e_{k-1} \rho (\text{sd} + k - 1) \\ & \quad \text{mkvec } k \\ \\ \text{code}_C [] \rho \text{sd} &= \text{code}_V [] \rho \text{sd} = \text{nil} \\ \\ \text{code}_C (e_1 : e_2) \rho \text{sd} &= \text{code}_V (e_1 : e_2) \rho \text{sd} = \text{code}_C e_1 \rho \text{sd} \\ & \quad \text{code}_C e_2 \rho (\text{sd} + 1) \\ & \quad \text{cons} \end{aligned}$$

25 Last Calls

A function application is called **last call** in an expression e if this application could deliver the value for e .

A last call usually is the **outermost** application of a defining expression.

A function definition is called **tail recursive** if all recursive calls are last calls.

Examples:

$r\ t\ (h : y)$ is a **last call** in `case x of [] → y; h : t → r t (h : y)`

$f\ (x - 1)$ is **not a last call** in `if x ≤ 1 then 1 else x * f (x - 1)`

Observation: Last calls in a function body need **no new** stack frame!



Automatic transformation of tail recursion into loops!!!

The code for a last call $l \equiv (e' e_0 \dots e_{m_1})$ inside a function f with k arguments must

1. allocate the arguments e_i and evaluate e' to a function (note: all this inside f 's frame!);
2. deallocate the local variables and the k consumed arguments of f ;
3. execute an `apply`.

```

codeV l ρ sd = codeC em-1 ρ sd
               codeC em-2 ρ (sd + 1)
               ...
               codeC e0 ρ (sd + m - 1)
               codeV e' ρ (sd + m)           // Evaluation of the function
               move r (m + 1)               // Deallocation of r cells
               apply

```

where $r = sd + k$ is the number of stack cells to deallocate.

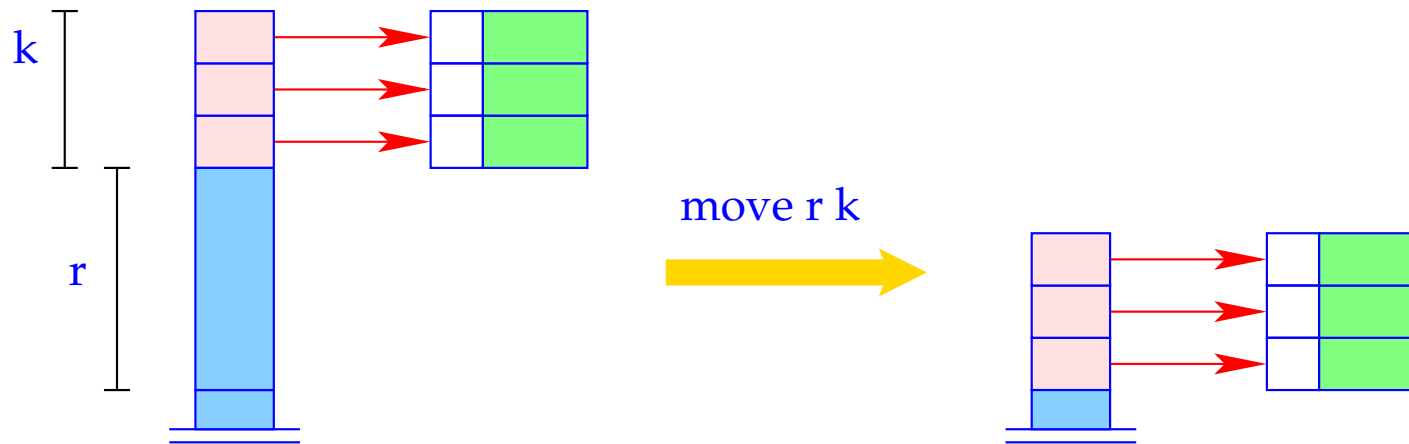
Example:

The body of the function

$$r = \mathbf{fn} \ x, y \Rightarrow \mathbf{case} \ x \ \mathbf{of} \ [] \rightarrow y; \ h : t \rightarrow r \ t \ (h : y)$$

0	targ 2	1	jump B	4	pushglob 0
0	pushloc 0			5	eval
1	eval	2	A: pushloc 1	5	move 4 3
1	tlist A	3	pushloc 4		apply
0	pushloc 1	4	cons		slide 2
1	eval	3	pushloc 1	1	B: return 2

Since the old stack frame is kept, **return 2** will only be reached by the direct jump at the end of the []-alternative.



```

SP = SP - k - r;
for (i=1; i ≤ k; i++)
    S[SP+i] = S[SP+i+r];
SP = SP + k;

```

The Translation of Logic Languages

26 The Language Proll

Here, we just consider the core language **Proll** (“Prolog-light” :-). In particular, we omit:

- arithmetic;
- the cut operator;
- self-modification of programs through **assert** and **retract**.

Example:

$\text{bigger}(X, Y) \leftarrow X = \textit{elephant}, Y = \textit{horse}$

$\text{bigger}(X, Y) \leftarrow X = \textit{horse}, Y = \textit{donkey}$

$\text{bigger}(X, Y) \leftarrow X = \textit{donkey}, Y = \textit{dog}$

$\text{bigger}(X, Y) \leftarrow X = \textit{donkey}, Y = \textit{monkey}$

$\text{is_bigger}(X, Y) \leftarrow \text{bigger}(X, Y)$

$\text{is_bigger}(X, Y) \leftarrow \text{bigger}(X, Z), \text{is_bigger}(Z, Y)$

? $\text{is_bigger}(\textit{elephant}, \textit{dog})$

A More Realistic Example:

$\text{app}(X, Y, Z) \leftarrow X = [], Y = Z$

$\text{app}(X, Y, Z) \leftarrow X = [H|X'], Z = [H|Z'], \text{app}(X', Y, Z')$

? $\text{app}(X, [Y, c], [a, b, Z])$

A More Realistic Example:

$\text{app}(X, Y, Z) \leftarrow X = [], Y = Z$

$\text{app}(X, Y, Z) \leftarrow X = [H|X'], Z = [H|Z'], \text{app}(X', Y, Z')$

? $\text{app}(X, [Y, c], [a, b, Z])$

Remark:

$[]$ \equiv the atom **empty list**

$[H|Z]$ \equiv **binary** constructor application

$[a, b, Z]$ \equiv shortcut for: $[a|[b|[Z|[]]]]$

A program p is constructed as follows:

$$\begin{aligned}t & ::= a \mid X \mid _ \mid f(t_1, \dots, t_n) \\g & ::= p(t_1, \dots, t_k) \mid X = t \\c & ::= p(X_1, \dots, X_k) \leftarrow g_1, \dots, g_r \\p & ::= c_1 \dots c_m ? g\end{aligned}$$

- A **term** t either is an atom, a variable, an anonymous variable or a constructor application.
- A **goal** g either is a literal, i.e., a predicate call, or a unification.
- A **clause** c consists of a **head** $p(X_1, \dots, X_k)$ with predicate name and list of formal parameters together with a **body**, i.e., a sequence of goals.
- A **program** consists of a sequence of clauses together with a single goal as **query**.

Procedural View of Proll programs:

goal	==	procedure call
predicate	==	procedure
body	==	definition
term	==	value
unification	==	basic computation step
binding of variables	==	side effect

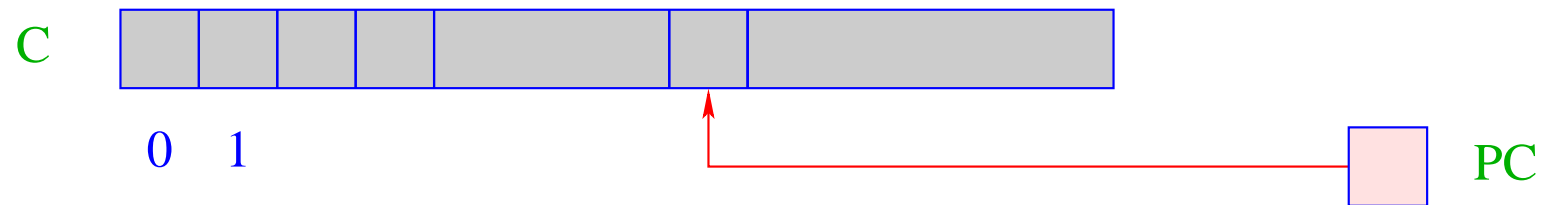
Note: Predicate calls ...

- ... do not have a return value.
- ... affect the caller through side effects only :-)
- ... may fail. Then the next definition is tried :-))

⇒ backtracking

27 Architecture of the WiM:

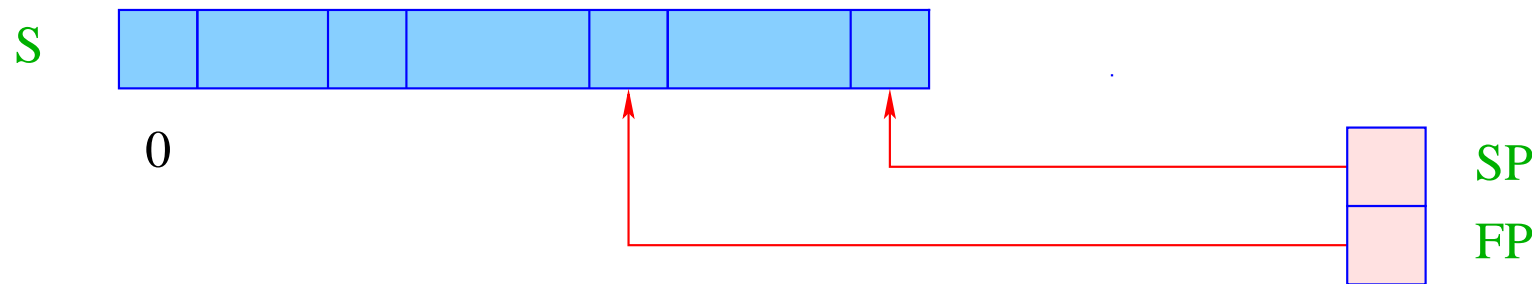
The Code Store:



C = Code store – contains WiM program;
every cell contains one instruction;

PC = Program Counter – points to the next instruction to executed;

The Runtime Stack:



S = Runtime **S**tack – every cell may contain a value or an address;

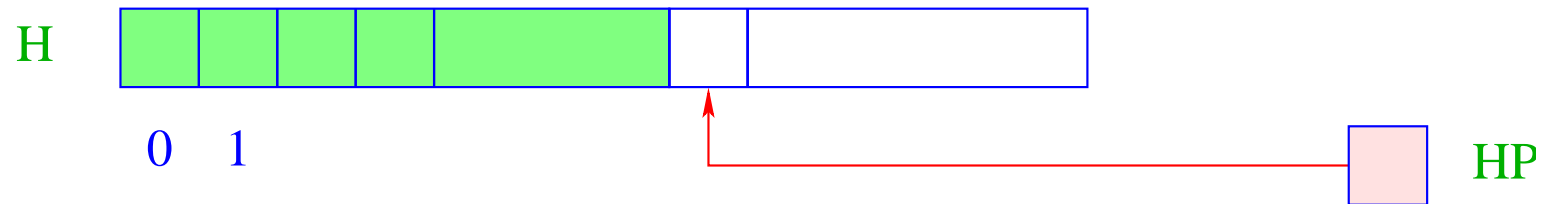
SP = **S**tack **P**ointer – points to the topmost occupied cell;

FP = **F**rame **P**ointer – points to the current stack frame.

Frames are created for predicate calls,

contain cells for each variable of the current clause

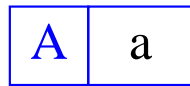
The Heap:



H = Heap for dynamicly constructed terms;

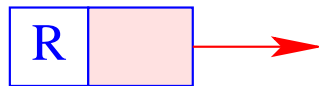
HP = Heap-Pointer – points to the first free cell;

- The heap is maintained like a **stack** as well :-)
- A **new**-instruction allocates an object in **H**.
- Objects are **tagged** with their types (as in the **MaMa**) ...



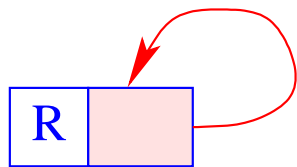
atom

1 cell



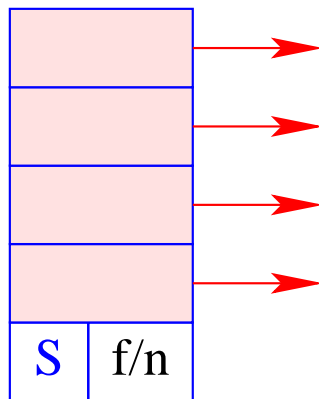
variable

1 cell



unbound variable

1 cell



structure

(n+1) cells

28 Construction of Terms in the Heap

Parameter terms of goals (calls) are constructed in the heap before passing.

Assume that the **address environment** ρ returns, for each clause variable X its address (relative to **FP**) on the stack. Then $\text{code}_A t \rho$ should ...

- construct (a presentation of) t in the heap; and
- return a reference to it on top of the stack.

Idea:

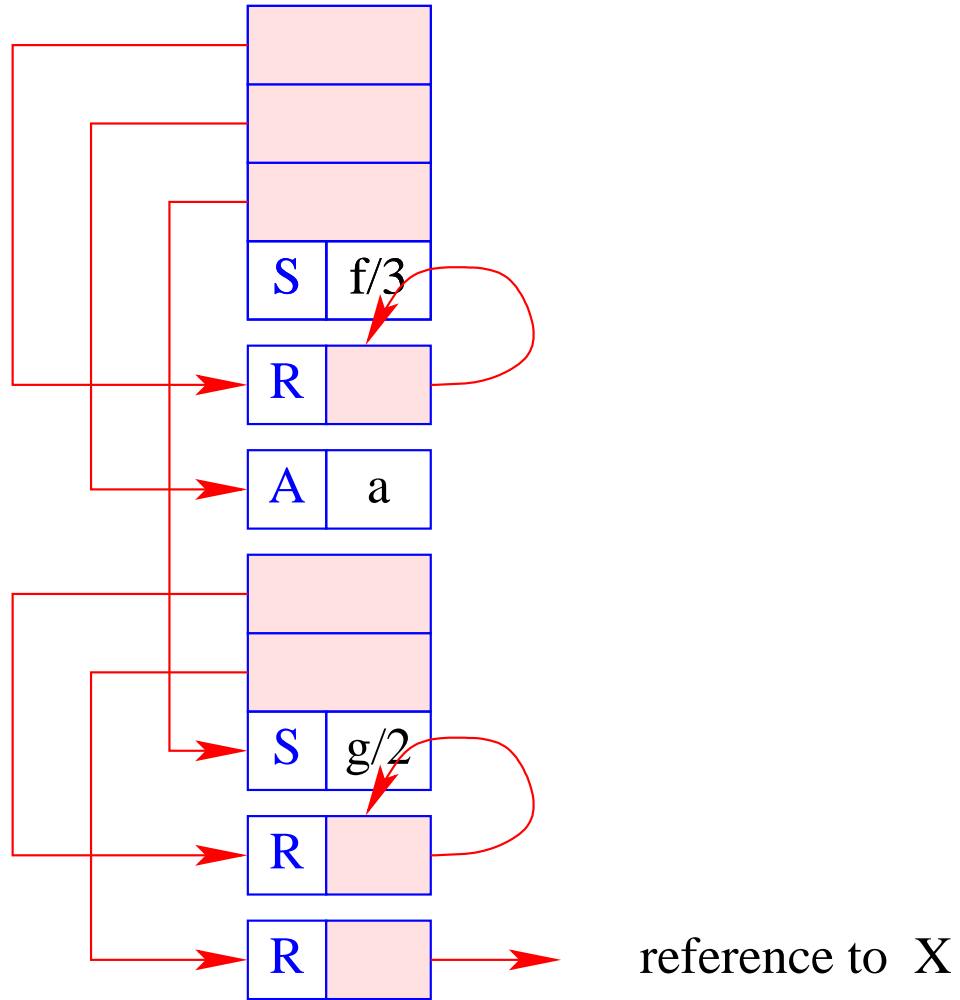
- Construct the tree during a **post-order** traversal of t
- with one instruction for each new node!

Example: $t \equiv f(g(X, Y), a, Z).$

Assume that X is **initialized**, i.e., $S[\text{FP} + \rho X]$ contains already a reference, Y and Z are not yet initialized.

Representing

$$t \equiv f(g(X, Y), a, Z) :$$



For a distinction, we mark occurrences of already initialized variables through **over-lining** (e.g. \bar{X}).

Note: Arguments are always initialized!

Then we define:

$\text{code}_A a \rho$	$=$	$\text{putatom } a$	$\text{code}_A f(t_1, \dots, t_n) \rho$	$=$	$\text{code}_A t_1 \rho$
$\text{code}_A X \rho$	$=$	$\text{putvar } (\rho X)$			\dots
$\text{code}_A \bar{X} \rho$	$=$	$\text{putref } (\rho X)$			$\text{code}_A t_n \rho$
$\text{code}_A _ \rho$	$=$	putanon			$\text{putstruct } f/n$

For a distinction, we mark occurrences of already initialized variables through over-lining (e.g. \bar{X}).

Note: Arguments are always initialized!

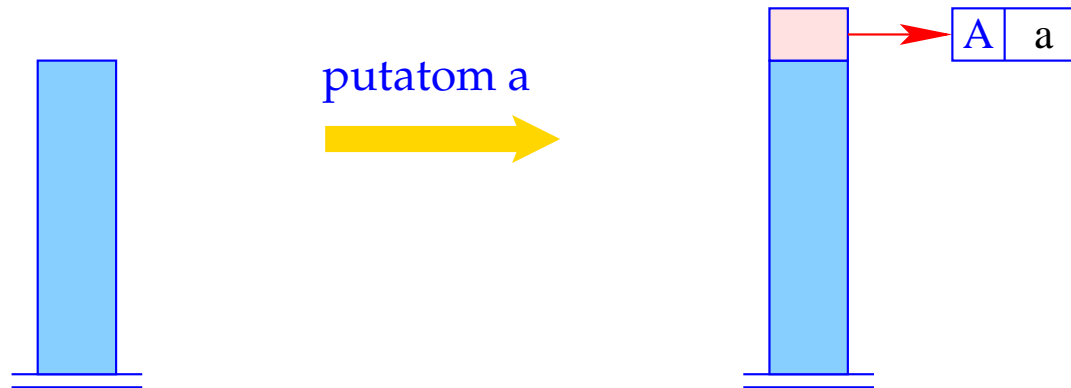
Then we define:

$$\begin{array}{ll} \text{code}_A a \rho & = \text{putatom } a & \text{code}_A f(t_1, \dots, t_n) \rho & = \text{code}_A t_1 \rho \\ \text{code}_A X \rho & = \text{putvar } (\rho X) & & \dots \\ \text{code}_A \bar{X} \rho & = \text{putref } (\rho X) & & \text{code}_A t_n \rho \\ \text{code}_A _ \rho & = \text{putanon} & & \text{putstruct } f/n \end{array}$$

For $f(g(X, Y), a, Z)$ and $\rho = \{X \mapsto 1, Y \mapsto 2, Z \mapsto 3\}$ this results in the sequence:

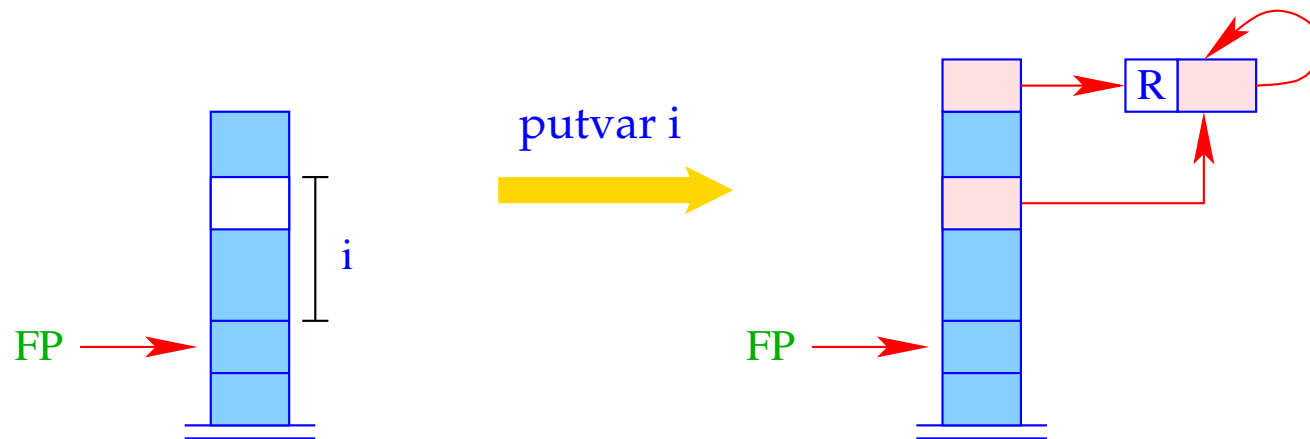
putref 1	putatom a
putvar 2	putvar 3
putstruct g/2	putstruct f/3

The instruction `putatom a` constructs an atom in the heap:



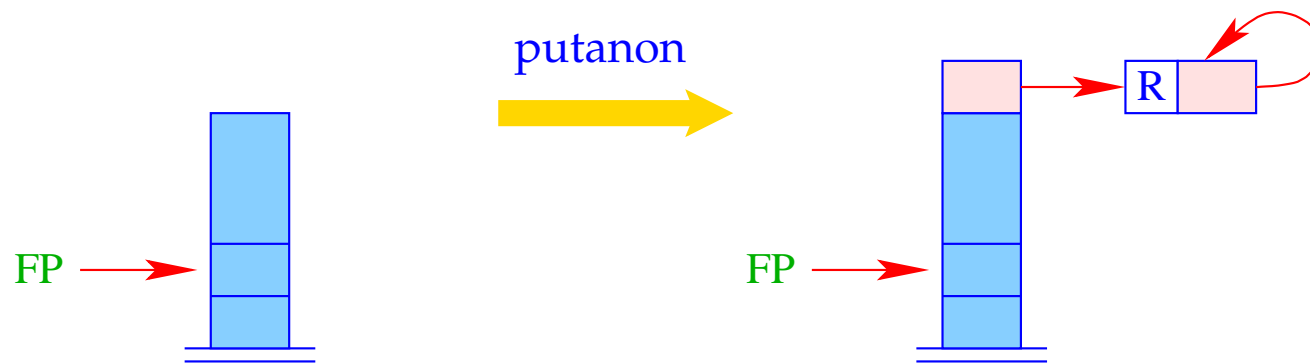
$SP++; S[SP] = \text{new } (A,a);$

The instruction `putvar i` introduces a new unbound variable and additionally initializes the corresponding cell in the stack frame:



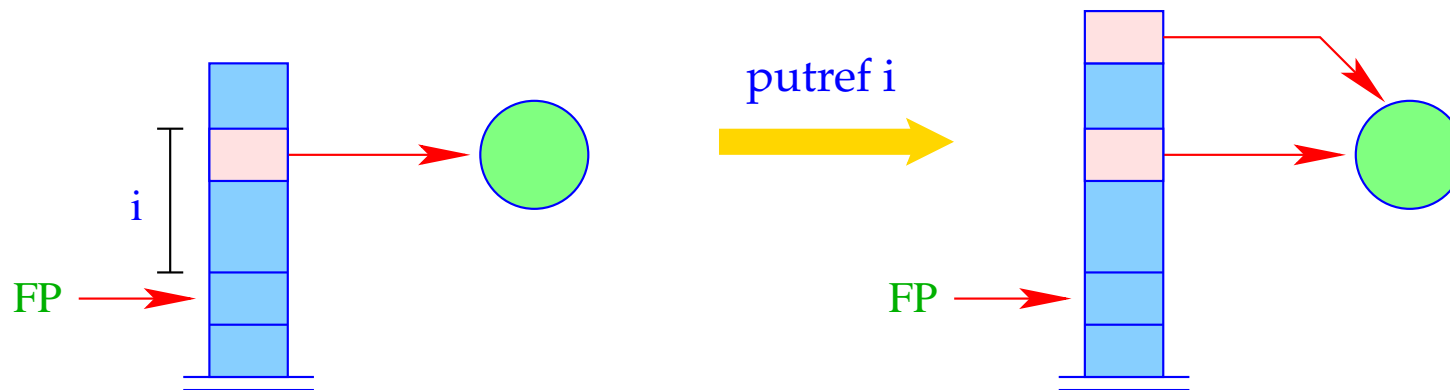
```
SP = SP + 1;  
S[SP] = new (R, HP);  
S[FP + i] = S[SP];
```

The instruction `putanon` introduces a new unbound variable but does not store a reference to it in the stack frame:



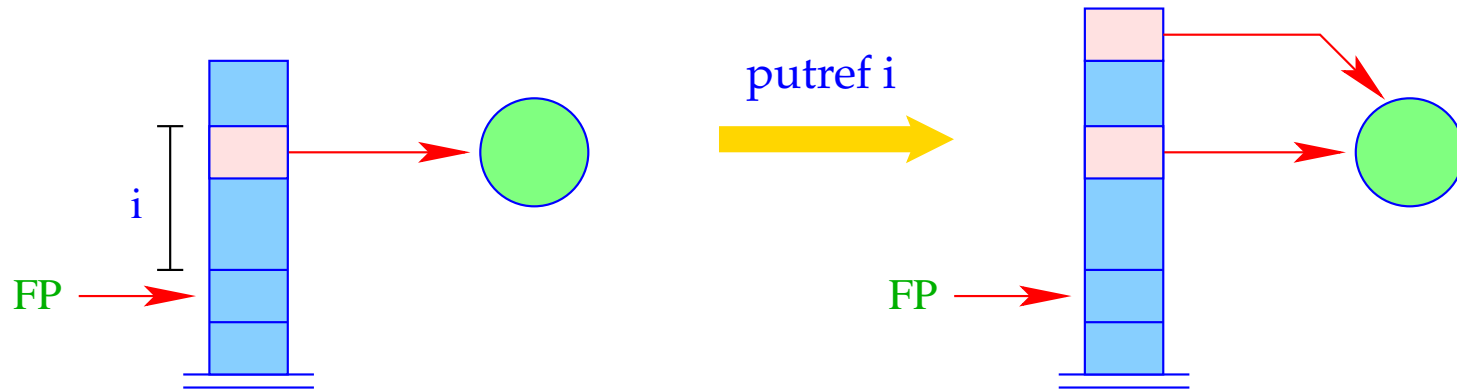
$SP = SP + 1;$
 $S[SP] = \text{new}(R, HP);$

The instruction `putref i` pushes the value of the variable onto the stack:



```
SP = SP + 1;  
S[SP] = deref S[FP + i];
```

The instruction `putref i` pushes the value of the variable onto the stack:

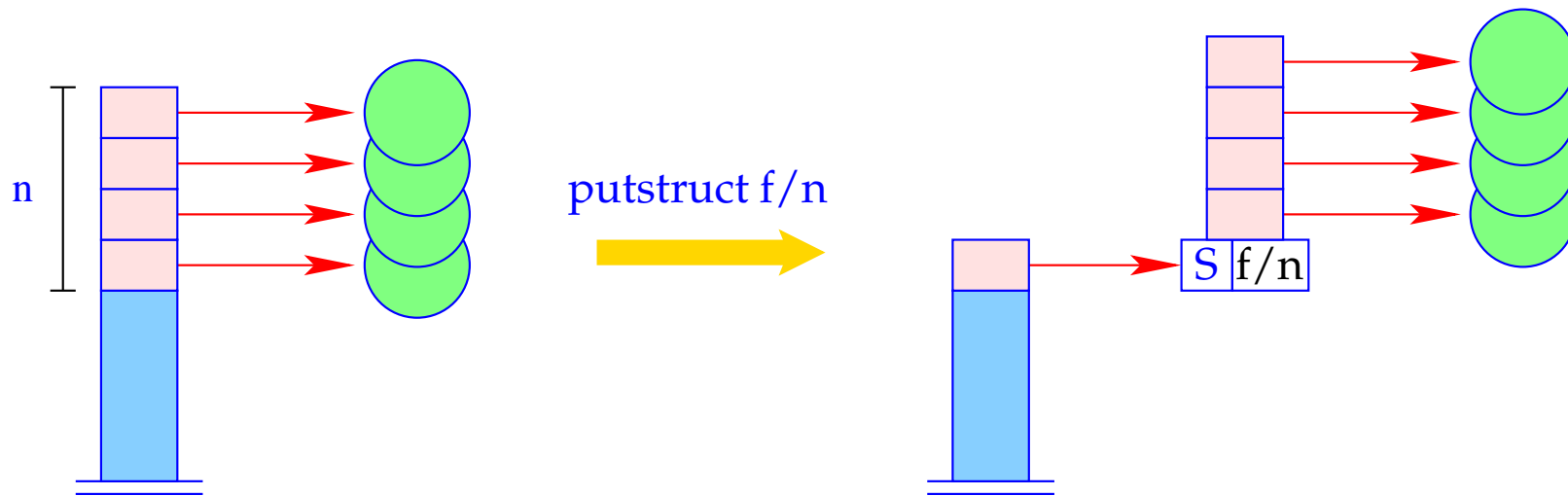


```
SP = SP + 1;  
S[SP] = deref S[FP + i];
```

The auxiliary function `deref` contracts **chains** of references:

```
ref deref (ref v) {  
    if (H[v]==(R,w) && v!=w) return deref (w);  
    else return v;  
}
```

The instruction `putstruct i` builds a constructor application in the heap:



```

v = new (S, f, n);
SP = SP - n + 1;
for (i=1; i<=n; i++)
    H[v + i] = S[SP + i - 1];
S[SP] = v;

```

Remarks:

- The instruction `putref i` does not just push the reference from $S[\text{FP} + i]$ onto the stack, but also dereferences it as much as possible
 \implies maximal contraction of reference chains.
- In constructed terms, references always point to **smaller** heap addresses.
Also otherwise, this will be often the case. Sadly enough, it cannot be **guaranteed** in general :-)

29 The Translation of Literals (Goals)

Idea:

- Literals are treated as **procedure calls**.
- We first allocate a stack frame.
- Then we construct the actual parameters (in the heap)
- ... and store references to these into the stack frame.
- Finally, we jump to the code for the procedure/predicate.

```
codeG p(t1, ..., tk) ρ =   mark B           // allocates the stack frame
                               codeA t1 ρ
                               ...
                               codeA tk ρ
                               call p/k         // calls the procedure p/k
B : ...
```

```

codeG p(t1, ..., tk) ρ =   mark B           // allocates the stack frame
                               codeA t1 ρ
                               ...
                               codeA tk ρ
                               call p/k         // calls the procedure p/k
                               B : ...

```

Example: $p(a, X, g(\bar{X}, Y))$ with $\rho = \{X \mapsto 1, Y \mapsto 2\}$

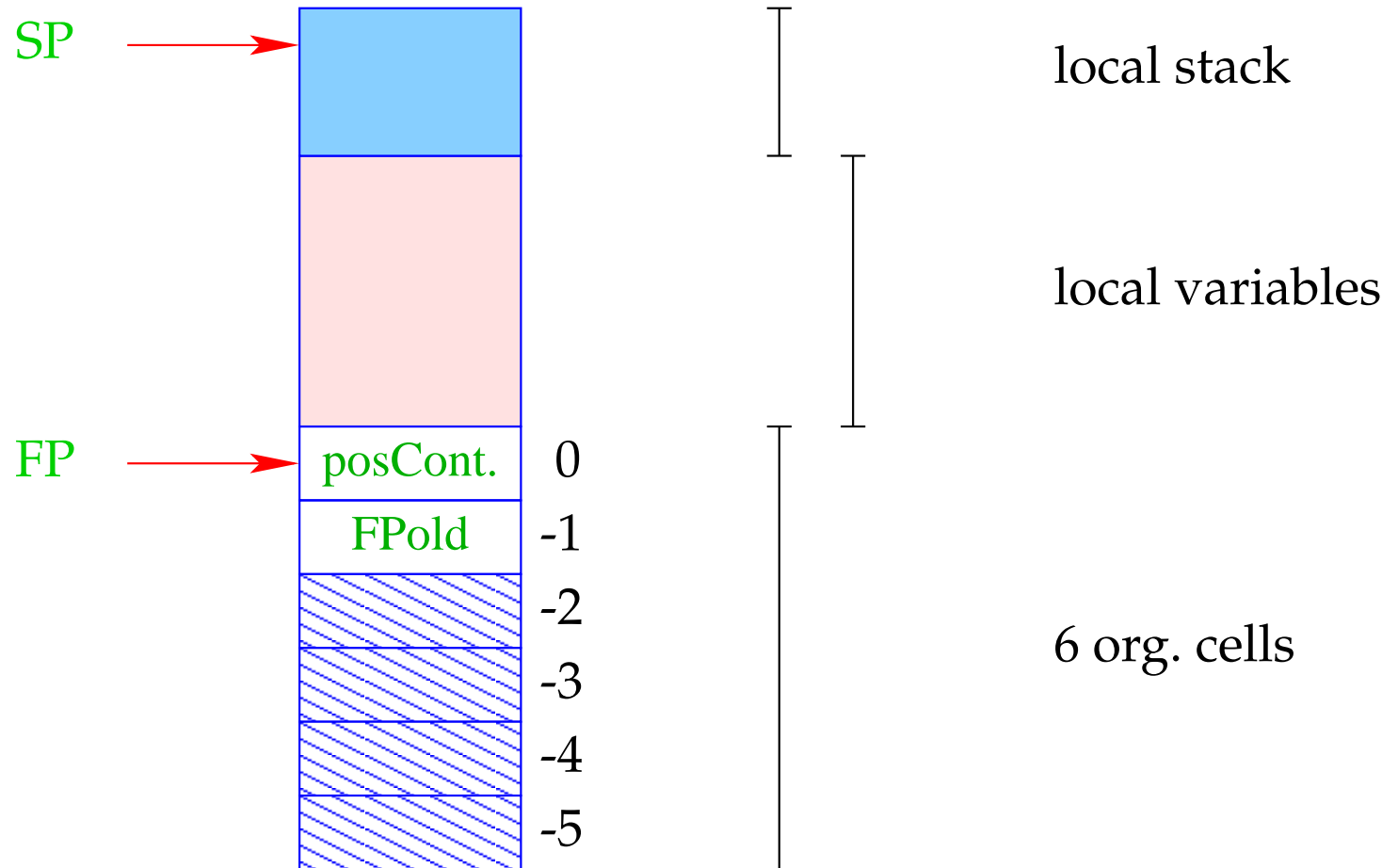
We obtain:

```

mark B           putref 1           call p/3
putatom a       putvar 2           B: ...
putvar 1        putstruct g/2

```

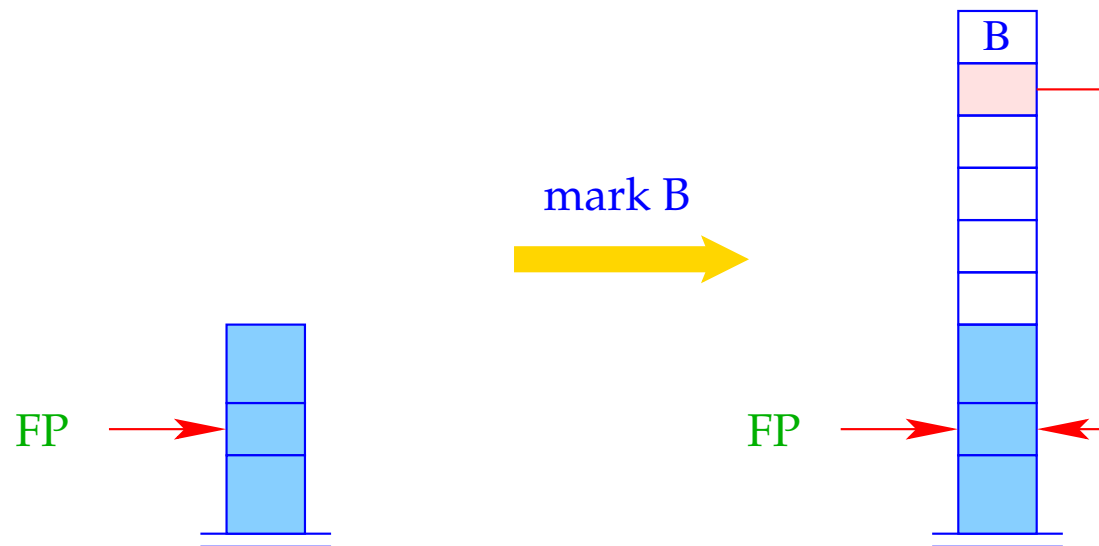
Stack Frame of the WiM:



Remarks:

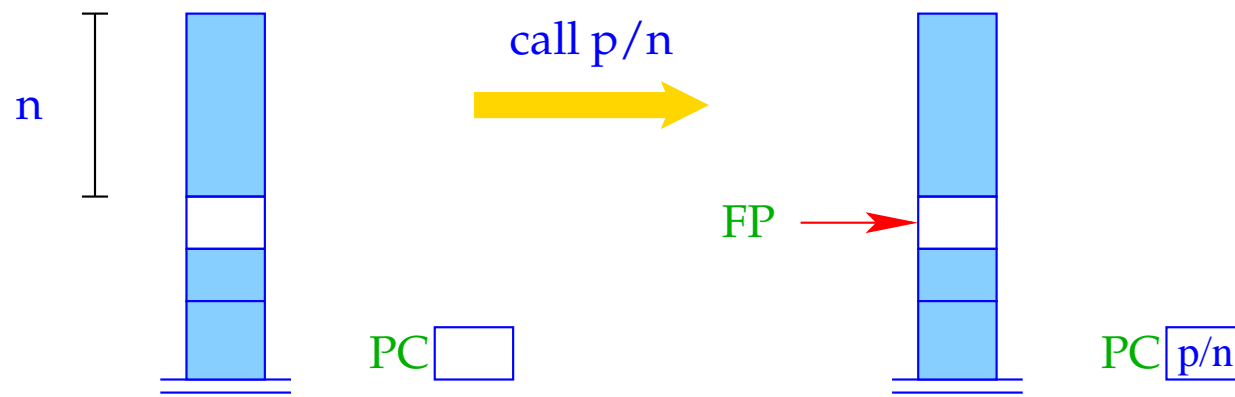
- The **positive** continuation address records where to continue after successful treatment of the goal.
- Additional organizational cells are needed for the implementation of **backtracking**
 \implies will be discussed at the translation of predicates.

The instruction `mark B` allocates a new stack frame:



```
SP = SP + 6;  
S[SP] = B; S[SP-1] = FP;
```

The instruction `call p/n` calls the n -ary predicate p :



$$\begin{aligned} \text{FP} &= \text{SP} - n; \\ \text{PC} &= p/n; \end{aligned}$$

30 Unification

Convention:

- By \tilde{X} , we denote an occurrence of X ;
it will be translated differently depending on whether the variable is initialized or not.
- We introduce the macro `put \tilde{X} ρ` :

`put X ρ` = `putvar (ρ X)`

`put $_$ ρ` = `putanon`

`put \bar{X} ρ` = `putref (ρ X)`

Let us translate the unification $\tilde{X} = t$.

Idea 1:

- Push a reference to (the binding of) X onto the stack;
- Construct the term t in the heap;
- Invent a new instruction implementing the unification :-)

Let us translate the unification $\tilde{X} = t$.

Idea 1:

- Push a reference to (the binding of) X onto the stack;
- Construct the term t in the heap;
- Invent a new instruction implementing the unification :-)

$$\begin{aligned} \text{code}_G (\tilde{X} = t) \rho &= \text{put } \tilde{X} \rho \\ &\quad \text{code}_A t \rho \\ &\quad \text{unify} \end{aligned}$$

Example:

Consider the equation:

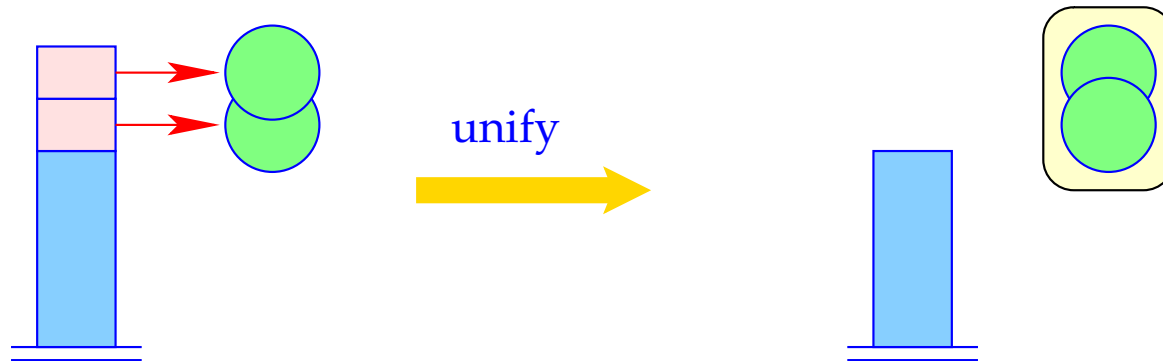
$$\bar{U} = f(g(\bar{X}, Y), a, Z)$$

Then we obtain for an address environment

$$\rho = \{X \mapsto 1, Y \mapsto 2, Z \mapsto 3, U \mapsto 4\}$$

putref 4	putref 1	putatom a	unify
	putvar 2	putvar 3	
	putstruct g/2	putstruct f/3	

The instruction `unify` calls the `run-time` function `unify()` for the topmost two references:



```
unify (S[SP-1], S[SP]);  
SP = SP-2;
```

The Function `unify()`

- ... takes two heap addresses.
For each call, we guarantee that these are **maximally de-referenced**.
- ... checks whether the two addresses are already **identical**.
If so, does nothing **:-)**
- ... binds **younger variables** (larger addresses) to **older variables** (smaller addresses);
- ... when binding a variable to a term, checks whether the variable occurs inside the term \implies **occur-check**;
- ... **records** newly created bindings;
- ... may **fail**. Then **backtracking** is initiated.

```

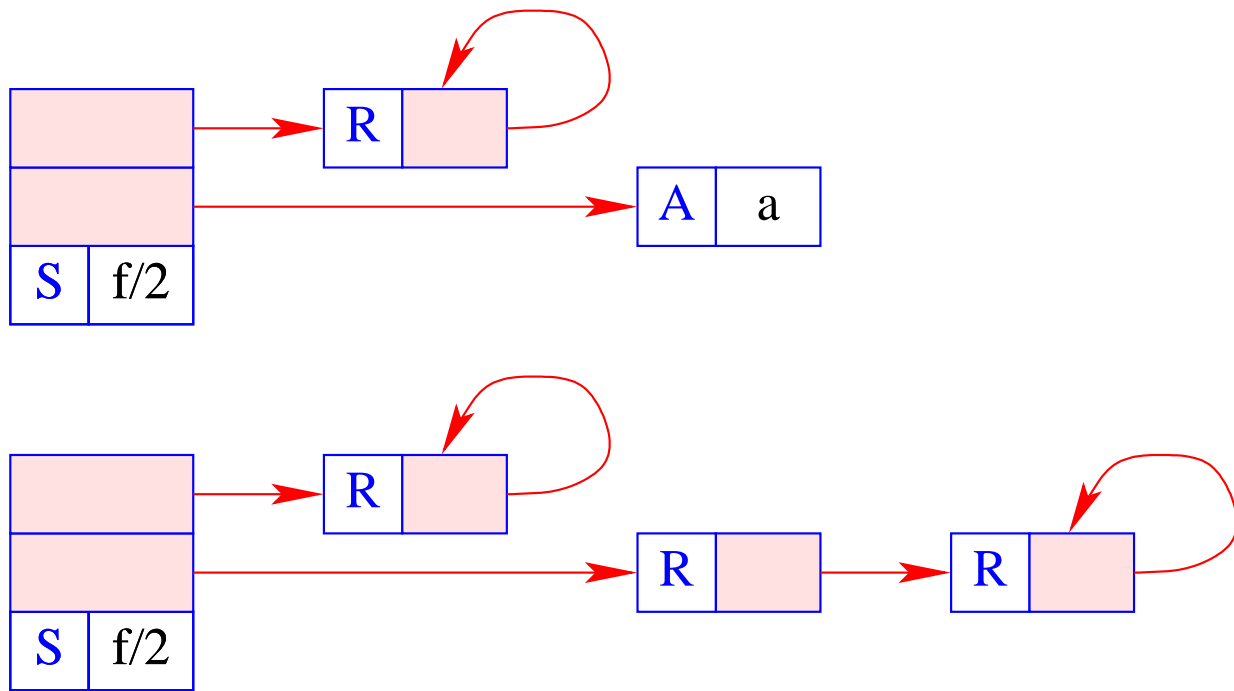
bool unify (ref u, ref v) {
    if (u == v) return true;
    if (H[u] == (R,_)) {
        if (H[v] == (R,_)) {
            if (u>v) {
                H[u] = (R,v); trail (u); return true;
            } else {
                H[v] = (R,u); trail (v); return true;
            }
        }
        } elseif (check (u,v)) {
            H[u] = (R,v); trail (u); return true;
        } else {
            backtrack(); return false;
        }
    }
}
...

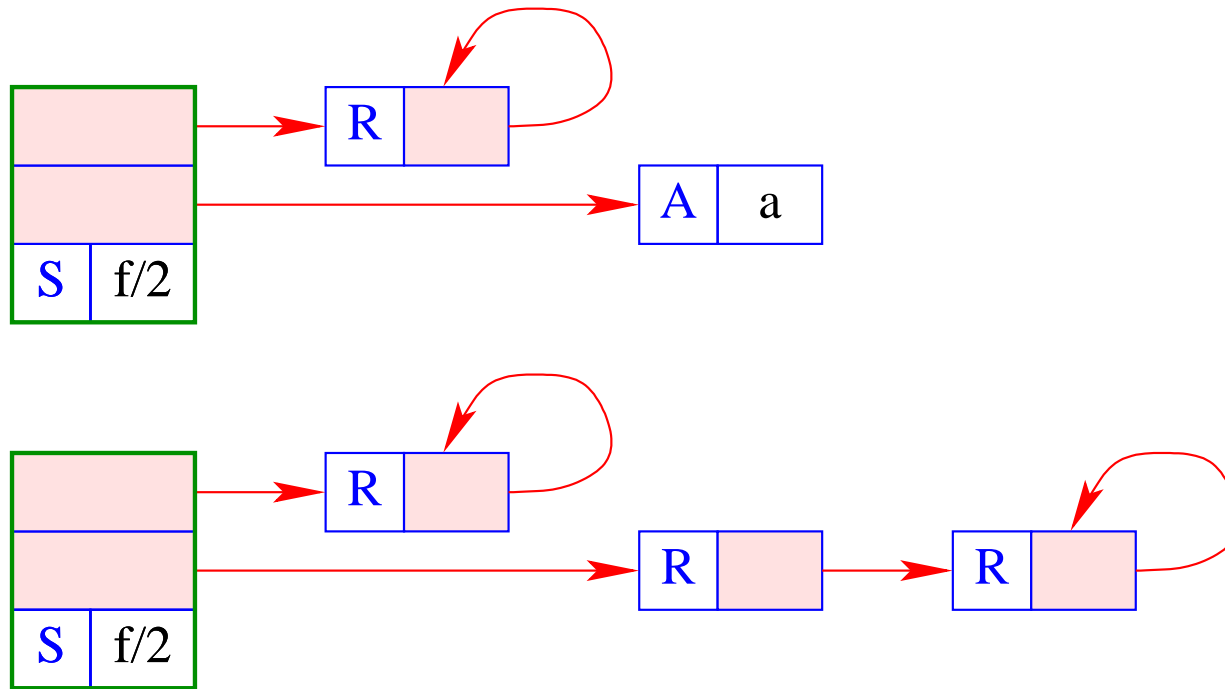
```

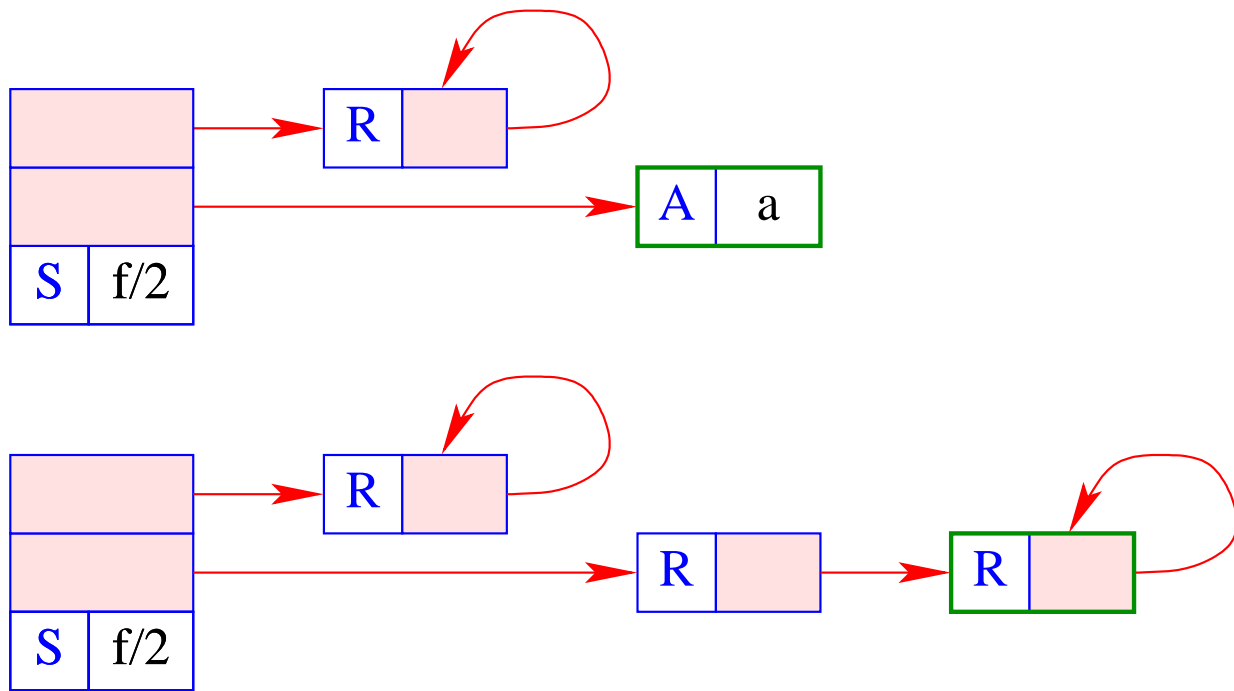
```

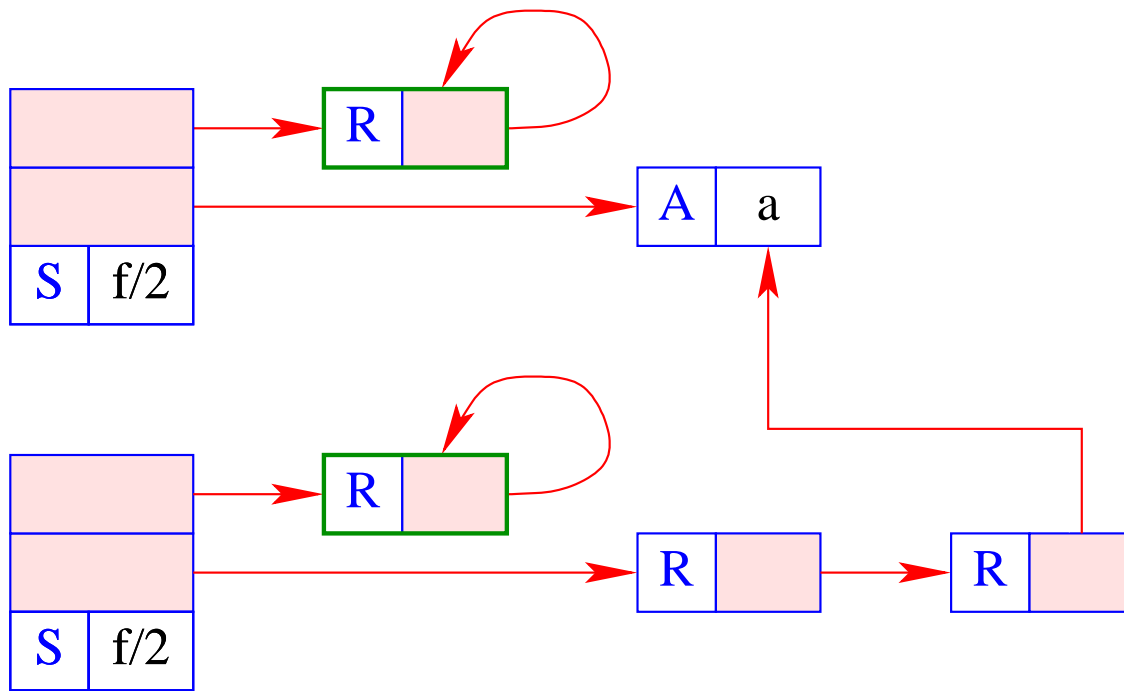
...
if ((H[v] == (R,_)) {
    if (check (v,u)) {
        H[v] = (R,u); trail (v); return true;
    } else {
        backtrack(); return false;
    }
}
if (H[u]==(A,a) && H[v]==(A,a))
    return true;
if (H[u]==(S, f/n) && H[v]==(S, f/n)) {
    for (int i=1; i<=n; i++)
        if(!unify (deref (H[u+i]), deref (H[v+i])) return false;
    return true;
}
backtrack(); return false;
}

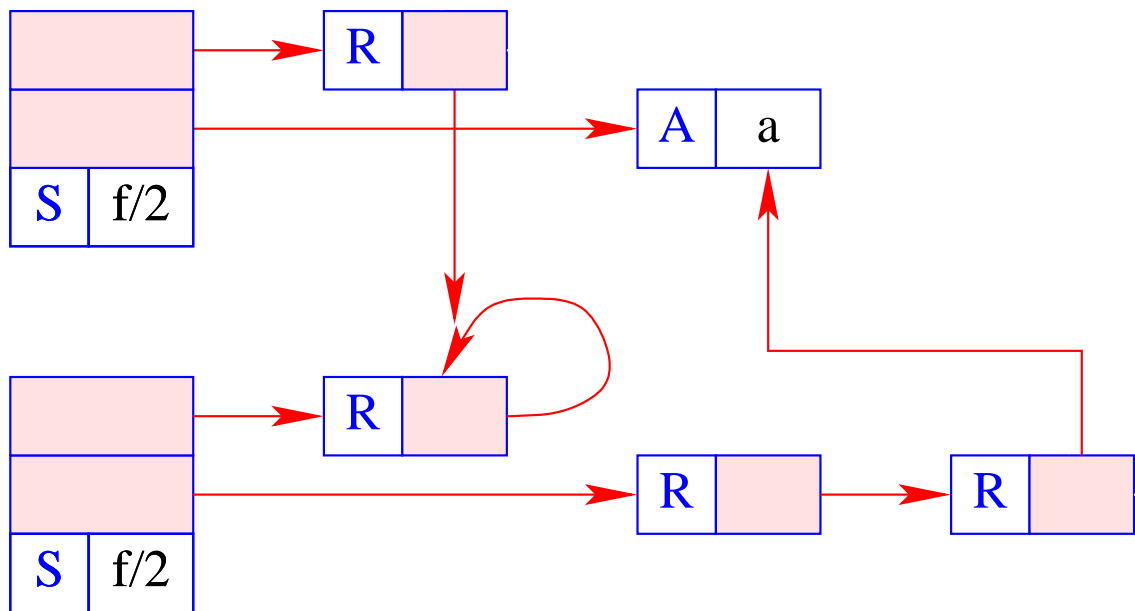
```











- The run-time function `trail()` **records** the a potential new binding.
- The run-time function `backtrack()` initiates **backtracking**.
- The auxiliary function `check()` performs the **occur-check**: it tests whether a variable (the first argument) **occurs inside** a term (the second argument).
- Often, this check is skipped, i.e.,

```
bool check (ref u, ref v) { return true; }
```

Otherwise, we could implement the run-time function `check()` as follows:

```
bool check (ref u, ref v) {
    if (u == v) return false;
    if (H[v] == (S, f/n)) {
        for (int i=1; i<=n; i++)
            if (!check(u, deref (H[v+i])))
                return false;
    }
    return true;
}
```

Discussion:

- The translation of an equation $\tilde{X} = t$ is very simple :-)
- Often the constructed cells immediately become garbage :-)

Idea 2:

- Push a reference to the run-time binding of the left-hand side onto the stack.
- Avoid to construct sub-terms of t whenever possible !
- Translate each node of t into an instruction which performs the unification with this node !!

Discussion:

- The translation of an equation $\tilde{X} = t$ is very simple :-)
- Often the constructed cells immediately become **garbage** :-)

Idea 2:

- Push a reference to the run-time binding of the left-hand side onto the stack.
- Avoid to construct sub-terms of t whenever possible !
- Translate each node of t into an instruction which performs the unification with this node !!

$$\text{code}_G (\tilde{X} = t) \rho = \text{put } \tilde{X} \rho \\ \text{code}_U t \rho$$

Let us first consider the unification code for atoms and variables only:

$\text{code}_U a \rho = \text{uatom } a$

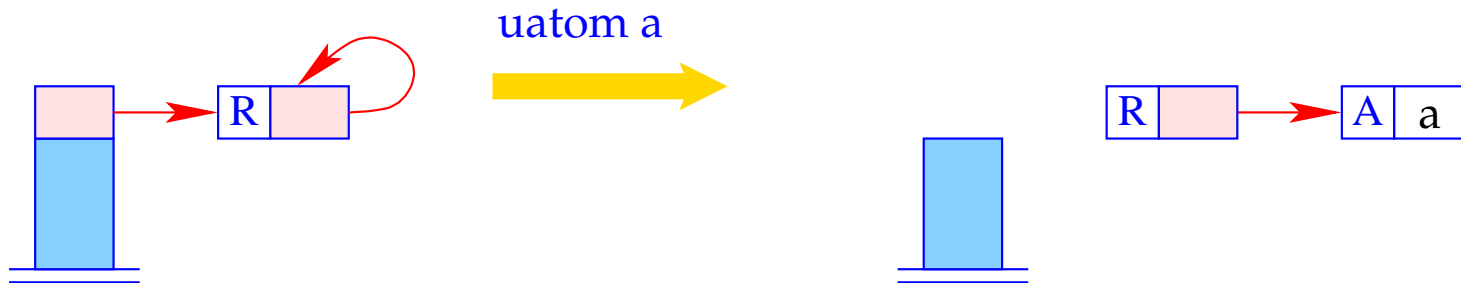
$\text{code}_U X \rho = \text{uvar } (\rho X)$

$\text{code}_U _ \rho = \text{pop}$

$\text{code}_U \bar{X} \rho = \text{uref } (\rho X)$

... // to be continued :-)

The instruction `uatom a` implements the unification with the atom `a`:



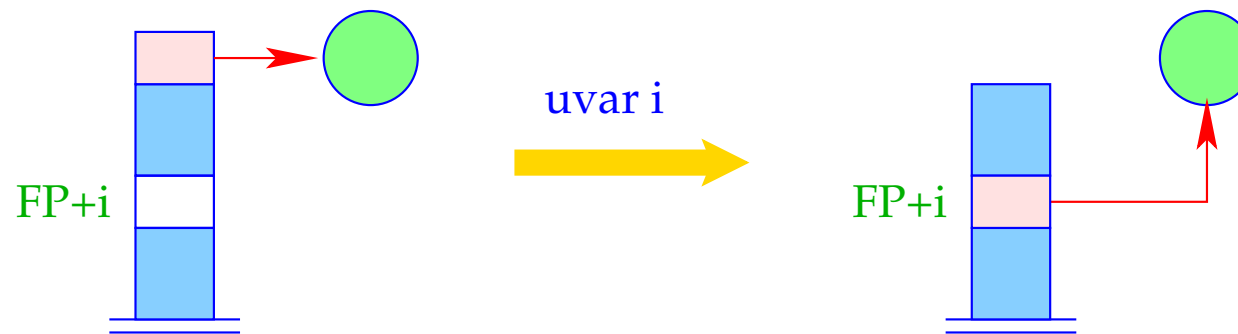
```

v = S[SP]; SP--;
switch (H[v]) {
case (A, a):    break;
case (R, _):    H[v] = (R, new (A, a));
                trail (v); break;
default:       backtrack();
}

```

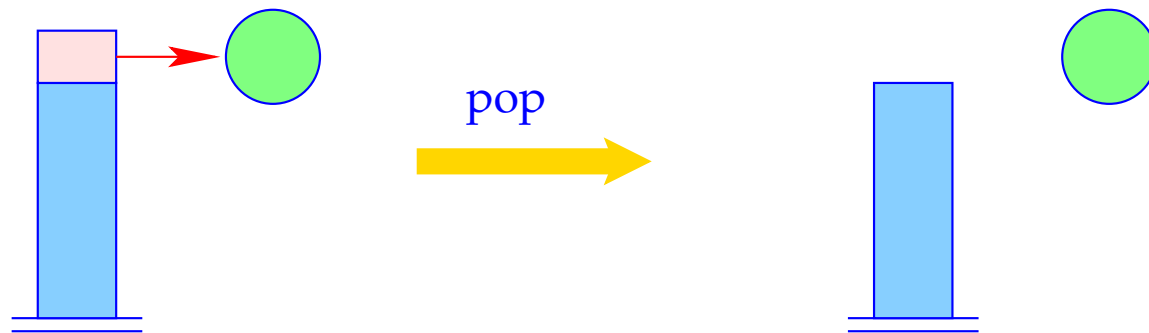
- The run-time function `trail()` records the a potential new binding.
- The run-time function `backtrack()` initiates backtracking.

The instruction `uvar i` implements the unification with an un-initialized variable:



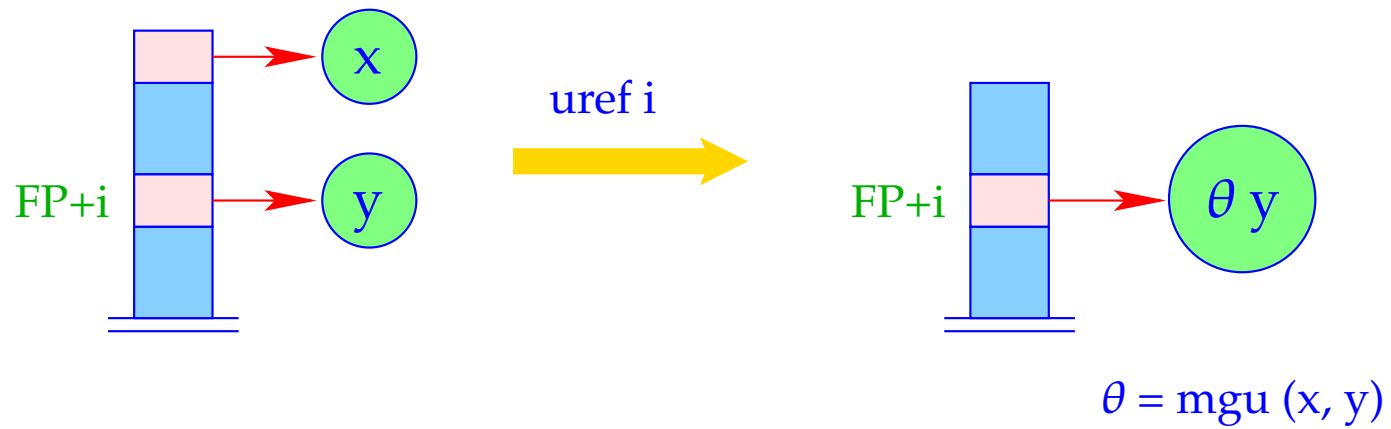
$$S[FP+i] = S[SP]; SP--;$$

The instruction `pop` implements the unification with an anonymous variable. It always succeeds :-)



`SP--;`

The instruction `uref i` implements the unification with an initialized variable:



```
unify (S[SP], deref (S[FP+i]));
SP--;
```

It is only here that the run-time function `unify()` is called :-)

- The unification code performs a **pre-order** traversal over t .
- In case, execution hits at an unbound variable, we **switch** from checking to building :-)

```

codeU f(t1, ..., tn) ρ =    ustruct f/n A                // test
                               son 1
                               codeU t1 ρ
                               ...
                               son n
                               codeU tn ρ
                               up B
A : check ivars(f(t1, ..., tn)) ρ    // occur-check
    codeA f(t1, ..., tn) ρ        // building !!
    bind                                // creation of bindings
B : ...

```

The Building Block:

Before constructing the new (sub-) term t' for the binding, we must exclude that it contains the variable X' on top of the stack !!!

This is the case iff **the binding** of no variable inside t' contains (a reference to) X' .

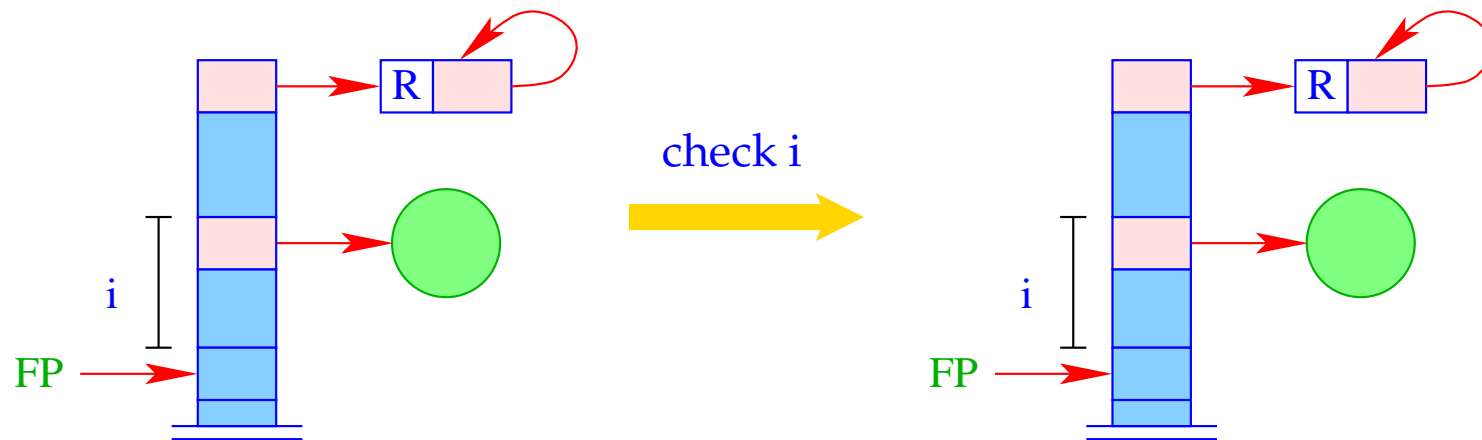
\implies $ivars(t')$ returns the set of **already initialized** variables of t .

\implies The macro **check** $\{Y_1, \dots, Y_d\} \rho$ generates the necessary tests on the variables Y_1, \dots, Y_d :

$$\begin{aligned} \text{check } \{Y_1, \dots, Y_d\} \rho &= \text{check } (\rho Y_1) \\ &\quad \text{check } (\rho Y_2) \\ &\quad \dots \\ &\quad \text{check } (\rho Y_d) \end{aligned}$$

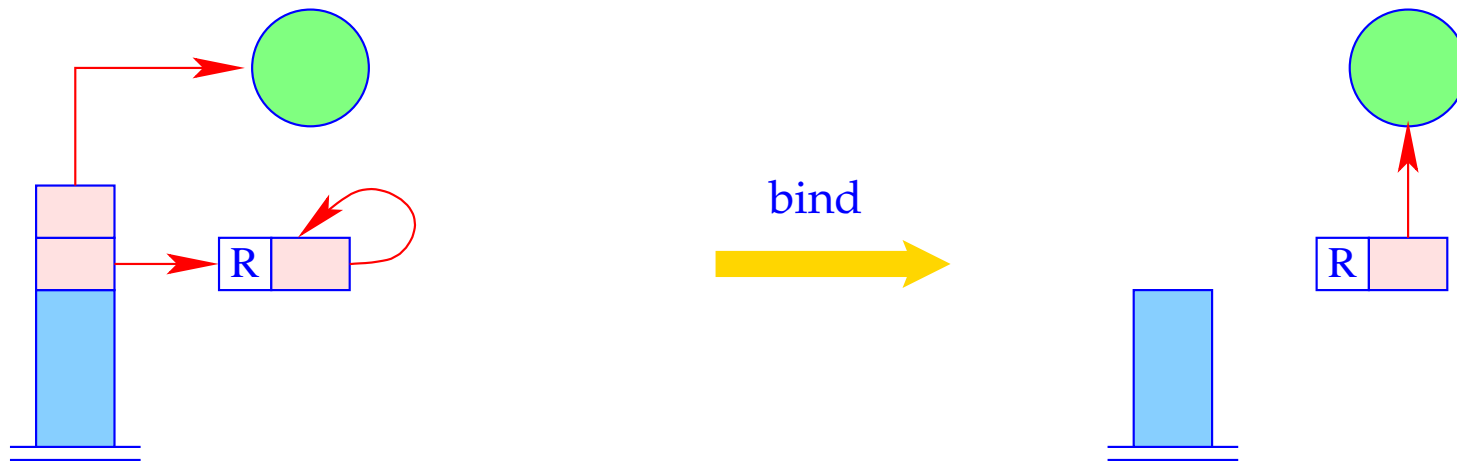
The instruction `check i` checks whether the (unbound) variable on top of the stack occurs inside the term bound to variable `i`.

If so, unification fails and **backtracking** is caused:



```
if (!check (S[SP], deref S[FP+i]))  
    backtrack();
```


The instruction `bind` terminates the building block. It binds the (unbound) variable to the constructed term:



```
H[S[SP-1]] = (R, S[SP]);  
trail (S[SP-1]);  
SP = SP - 2;
```

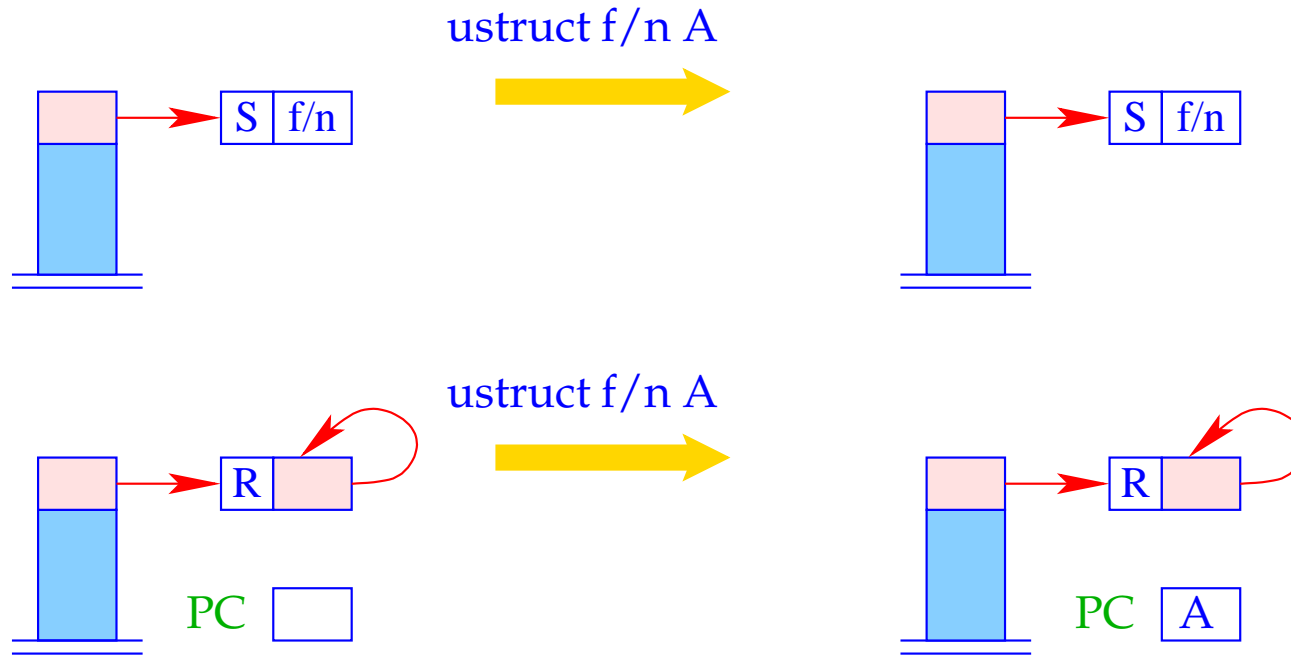
The Pre-Order Traversal:

- First, we **test** whether the topmost reference is an unbound variable. If so, we jump to the building block.
- Then we compare the root node with the constructor **f/n**.
- Then we **recursively descend** to the children.
- Then we **pop** the stack and proceed behind the unification code:

Once again the unification code for constructed terms:

```
codeU f(t1, ..., tn) ρ =    ustruct f/n A           // test
                               son 1           // recursive descent
                               codeU t1 ρ
                               ...
                               son n           // recursive descent
                               codeU tn ρ
                               up B            // ascent to father
A : check ivars(f(t1, ..., tn)) ρ
    codeA f(t1, ..., tn) ρ
    bind
B : ...
```

The instruction `ustruct i` implements the test of the root node of a structure:



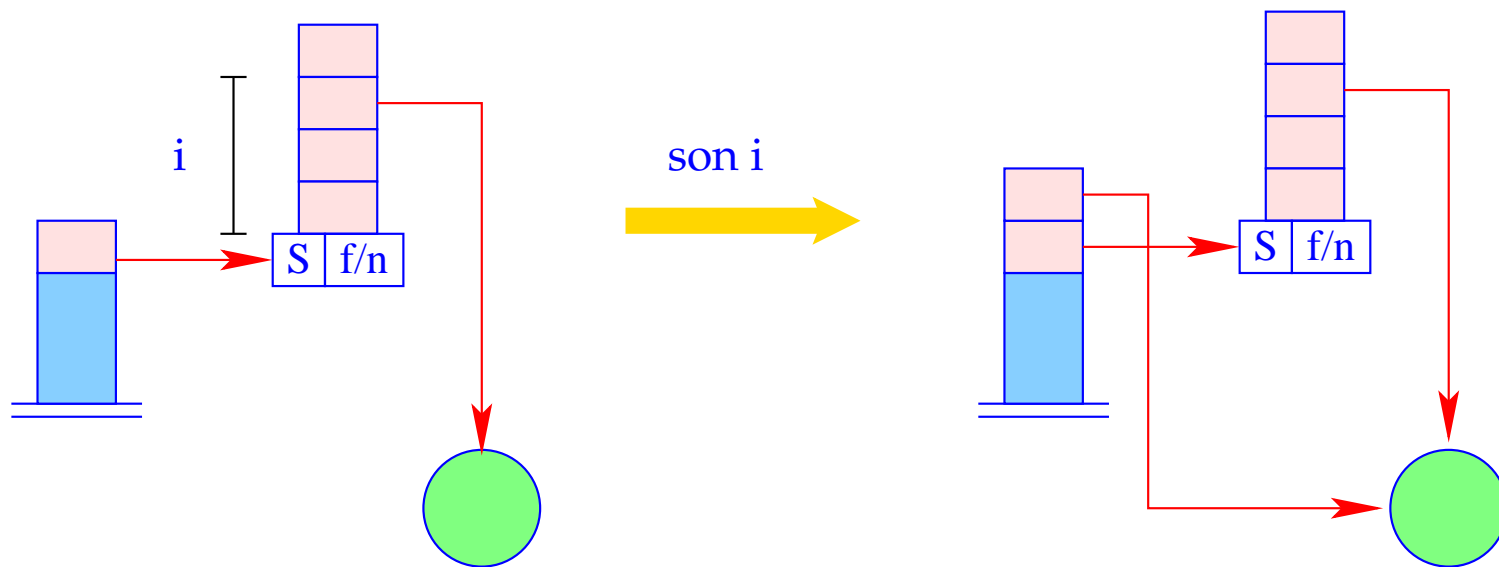
```

switch (H[S[SP]]) {
  case (S, f/n):  break;
  case (R, _):   PC = A; break;
  default:      backtrack();
}

```

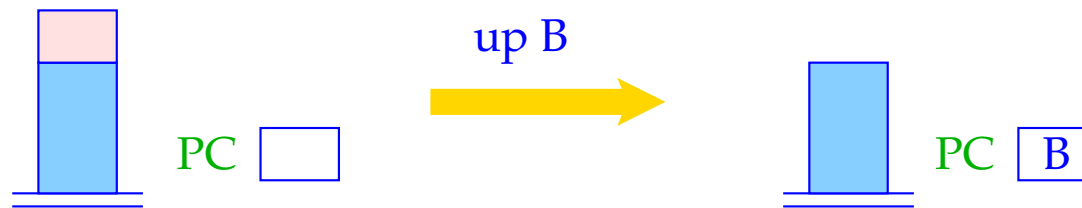
... the argument reference is **not yet** popped :-)

The instruction `son i` pushes the (reference to the) i -th sub-term from the structure pointed at from the topmost reference:



$S[SP+1] = \text{deref}(H[S[SP]+i]); SP++;$

It is the instruction `up B` which finally pops the reference to the structure:



`SP--; PC = B;`

The continuation address `B` is the next address after the `build`-section.

Example:

For our example term $f(g(\bar{X}, Y), a, Z)$ and
 $\rho = \{X \mapsto 1, Y \mapsto 2, Z \mapsto 3\}$ we obtain:

ustruct f/3 A_1	up B_2	B_2 :	son 2	putvar 2	
son 1			uatom a	putstruct g/2	
ustruct g/2 A_2	A_2 :	check 1	son 3	putatom a	
son 1		putref 1	uvar 3	putvar 3	
uref 1		putvar 2	up B_1	putstruct f/3	
son 2		putstruct g/2	A_1 :	check 1	bind
uvar 2		bind	putref 1	B_1 :	...

Code size can grow quite considerably — for **deep** terms. In practice, though, deep terms are “rare” :-)

31 Clauses

Clausal code must

- **allocate** stack space for locals;
- **evaluate** the body;
- **free** the stack frame (whenever possible :-)

Let r denote the clause: $p(X_1, \dots, X_k) \leftarrow g_1, \dots, g_n.$

Let $\{X_1, \dots, X_m\}$ denote the set of locals of r and ρ the address environment:

$$\rho X_i = i$$

Remark: The first k locals are always the **formals** :-)

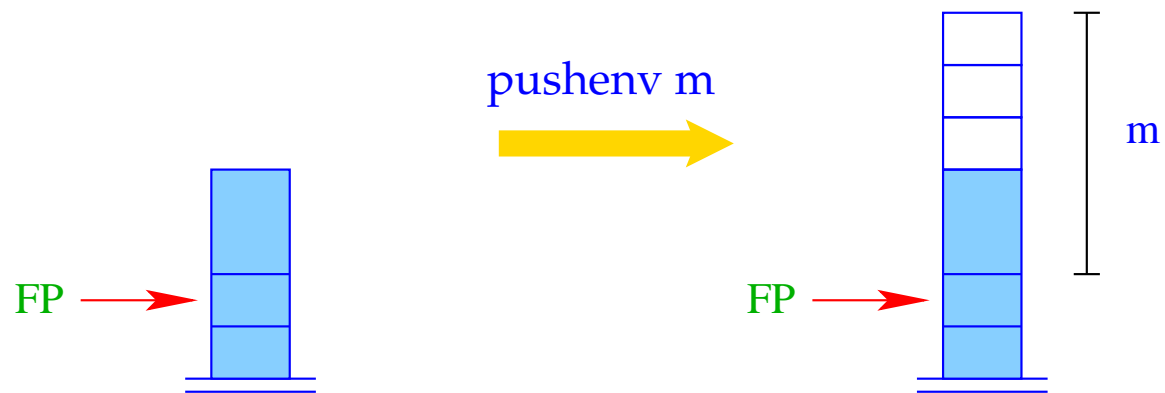
Then we translate:

```
codeC r = pushenv m           // allocates space for locals
          codeG g1 ρ
          ...
          codeG gn ρ
          popenv
```

The instruction `popenv` restores `FP` and `PC` and `tries to pop` the current stack frame.

It should succeed whenever program execution will never return to this stack frame :-)

The instruction `pushenv m` sets the stack pointer:



$$SP = FP + m;$$

Example:

Consider the clause r :

$$a(X, Y) \leftarrow f(\bar{X}, X_1), a(\bar{X}_1, \bar{Y})$$

Then `codeC r` yields:

pushenv 3

mark A

A: mark B

B: popenv

putref 1

putref 3

putvar 3

putref 2

call f/2

call a/2

32 Predicates

A predicate q/k is defined through a sequence of clauses $rr \equiv r_1 \dots r_f$.

The translation of q/k provides the translations of the individual clauses r_i .

In particular, we have for $f = 1$:

$$\text{code}_P rr = \text{code}_C r_1$$

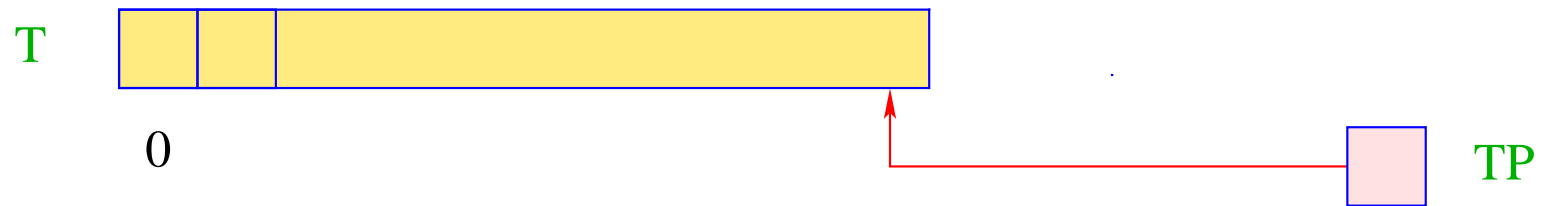
If q/k is defined through several clauses, the first alternative must be tried.

On failure, the next alternative must be tried

\implies backtracking :-)

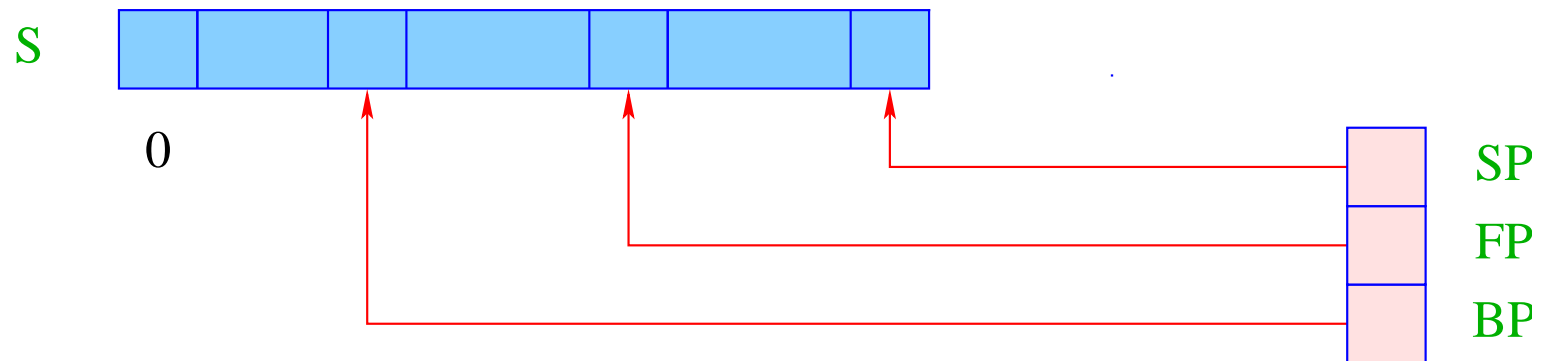
32.1 Backtracking

- Whenever unification fails, we call the run-time function `backtrack()`.
- The goal is to **roll back** the whole computation to the (**dynamically :-**) latest goal where another clause can be chosen \implies the last **backtrack point**.
- In order to undo intermediate variable bindings, we always have recorded new bindings with the run-time function `trail()`.
- The run-time function `trail()` stores variables in the data-structure **trail**:



TP == Trail Pointer
points to the topmost occupied Trail cell

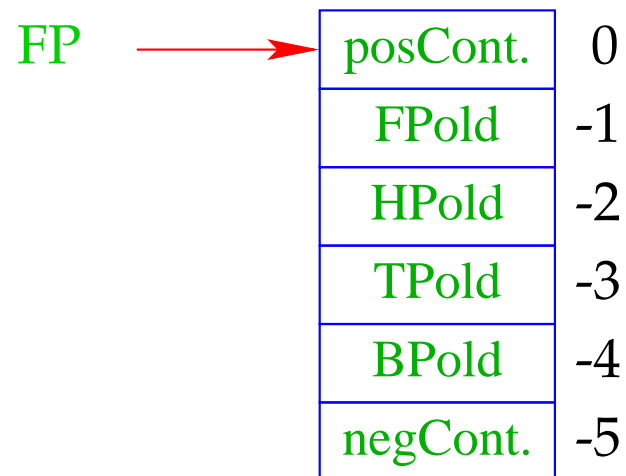
The current stack frame where backtracking should return to is pointed at by the extra register **BP**:



A **backtrack point** is stack frame to which program execution possibly returns.

- We need the code address for trying the **next** alternative (**negative continuation address**);
- We save the old values of the registers **HP**, **TP** and **BP**.
- **Note:** The **new BP** will receive the value of the current **FP** :-)

For this purpose, we use the corresponding four organizational cells:



For more comprehensible notation, we thus introduce the macros:

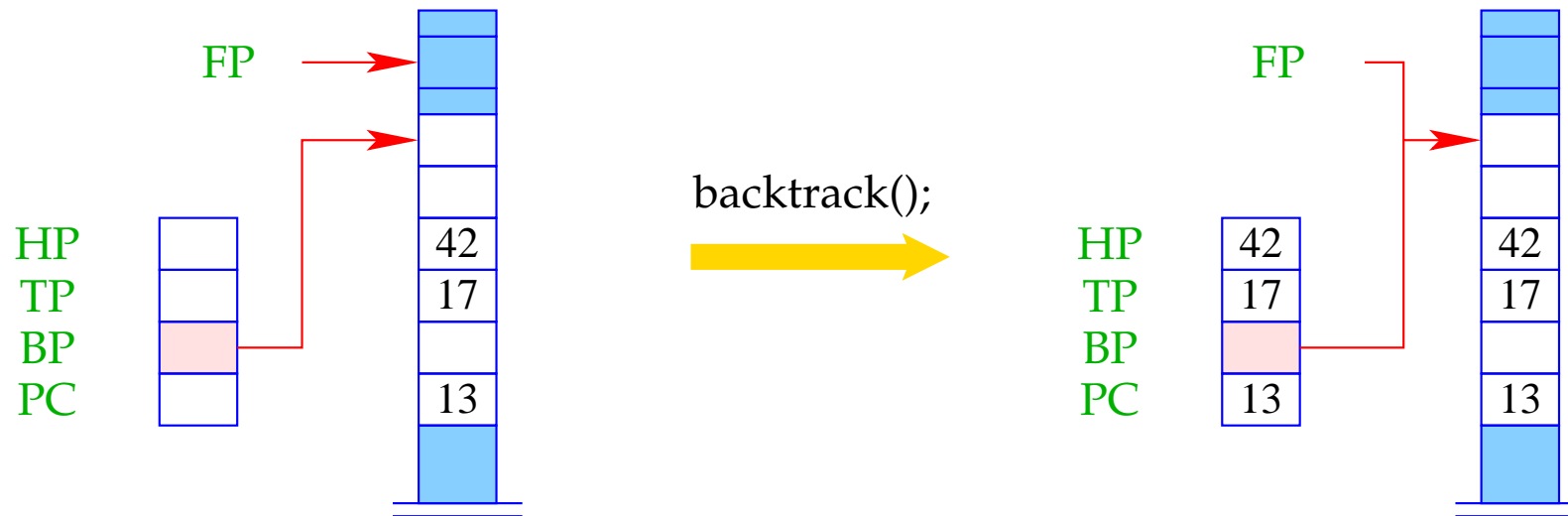
$$\begin{aligned}\text{posCont} &\equiv S[\text{FP}] \\ \text{FPold} &\equiv S[\text{FP} - 1] \\ \text{HPold} &\equiv S[\text{FP} - 2] \\ \text{TPold} &\equiv S[\text{FP} - 3] \\ \text{BPold} &\equiv S[\text{FP} - 4] \\ \text{negCont} &\equiv S[\text{FP} - 5]\end{aligned}$$

for the corresponding addresses.

Remark:

Occurrence on the **left** \equiv saving the register
Occurrence on the **right** \equiv restoring the register

Calling the run-time function `void backtrack()` yields:



```
void backtrack() {  
    FP = BP; HP = HPold;  
    reset (TPold, TP);  
    TP = TPold; PC = negCont;  
}
```

where the run-time function `reset()` undoes the bindings of variables established **since** the backtrack point.

32.2 Resetting Variables

Idea:

- The variables which have been created since the last backtrack point can be removed together with their bindings by popping the heap !!! :-)
- This works fine if **younger** variables always point to **older** objects.
- Bindings of **old** variables to younger objects, though, must be reset **manually** :-(
- These are therefore recorded in the trail.

Functions `void trail(ref u)` and `void reset (ref y, ref x)` can thus be implemented as:

```
void trail (ref u) {
    if (u < S[BP-2]) {
        TP = TP+1;
        T[TP] = u;
    }
}

void reset (ref x, ref y) {
    for (ref u=y; x<u; u--)
        H[T[u]] = (R,T[u]);
}
```

Here, `S[BP-2]` represents the heap pointer when creating the last backtrack point.

32.3 Wrapping it Up

Assume that the predicate q/k is defined by the clauses r_1, \dots, r_f ($f > 1$).

We provide code for:

- **setting** up the backtrack point;
- successively **trying** the alternatives;
- **deleting** the backtrack point.

This means:

```

codeP rr = q/k : setbtp
                try A1
                ...
                try Af-1
                delbtp
                jump Af
A1 : codeC r1
    ...
Af : codeC rf

```

Note:

- We delete the backtrack point **before** the last alternative **:-)**
- We **jump** to the last alternative — never to return to the present frame **:-))**

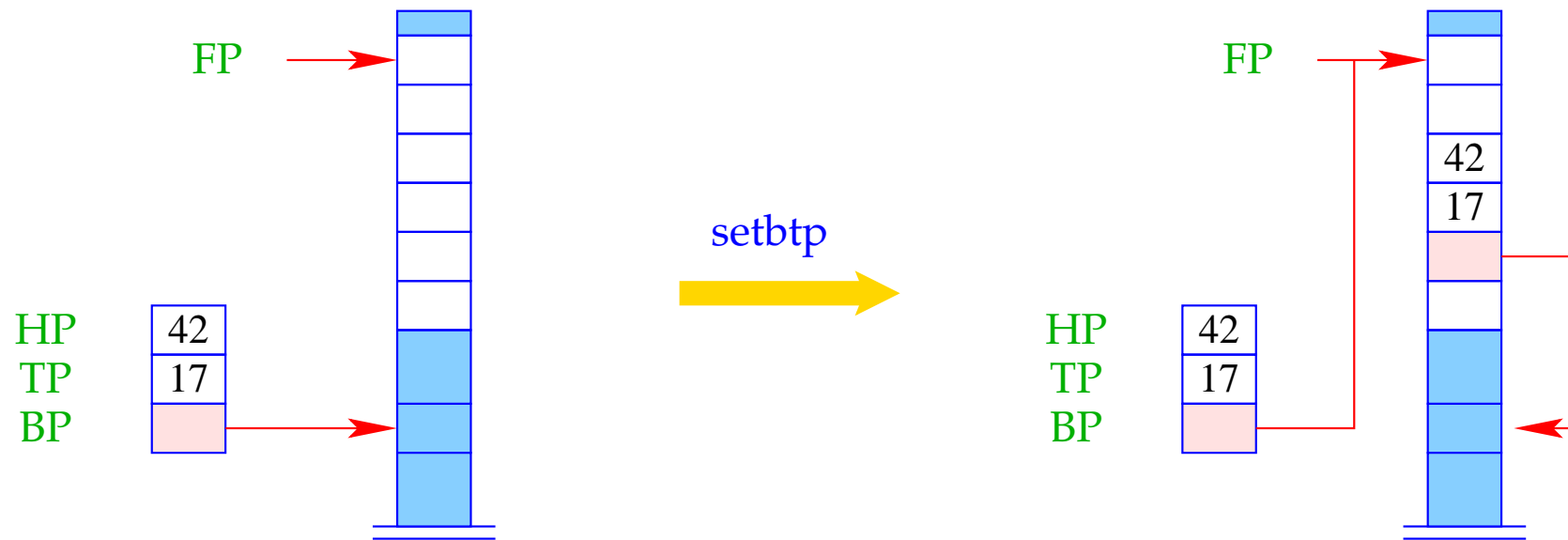
Example:

$$\begin{aligned} s(X) &\leftarrow t(\bar{X}) \\ s(X) &\leftarrow \bar{X} = a \end{aligned}$$

The translation of the predicate s yields:

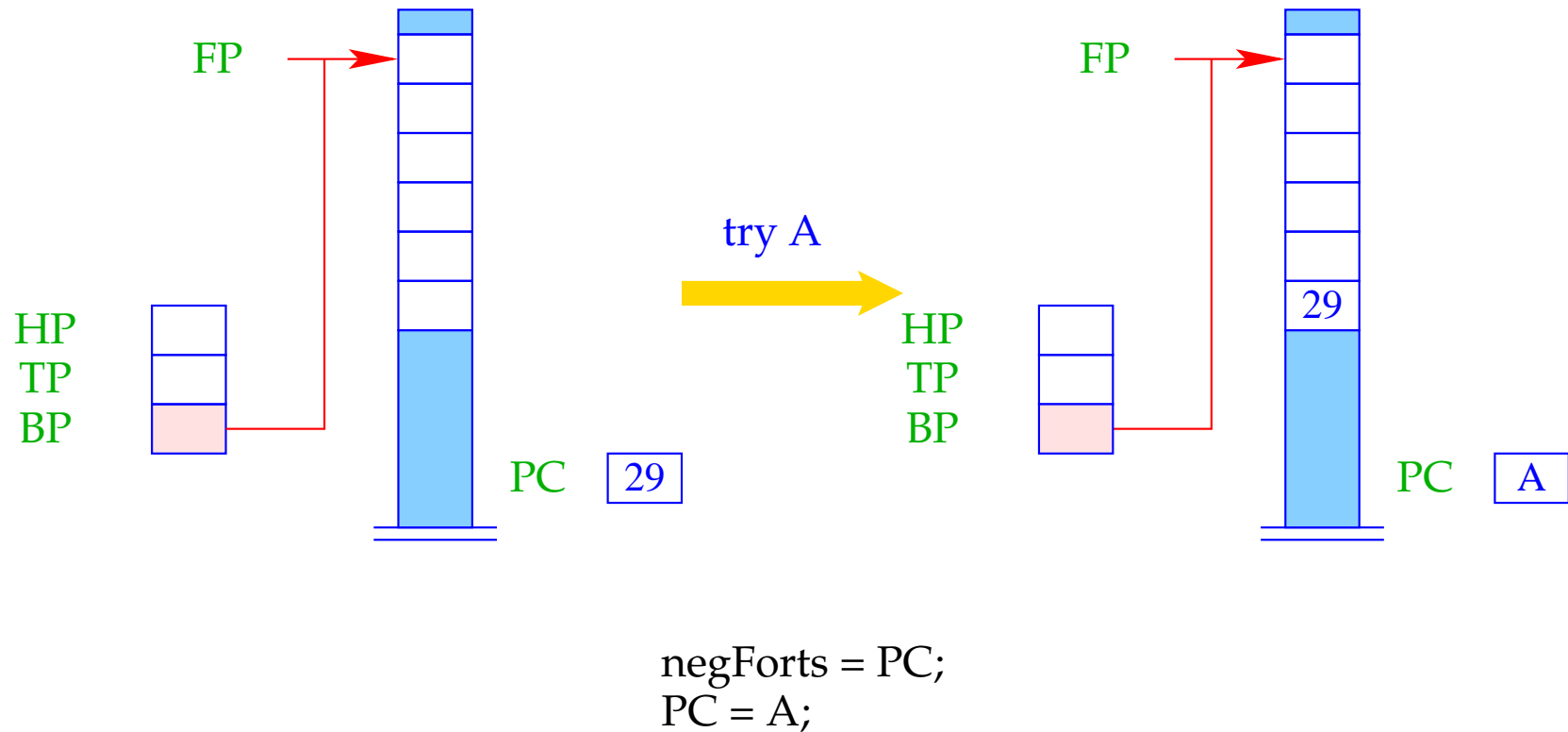
$s/1$:	setbtp	A:	pushenv 1	B:	pushenv 1
	try A		mark C		putref 1
	delbtp		putref 1		uatom a
	jump B		call t/1		popenv
		C:	popenv		

The instruction `setbtp` saves the registers `HP`, `TP`, `BP`:

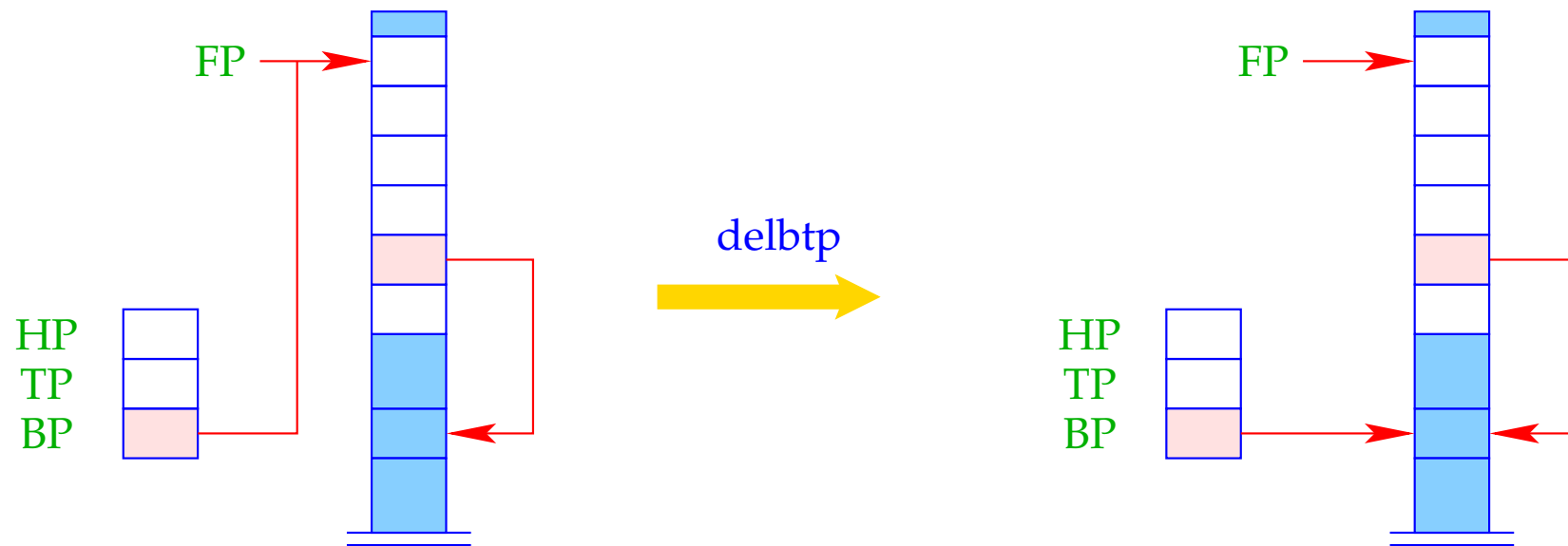


HPold = HP;
TPold = TP;
BPold = BP;
BP = FP;

The instruction `try A` tries the alternative at address `A` and updates the negative continuation address to the current `PC`:



The instruction `delbtp` restores the old backtrack pointer:



$BP = BP_{old};$

32.4 Popping of Stack Frames

Recall the translation scheme for clauses:

$$\begin{aligned} \text{code}_C r &= \text{pushenv } m \\ &\quad \text{code}_G g_1 \rho \\ &\quad \dots \\ &\quad \text{code}_G g_n \rho \\ &\quad \text{popenv} \end{aligned}$$

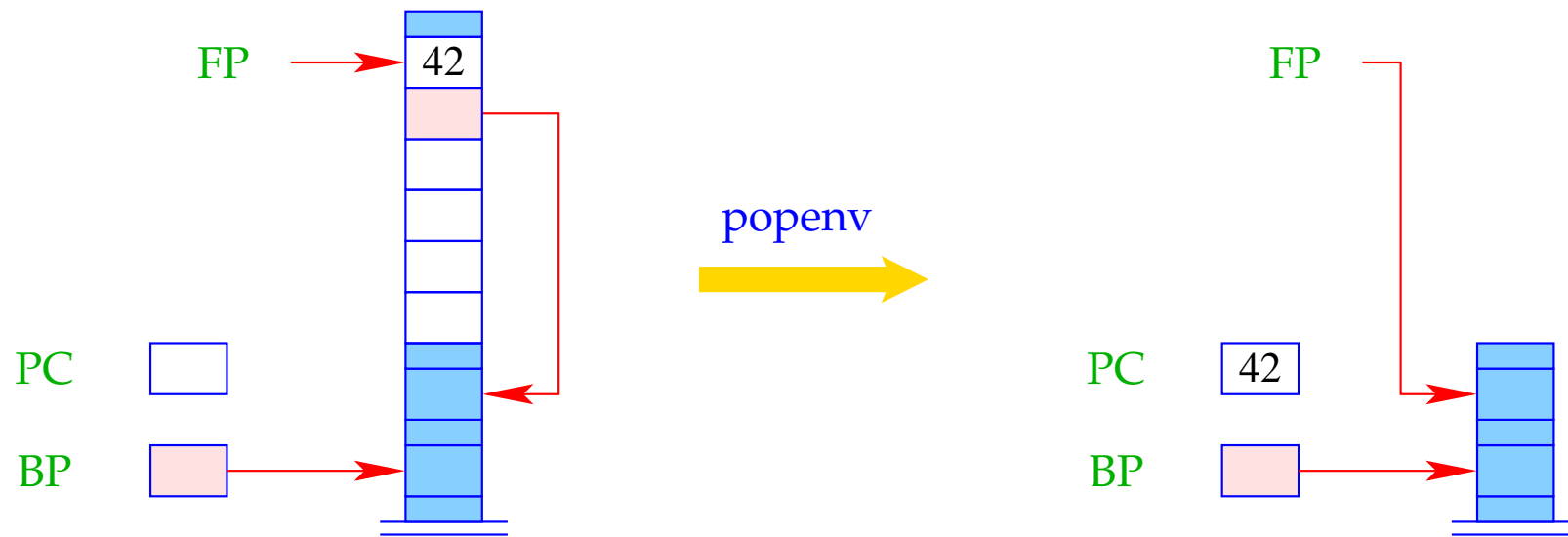
The present stack frame can be **popped** ...

- if the applied clause was the **last** (or **only**); and
- if all goals in the body are definitely **finished**.

\implies the backtrack point is **older** :-)

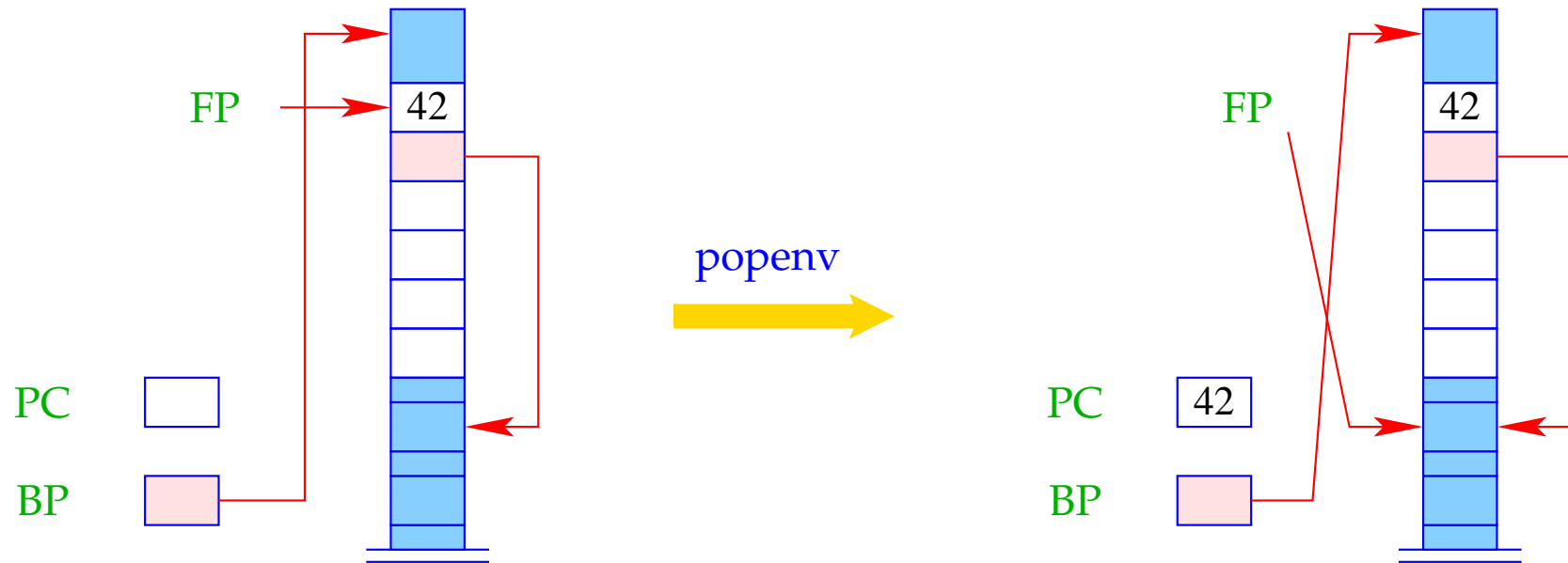
\implies **FP > BP**

The instruction `popenv` restores the registers `FP` and `PC` and possibly pops the stack frame:



```
if (FP > BP) SP = FP - 6;  
PC = posCont;  
FP = FPold;
```

Warning: `popenv` may fail to de-allocate the frame !!!



```

if (FP > BP) SP = FP - 6;
PC = posCont;
FP = FPold;

```

If popping the stack frame fails, new data are allocated on top of the stack. When returning to the frame, the locals still can be accessed through the **FP** :-))

33 Queries and Programs

The translation of a program: $p \equiv rr_1 \dots rr_h ? g$
consists of:

- an instruction `no` for failure;
- code for evaluating the query `g`;
- code for the predicate definitions rr_i .

Preceding query evaluation:

- \implies initialization of registers
- \implies allocation of space for the globals

Succeeding query evaluation:

- \implies returning the values of globals

```

code p =      init A
              pushenv d
              codeG g ρ
              halt d
            A: no
              codeP rr1
              ...
              codeP rrh

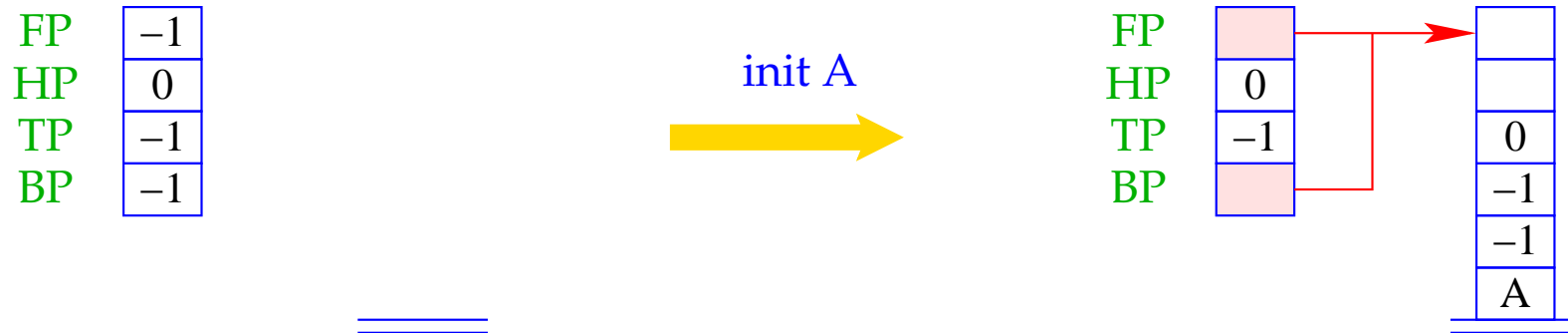
```

where $free(g) = \{X_1, \dots, X_d\}$ and ρ is given by $\rho X_i = i$.

The instruction `halt d ...`

- ... terminates the program execution;
- ... returns the bindings of the d globals;
- ... causes backtracking — if demanded by the user :-)

The instruction `init A` is defined by:



BP = FP = SP = 5;
S[0] = A;
S[1] = S[2] = -1;
S[3] = 0;
BP = FP;

At address "A" for a failing goal we have placed the instruction `no` for printing `no` to the standard output and halt :-)

The Final Example:

$$\begin{array}{lll}
 t(X) \leftarrow \bar{X} = b & q(X) \leftarrow s(\bar{X}) & s(X) \leftarrow \bar{X} = a \\
 p \leftarrow q(X), t(\bar{X}) & s(X) \leftarrow t(\bar{X}) & ? \quad p
 \end{array}$$

The translation yields:

	init N		popenv	q/1:	pushenv 1	E:	pushenv 1
	pushenv 0	p/0:	pushenv 1		mark D		mark G
	mark A		makr B		putref 1		putref 1
	call p/0		putvar 1		call s/1		call t/1
A:	halt 0		call q/1	D:	popenv	G:	popenv
N:	no	B:	mark C	s/1:	setbtp	F:	pushenv 1
t/1:	pushenv 1		putref 1		try E		putref 1
	putref 1		call t/1		delbtp		uatom a
	uatom b	C:	popenv		jump F		popenv

34 Last Call Optimization

Consider the `app` predicate from the beginning:

$$\text{app}(X, Y, Z) \leftarrow X = [], Y = Z$$
$$\text{app}(X, Y, Z) \leftarrow X = [H|X'], Z = [H|Z'], \text{app}(X', Y, Z')$$

We observe:

- The recursive call occurs in the **last** goal of the clause.
- Such a goal is called **last call**.

\implies we try to evaluate it in the **current** stack frame !!!

\implies after (successful) completion, we will not return to the current caller !!!

Consider a clause r :
 with m locals where
 code_G :

$p(X_1, \dots, X_k) \leftarrow g_1, \dots, g_n$
 $g_n \equiv q(t_1, \dots, t_h)$. The interplay between code_C and

$\text{code}_C r =$

```

pushenv m
code_G g1 ρ
...
code_G gn-1 ρ
mark B
code_A t1 ρ
...
code_A th ρ
call q/h
B : popenv
  
```

Replacement:	mark B	\implies	lastmark
	call q/h; popenv	\implies	lastcall q/h m

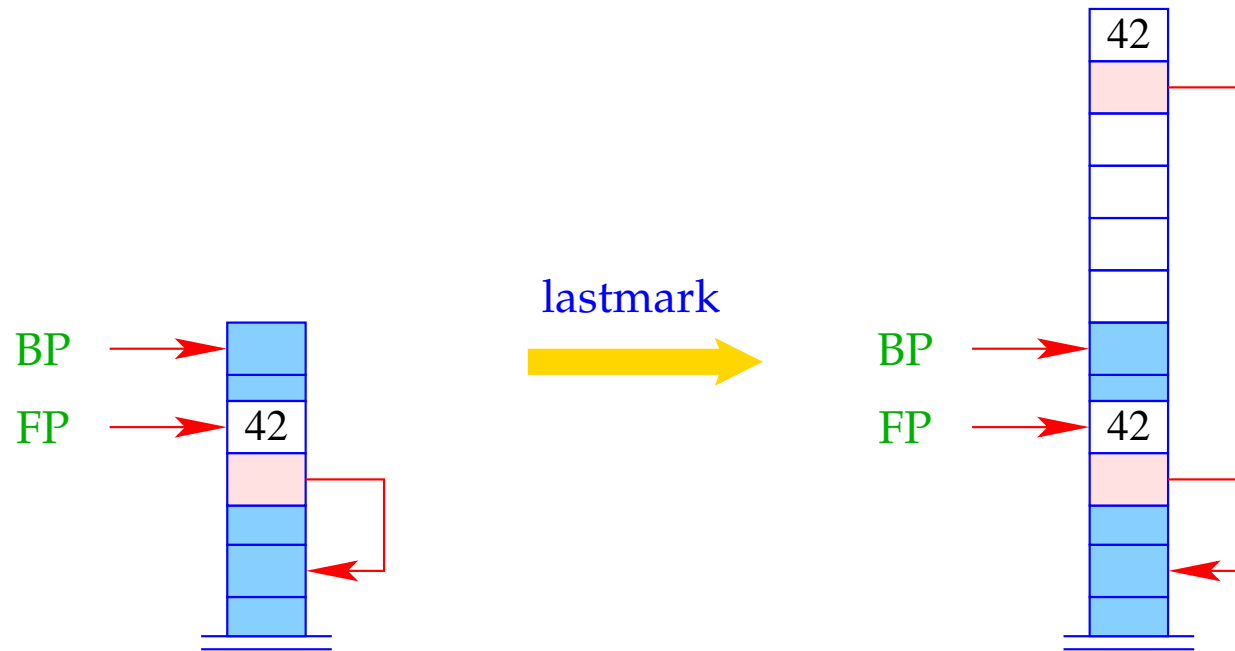
Consider a clause r : $p(X_1, \dots, X_k) \leftarrow g_1, \dots, g_n$
 with m locals where $g_n \equiv q(t_1, \dots, t_h)$. The interplay between code_C and code_G :

$\text{code}_C r =$ $\text{pushenv } m$
 $\text{code}_G g_1 \rho$
 ...
 $\text{code}_G g_{n-1} \rho$
 lastmark
 $\text{code}_A t_1 \rho$
 ...
 $\text{code}_A t_h \rho$
 $\text{lastcall } q/h \ m$

Replacement:	$\text{mark } B$	\implies	lastmark
	$\text{call } q/h; \text{popenv}$	\implies	$\text{lastcall } q/h \ m$

If the current clause is not **last** or the g_1, \dots, g_{n-1} have created backtrack points, then **FP** \leq **BP** :-)

Then **lastmark** creates a new frame but stores a reference to the **predecessor**:



```

if (FP  $\leq$  BP) {
    SP = SP + 6;
    S[SP] = posCont; S[SP-1] = FPold;
}

```

If **FP** $>$ **BP** then **lastmark** does nothing :-)

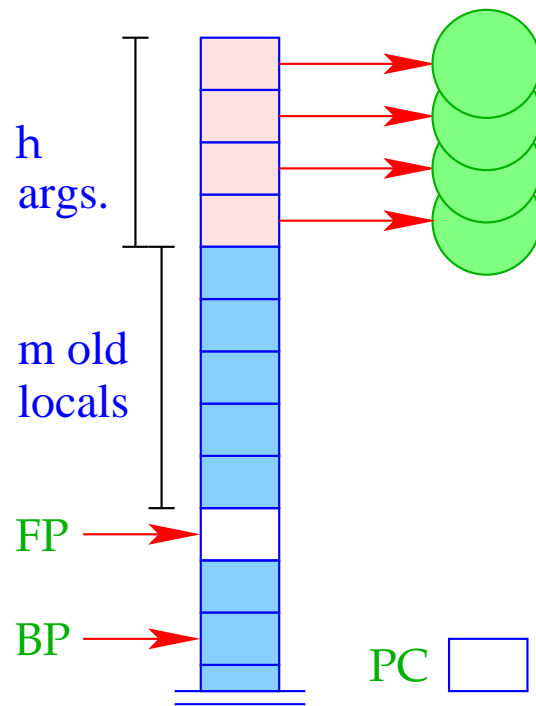
If $FP \leq BP$, then `lastcall q/h m` behaves like a normal `call q/h`.

Otherwise, the current stack frame is re-used. This means that:

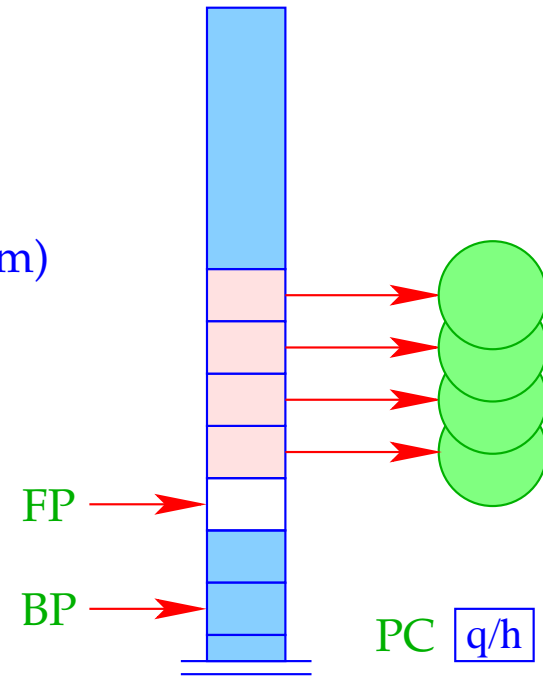
- the cells $S[FP+1], S[FP+2], \dots, S[FP+h]$ receive the new values and
- `q/h` can be jumped to `:-)`

```
lastcall q/h m = if (FP ≤ BP) call q/h;
                else {
                    move m h;
                    jump q/h;
                }
```

The difference between the old and the new addresses of the parameters `m` just equals the number of the `local variables` of the current clause `:-))`



lastcall (q/h,m)



Example:

Consider the clause:

$$a(X, Y) \leftarrow f(\bar{X}, X_1), a(\bar{X}_1, \bar{Y})$$

The last-call optimization for `codeC r` yields:

	mark A	A:	lastmark
pushenv 3	putref 1		putref 3
	putvar 3		putref 2
	call f/2		lastcall a/2 3

Example:

Consider the clause:

$$a(X, Y) \leftarrow f(\bar{X}, X_1), a(\bar{X}_1, \bar{Y})$$

The last-call optimization for `codeC r` yields:

	mark A	A: lastmark
pushenv 3	putref 1	putref 3
	putvar 3	putref 2
	call f/2	lastcall a/2 3

Note:

If the clause is **last** and the last literal is the **only one**, we can skip **lastmark** and can replace **lastcall q/h m** with the sequence **move m n; jump p/n :-))**

Example:

Consider the **last** clause of the `app` predicate:

$$\text{app}(X, Y, Z) \leftarrow \bar{X} = [H|X'], \bar{Z} = [\bar{H}|Z'], \text{app}(\bar{X}', \bar{Y}, \bar{Z}')$$

Here, the last call is the **only one** :-). Consequently, we obtain:

A: pushenv 6			uref 4	bind
putref 1	B: putvar 4		son 2	E: putref 5
ustruct [[]]/2 B	putvar 5		uvar 6	putref 2
son 1	putstruct [[]]/2		up E	putref 6
uvar 4	bind	D: check 4		move 6 3
son 2	C: putref 3	putref 4		jump app/3
uvar 5	ustruct [[]]/2 D	putvar 6		
up C	son 1	putstruct [[]]/2		

35 Trimming of Stack Frames

Idea:

- Order local variables according to their **life times**;
- Pop the **dead** variables — if possible **:-}**

35 Trimming of Stack Frames

Idea:

- Order local variables according to their **life times**;
- Pop the **dead** variables — if possible **:-}**

Example:

Consider the clause:

$$a(X, Z) \leftarrow p_1(\bar{X}, X_1), p_2(\bar{X}_1, X_2), p_3(\bar{X}_2, X_3), p_4(\bar{X}_3, \bar{Z})$$

35 Trimming of Stack Frames

Idea:

- Order local variables according to their **life times**;
- Pop the **dead** variables — if possible **:-}**

Example:

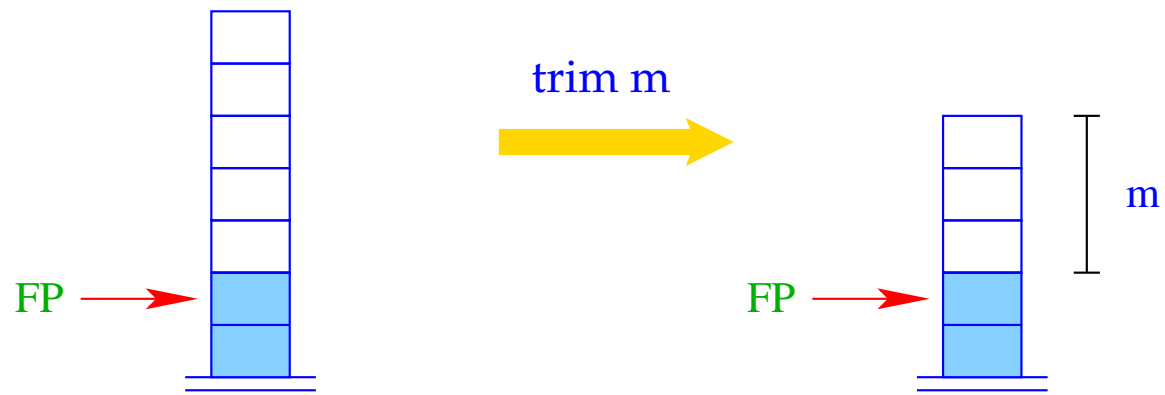
Consider the clause:

$$a(X, Z) \leftarrow p_1(\bar{X}, X_1), p_2(\bar{X}_1, X_2), p_3(\bar{X}_2, X_3), p_4(\bar{X}_3, \bar{Z})$$

After the query $p_2(\bar{X}_1, X_2)$, variable X_1 is dead.

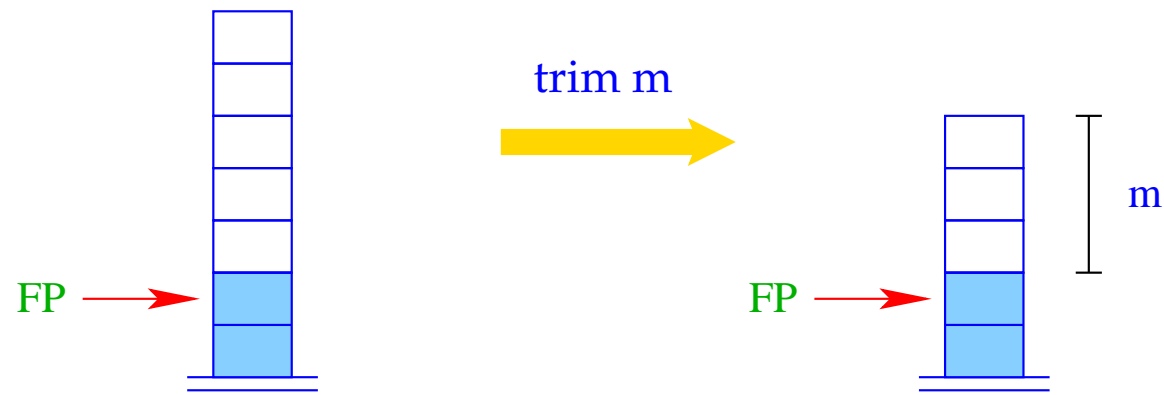
After the query $p_3(\bar{X}_2, X_3)$, variable X_2 is dead **:-)**

After every non-last goal with dead variables, we insert the instruction `trim` :



```
if (FP ≥ BP)
    SP = FP + m;
```

After every non-last goal with dead variables, we insert the instruction `trim` :



```
if (FP ≥ BP)
    SP = FP + m;
```

The dead locals can only be popped if no new backtrack point has been allocated :-)

Example (continued):

$$a(X, Z) \leftarrow p_1(\bar{X}, X_1), p_2(\bar{X}_1, X_2), p_3(\bar{X}_2, X_3), p_4(\bar{X}_3, \bar{Z})$$

Ordering of the variables:

$$\rho = \{X \mapsto 1, Z \mapsto 2, X_3 \mapsto 3, X_2 \mapsto 4, X_1 \mapsto 5\}$$

The resulting code:

pushenv 5	A:	mark B	mark C	lastmark
mark A		putref 5	putref 4	putref 3
putref 1		putvar 4	putvar 3	putref 2
putvar 5		call p ₂ /2	call p ₃ /2	lastcall p ₄ /2 3
call p ₁ /2	B:	trim 4	C:	trim 3

36 Clause Indexing

Observation:

Often, predicates are implemented by case distinction on the first argument.

- ⇒ Inspecting the first argument, many alternatives can be excluded :-)
- ⇒ Failure is earlier detected :-)
- ⇒ Backtrack points are earlier removed. :-))
- ⇒ Stack frames are earlier popped :-)))

Example: The app-predicate:

$$\text{app}(X, Y, Z) \leftarrow X = [], Y = Z$$
$$\text{app}(X, Y, Z) \leftarrow X = [H|X'], Z = [H|Z'], \text{app}(X', Y, Z')$$

- If the root constructor is $[]$, only the first clause is applicable.
- If the root constructor is $[[]]$, only the second clause is applicable.
- Every other root constructor should **fail !!**
- Only if the first argument equals an unbound variable, both alternatives must be tried **;-)**

Idea:

- Introduce separate try chains for every possible constructor.
- Inspect the root node of the first argument.
- Depending on the result, perform an **indexed** jump to the appropriate try chain.

Assume that the predicate p/k is defined by the sequence rr of clauses $r_1 \dots r_m$.

Let **tchains** rr denote the sequence of try chains as built up for the root constructors occurring in unifications $X_1 = t$.

Example:

Consider again the `app`-predicate, and assume that the code for the two clauses start at addresses A_1 and A_2 , respectively.

Then we obtain the following four `try chains`:

<code>VAR:</code>	<code>setbtp</code>	<code>// variables</code>	<code>NIL:</code>	<code>jump A₁</code>	<code>// atom []</code>
	<code>try A₁</code>				
	<code>delbtp</code>		<code>CONS:</code>	<code>jump A₂</code>	<code>// constructor []</code>
	<code>jump A₂</code>		<code>ELSE:</code>	<code>fail</code>	<code>// default</code>

Example:

Consider again the `app`-predicate, and assume that the code for the two clauses start at addresses A_1 and A_2 , respectively.

Then we obtain the following four **try chains**:

<code>VAR:</code>	<code>setbtp</code>	<code>// variables</code>	<code>NIL:</code>	<code>jump A₁</code>	<code>// atom []</code>
	<code>try A₁</code>				
	<code>delbtp</code>		<code>CONS:</code>	<code>jump A₂</code>	<code>// constructor [[]]</code>
	<code>jump A₂</code>				
			<code>ELSE:</code>	<code>fail</code>	<code>// default</code>

The new instruction `fail` takes care of any constructor besides `[]` and `[[]]` ...

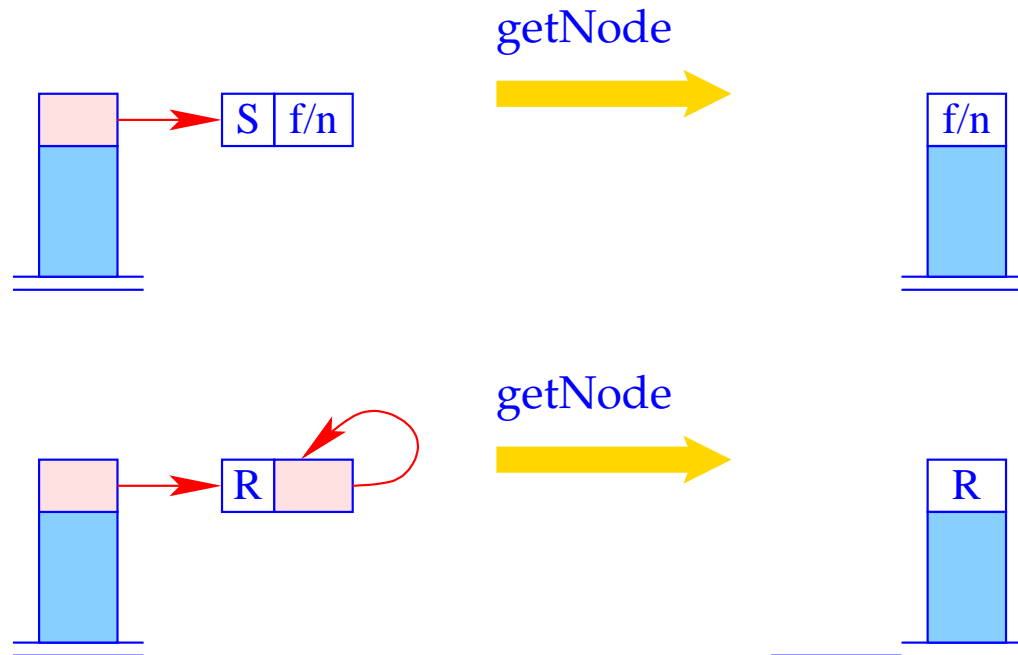
`fail` = `backtrack()`

It directly triggers **backtracking** :-)

Then we generate for a predicate p/k :

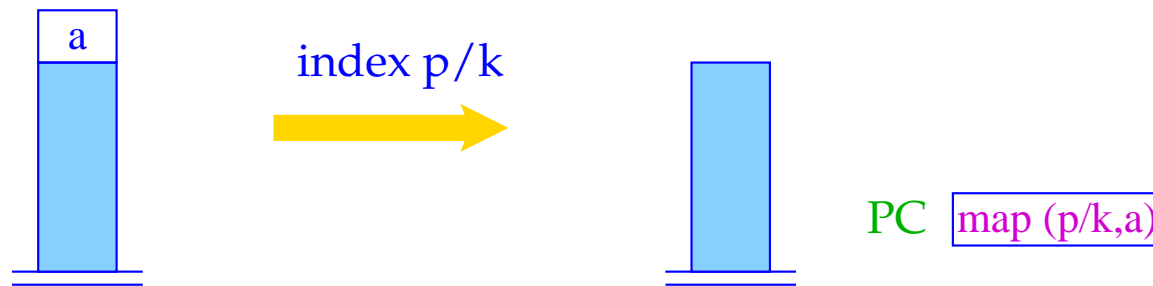
```
codep rr    =      putref 1
                   getNode // extracts the root label
                   index p/k // jumps to the try block
                   tchains rr
A1 : codeC r1
           ...
Am : codeC rm
```

The instruction `getNode` returns "R" if the pointer on top of the stack points to an unbound variable. Otherwise, it returns the content of the heap object:



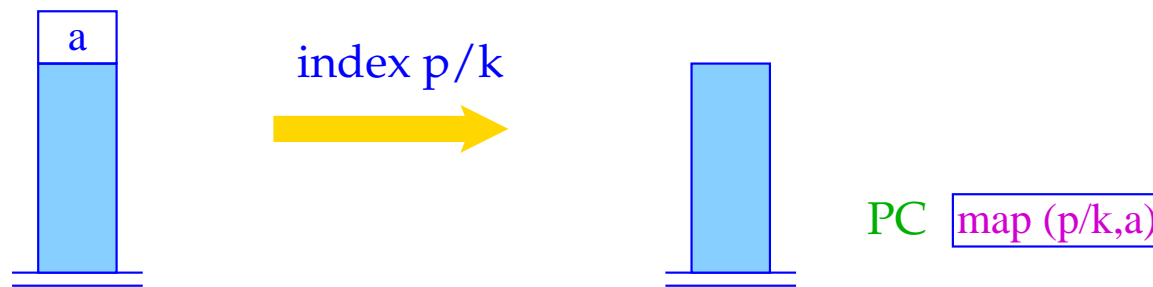
```
switch (H[S[SP]]) {  
  case (S, f/n):  S[SP] = f/n; break;  
  case (A,a):    S[SP] = a; break;  
  case (R,_):   S[SP] = R;  
}
```

The instruction `index p/k` performs an indexed jump to the appropriate try chain:



```
PC = map (p/k,S[SP]);  
SP--;
```


The instruction `index p/k` performs an indexed jump to the appropriate try chain:



```
PC = map (p/k,S[SP]);  
SP--;
```

The function `map()` returns, for a given predicate and node content, the start address of the appropriate try chain :-)

It typically is defined through some hash table :-))

37 Extension: The Cut Operator

Realistic Prolog additionally provides an operator “!” (cut) which explicitly allows to prune the search space of backtracking.

Example:

$$\begin{aligned} \text{branch}(X, Y) &\leftarrow p(X), !, q_1(X, Y) \\ \text{branch}(X, Y) &\leftarrow q_2(X, Y) \end{aligned}$$

Once the queries before the cut have succeeded, the choice is committed:

Backtracking will return only to backtrack points preceding the call to the left-hand side ...

The Basic Idea:

- We restore the `oldBP` from our current stack frame;
- We pop all stack frames on top of the local variables.

Accordingly, we translate the cut into the sequence:

```
prune  
pushenv m
```

where `m` is the number of (still used) local variables of the clause.

Example:

Consider our example:

$$\begin{aligned} \text{branch}(X, Y) &\leftarrow p(X), !, q_1(X, Y) \\ \text{branch}(X, Y) &\leftarrow q_2(X, Y) \end{aligned}$$

We obtain:

setbtp	A:	pushenv 2	C:	prune	lastmark	B:	pushenv 2
try A		mark C		pushenv 2	putref 1		putref 2
delbtp		putref 1			putref 2		putref 2
jump B		call p/1			lastcall q ₁ /2 2		move 2 2
							jump q ₂ /2

Example:

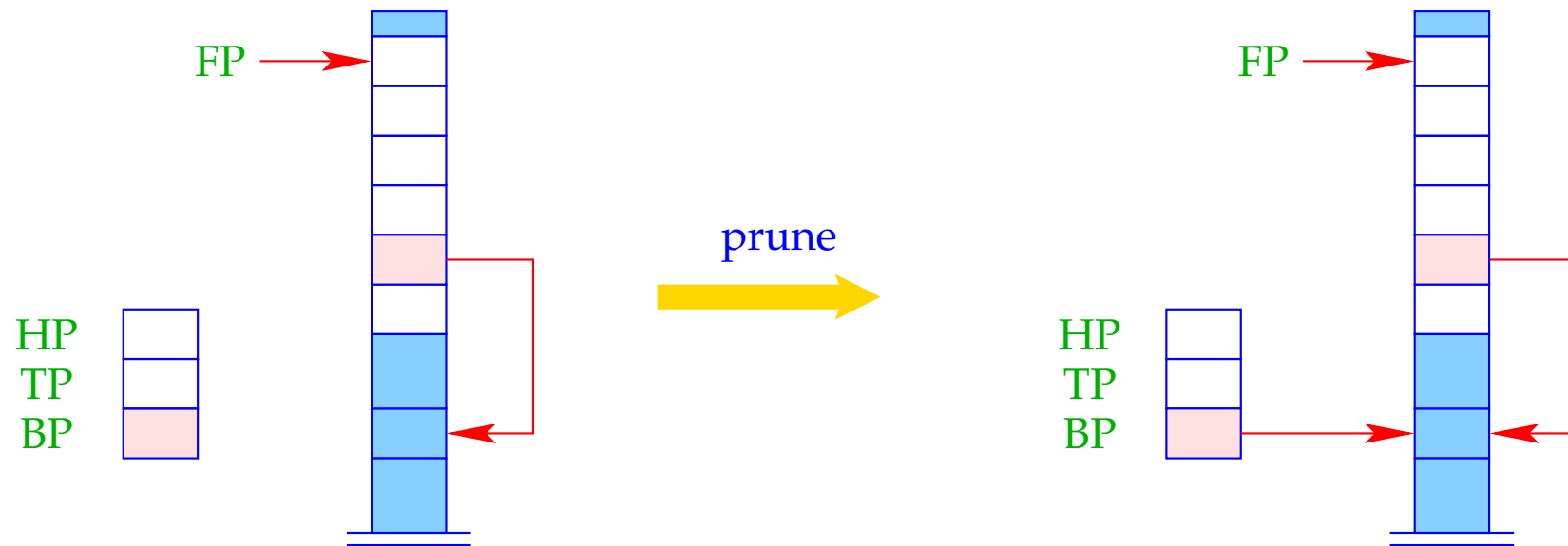
Consider our example:

$$\begin{aligned}\text{branch}(X, Y) &\leftarrow p(X), !, q_1(X, Y) \\ \text{branch}(X, Y) &\leftarrow q_2(X, Y)\end{aligned}$$

In fact, an **optimized** translation even yields here:

setbtp	A:	pushenv 2	C:	prune	putref 1	B:	pushenv 2
try A		mark C		pushenv 2	putref 2		putref 1
delbtp		putref 1			move 2 2		putref 2
jump B		call p/1			jump q ₁ /2		move 2 2
							jump q ₂ /2

The new instruction `prune` simply restores the backtrack pointer:



$BP = BP_{old};$

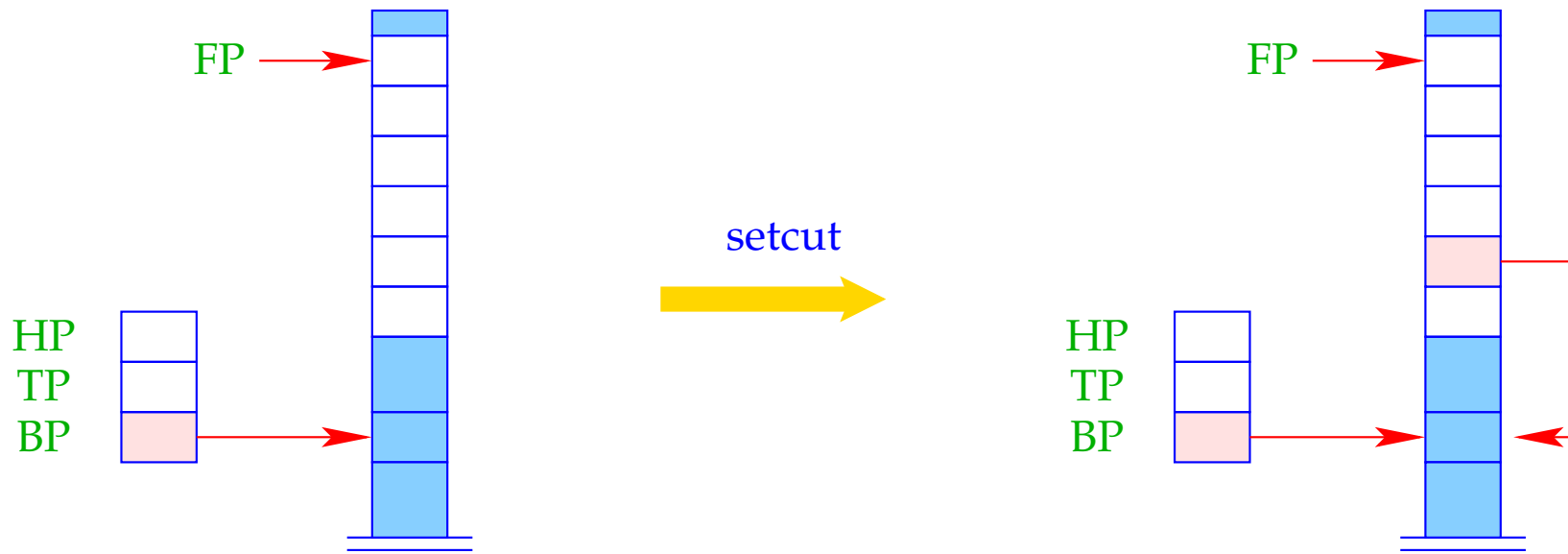
Problem:

If a clause is **single**, then (at least so far **;-)** we have not stored the old **BP** inside the stack frame **:-(**



For the cut to work also with **single-clause** predicates or try chains of length 1, we insert an extra instruction **setcut** before the clausal code (or the jump):

The instruction `setcut` just stores the current value of `BP`:



`BPold = BP;`

The Final Example: Negation by Failure

The predicate `notP` should succeed whenever `p` fails (and vice versa :-)

```
notP(X) ← p(X), !, fail
notP(X) ←
```

where the goal `fail` never succeeds. Then we obtain for `notP` :

```
setbtp   A:   pushenv 1   C:   prune       B:   pushenv 1
try A    mark C
delbtp   putref 1       fail
jump B   call p/1      popenv
```

38 Garbage Collection

- Both during execution of a **MaMa**- as well as a **WiM**-programs, it may happen that some objects can no longer be reached through references.
- Obviously, they cannot affect the further program execution. Therefore, these objects are called **garbage**.
- Their storage space should be freed and reused for the creation of other objects.

Warning:

The **WiM** provides some kind of heap de-allocation. This, however, only frees the storage of **failed alternatives** !!!

Operation of a stop-and-copy-Collector:

- Division of the heap into two parts, the **to-space** and the **from-space** — which, after each collection flip their roles.
- Allocation with **new** in the current **from-space**.
- In case of memory exhaustion, call of the collector.

The Phases of the Collection:

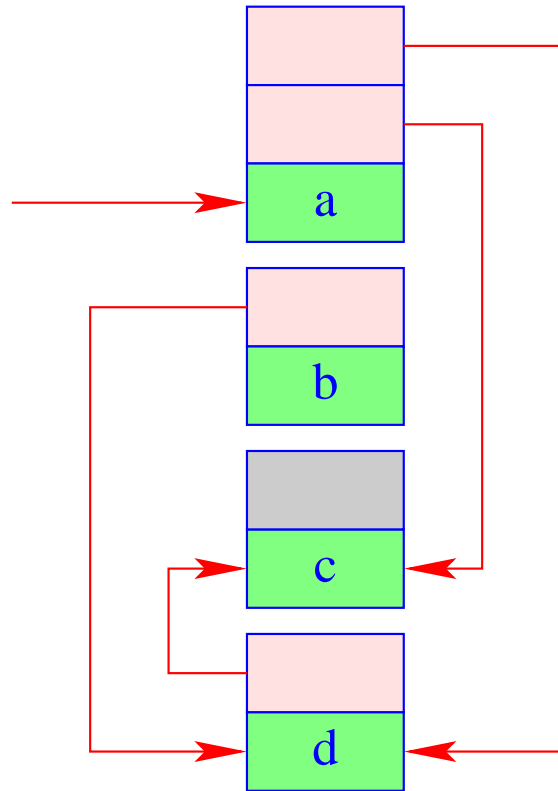
1. Marking of all reachable objects in the **from-space**.
2. Copying of all marked objects into the **to-space**.
3. Correction of references.
4. Exchange of **from-space** and **to-space**.

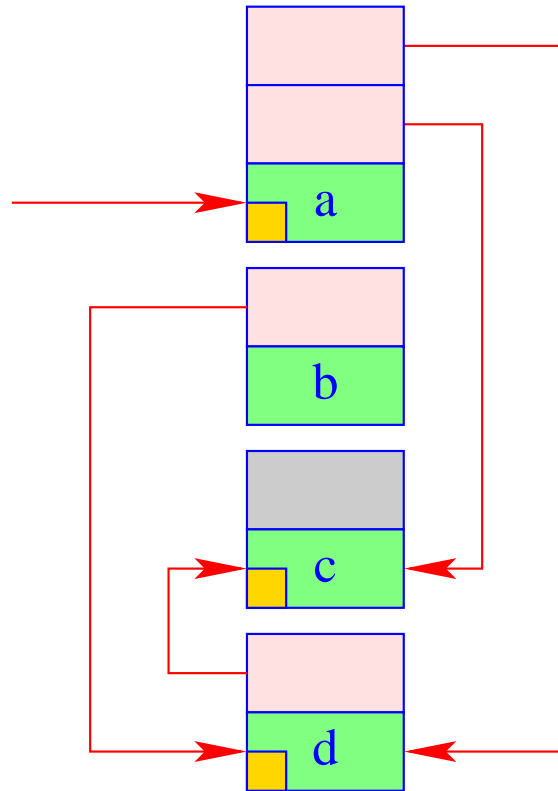
(1) **Mark:** Detection of **live** objects:

- all references in the stack point to live objects;
- every reference of a live object points to a live object.



Graph Reachability

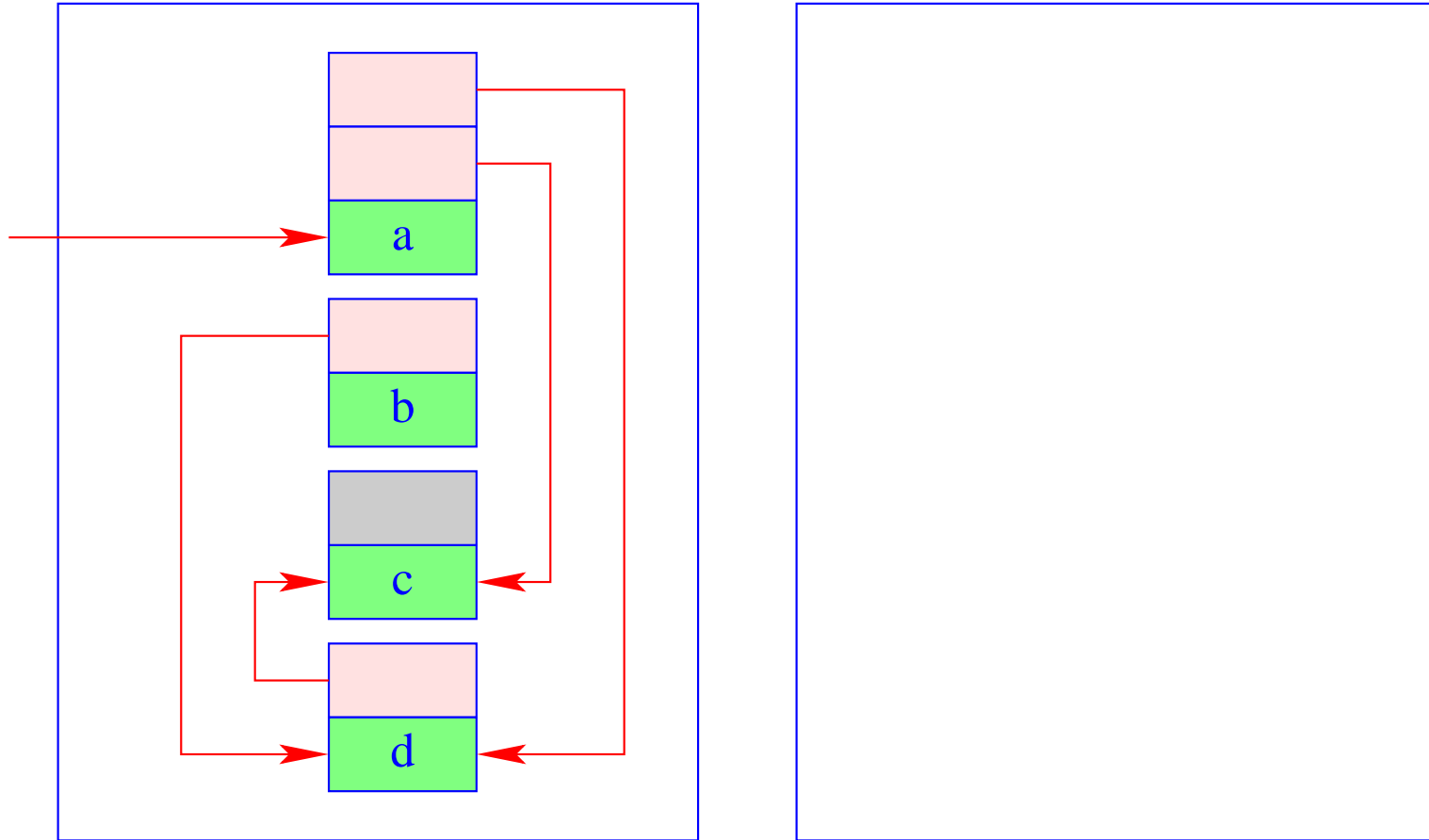


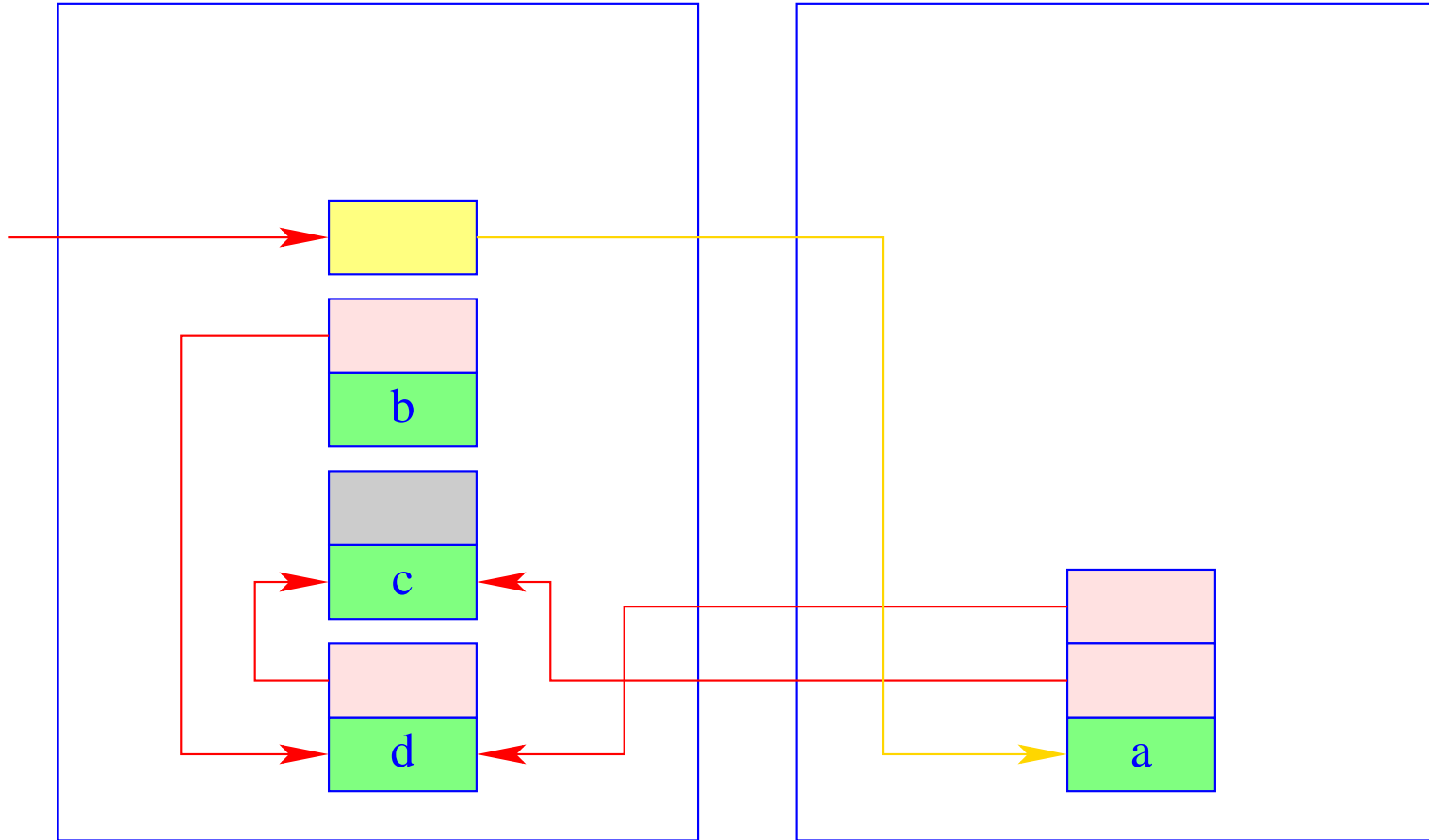


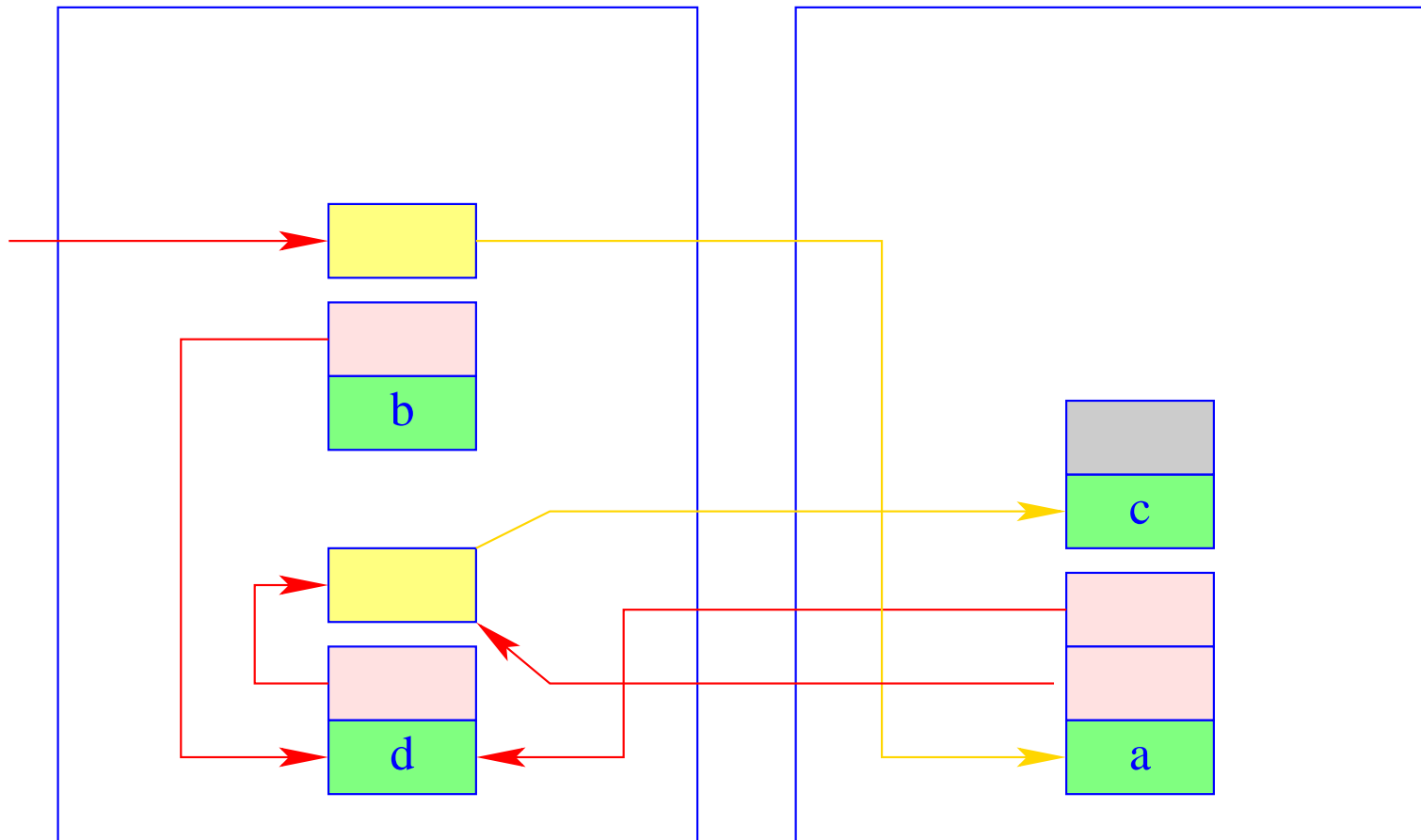
- (2) **Copy:** Copying of all live objects from the current **from-space** into the current **to-space**. This means for every detected object:
- Copying the object;
 - Storing a forward reference to the new place at the old place :-)

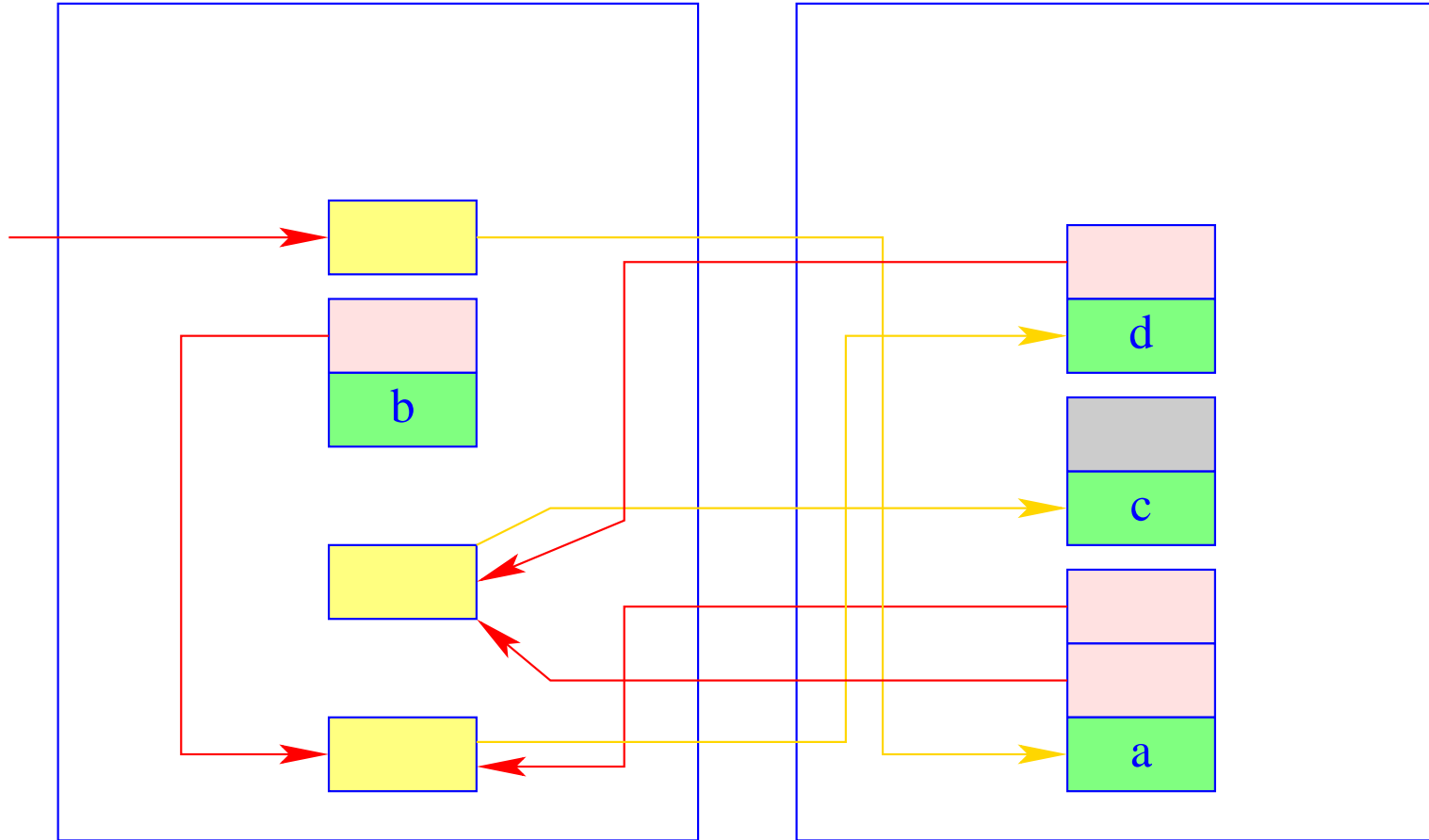


all references of the copied objects point to the forward references in the **from-space**.

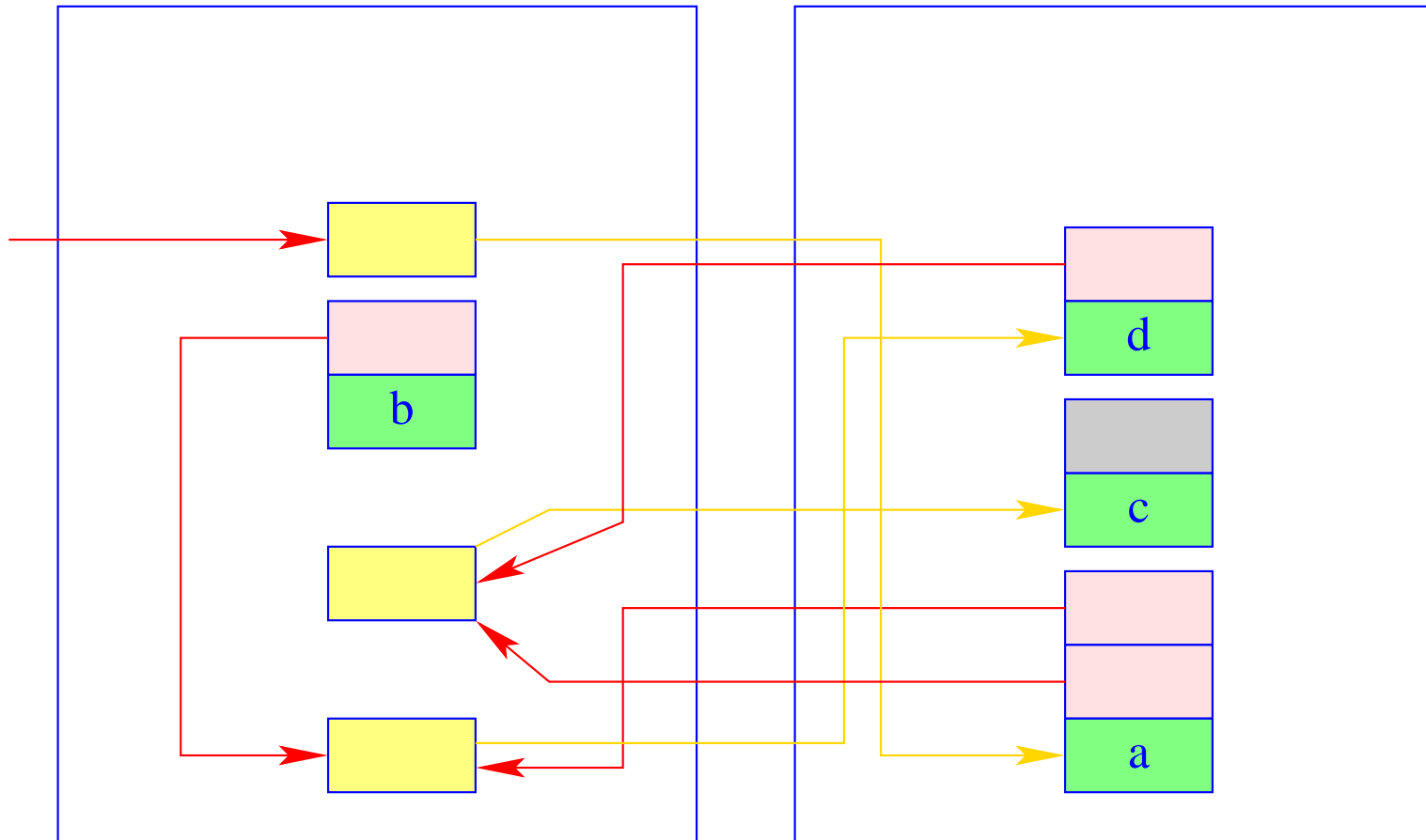


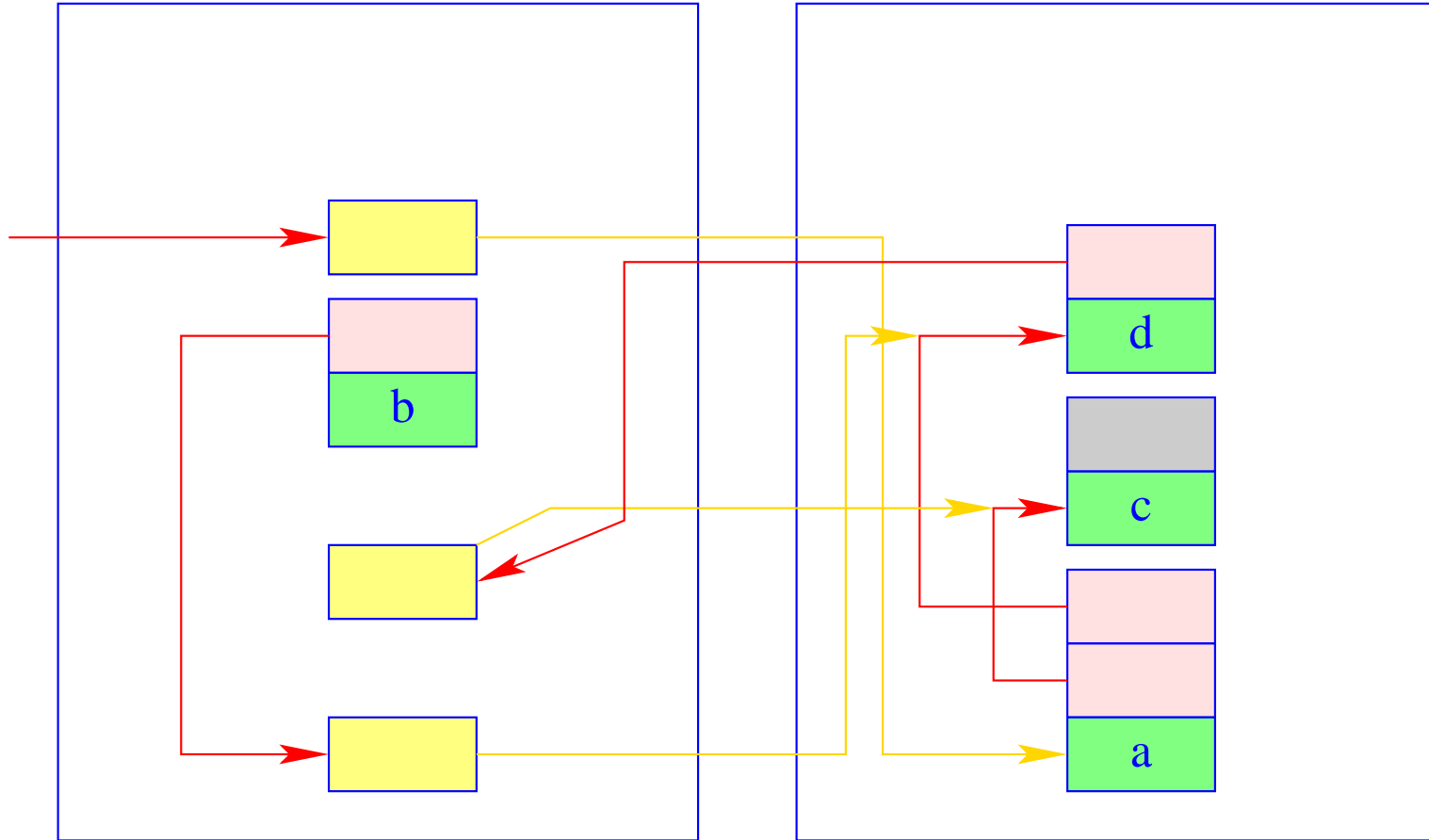


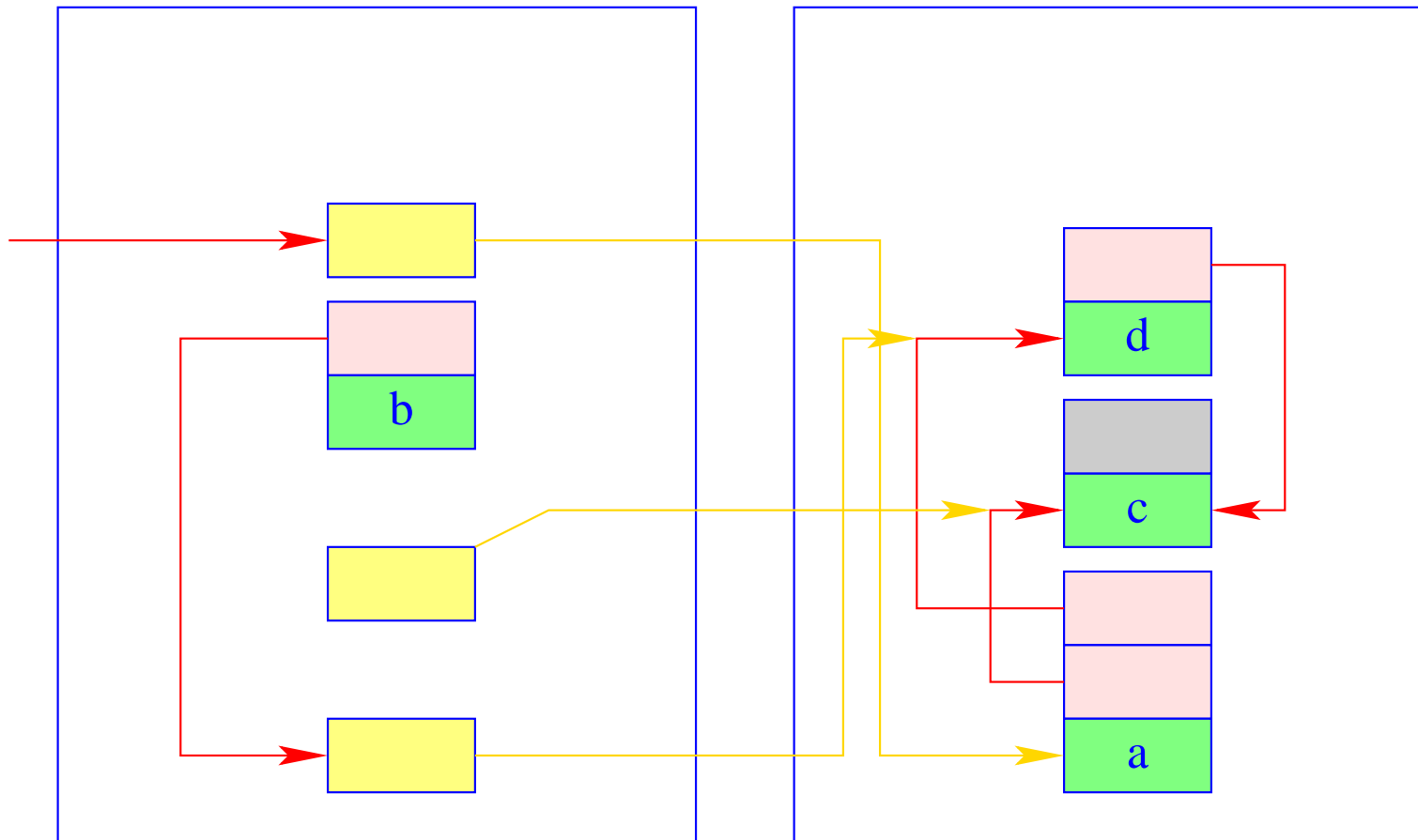


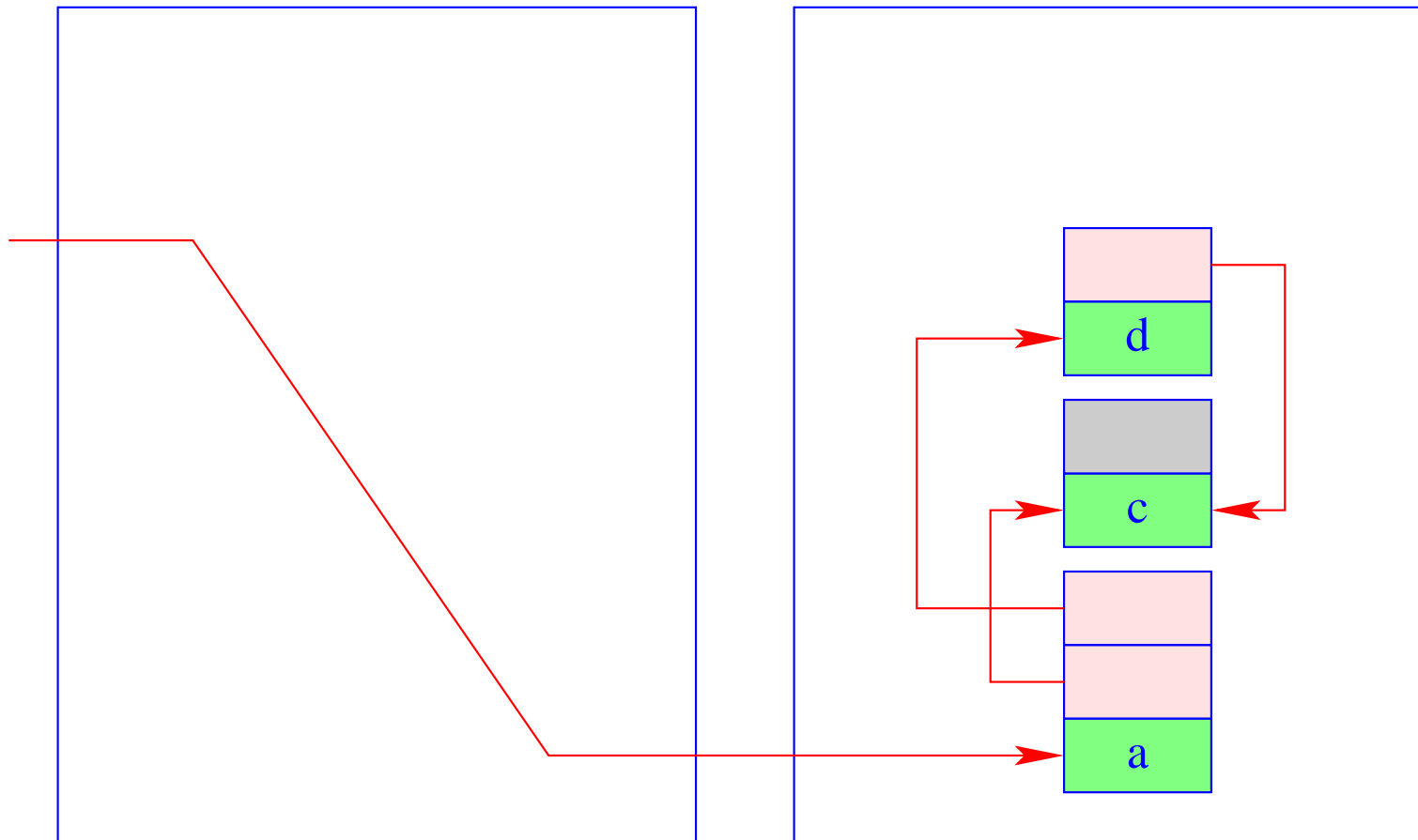


(3) Traversing of the **to-space** in order to correct the references.

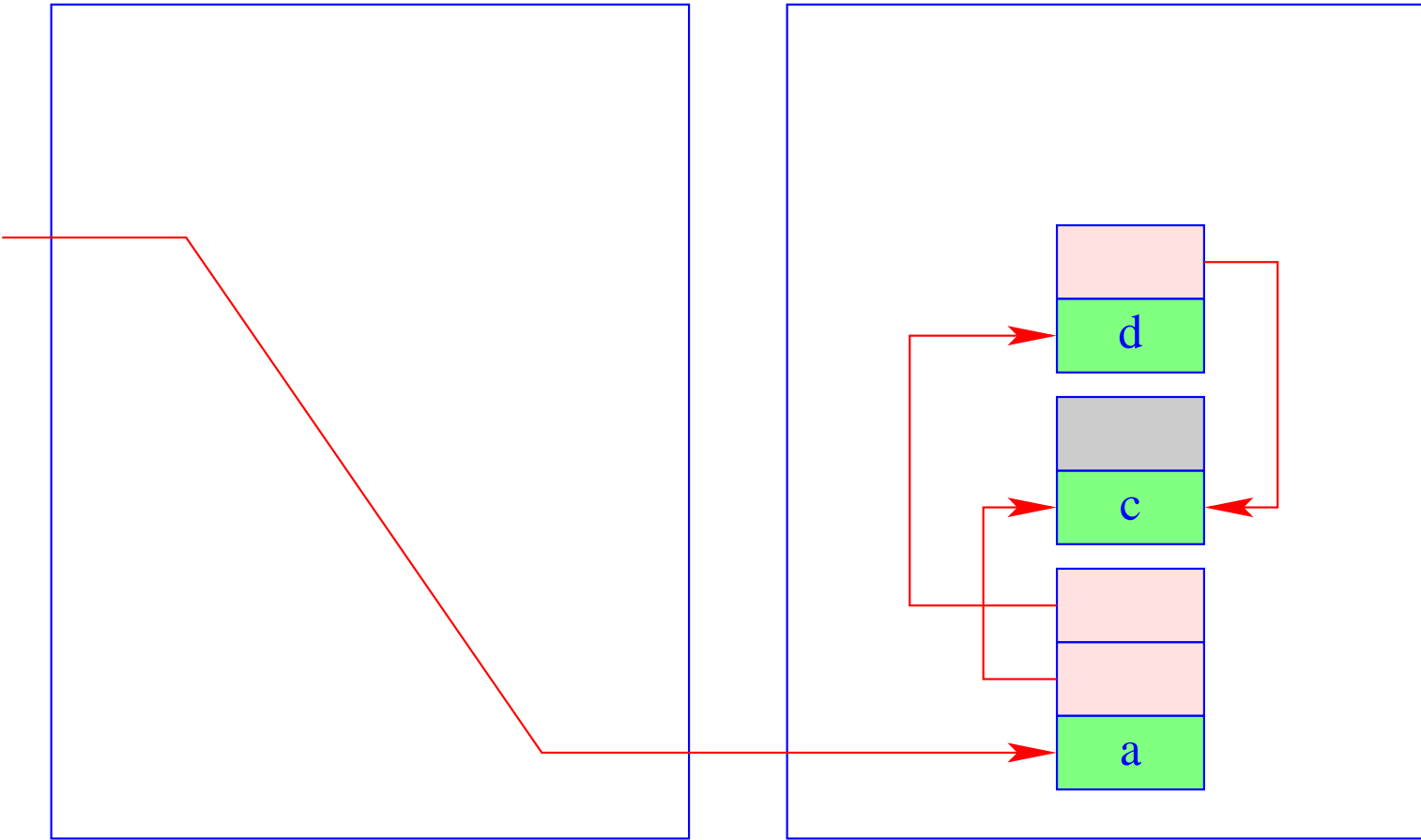


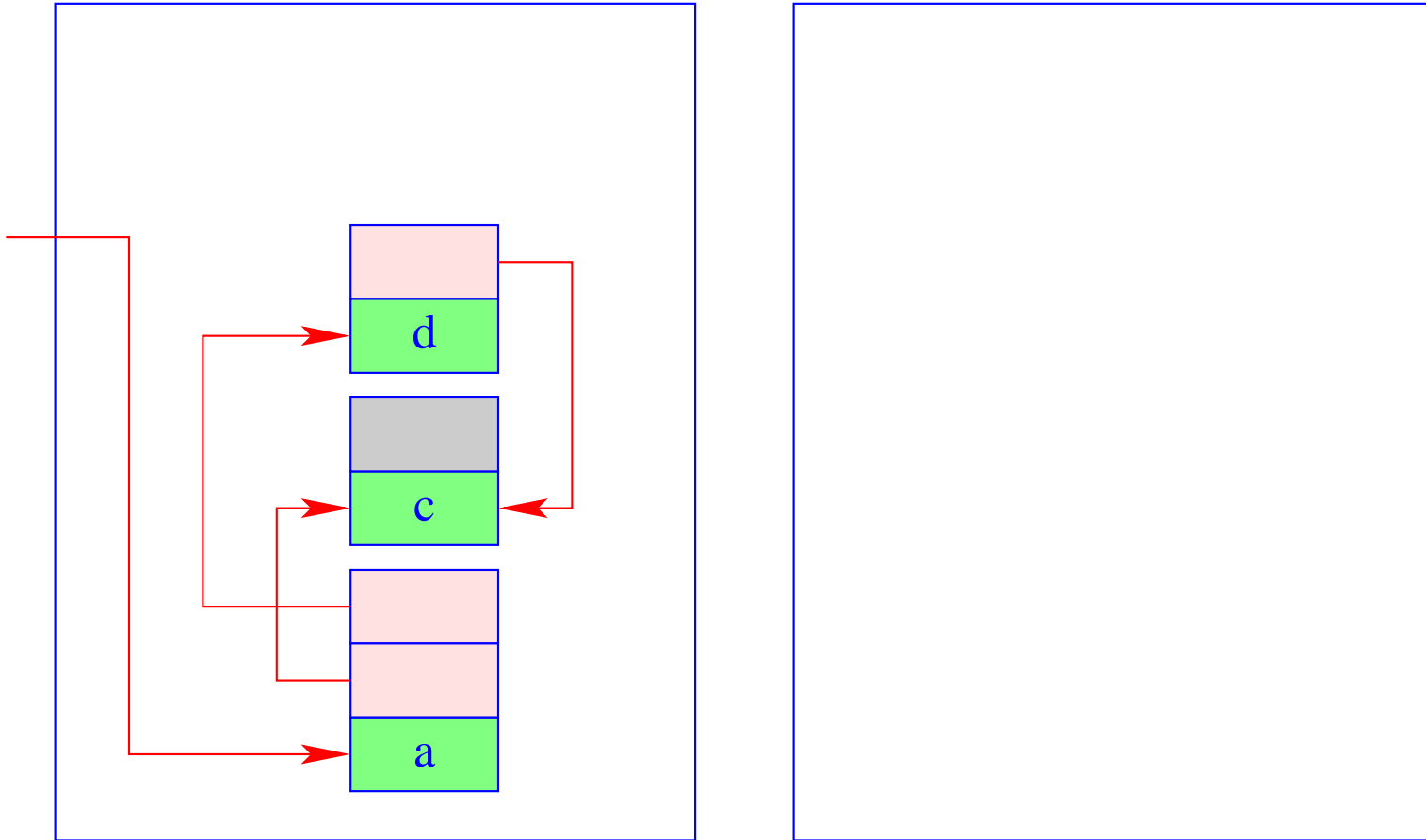






(4) Exchange of to-space and from-space.



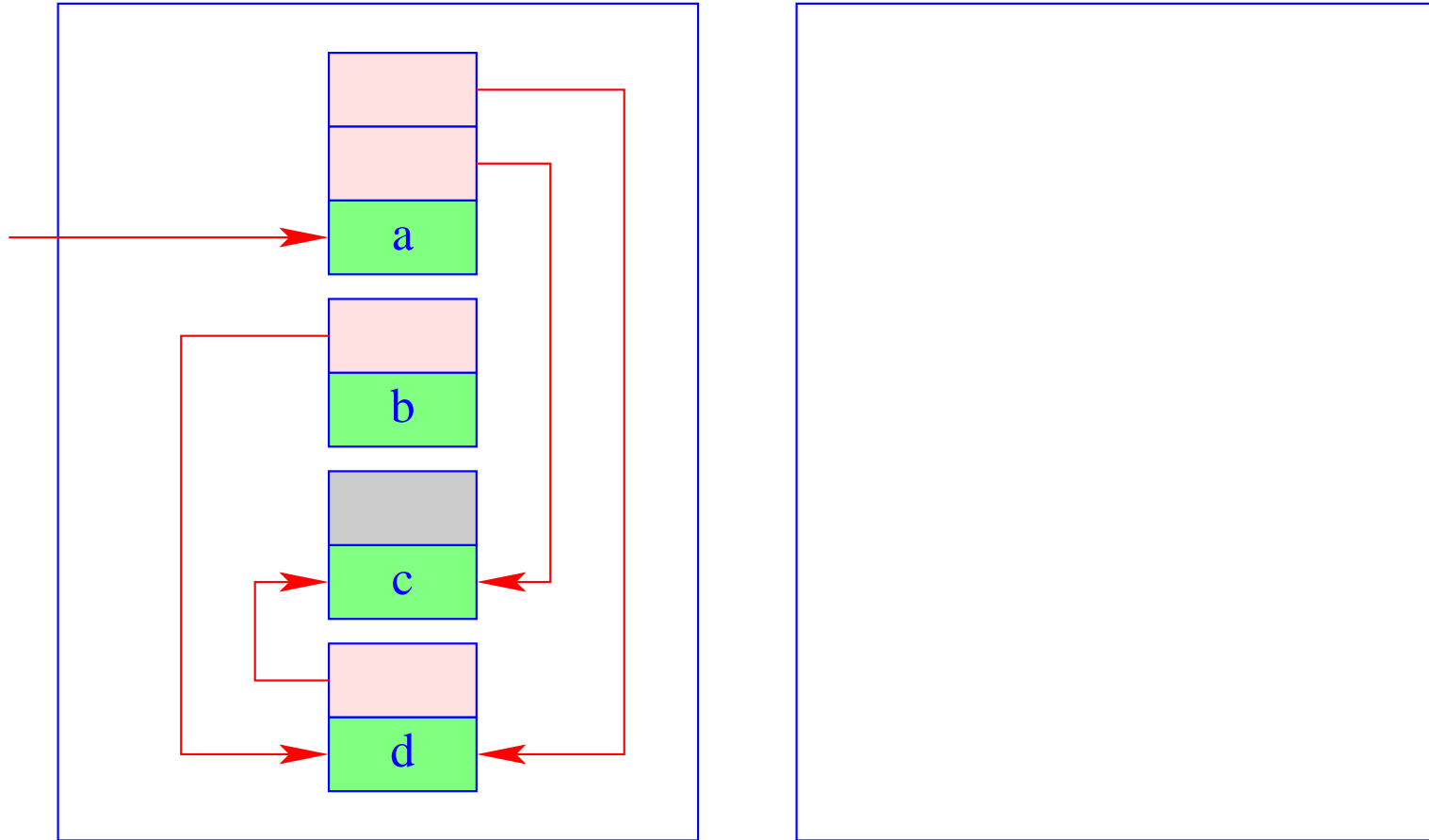


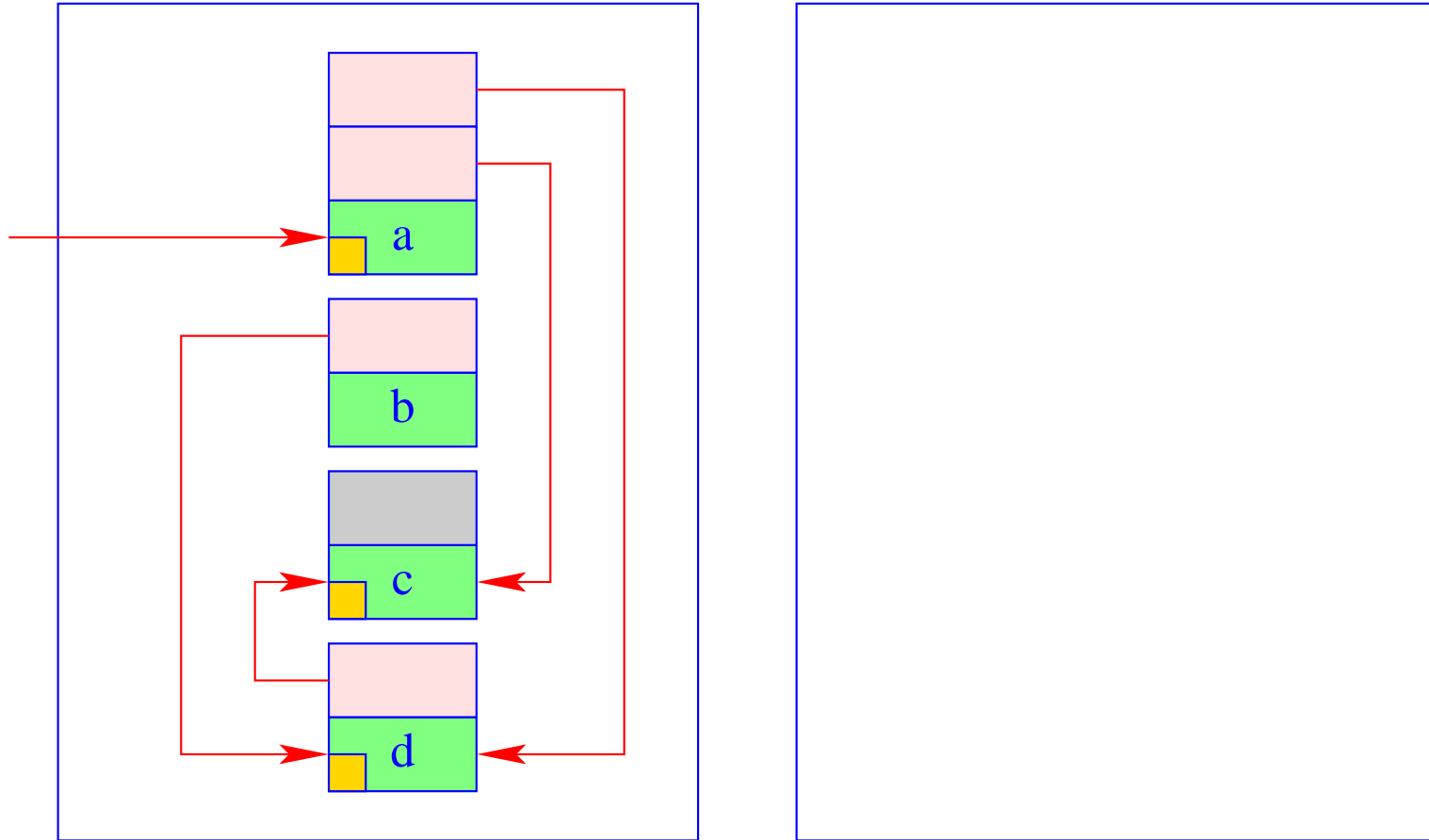
Warning:

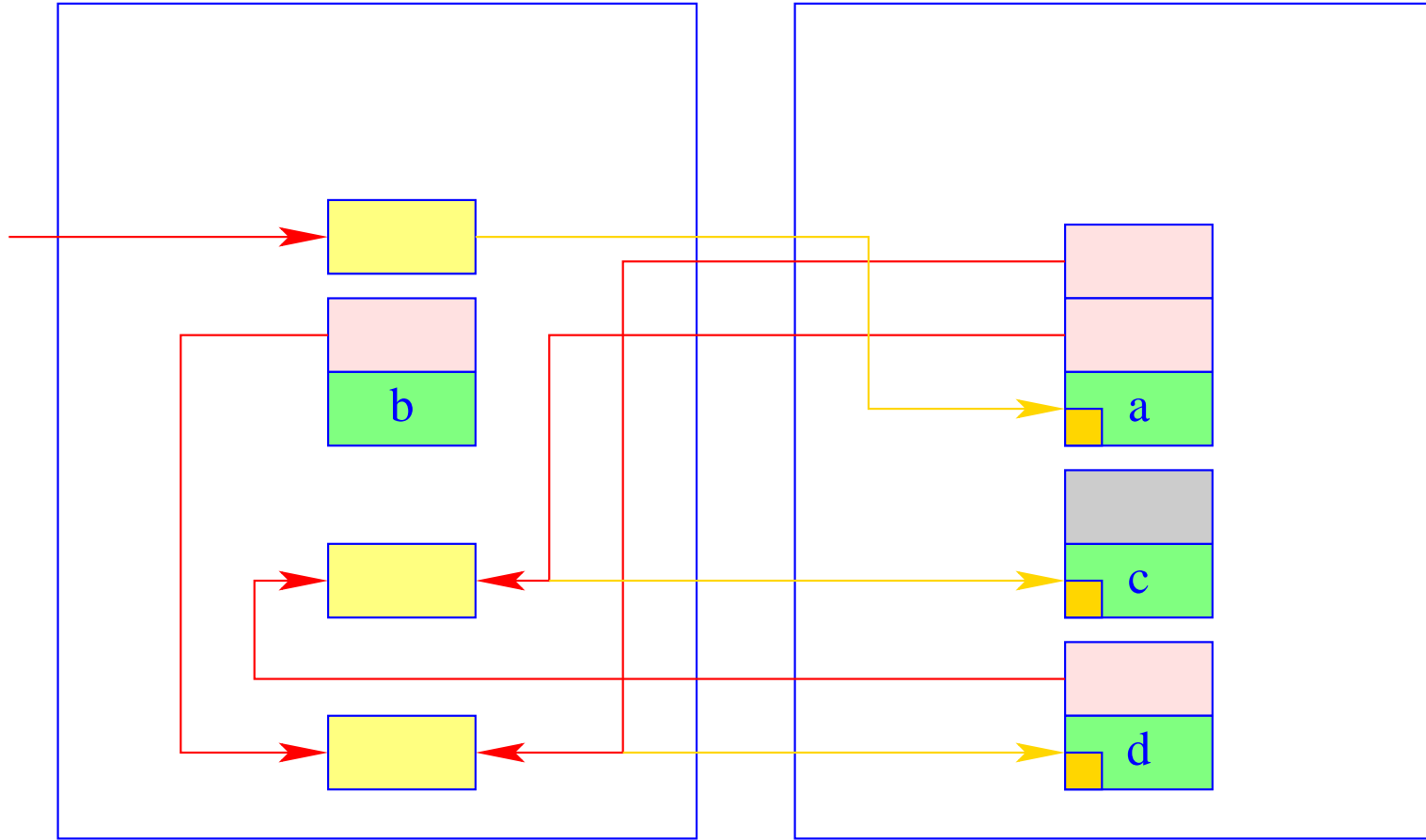
The garbage collection of the **WiM** must **harmonize** with backtracking.

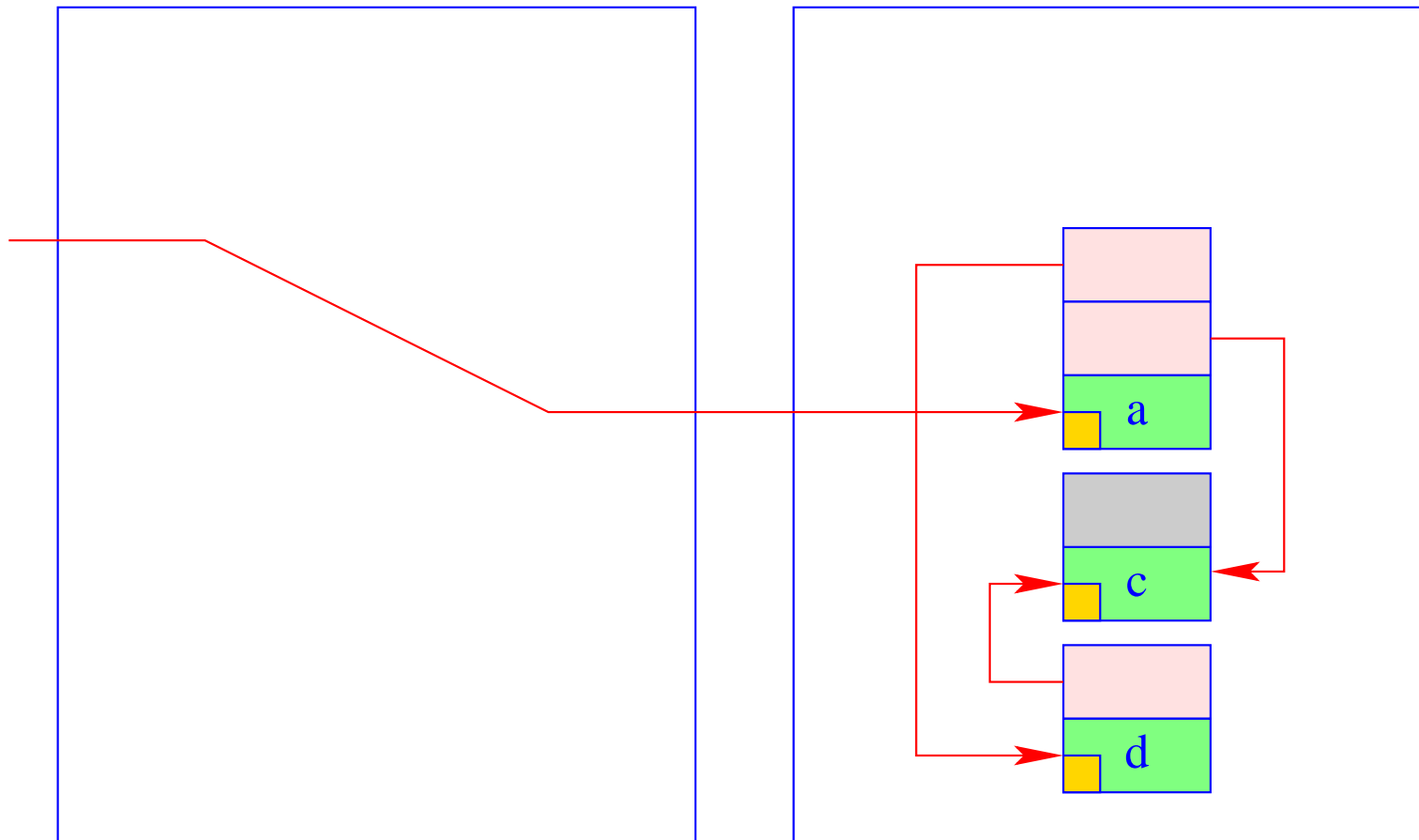
This means:

- The relative position of heap objects must not change during copying :-!!
- The heap references in the trail must be updated to the new positions.
- If heap objects are collected which have been created before the last backtrack point, then also the heap pointers in the stack must be updated.









Threads

39 The Language ThreadedC

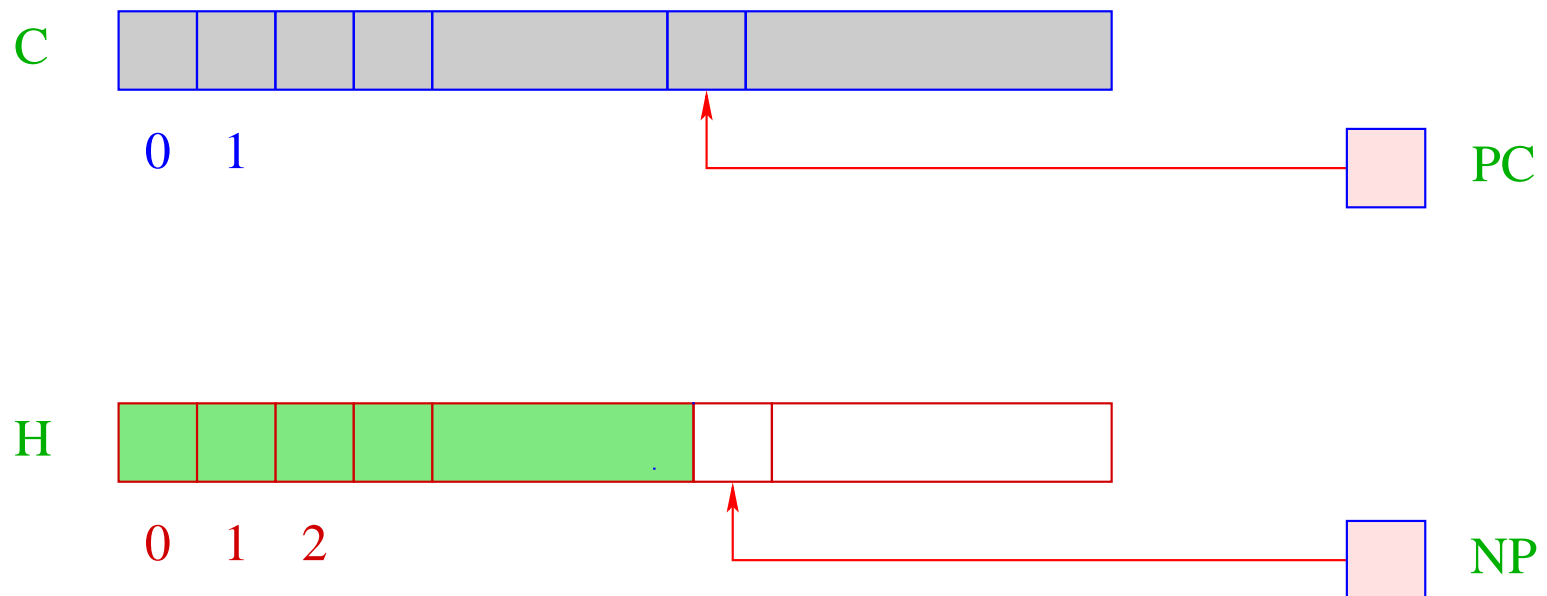
We extend C by a simple thread concept. In particular, we provide functions for:

- generating new threads: `create()`;
- terminating a thread: `exit()`;
- waiting for termination of a thread: `join()`;
- mutual exclusion: `lock()`, `unlock()`; ...

In order to enable a parallel program execution, we **extend** the abstract machine (what else? :-)

40 Storage Organization

All threads share the same common code store and heap:

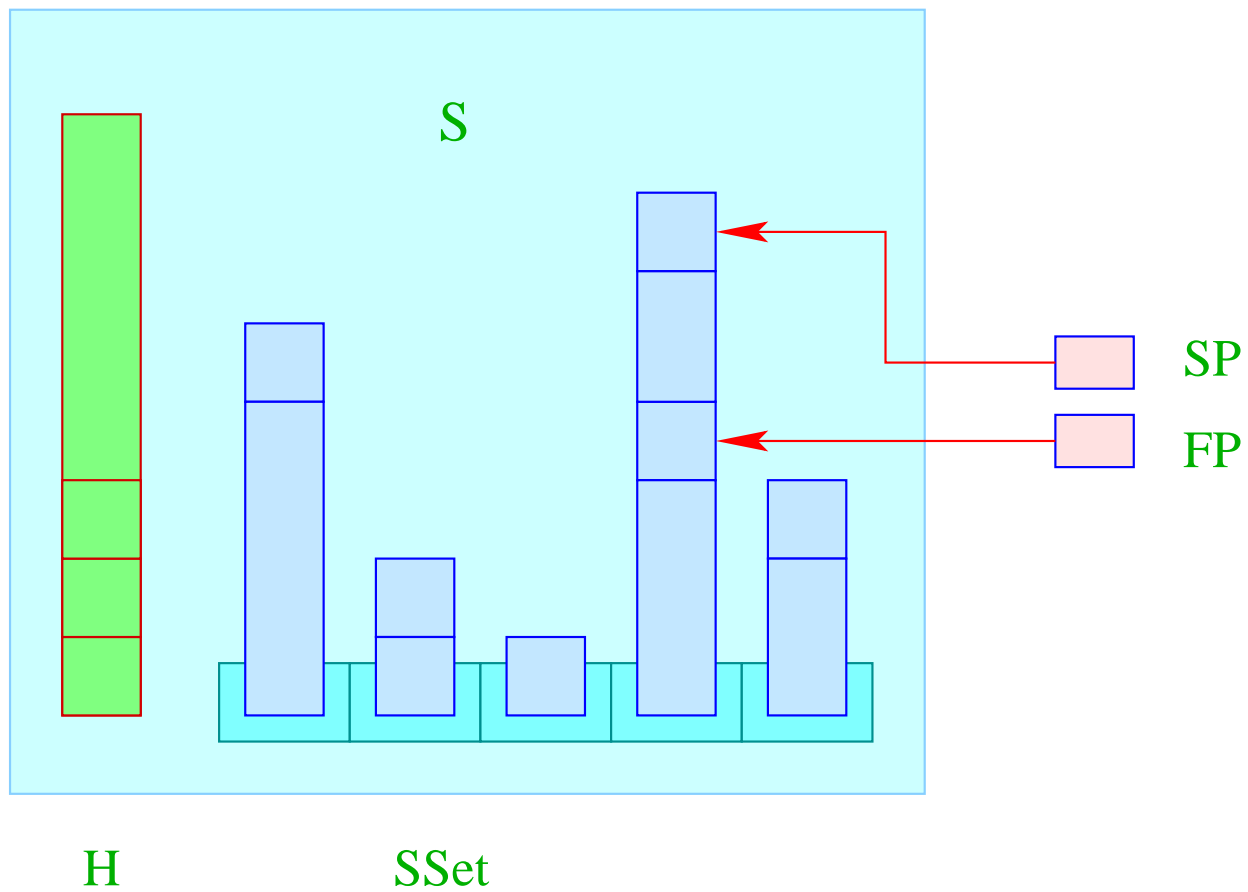


... similar to the **CMa**, we have:

- C** = **C**ode Store – contains the **CMa** program;
every cell contains one instruction;
- PC** = **P**rogram-**C**ounter – points to the next executable instruction;
- H** = **H**eap –
every cell may contain a base value or an address;
the **globals** are stored at the bottom;
- NP** = **N**ew-**P**ointer – points to the **first free** cell.

For a simplification, we assume that the heap is stored in a separate segment.
The function `malloc()` then fails whenever **NP** exceeds the topmost border.

Every thread on the other hand needs its **own stack**:



In contrast to the **CMa**, we have:

- SSet** = **Set** of **S**tacks – contains the stacks of the threads;
every cell may contain a base value of an address;
- S** = common address space for heap and the stacks;
- SP** = **S**tack-**P**ointer – points to the **current** topmost occupied stack cell;
- FP** = **F**rame-**P**ointer – points to the **current** stack frame.

Warning:

- If all references pointed into the heap, we could use separate address spaces for each stack.
Besides **SP** and **FP**, we would have to record the number of the current stack :-)
- In the case of **C**, though, we must assume that all storage regions live within the same address space — only at different locations :-)
SP and **FP** then uniquely identify storage locations.
- For simplicity, we omit the extreme-pointer **EP**.

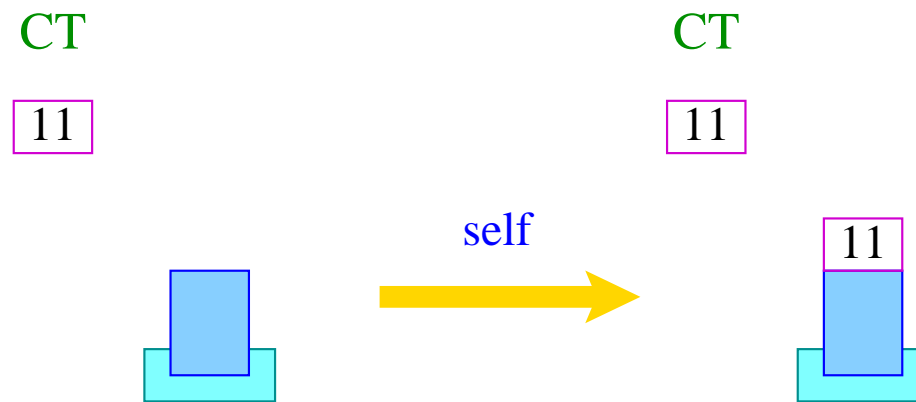
41 The Ready-Queue

Idea:

- Every thread has a unique number `tid`.
- A table `TTab` allows to determine for every `tid` the corresponding thread.
- At every point in time, there can be several `executable` threads, but only one `running` thread (per processor :-)
- the `tid` of the currently running thread is cept in the register `CT` (`C`urrent `T`hread).
- The function: `tid self ()` returns the `tid` of the current thread.
Accordingly:

`codeR self () ρ = self`

... where the instruction `self` pushes the content of the register `CT` onto the (current) stack:

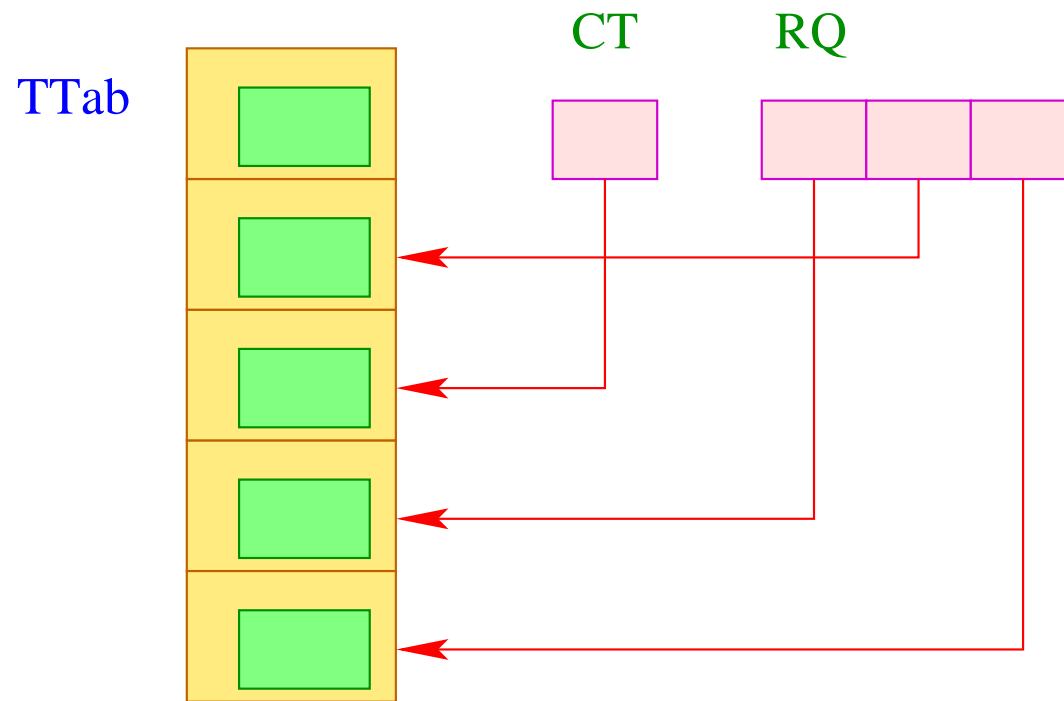


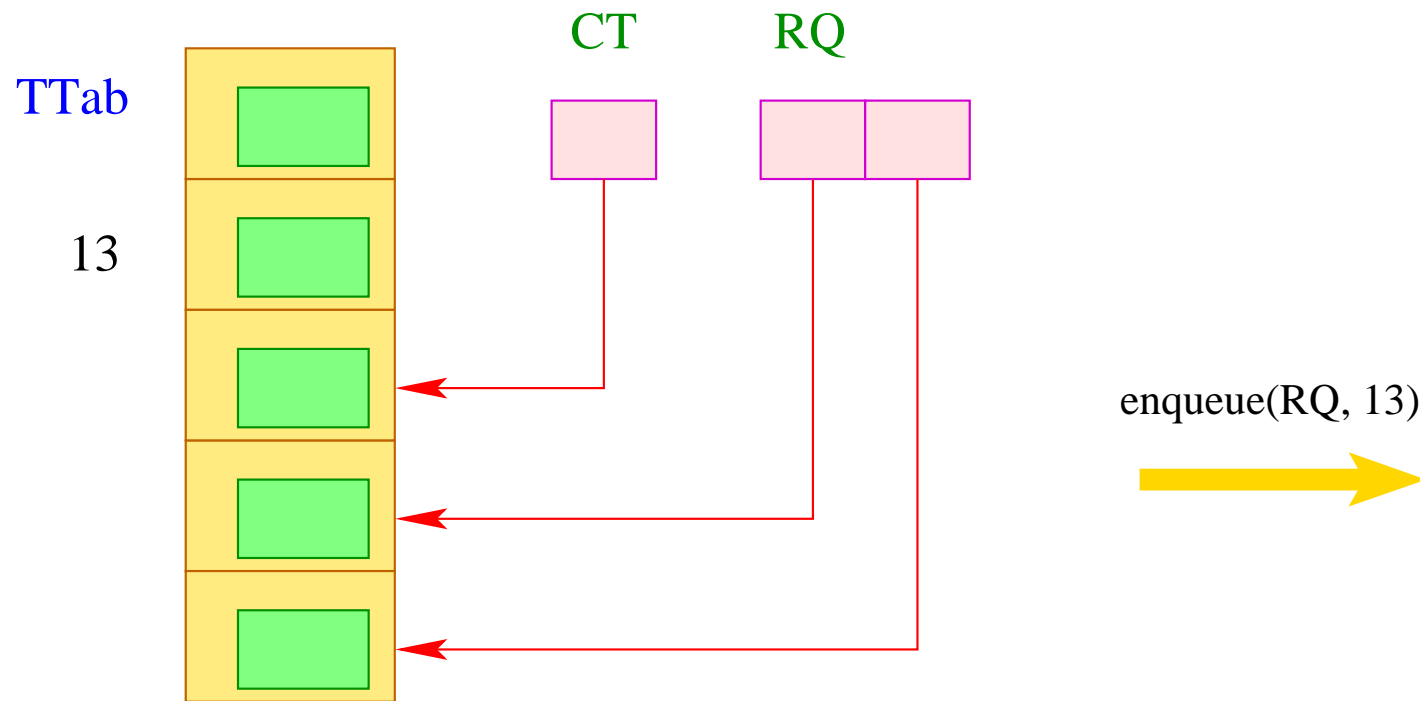
$S[SP++] = CT;$

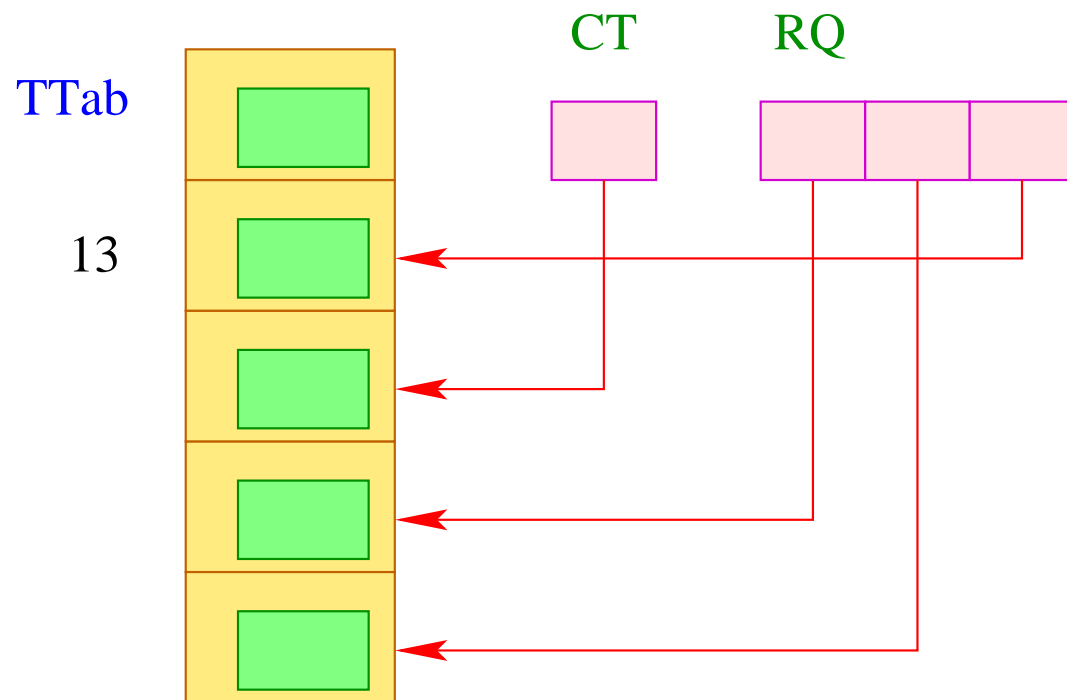
- The remaining executable threads (more precisely, their `tid`'s) are maintained in the queue `RQ` (`Ready-Queue`).
- For queues, we need the functions:

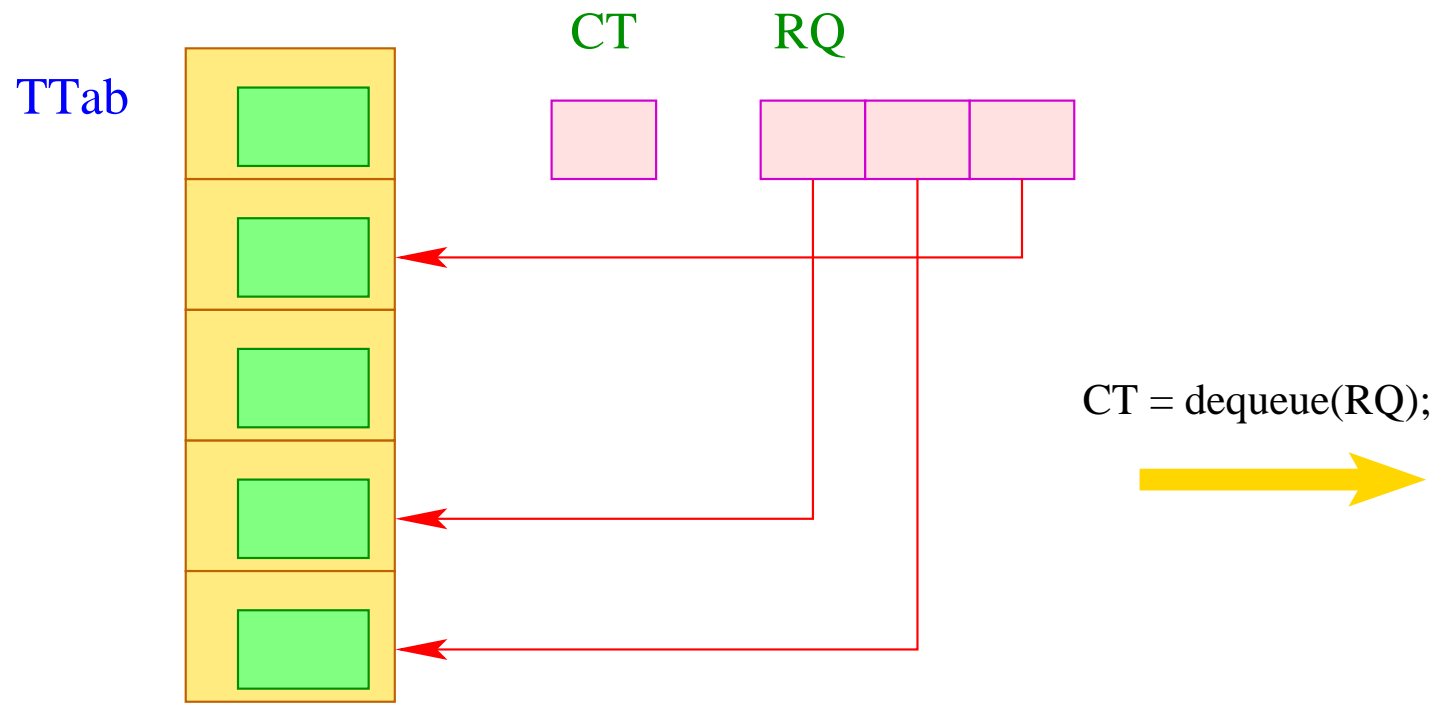
```
void enqueue (queue q, tid t),  
tid dequeue (queue q)
```

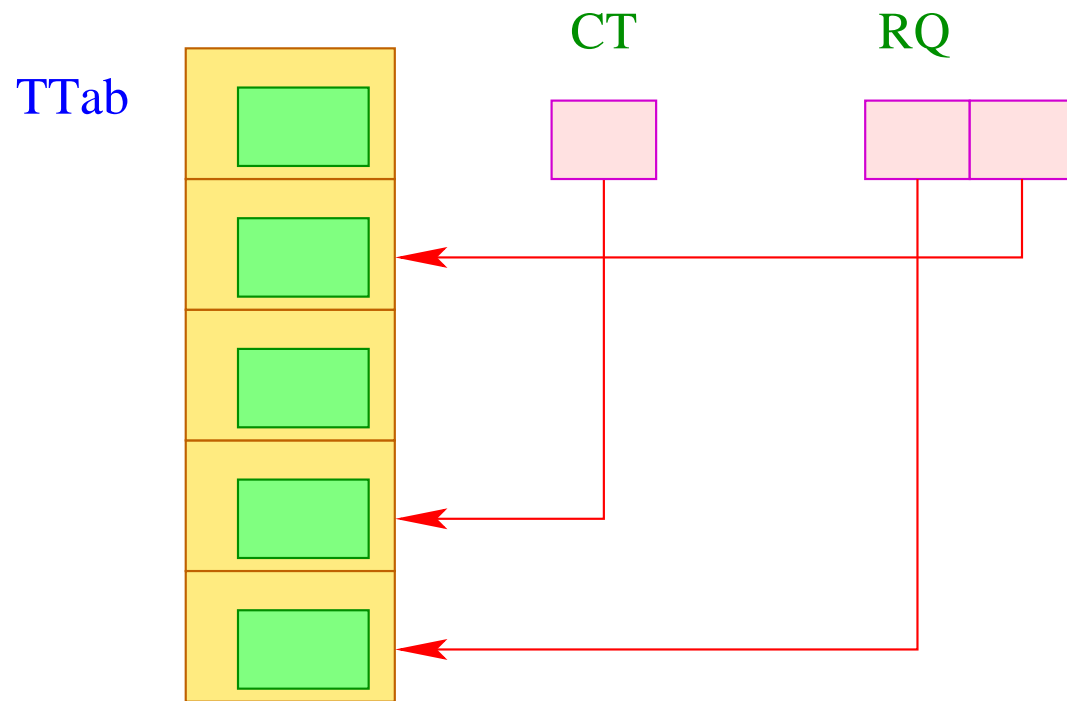
which insert a `tid` into a queue and return the first one, respectively ...





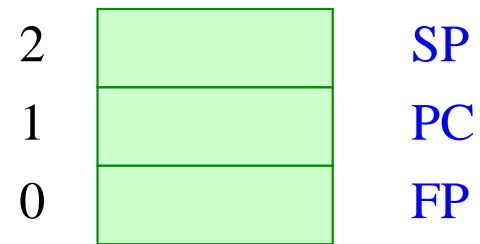






If a call to `dequeue ()` failed, it returns a value < 0 :-)

The thread table must contain for every thread, all information which is needed for its execution. In particular it consists of the registers **PC**, **SP** und **FP**:



Interrupting the current thread therefore requires to save these registers:

```
void save () {  
    TTab[CT][0] = FP;  
    TTab[CT][1] = PC;  
    TTab[CT][2] = SP;  
}
```

Analogously, we **restore** these registers by calling the function:

```
void restore () {  
    FP = TTab[CT][0];  
    PC = TTab[CT][1];  
    SP = TTab[CT][2];  
}
```

Thus, we can realize an instruction **yield** which causes a **thread-switch**:

```
tid ct = dequeue ( RQ );  
if (ct ≥ 0) {  
    save (); enqueue ( RQ, CT );  
    CT = ct;  
    restore ();  
}
```

Only if the ready-queue is **non-empty**, the current thread is replaced **:-)**

42 Switching between Threads

Problem:

We want to give each executable thread a fair chance to be completed.



- Every thread must former or later be scheduled for running.
- Every thread must former or later be interrupted.

Possible Strategies:

- Thread switch only at explicit calls to a function `yield()` :-(
• Thread switch after `every` instruction \implies too expensive :-(
• Thread switch after a `fixed number` of steps \implies we must install a counter and execute `yield` at dynamically chosen points :-(

We insert thread switches at selected program points ...

- at the **beginning** of function bodies;
- before every jump whose target does not exceed the current **PC** ...

⇒ rare :-))

The modified scheme for loops $s \equiv \mathbf{while} (e) s$ then yields:

```
code  $s \rho$  = A : codeR  $e \rho$   
                jumpz B  
                code  $s \rho$   
                yield  
                jump A  
B : ...
```


Note:

- **If-then-else**-Statements do not necessarily contain thread switches.
- **do-while**-Loops require a thread switch at the end of the condition.
- Every loop should contain (at least) one thread switch :-)
- Loop-Unrolling reduces the number of thread switches.
- At the translation of **switch**-statements, we created a jump table **behind** the code for the alternatives. Nonetheless, we can avoid thread switches here.
- At **freely programmed** uses of **jumpi** as well as **jumpz** we should also insert thread switches **before** the jump (or at the jump target).
- If we want to reduce the number of executed thread switches even further, we could switch threads, e.g., only at every 100th call of **yield** ...

43 Generating New Threads

We assume that the expression: $s \equiv \mathbf{create} (e_0, e_1)$ first evaluates the expressions e_i to the values f, a and then creates a new thread which computes $f(a)$.

If thread creation fails, s returns the value -1 .

Otherwise, s returns the new thread's **tid**.

Tasks of the Generated Code:

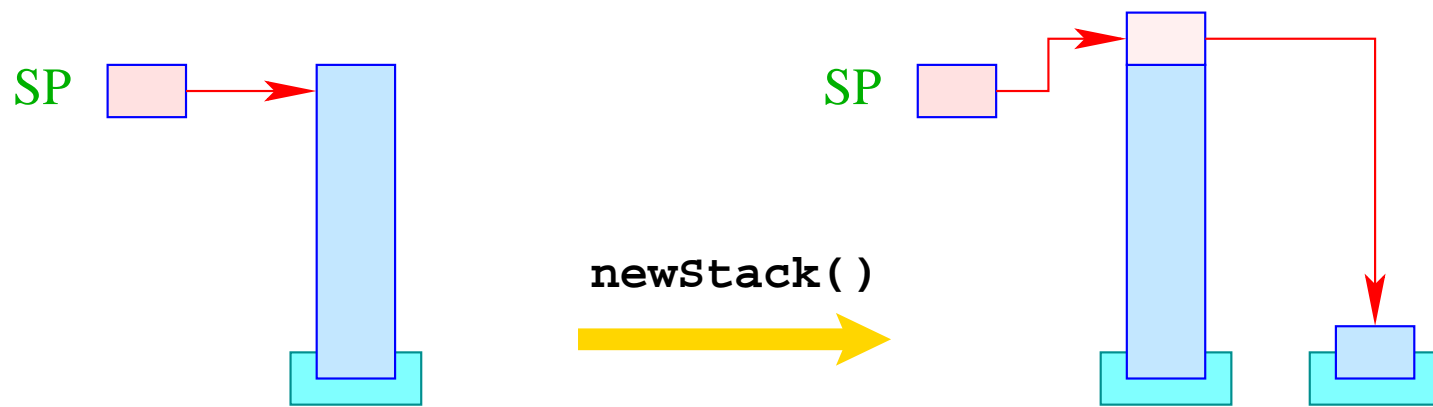
- Evaluation of the e_i ;
- Allocation of a new run-time stack together with a stack frame for the evaluation of $f(a)$;
- Generation of a new **tid**;
- Allocation of a new entry in the **TTab**;
- Insertion of the new **tid** into the ready-queue.

The translation of s then is quite simple:

$$\text{code}_R s \rho = \begin{array}{l} \text{code}_R e_0 \rho \\ \text{code}_R e_1 \rho \\ \text{initStack} \\ \text{initThread} \end{array}$$

where we assume the argument value occupies 1 cell :-)

For the implementation of `initStack` we need a run-time function `newStack()` which returns a pointer onto the first element of a new stack:



If the creation of a new stack fails, the value 0 is returned.



```

newStack();
if (S[SP]) {
    S[S[SP]+1] = -1;
    S[S[SP]+2] = f;
    S[S[SP]+3] = S[SP-1];
    S[SP-1] = S[SP]; SP--
}
else S[SP = SP - 2] = -1;

```

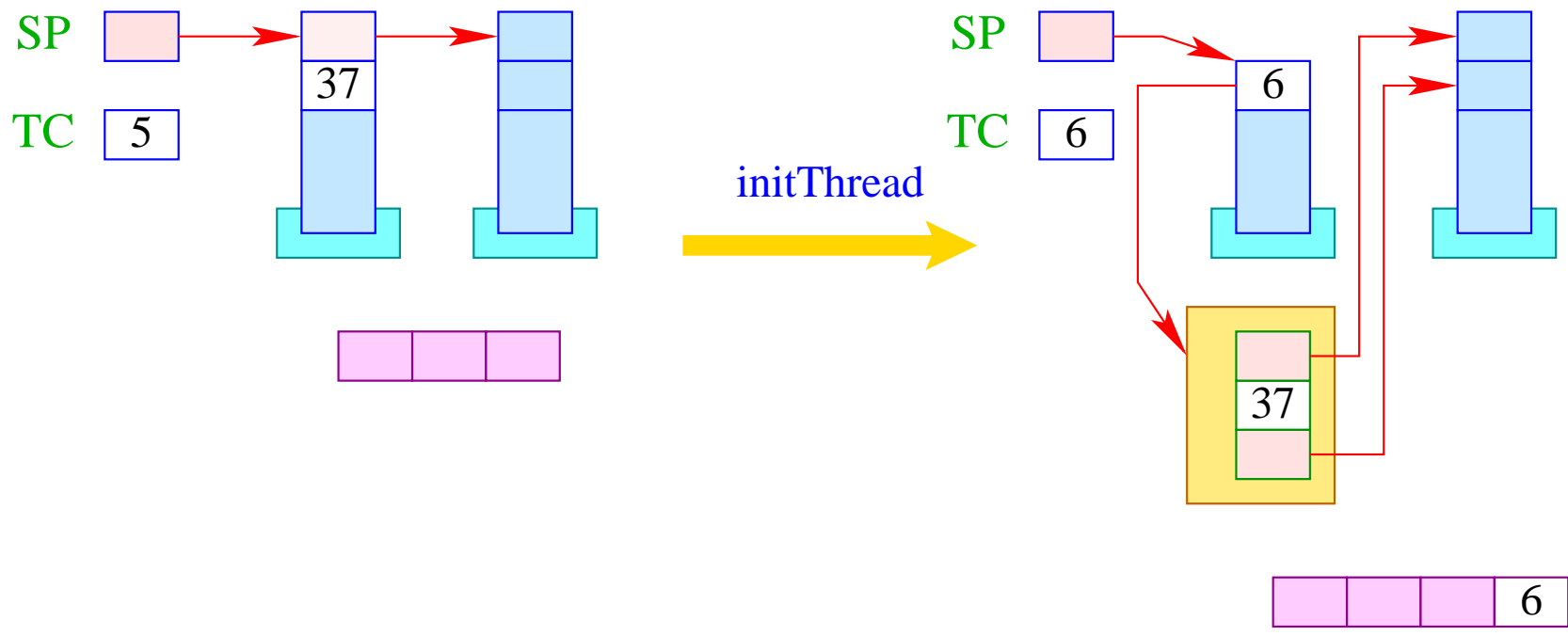
Note:

- The continuation address `f` points to the (fixed) code for the termination of threads.
- Inside the stack frame, we no longer allocate space for the `EP` \implies the return value has relative address -2 .
- The bottom stack frame can be identified through `FPold = -1` $:-)$

In order to create `new` thread ids, we introduce a new register `TC` (Thread Count).

Initially, `TC` has the value 0 (corresponds to the `tid` of the initial thread).

Before thread creation, `TC` is incremented by 1.



```
if (S[SP] ≥ 0) {  
    tid = ++TCount;  
    TTab[tid][0] = S[SP]-1;  
    TTab[tid][1] = S[SP-1];  
    TTab[tid][2] = S[SP];  
    S[--SP] = tid;  
    enqueue( RQ, tid );  
}
```


44 Terminating Threads

Termination of a thread (usually `:-`) returns a value. There are two (regular) ways to terminate a thread:

1. The initial function call has terminated. Then the return value is the return value of the call.
2. The thread executes the statement `exit (e)`; Then the return value equals the value of `e`.

Warning:

- We want to return the return value in the bottom stack cell.
- `exit` may occur arbitrarily deeply nested inside a recursion. Then we de-allocate all stack frames ...
- ... and jump to the terminal treatment of threads at address `f` .

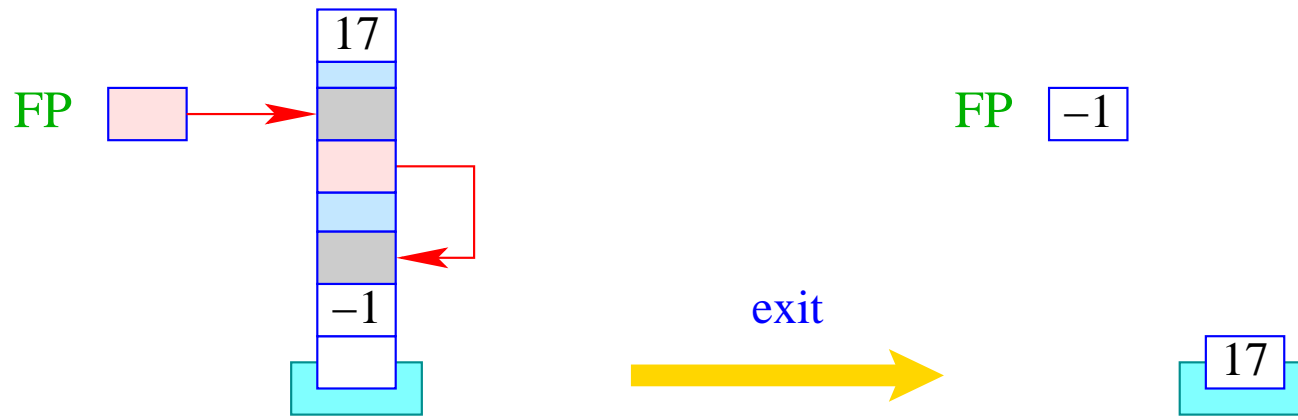
Therefore, we translate:

```
code exit (e); ρ = codeR e ρ
                  exit
                  term
                  next
```

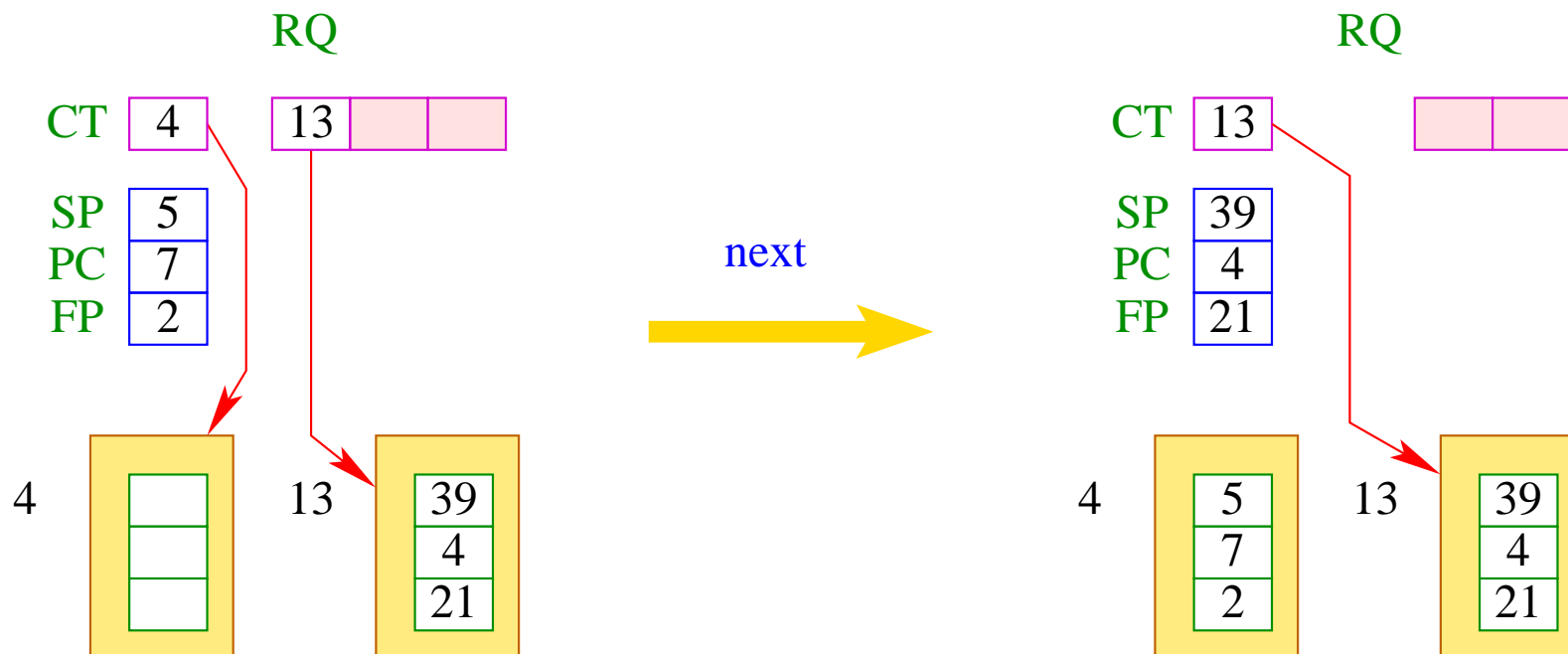
The instruction `term` is explained later :-)

The instruction `exit` successively pops all stack frames:

```
result = S[SP];
while (FP ≠ -1) {
    SP = FP-2;
    FP = S[FP-1];
}
S[SP] = result;
```



The instruction `next` activates the next executable thread:
in contrast to `yield` the current thread is **not** inserted into `RQ`.



Ist die Schlange **RQ** leer, wird zusätzlich If the queue **RQ** is empty, we additionally terminate the whole program:

```
if (0 > ct = dequeue( RQ )) halt;  
else {  
    save ();  
    CT = ct;  
    restore ();  
}
```

45 Waiting for Termination

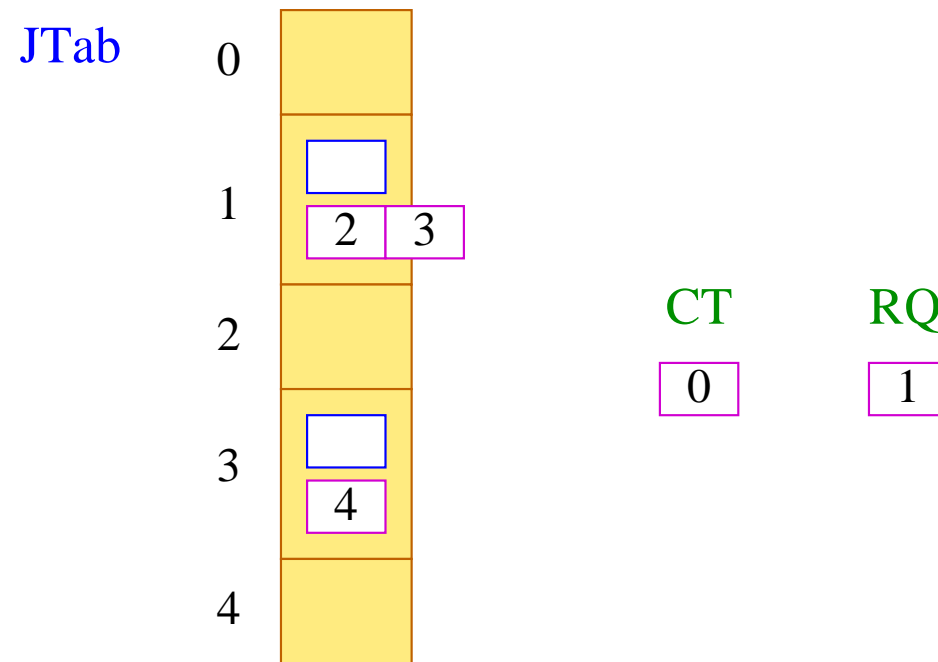
Occasionally, a thread may only continue with its execution, if some other thread has terminated. For that, we have the expression `join (e)` where we assume that e evaluates to a thread id `tid`.

- If the thread with the given `tid` is already terminated, we return its return value.
- If it is not yet terminated, we interrupt the current thread execution.
- We insert the current thread into the queue of threads already waiting for the termination.

We save the current registers and switch to the next executable thread.

- Thread waiting for termination are maintained in the table `JTab`.
- There, we also store the return values of threads `:-)`

Example:



Thread 0 is running, thread 1 could run, threads 2 and 3 wait for the termination of 1, and thread 4 waits for the termination of 3.

Thus, we translate:

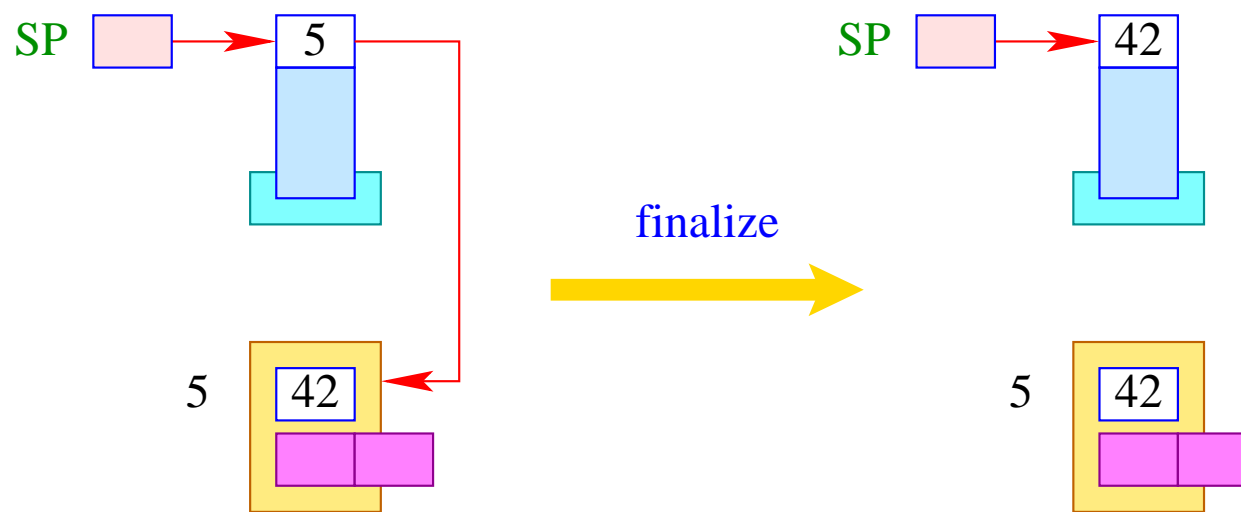
$$\text{code}_R \text{ join } (e) \rho = \text{code}_R e \rho$$

`join`
`finalize`

... where the instruction `join` is defined by:

```
tid = S[SP];
if (TTab[tid][1] ≥ 0) {
    enqueue ( JTab[tid], CT );
    next
}
```


... accordingly:



$S[SP] = JTab[tid][1];$

The instruction sequence:

`term`

`next`

is executed before a thread is terminated.

Therefore, we store them at the location `f`.

The instruction `next` switches to the next executable thread. Before that, though,

- ... the last stack frame must be popped and the result be stored in the table `JTab` ;
- ... the thread must be marked as terminated, e.g., by additionally setting the `PC` to `-1`;
- ... all threads must be notified which have waited for the termination.

For the instruction `term` this means:

```
PC = -1;  
JTab[CT][1] = S[SP];  
freeStack(SP);  
while (0 ≤ tid = dequeue ( JTab[CT][0] ))  
    enqueue ( RQ, tid );
```

The run-time function `freeStack (int adr)` removes the (one-element) stack at the location `adr` :



46 Mutual Exclusion

A **mutex** is an (abstract) datatype (in the heap) which should allow the programmer to dedicate exclusive access to a shared resource (**mutual exclusion**).

The datatype supports the following operations:

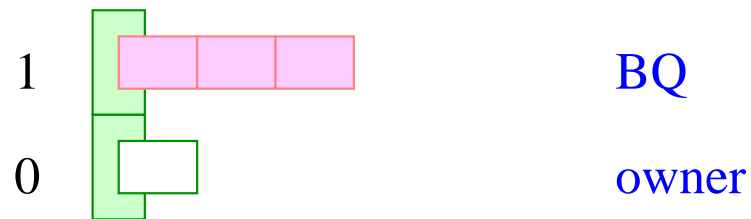
Mutex * newMutex ();	—	creates a new mutex;
void lock (Mutex *me);	—	tries to acquire the mutex;
void unlock (Mutex *me);	—	releases the mutex;

Warning:

A thread is only allowed to release a mutex if it has owned it beforehand :-)

A mutex me consists of:

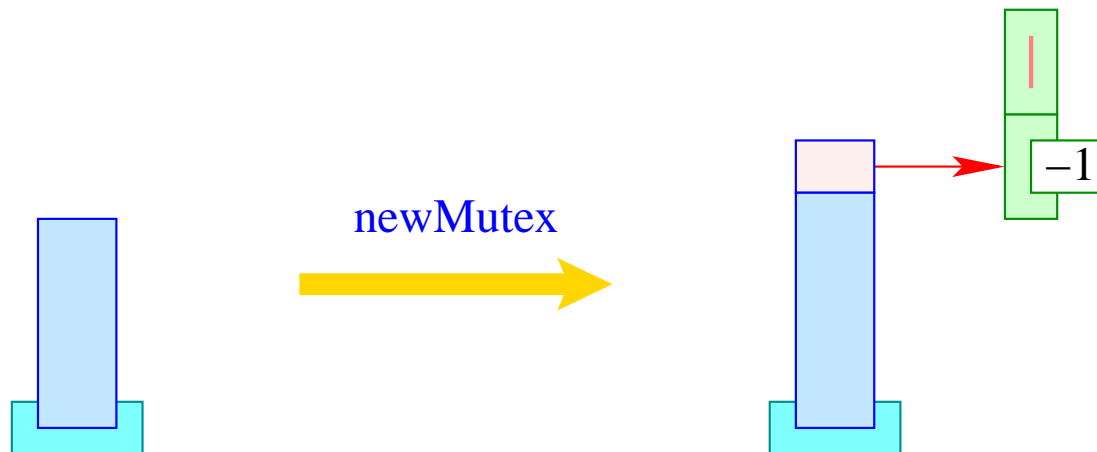
- the tid of the current owner (or -1 if there is no one);
- the queue **BQ** of **blocked** threads which want to acquire the mutex.



Then we translate:

$$\text{code}_R \text{ newMutex } () \rho = \text{newMutex}$$

where:

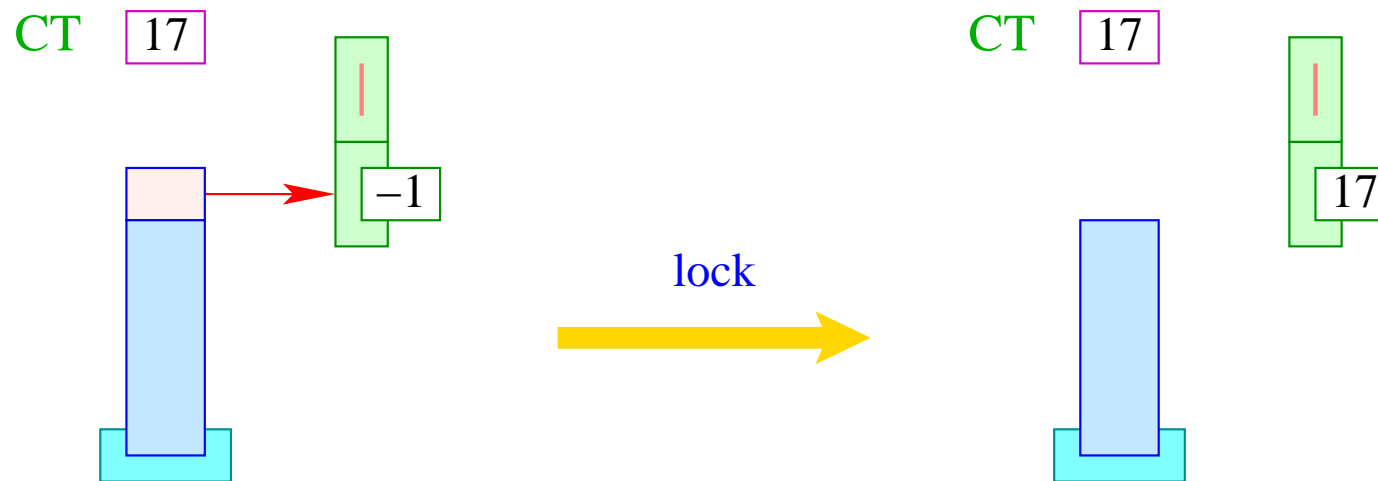


Then we translate:

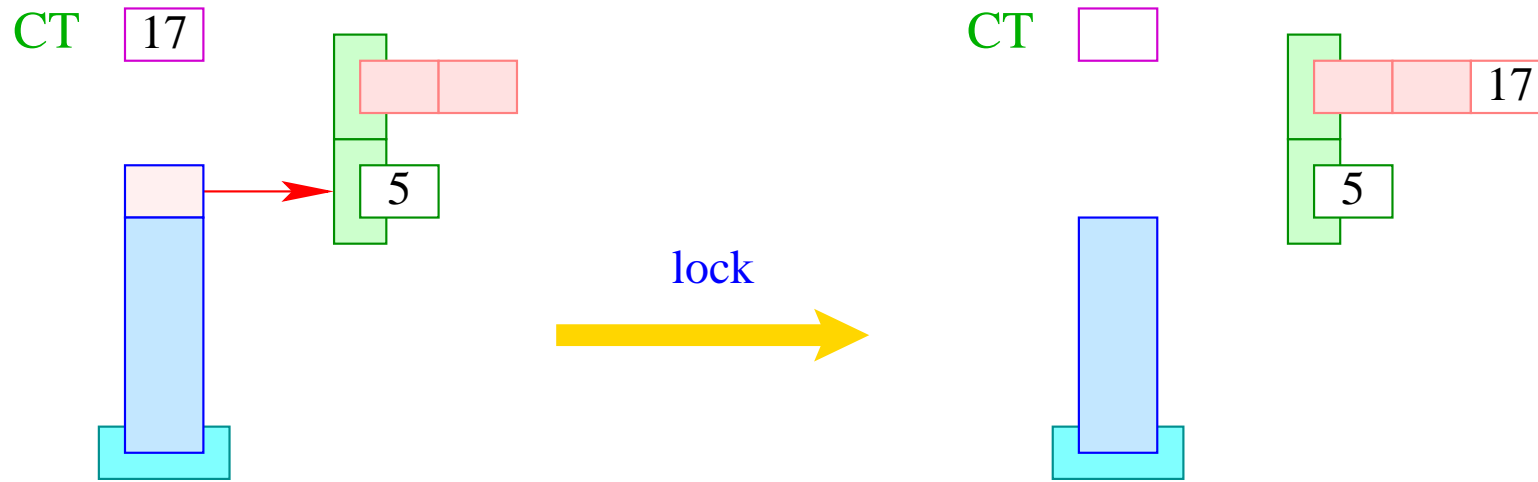
$$\text{code lock}(e); \rho = \text{code}_R e \rho$$

lock

where:



If the mutex is already owned by someone, the current thread is interrupted:

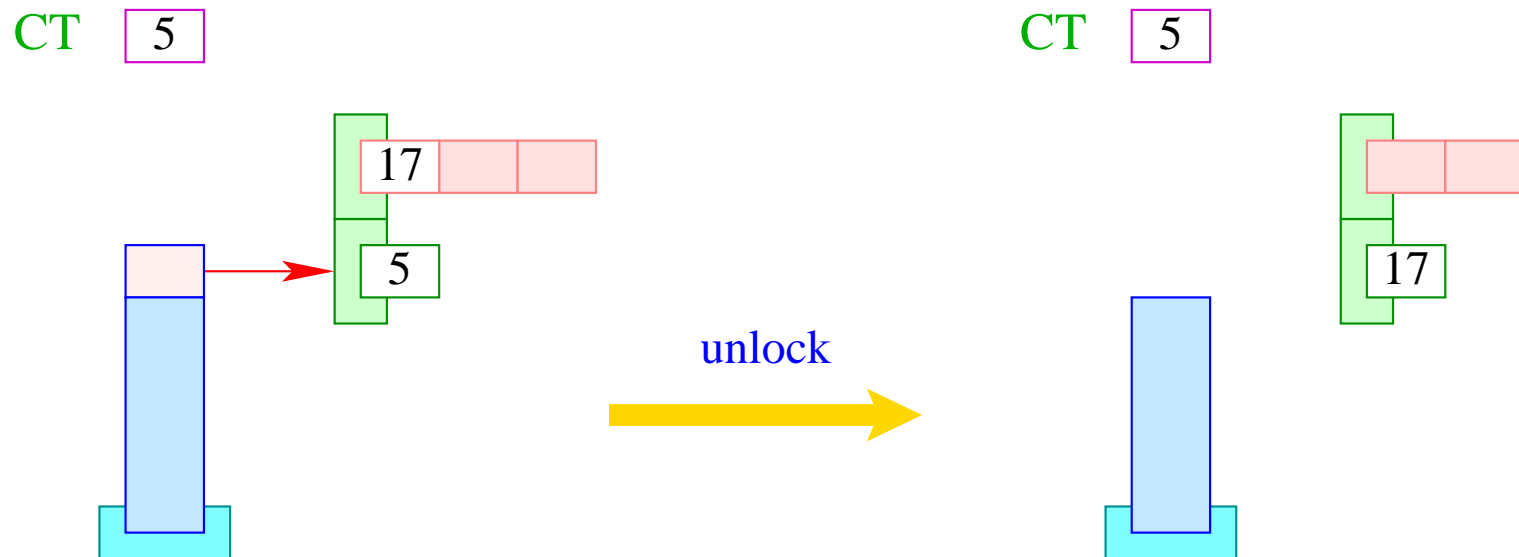


```
if (S[S[SP]] < 0) S[S[SP--]] = CT;  
else {  
    enqueue ( S[SP--]+1, CT );  
    next;  
}
```

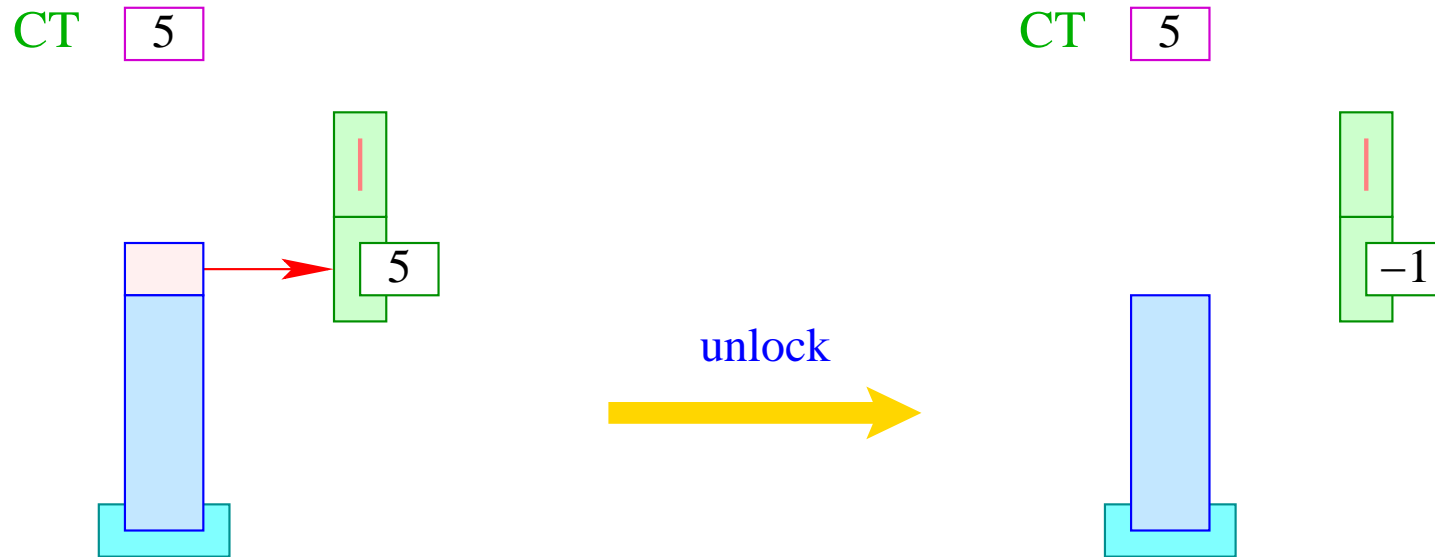

Accordingly, we translate:

$$\text{code unlock } (e); \rho = \text{code}_R e \rho \\ \text{unlock}$$

where:



If the queue **BQ** is empty, we release the mutex:



```

if (S[S[SP]]  $\neq$  CT) Error ("Illegal unlock!");
if (0 > tid = dequeue ( S[SP]+1)) S[S[SP--]] = -1;
else {
    S[S[SP--]] = tid;
    enqueue ( RQ, tid );
}

```

47 Waiting for Better Weather

It may happen that a thread owns a mutex but must wait until some extra condition is true.

Then we want the thread to remain in-active until it is told otherwise.

For that, we use **condition variables**. A condition variable consists of a queue **WQ** of waiting threads :-)



For condition variables, we introduce the functions:

CondVar * newCondVar ();	— creates a new condition variable;
void wait (CondVar * cv), Mutex * me);	— enqueues the current thread;
void signal (CondVar * cv);	— re-animates one waiting thread;
void broadcast (CondVar * cv);	— re-animates all waiting threads.

Then we translate:

$$\text{code}_R \text{ newCondVar } () \rho = \text{newCondVar}$$

where:

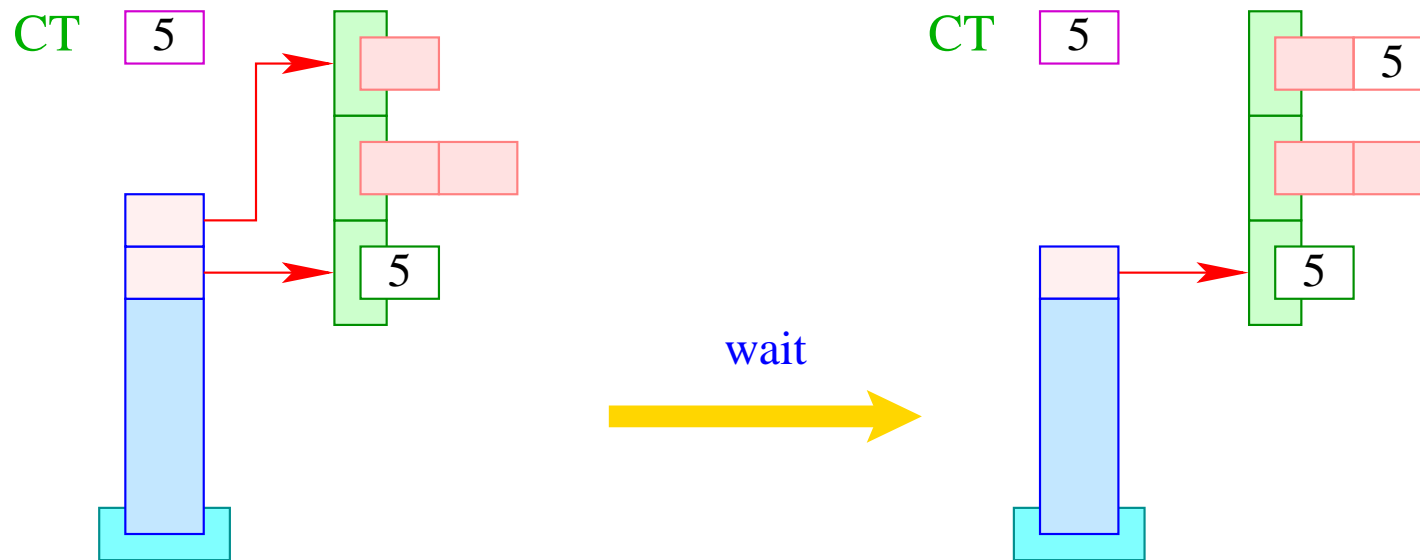


After enqueueing the current thread, we release the mutex. After re-animation, though, we must acquire the mutex again.

Therefore, we translate:

$$\text{code wait } (e_0, e_1); \rho = \begin{array}{l} \text{code}_R e_1 \rho \\ \text{code}_R e_0 \rho \\ \text{wait} \\ \text{dup} \\ \text{unlock} \\ \text{next} \\ \text{lock} \end{array}$$

where ...



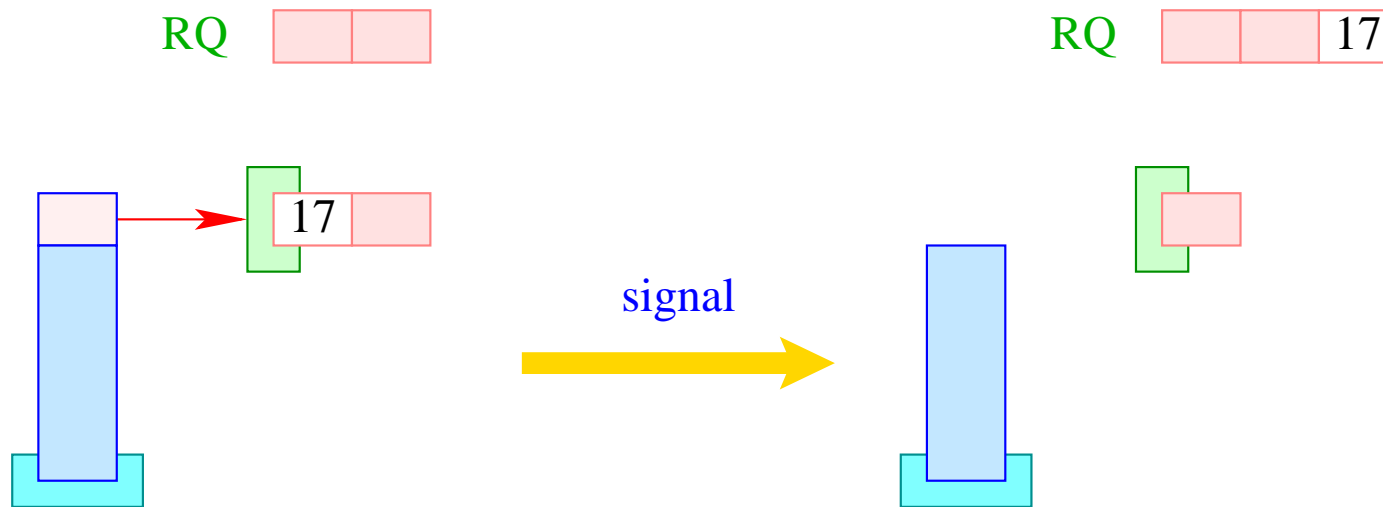
```

if (S[S[SP-1]] ≠ CT) Error ("Illegal wait!");
enqueue ( S[SP], CT ); SP--;

```

Accordingly, we translate:

`code signal (e); ρ = codeR e ρ`
`signal`



```
if (0 ≤ tid = dequeue ( S[SP]))  
    enqueue ( RQ, tid );  
SP--;
```


Analogously:

`code broadcast (e); ρ = codeR e ρ`
`broadcast`

where the instruction `broadcast` enqueues all threads from the queue `WQ` into the ready-queue `RQ` :

```
while (0 ≤ tid = dequeue ( S[SP]))
    enqueue ( RQ, tid );
SP--;
```

Warning:

The re-animated threads are not `blocked` !!!

When they become running, though, they first have to acquire their mutex :-)

48 Example: Semaphores

A semaphore is an abstract datatype which controls the access of a bounded number of (identical) resources.

Operations:

- `Sema * newSema (int n)` — creates a new semaphore;
- `void Up (Sema * s)` — increases the number of free resources;
- `void Down (Sema * s)` — decreases the number of available resources.

Therefore, a semaphore consists of:

- a **counter** of type **int**;
- a mutex for synchronizing the semaphore operations;
- a condition variable.

```
typedef struct {  
    Mutex * me;  
    CondVar * cv;  
    int count;  
} Sema;
```

```
Sema * newSema (int n) {  
    Sema * s;  
    s = (Sema *) malloc (sizeof (Sema));  
    s→me = newMutex ();  
    s→cv = newCondVar ();  
    s→count = n;  
    return (s);  
}
```

The translation of the body amounts to:

alloc 1	newMutex	newCondVar	loadr 1	loadr 2
loadc 3	loadr 2	loadr 2	loadr 2	storer -2
new	store	loadc 1	loadc 2	return
storer 2	pop	add	add	
pop		store	store	
		pop	pop	

The function `Down()` **decrements** the counter.

If the counter becomes negative, `wait` is called:

```
void Down (Sema * s) {  
    Mutex *me;  
    me = s→me;  
    lock (me);  
    s→count--;  
    if (s→count < 0) wait (s→cv,me);  
    unlock (me);  
}
```

The translation of the body amounts to:

alloc 1	loadc 2	add		loadc 1
loadr 1	add	store		add
load	load	loadc 0		load
storer 2	loadc 1	le		wait
lock	sub	jumpz A	A:	loadr 2
	loadr 1	loadr 2		unlock
loadr 1	loadc 2	loadr 1		return

The function `Up()` increments the counter again.

If it is afterwards **not yet positive**, there still must exist waiting threads. One of these is sent a signal:

```
void Up (Sema * s) {  
    Mutex *me;  
    me = s→me;  
    lock (me);  
    s→count++;  
    if (s→count ≤ 0) signal (s→cv);  
    unlock (me);  
}
```


The translation of the body amounts to:

alloc 1	loadc 2	add	loadc 1
loadr 1	add	store	add
load	load	loadc 0	load
storer 2	loadc 1	le	signal
lock	add	jumpz A	A: loadr 2
	loadr 1		unlock
loadr 1	loadc 2	loadr 1	return

49 Stack-Management

Problem:

- All threads live within the same storage.
- Every thread requires its own stack (at least conceptually).

1. Idea:

Allocate for each new thread a **fixed amount** of storage space.



Then we implement:

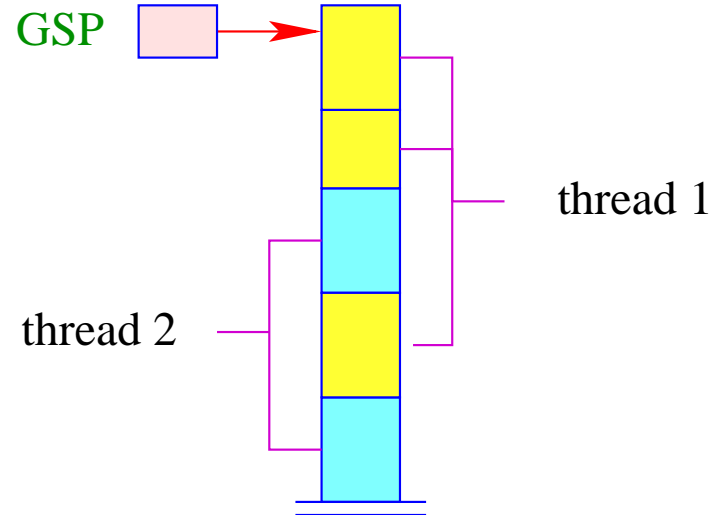
```
void *newStack() { return malloc(M); }  
void freeStack(void *adr) { free(adr); }
```

Problem:

- Some threads consume much, some only little stack space.
- The necessary space is statically typically unknown :-)

2. Idea:

- Maintain all stacks in one joint **Frame-Heap** **FH** :-)
- Take care that the space inside the stack frame is sufficient at least for the current function call.
- A global stack-pointer **GSP** points to the overall topmost stack cell ...



Allocation and de-allocation of a stack frame makes use of the run-time functions:

```
int newFrame(int size) {  
    int result = GSP;  
    GSP = GSP+size;  
    return result;  
}
```

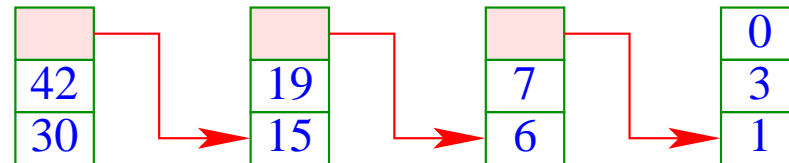
```
void freeFrame(int sp, int size);
```

Warning:

The de-allocated block may reside inside the stack :-)



We maintain a list of freed stack blocks :-)

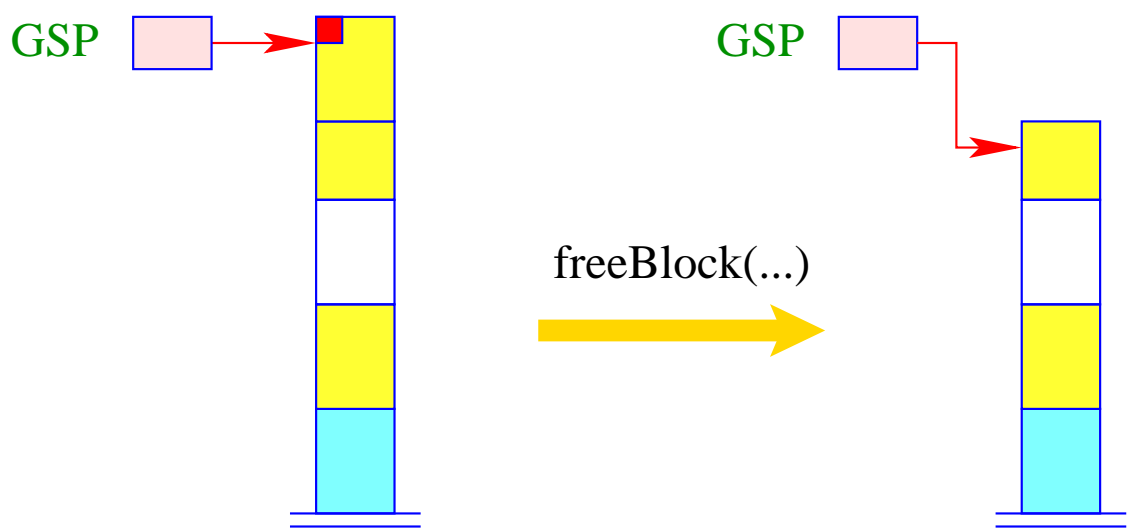


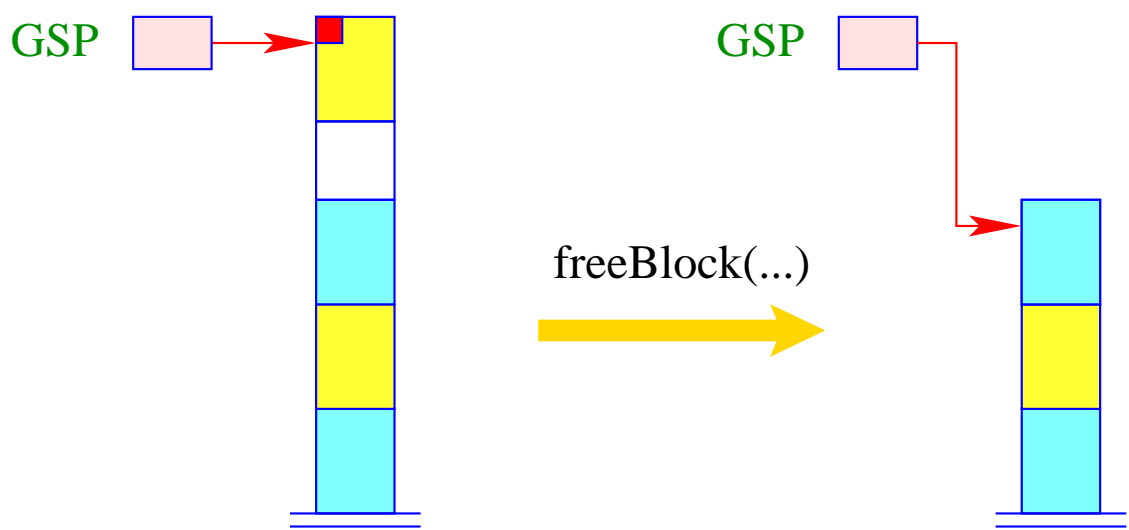
This list supports a function

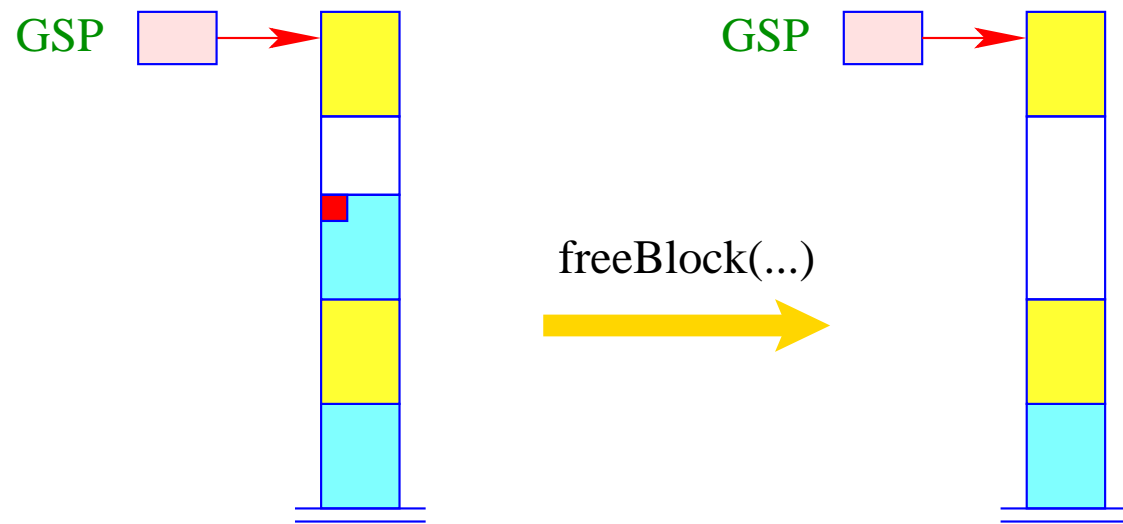
```
void insertBlock(int max, int min)
```

which allows to free single blocks.

- If the block is on top of the stack, we pop the stack immediately;
- ... together with the blocks below – given that these have already been marked as de-allocated.
- If the block is inside the stack, we merge it with neighbored free blocks:





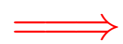


Approach:

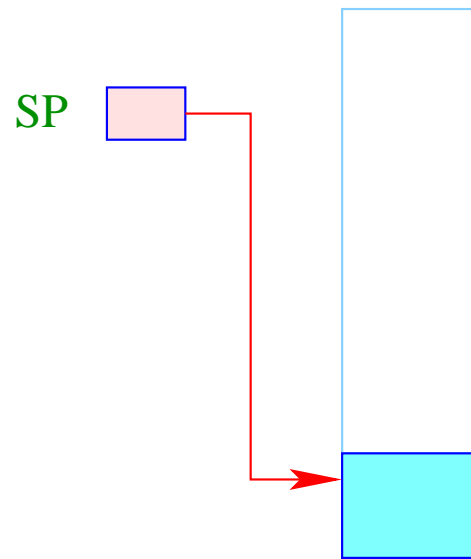
We allocate a fresh block for every function call ...

Problem:

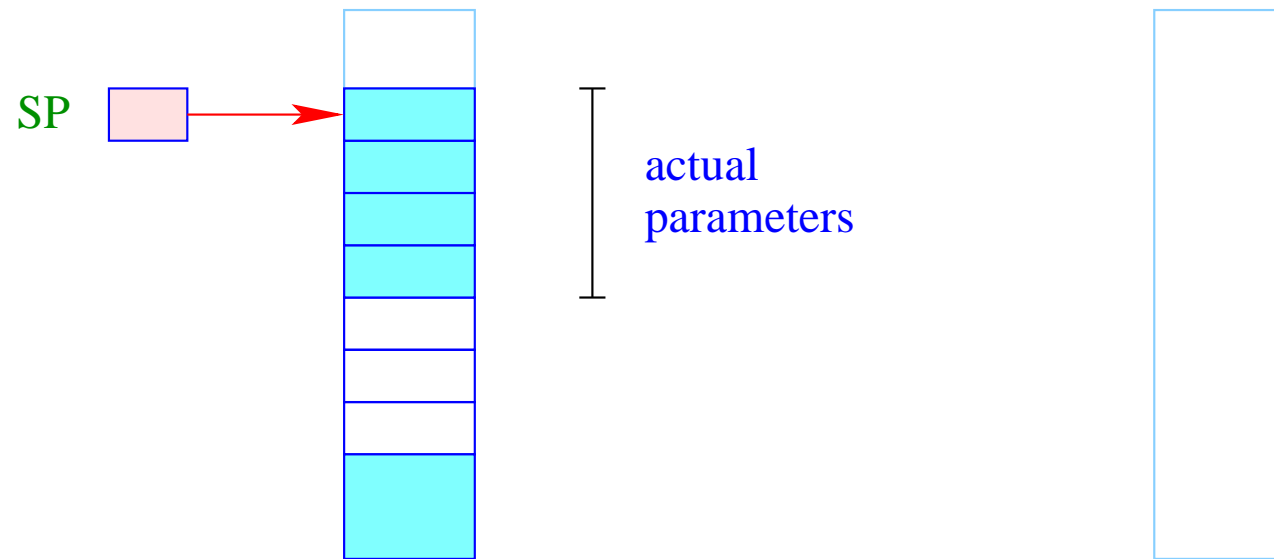
When ordering the block **before** the call, we do not yet know the space consumption of the called function :-)



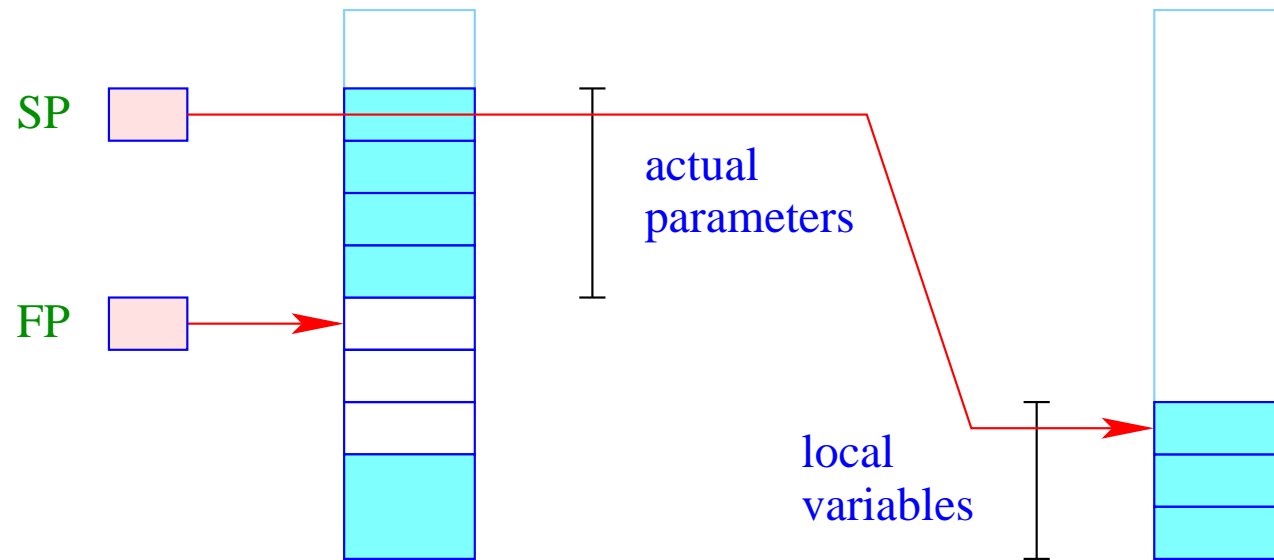
We order the new block **after** entering the function body!



Organisational cells as well as actual parameters must be allocated inside the old block ...



When entering the new function, we now allocate the new block ...



In particular, the **local** variables reside in the new block ...



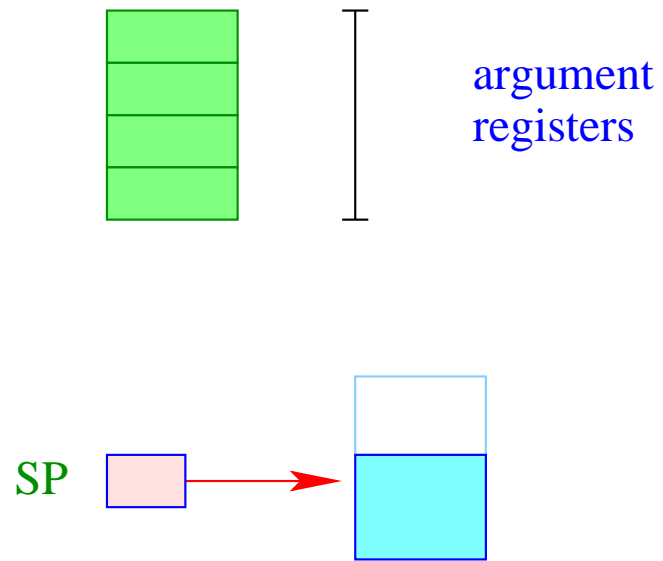
We address ...

- the formal parameters **relatively** to the frame-pointer;
- the local variables **relatively** to the stack-pointer :-)

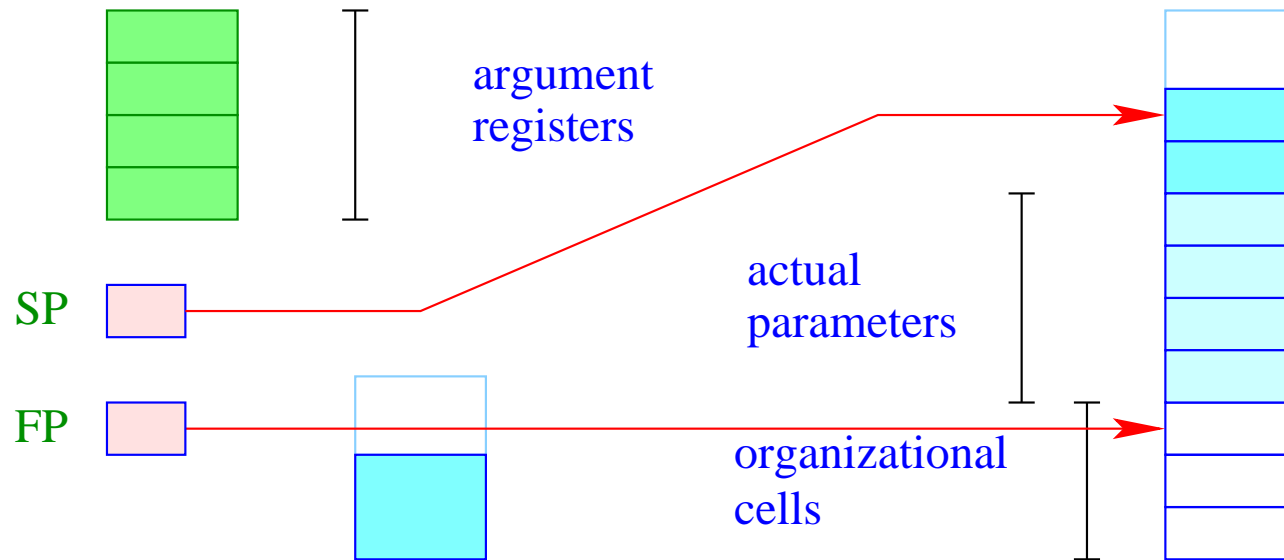


We must re-organize the complete code generation ... :-)

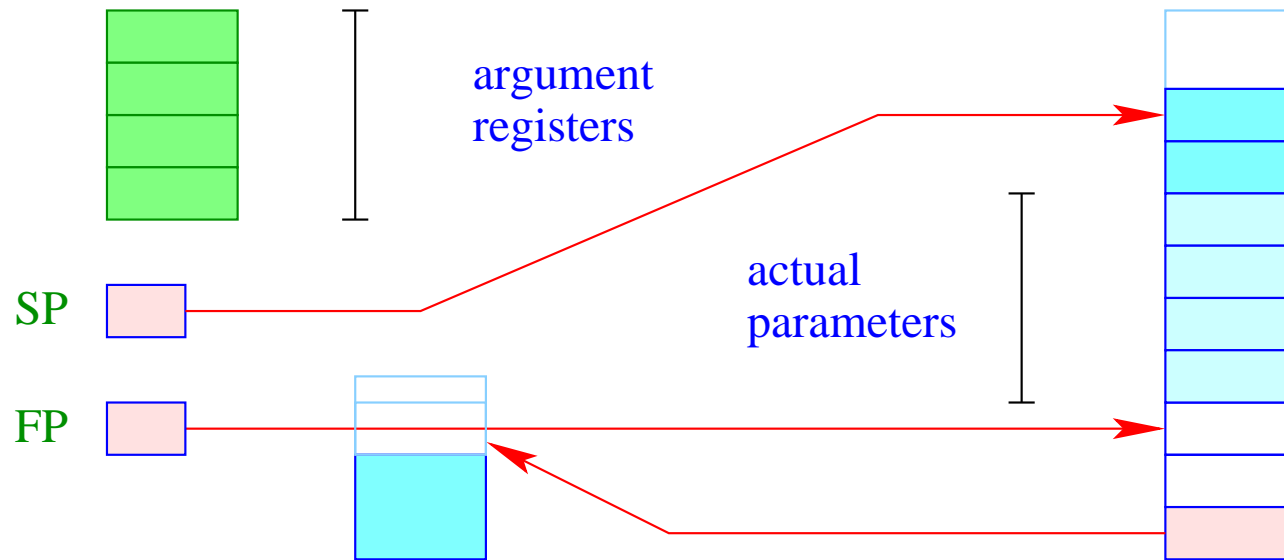
Alternative: Passing of parameters in registers ... :-)



The values of the actual parameters are determined **before** allocation of the new stack frame.



The **complete** frame is allocated inside the new block – plus the space for the current parameters.



Inside the new block, though, we must store the old **SP** (possibly +1) in order to correctly return the result ... :-)

3. Idea: Hybrid Solution

- For the first k threads, we allocate a separate stack area.
- For all further threads, we successively use one of the existing ones !!!



- For few threads extremely **simple** and **efficient**;
- For many threads **amortized** storage usage :-))