

# Vervolg: compilatie van een imperatieve taal

## Nadruk op de HL features:

- Control flow
- Datatypes, inclusief arrays en structs
- Procedures + parameters
- Scoping

## Centrale vraag:

**Waar vinden de VM instructies hun data?**

In het voorbeeld was dat triviaal omdat ook de doeltaal met symbolische variabelen werkte

# Compilatie van een imperatieve taal (Pascal) naar een VM (P-machine)

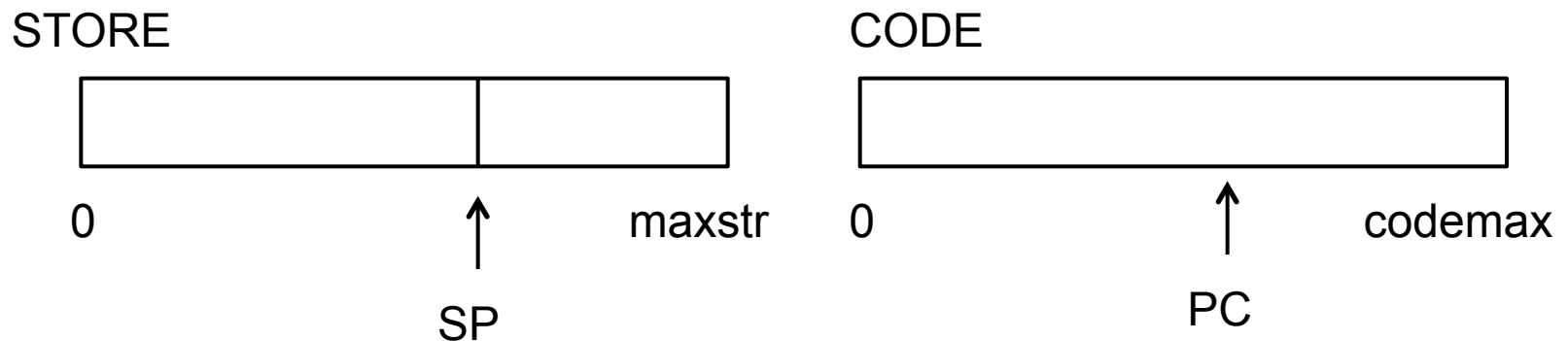
## Features van een imperatieve taal:

- Variabelen
- State = mapping van variabelen op waarden
- State wordt gewijzigd d.m.v. **assignments**
- Expressions worden geëvalueerd
- Variabelen hebben een scope
- Complexe datastructuren (arrays,...) + types
- Expliciete specificatie van control flow (if-then-else, while, ...)
- Procedures, met locale variabelen, parameters,...
- Variabelen kunnen meerdere incarnaties hebben (meerdere versies die tegelijk bestaan)

Architectuur van de P (Pascal) Machine – een VM op maat van de te vertalen taal. Vormt ook de definitie van de betekenis van die taal

Definitie: organisatie van het geheugen + instructies

Layout van het geheugen: een deel voor data, een deel voor code



Main loop

```
do
  PC := PC + 1
  execute instruction at location CODE[PC - 1]
od
```

PC eerst aanpassen maakt de behandeling van jumps wat eenvoudiger

# De instructies van de P-Machine

- P-machine instructies zijn macro's bestaande uit meer elementaire instructies voor de P-machine. We introduceren ze stapsgewijs, voor de verschillende instructies van Pascal.
- We maken een onderscheid tussen **L-values** (locaties, waaraan waarden toegekend kunnen worden) en **R-values** (e.g. waarde van een expressie). Afhankelijk van de context zal een variabele dus omgezet moeten worden in ofwel een L-value , ofwel een R-value.

Formeel wordt de vertaling gedefiniëerd met behulp van de *code*-functies; die zetten hun argument (een element van de brontaal) om in een rij P-machine instructies.

Van de code-functies definiëren we meerdere varianten, te gebruiken naargelang de context: *code<sub>L</sub>*, *code<sub>R</sub>* en *code*, voor het genereren van code die bij uitvoering ofwel een L-value of een R-value moet produceren, of voor het geval dat irrelevant is. Later hebben we nog een paar andere varianten nodig

We behandelen eerst expressies en assignments: we geven de eerst de relevante P-machine instructies en dan de code-functies.

De instructies:

1. Rekenkundige operatoren
2. Logische connectieven
3. Relationele operatoren
4. Transport van en naar het geheugen

Instructies voor  
evaluatie van  
expressies op de  
stack

i: integer  
r: real  
b: boolean  
a: adres  
N: numeriek type  
T: willekeurig type

Instr	Meaning	Cond.	Result
<b>add</b> <i>N</i>	$STORE[SP-1] := STORE[SP-1] +_N STORE[SP];$ $SP := SP-1$	( <i>N</i> , <i>N</i> )	( <i>N</i> )
<b>sub</b> <i>N</i>	$STORE[SP-1] := STORE[SP-1] -_N STORE[SP];$ $SP := SP-1$	( <i>N</i> , <i>N</i> )	( <i>N</i> )
<b>mul</b> <i>N</i>	$STORE[SP-1] := STORE[SP-1] *_N STORE[SP];$ $SP := SP-1$	( <i>N</i> , <i>N</i> )	( <i>N</i> )
<b>div</b> <i>N</i>	$STORE[SP-1] := STORE[SP-1] /_N STORE[SP];$ $SP := SP-1$	( <i>N</i> , <i>N</i> )	( <i>N</i> )
<b>neg</b> <i>N</i>	$STORE[SP] := -_N STORE[SP]$	( <i>N</i> )	( <i>N</i> )
<b>and</b>	$STORE[SP-1] := STORE[SP-1] \text{ and } STORE[SP];$ $SP := SP-1$	(b, b)	(b)
<b>or</b>	$STORE[SP-1] := STORE[SP-1] \text{ or } STORE[SP];$ $SP := SP-1$	(b, b)	(b)
<b>not</b>	$STORE[SP] := \text{not } STORE[SP]$	(b)	(b)
<b>equ</b> <i>T</i>	$STORE[SP-1] := STORE[SP-1] =_T STORE[SP];$ $SP := SP-1$	( <i>T</i> , <i>T</i> )	(b)
<b>geq</b> <i>I</i>	$STORE[SP-1] := STORE[SP-1] \geq_I STORE[SP];$ $SP := SP-1$	( <i>T</i> , <i>T</i> )	(b)
<b>leq</b> <i>T</i>	$STORE[SP-1] := STORE[SP-1] \leq_T STORE[SP];$ $SP := SP-1$	( <i>T</i> , <i>T</i> )	(b)
<b>les</b> <i>I</i>	$STORE[SP-1] := STORE[SP-1] <_I STORE[SP];$ $SP := SP-1$	( <i>T</i> , <i>T</i> )	(b)
<b>grt</b> <i>I</i>	$STORE[SP-1] := STORE[SP-1] >_I STORE[SP];$ $SP := SP-1$	( <i>T</i> , <i>T</i> )	(b)
<b>neq</b> <i>T</i>	$STORE[SP-1] := STORE[SP-1] \neq_T STORE[SP];$ $SP := SP-1$	( <i>T</i> , <i>T</i> )	(b)

Types van de  
elementen  
boven op de  
stack

(*N*,*N*): twee keer  
**hetzelfde**  
numerieke type



Instructions: **load** (van store naar stack)  
**store** (van stack naar store)

**Table 2.2** Store and load instructions. **ldo** loads from a location given by an absolute address, **ldc** loads a constant given in the instruction, **ind** loads indirectly using the highest stack location, **sro** stores in a location given by an absolute address, and **sto** stores in a location addressed by the second-highest location on the stack.

Instr.	Meaning	Cond.	Result
<b>ldo</b> $T\ q$	$SP := SP + 1;$ $STORE[SP] := STORE[q]$	$q \in [0, maxstr]$	$(T)$
<b>ldc</b> $T\ q$	$SP := SP + 1;$ $STORE[SP] := q$	$Type(q) = T$	$(T)$
<b>ind</b> $T$	$STORE[SP] := STORE[STORE[SP]]$	(a)	$(T)$
<b>sro</b> $T\ q$	$STORE[q] := STORE[SP];$ $SP := SP - 1$	$(T)$ $q \in [0, maxstr]$	
<b>sto</b> $T$	$STORE[STORE[SP - 1]] := STORE[SP];$ $SP := SP - 2$	(a, $T$ )	

## Definitie van de *code*-functie(s) voor de assignment instructie.

Omdat *code* code moet produceren waarin variabelen omgezet zijn in adressen, heeft die functie niet 1 maar **2** parameters:

- het te vertalen programmadeel
- een **environment**  $p$ ; dit is een functie met voor elke variabele het corresponderende (relatieve) adres.  $p$  komt overeen met de symbooltabel, maar dan wel uitgebreid (type + adres, voor elke variabele)

Schema voor assignment  $x := y$  voor simpele variabelen  $x$  en  $y$  :

Bereken L-value van $x$
Bereken R-value van $y$
<b>sto</b> $i$

## Code voor expressies en assignment

Uitvoering plaatst het  
resultaat op de stack, en  
laat de stack voor de  
rest ongewijzigd

Notatie: schrijf  $f \ x \ y$   
ipv  $f(x,y)$

**Table 2.3** The compilation of assignments.

Function	Condition
$code_R(e_1 = e_2) \rho$	$= code_R \ e_1 \ \rho; \ code_R \ e_2 \ \rho; \text{equ } T$
$code_R(e_1 \neq e_2) \rho$	$= code_R \ e_1 \ \rho; \ code_R \ e_2 \ \rho; \text{neq } T$
$\vdots$	
$code_R(e_1 + e_2) \rho$	$= code_R \ e_1 \ \rho; \ code_R \ e_2 \ \rho; \text{add } N$
$code_R(e_1 - e_2) \rho$	$= code_R \ e_1 \ \rho; \ code_R \ e_2 \ \rho; \text{sub } N$
$code_R(e_1 * e_2) \rho$	$= code_R \ e_1 \ \rho; \ code_R \ e_2 \ \rho; \text{mul } N$
$code_R(e_1 / e_2) \rho$	$= code_R \ e_1 \ \rho; \ code_R \ e_2 \ \rho; \text{div } N$
$code_R(-e) \rho$	$= code_R \ e \ \rho; \text{neg } N$
$code_R \ x \ \rho$	$= code_L \ x \ \rho; \text{ind } T$
$code_R \ c \ \rho$	$= \text{ldc } T \ c$
$code(x := e) \rho$	$= code_L \ x \ \rho; \ code_R \ e \ \rho; \text{sto } T$
$code_L \ x \ \rho$	$= \text{ldc } a \ \rho(x)$

code<sub>L</sub> voor x  
code<sub>R</sub> voor e

Stel dat we in een programma met integer variabelen  $a, b, c$  het assignment  $a := (b + (b \times c))$  willen vertalen, en dat voor  $\rho$  geldt:  $\rho(a) = 5$ ,  $\rho(b) = 6$ ,  $\rho(c) = 7$ . Dan is de vertaling dus het volgende:

Stel dat we in een programma met integer variabelen  $a, b, c$  het assignment  $a := (b + (b \times c))$  willen vertalen, en dat voor  $\rho$  geldt:  $\rho(a) = 5$ ,  $\rho(b) = 6$ ,  $\rho(c) = 7$ . Dan is de vertaling dus het volgende:

```

code( $a := (b + (b * c))$ )  $\rho$ 
= codeL  $a$   $\rho$ ; codeR( $b + (b * c)$ )  $\rho$ ; sto  $i$ 
= ldc  $a$   $5$ ; codeR( $b + (b * c)$ )  $\rho$ ; sto  $i$ 
= ldc  $a$   $5$ ; codeR( $b$ )  $\rho$ ; codeR( $b * c$ )  $\rho$ ; add  $i$ ; sto  $i$ 
= ldc  $a$   $5$ ; ldc  $a$   $6$ ; ind  $i$ ; codeR( $b * c$ )  $\rho$ ; add  $i$ ; sto  $i$ 
= ldc  $a$   $5$ ; ldc  $a$   $6$ ; ind  $i$ ; codeR( $b$ )  $\rho$ ; codeR( $c$ )  $\rho$ ; mul  $i$ ; add  $i$ ; sto  $i$ 
= ldc  $a$   $5$ ; ldc  $a$   $6$ ; ind  $i$ ; ldc  $a$   $6$ ; ind  $i$ ; codeR( $c$ )  $\rho$ ; mul  $i$ ; add  $i$ ; sto  $i$ 
= ldc  $a$   $5$ ; ldc  $a$   $6$ ; ind  $i$ ; ldc  $a$   $6$ ; ind  $i$ ; ldc  $a$   $7$ ; ind  $i$ ; mul  $i$ ; add  $i$ ; sto  $i$ 

```

□

Control flow:

Conditionele en iteratieve statements, rijen van statements

**Table 2.4** Conditional and unconditional branch.

Instr	Meaning	Comments	Cond.	Result
<b>ujp</b> $q$	$PC := q$	Unconditional branch	$q \in [0, codemax]$	
<b>fjp</b> $q$	if $STORE[SP] = false$ then $PC := q$ fi $SP := SP - 1$	Conditional branch	(b) $q \in [0, codemax]$	

Control flow:

Conditionele en iteratieve statements, rijen van statements

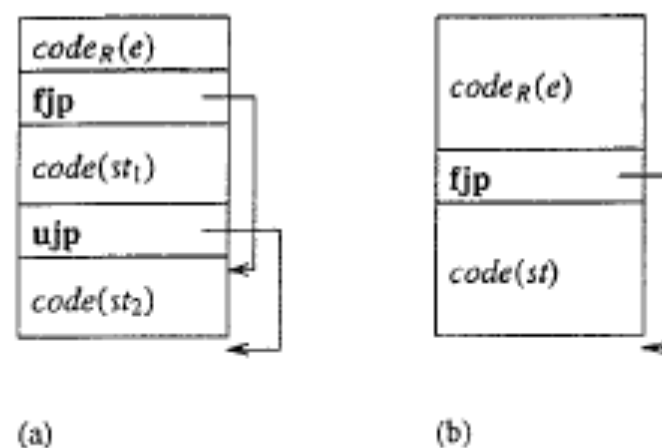
Table 2.4 Conditional and unconditional branch.

Instr	Meaning	Comments	Cond.	Result
<b>ujp</b> $q$	$PC := q$	Unconditional branch	$q \in [0, codemax]$	
<b>fjp</b> $q$	if $STORE[SP] = false$ then $PC := q$ fi $SP := SP - 1$	Conditional branch	(b) $q \in [0, codemax]$	

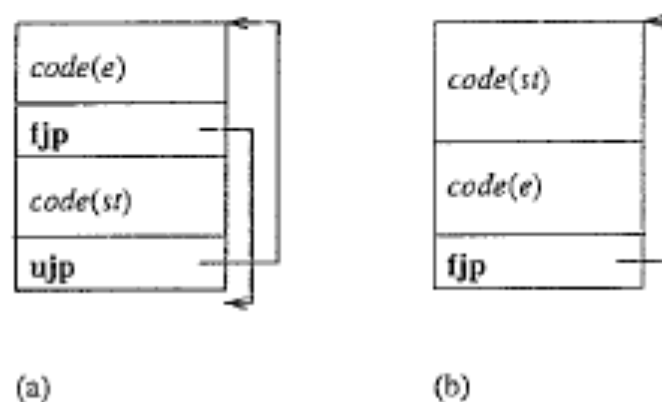
*code* (if  $e$  then  $st_1$  else  $st_2$  fi)  $p =$

*code* <sub>$R$</sub>   $e$   $p$ ; **fjp**  $l_1$ ; *code*  $st_1$   $p$ ; **ujp**  $l_2$ ;  $l_1$ : *code*  $st_2$   $p$ ;  $l_2$ :

$l_1, l_2$  : labels



**Figure 2.2** Instruction sequences for conditional statements. (a) Bidirectional, (b) unidirectional.



**Figure 2.3** Code generation for loops, (a) **while**  $e$  do  $st$  do and (b) **repeat**  $st$  until  $e$ .



# Switch / case statement

case e of

0:  $st_0$  ;

1:  $st_1$  ;

...

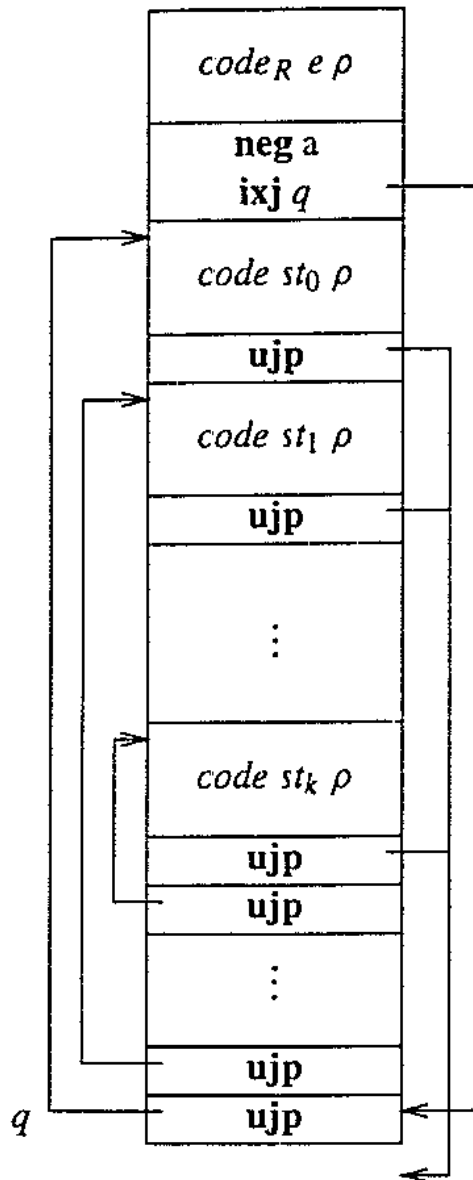
k:  $st_k$  ;

end

de waarde van e bepaalt welke  
van de  $st_i$  gekozen wordt

P-instructie: **ixj**

Instr	Meaning	Cond.	Result
<b>ixj</b> $q$	$PC := STORE[SP] + q;$ $SP := SP - 1$	(i)	



leaves the value of the selector  
on the top of the stack

indexed branch to the end of the  
branching table

branch to the end of the statement

Om k opeenvolgende doel-adressen  
te hebben voor de jumps

branching table

Om de parameter  $q$  te kunnen invullen in dit schema hebben we de lengte van de stukken  $\text{code } st_i$   $\rho$  nodig. Die is pas na hun generatie bekend, dus die  $q$  wordt pas achteraf ingevuld.

Er is ook nog een check nodig op de range waarin de waarde van  $e$  mag liggen.

Vorm: **check**  $0 \leq k \leq q$

**check 0 k q**

**dup**

dupliceer de top van de stack

**ldc 0**

**geq a**

**fjp r**

**dup**

**ldc k**

**leq a**

**fjp r**

**ixj q**

**r: pop**

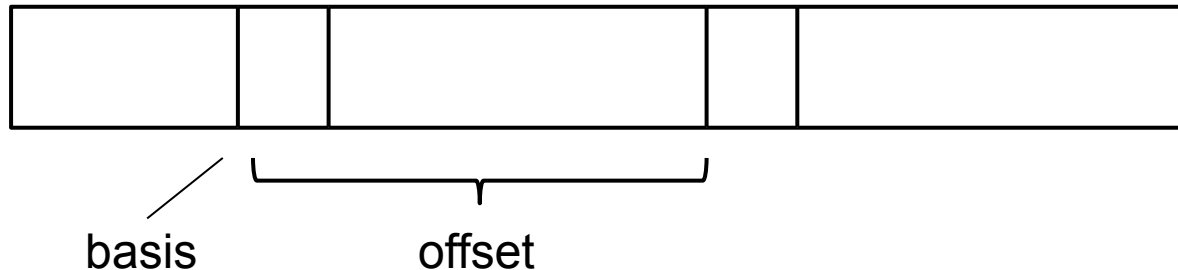
**ldc k**

k gebruikt als default

**ixj q**

# Geheugenallocatie: variabelen van een simpel type

$p$  moet voor elke varabele een numeriek adres geven. Gewoonlijk is dat een relatief adres, dus een offset t.o.v. een basisadres.  $p$  moet uiteraard “at compile time” berekend worden, de adressen worden gebruikt bij uitvoering, dus “at runtime”.



We zullen aannemen dat locaties 0...4 gereserveerd zijn (zie verder: layout van de stack frames)

# Geheugenallocatie: variabelen van een simpel type

Voor simpele types: assigneer aaneensluitende geheugenplaatsen

Voor een lijst van declaraties `var  $n_1:t_1$ ; ...  $n_k:t_k$ ;`

Definiëer  $\rho$  door:  $\rho(n_i) = i + 4$  of beter, rekening houdend met de grootte van de types:

$$\rho(n_i) = \begin{cases} 5 & \text{als } i = 1 \\ 5 + \sum_{j=1}^{i-1} \text{size}(t_j) & \text{anders} \end{cases}$$

waar  $\text{size}(t_j)$  het aantal plaatsen is nodig voor de voorstelling van een element van type  $t_j$ .

# Geheugenallocatie: statische arrays

Probleem: bij gebruik van een component, bv `a[e]`, is de waarde van `e` maar bekend bij uitvoering. Ze kan bovendien variëren, bv omdat die `a[e]` een lus voorkomt:

```
while ... do
    ...
    x := a[e]
    ...
od
```

Het adres waarnaar `a[e]` verwijst is **dynamische** informatie, m.a.w. informatie die slechts bekend is **at runtime** (anders spreekt men van **statische** informatie, die is bekend **at compile time**)

# Geheugenallocatie: statische arrays

Voorbeeld: `var a : array[-5..5,1..9] of integer`

Afspraak: we gebruiken **row-major** ordening: de elementen  
Van de array staan in de volgorde

$a[-5,1], a[-5,2], \dots, a[-5,9],$   
 $a[-4,1], a[-4,2], \dots, a[-4,9],$   
 $\dots$   
 $a[5,1], a[5,2], \dots, a[5,9]$

m.a.w. de rijen na elkaar. Anders gezegd, als we de elementen doorlopen varieert de meest rechtse index het snelst.



# Geheugenallocatie: statische arrays

Algemeen geval, met k dimensies:

- `var b : array [u1, ... , o1, ... , uk, ... , ok] of integer`
- dus op `b[i1, ..., ij, oj+1, ..., ok]` volgt `b[i1, ..., ij+1, uj+1, ..., uk]`
- aantal elementen in dimensie i :  $(o_i - u_i + 1)$
- `size(t)` for dit array type:  $\prod_1^k (o_i - u_i + 1)$  (als integer size 1 heeft)
- Voor een rij van declaraties: zoals voordien,

$$\rho(n_i) = \begin{cases} 5 & \text{als } i = 1 \\ 5 + \sum_{j=1}^{i-1} \text{size}(t_j) & \text{anders} \end{cases}$$

# Geheugenallocatie: statische arrays

Voorbeeld: var a : array[-5..5,1..9] of integer

a[-5,1], a[-5,2], ... , a[-5,9],  
a[-4,1], a[-4,2], ... , a[-4,9],  
...  
a[5,1], a[5,2], ... , a[5,9]

Het element a[-2,6] staat op plaats (we tellen vanaf 0)

$$(-2 - (-5)) \times (9 - 1 + 1) + 6 - 1 = (-2 + 5) \times 9 + 6 - 1 = 32$$



aantal elementen in een rij (= aantal kolommen)

# Geheugenallocatie: statische arrays

Toekenning van 0 aan de eerste component van een array a:

ldc a p(a);

ldc i 0;

sto i;

$p(a)$  is het adres van het **eerste** element

waar  $p(a)$  bekend is bij compilatie.

Maar wat met  $a[i,j]:=0$  ? De waarden  $\hat{i}, \hat{j}$  van  $i$  en  $j$  zijn alleen bekend bij uitvoering. Die waarden zijn dynamische informatie. We moeten dus instructies genereren die at runtime de waarde

$$r = (\hat{i} - (-5)) \times (9 - 1 + 1) + \hat{j} - 1 = (\hat{i} + 5) \times 9 + \hat{j} - 1$$

berekenen. De component staat dan op adres  $p(a) + r$

# Geheugenallocatie: statische arrays

Het algemene geval:

$$r = (\hat{I}_1 - u_1) \times \text{size}(\text{array}[u_2, \dots, o_2, \dots, u_k, \dots, o_k] \text{ of integer}) + \\ (\hat{I}_2 - u_2) \times \text{size}(\text{array}[u_3, \dots, o_3, \dots, u_k, \dots, o_k] \text{ of integer}) +$$

...

$$(\hat{I}_{k-1} - u_{k-1}) \times \text{size}(\text{array}[u_k, \dots, o_k] \text{ of integer}) + \\ (\hat{I}_k - u_k)$$

laat  $d_i = o_i - u_i + 1$  (range in de i-de dimensie)

We kunnen  $r$  opsplitsen in een deel dat bekend is bij compilatie, en de rest.

$$\begin{aligned}
 r = & (\bar{t}_1 \times d_2 \times d_3 \times \dots \times d_k + \bar{t}_2 \times d_3 \times d_4 \times \dots \times d_k + \\
 & \dots + \bar{t}_{k-1} \times d_k + \bar{t}_k) \\
 & - (u_1 \times d_2 \times d_3 \times \dots \times d_k + u_2 \times d_3 \times d_4 \times \dots \times d_k + \\
 & \dots + u_{k-1} \times d_k + u_k)
 \end{aligned}$$

Het deel dat bekend is bij compilatie (hier het tweede dus) kan uiteraard dan ook al berekend worden, en hoeft niet herberekend te worden elke keer een component gebruikt wordt (zie verder, dynamische arrays).

## Aangepaste P-machine instructie

Instr.	Meaning	Cond.	Results
<b>ixa</b> $q$	$STORE[SP - 1] := STORE[SP - 1] +$ $STORE[SP] * q$ $SP := SP - 1$	(a,i)	(a)

Berekening van geïndexeerd adres.  $STORE[SP-1]$  bevat een start adres,  $STORE[SP]$  bevat de index van de geselecteerde subarray en  $q$  bevat de grootte van de subarray. In het eerste voorbeeld dus 9.

## Aangepaste P-machine instructies

Instr.	Meaning	Cond.	Result
<b>inc</b> $T\ q$	$STORE[SP] := STORE[SP] + q$	$(T)$ and type $(q) = i$	$(T)$
<b>dec</b> $T\ q$	$STORE[SP] := STORE[SP] - q$	$(T)$ and type $(q) = i$	$(T)$

Om de rij indices te compileren gebruiken we de functie `codel` met als tweede parameter de grootte van de componenten,  $g$ . Het startadres van array  $c$  is  $\rho(c)$ .

Verder is  $d^{(j)} = d_{j+1} \times d_{j+2} \times \dots \times d_k$  (met  $d^{(k)} = 1$ ) en  
 $d = u_1 \times d^{(1)} + u_2 \times d^{(2)} + \dots + u_k \times d^{(k)}$

$$\begin{aligned} \text{code}_L\ c[i_1, \dots, i_k]\ \rho &= \text{ldc}\ a\ \rho(c); \text{code}_I\ [i_1, \dots, i_k]\ g\ \rho \\ \text{code}_I\ [i_1, \dots, i_k]\ g\ \rho &= \text{code}_R\ i_1\ \rho; \text{ixa}\ g \cdot d^{(1)}; \\ &\quad \text{code}_R\ i_2\ \rho; \text{ixa}\ g \cdot d^{(2)}; \\ &\quad \vdots \\ &\quad \text{code}_R\ i_k\ \rho; \text{ixa}\ g \cdot d^{(k)}; \\ &\quad \text{dec}\ a\ g \cdot d; \end{aligned}$$



Met “rangecheck”: test of de waarden van de indices binnen de grenzen blijven.

$$\begin{aligned}
 code_I [i_1, \dots, i_k] \text{ arr } \rho &= code_R i_1 \rho; \mathbf{chk} \ u_1 \ o_1; \mathbf{ixa} \ g \cdot d^{(1)}; \\
 &\quad code_R i_2 \rho; \mathbf{chk} \ u_2 \ o_2; \mathbf{ixa} \ g \cdot d^{(2)}; \\
 &\quad \vdots \\
 &\quad code_R i_k \rho; \mathbf{chk} \ u_k \ o_k; \mathbf{ixa} \ g \cdot d^{(k)}; \\
 &\quad \mathbf{dec} \ a \ g \cdot d;
 \end{aligned}$$

where  $arr = (g; u_1, o_1, \dots, u_n, o_n)$ .

$g$ : componentgrootte       $u_i, o_i$ : onder- en bovengrens

Instr.	Meaning	Cond.	Result
<b>chk</b> $p \ q$	<b>if</b> ( $STORE[SP] < p$ ) <b>or</b> ( $STORE[SP] > q$ ) <b>then</b> <i>error</i> ('value out of range') <b>fi</b>	(i,i)	(i)

Voorbeeld:     var i,j: integer;  
                  a: array[-5..5,1..9] of integer  
waar  $\rho(i)=5, \rho(j)=6$  en  $\rho(a)=7$

*Code* (a[i+1,j]:=0)  $\rho =$

ldc a 7	
ldc a 5	
ind i 5	
ldc i 1	
add i	
chk -5 5	
ixa 9	
ldc a 6	dec a -44
ind i	ldc i 0
chk 1 9	sto i
ixa 1	