

PROCESADORES VECTORIALES SEGMENTADOS

Víctor Viñals Yúfera y Jesús Alastruey Benedé

Multiprocesadores

Curso 2019-20

3º Grado en Ingeniería Informática

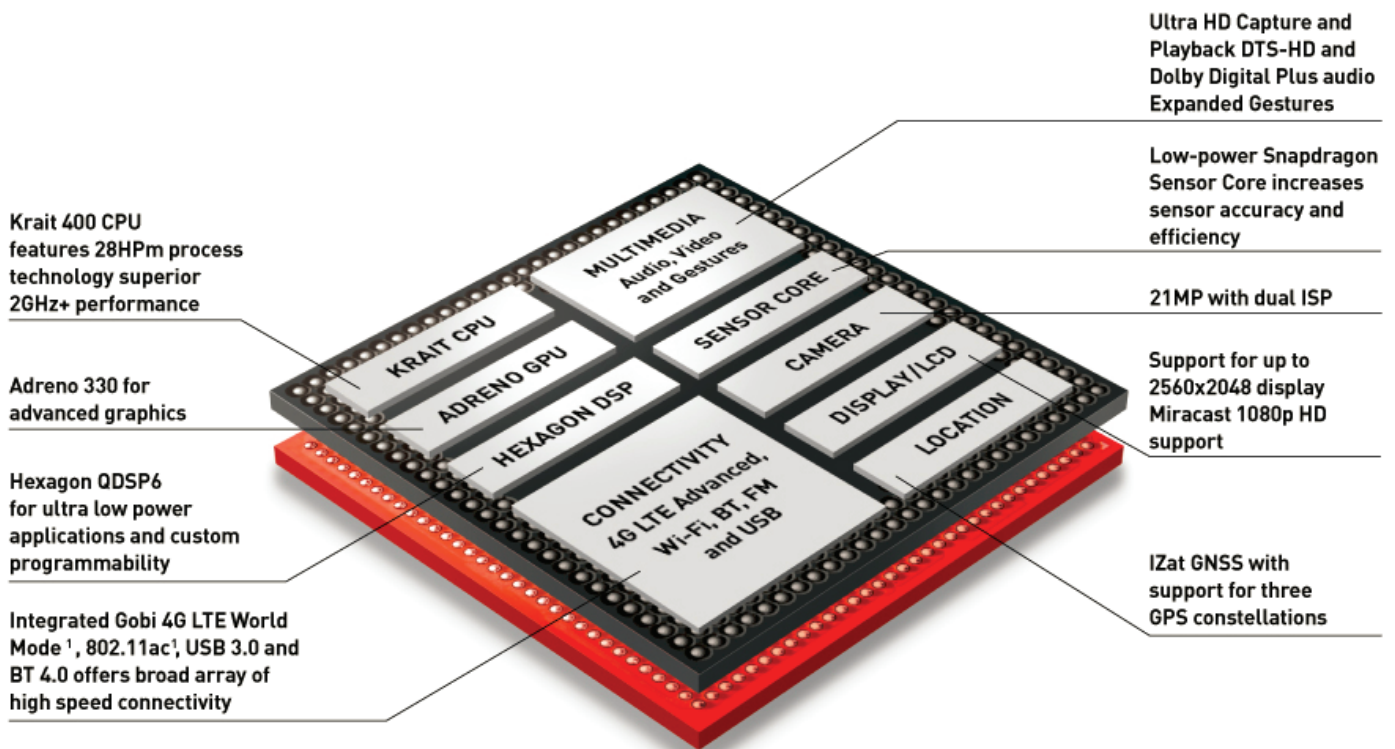
Especialidad Ingeniería de Computadores

- ❑ Basado en los apuntes de la asignatura “Procesadors Vectorials” de la profesora Montserrat Peirón Guardia, Facultat d’Informàtica de Barcelona (FIB), UPC.
- ❑ Algunas transparencias tomadas de cursos de los profesores José María LLabería y Mateo Valero Cortés, FIB, UPC.
- ❑ Basada en la asignatura “Fundamentos de Arquitecturas Paralelas”, de Ingeniería informática, impartida desde el Curso 1995-96 hasta el 2012-13.



OBJETIVOS

- (1) ↓ tiempo de ejecución de una aplicación
 - (2) ↑ productividad múltiples usuarios
 - (3) ↑ productividad aplicaciones multihilo
- servidores web, bases de datos, ...
- Consecuencias
- Hardware
 - Compilación
 - Software
- (4) Tolerancia a fallos: p.ej. sistemas navegación de un avión
 - (5) Simplificar componentes y especializar función.
Por ejemplo: sistemas empotrados en-chip (teléfono móvil)



<http://www.qualcomm.com/snapdragon/processors/800>
Snapdragon 800 series

CONSECUENCIAS

HARDWARE para (1, 2, 3)

- ❑ Combinación de:
 - ✓ ejecutar varias instrucciones por ciclo → ILP
 - ✓ la misma instrucción opera sobre varios datos → SIMD
 - ✓ un procesador ejecuta varios hilos (*threads*) → MT
 - ✓ muchos procesadores interconectados → MIMD
- ❑ Mucha memoria accesible desde los procesadores con gran ancho de banda
- ❑ Mucho disco accesible desde la memoria con gran ancho de banda (entrada/salida)

AMD EPYC Rome

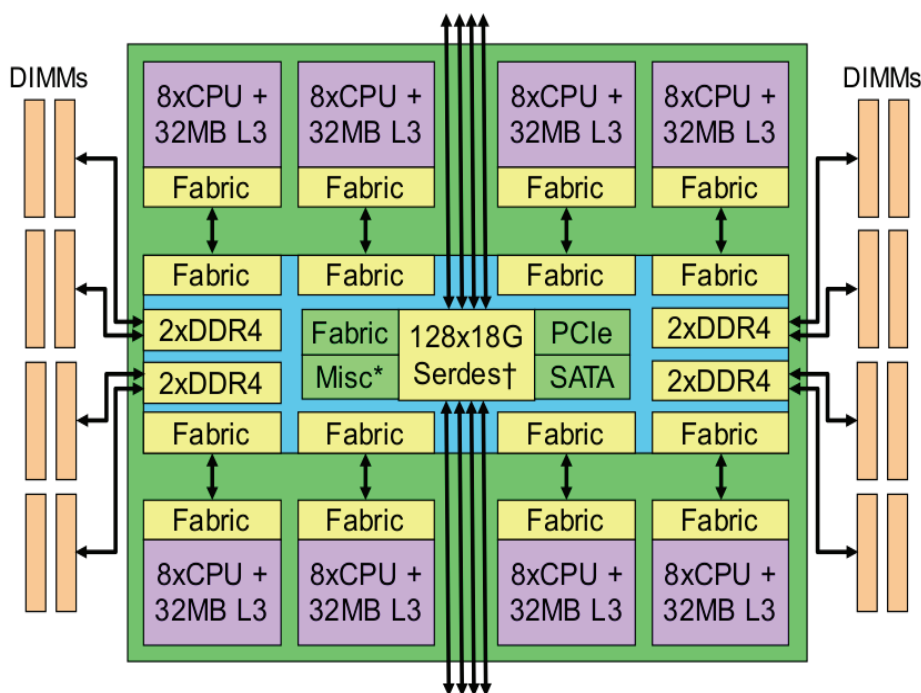


Figura: Linley Gwennap. AMD Rome Ruins Intel Hegemony. Microprocessor Report. Agosto 2019.

CONSECUENCIAS

COMPILACIÓN

- ❑ Extracción automática de paralelismo a partir de códigos secuenciales

- ✓ vectorización → SIMD
- ✓ paralelización → MIMD

Por ejemplo: <<https://godbolt.org/z/B67Rxu>>

SOFTWARE: modelos de Programación Paralela

- ❑ Paralelismo vectorial: FORTRAN 90
- ❑ Paralelismo de datos: p.ej. High Performance FORTRAN
- ❑ Single-program, multiple-data (SPMD): Co-Array FORTRAN
- ❑ Memoria compartida:
 - ✓ OpenMP <<http://www.openmp.org>>
 - ✓ Pthreads (ANSI/IEEE POSIX std. 1003.1)
 - ✓ Java
 - ✓ Intel Threading Building Blocks (TBB) en C++
 - ✓ C11, C++11, ...
- ❑ Paso de mensajes:
 - ✓ MPI
 - ✓ ...

**OBJETIVO (1): ↓ Tex de una aplicación
Supercomputación**

- ❑ Aplicaciones numéricas
 - ✓ predicción meteorológica, dinámica de fluidos, dinámica molecular, alineamiento genético
 - ✓ simulación aerodinámica
 - ala, estacionario: 10^{18} operaciones doble precisión
 - ala, turbulencia: 10^{20} DP FLOPs
 - avión, turbulencia: 10^{23} DP FLOPs
 - ✓ ¡Y muchas más!, ver p.ej.
<<http://www.bsc.es/index.php>>
computer, earth and life sciences
- ❑ Estructuras de datos: grandes matrices densas o dispersas
- ❑ Tipos de datos: generalmente números reales (coma flotante), 32/64 bits, IEEE 754
- ❑ “Pocos” bucles con muchas iteraciones
- ❑ Tiempo de ejecución limitado
por el cálculo *compute-bound*
o por el acceso a memoria *memory-bound*
- ❑ Paralelismo de datos
= **Supercomputadores numéricos**

CÓDIGO

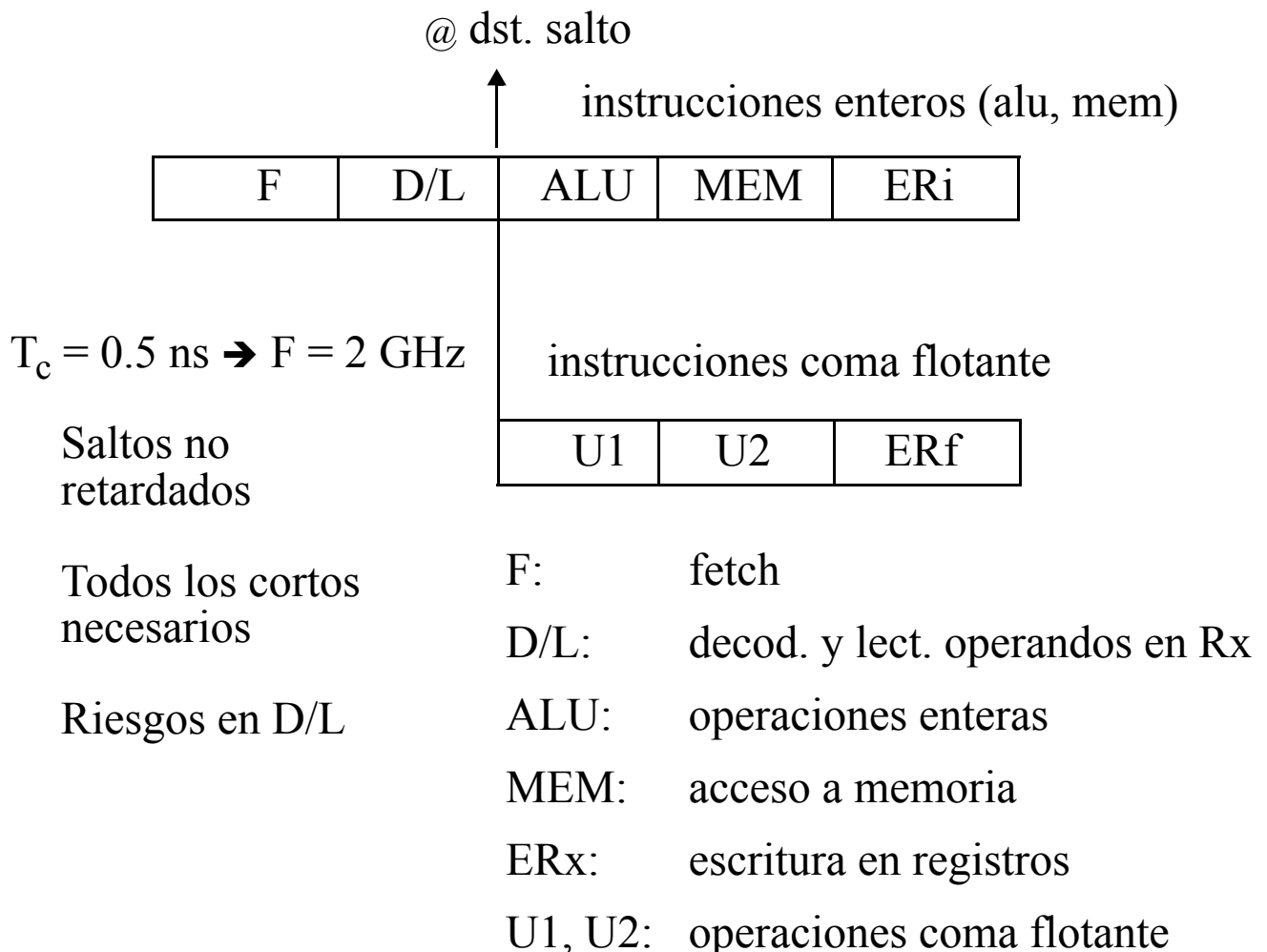
```
Real*8 C(max), A(max), B(max)
```

```
DO I = 1, max
```

```
    C(I) = A(I) + B(I)
```

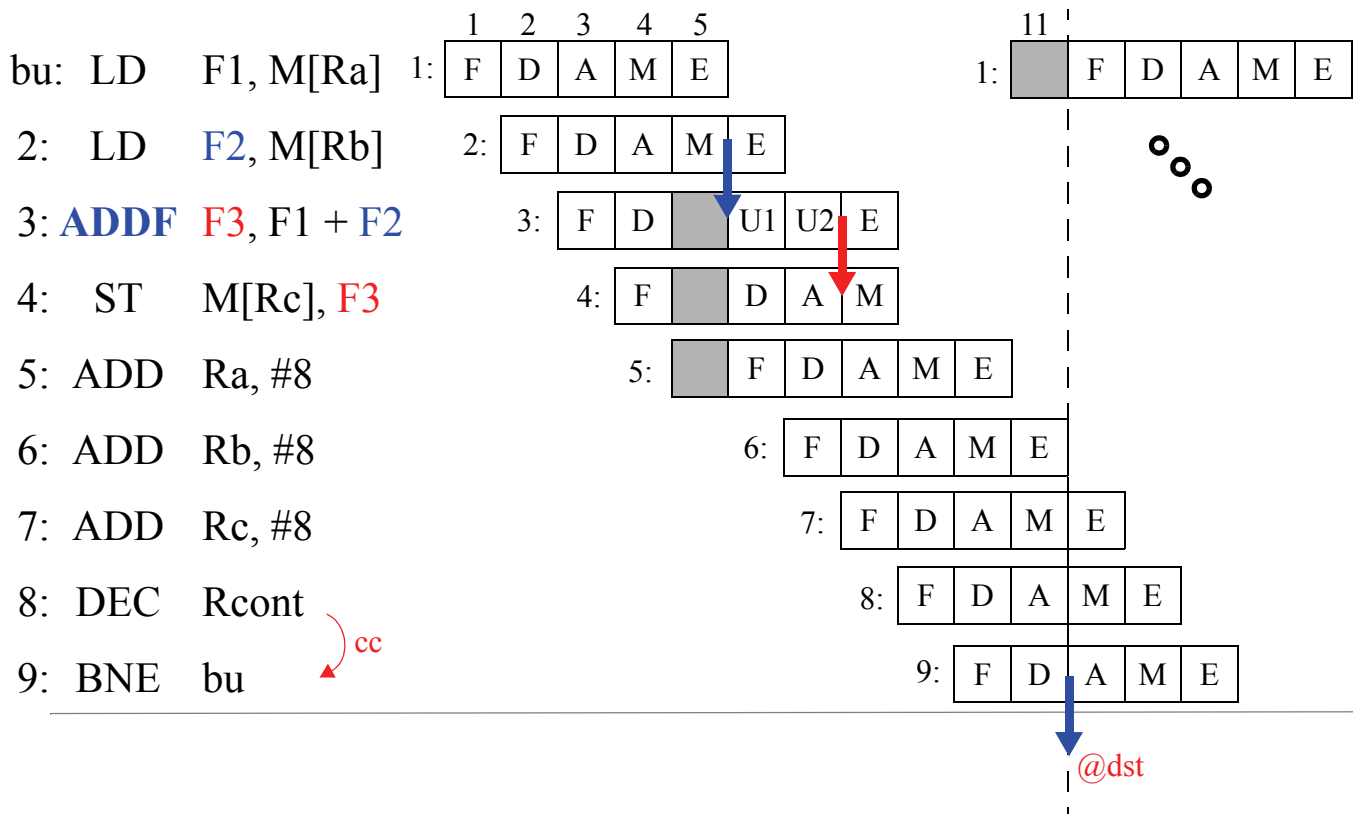
```
ENDDO
```

SEGMENTACIÓN



PROCESADOR SEGMENTADO

Rcont = max; Ra, Rb, Rc = &A, &B, &C



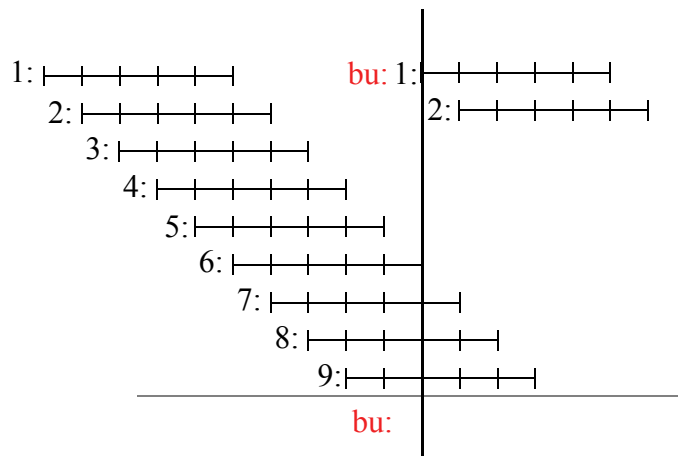
CÓDIGO Y DATOS cargados en caches

- ❑ Ciclos / iteración =
- ❑ CPF(ciclos/flop) =
- ❑ TPF(ns/flop) =
- ❑ R(MFLOPS) =

REORDENANDO CÓDIGO

- **bu:** LD F1, M[Ra]
- 2: ADD Ra, #8
- 3: LD F2, M[Rb]
- 4: ADD Rb, #8
- ➔ 5: **ADDF** F3, F1, F2
- 6: ST M[Rc], F3
- 7: ADD Rc, #8
- 8: DEC Rcont
- 9: BNE **bu**

[--]



- Un **flop/iteración** (instr. ➔)

TPF = 5 ns, R = 200 MFLOPS

rendimiento máximo: 2 GFLOPS

- Sobrecarga (*overhead*) del bucle:

lectura y escritura de operandos (instr. ○)

control de bucle y punteros (instr. □)

DESENROLLANDO D = 4

- bu: LD F10, M[Ra]
- 2: LD F11, M[Ra+8]
- 3: LD F12, M[Ra+16]
- 4: LD F13, M[Ra+24]
- 5: LD F20, M[Rb]
- 6: LD F21, M[Rb+8]
- 7: LD F22, M[Rb+16]
- 8: LD F23, M[Rb+24]

➔ 9: **ADDF** F30, F10, F20

➔ 10: **ADDF** F31, F11, F21

➔ 11: **ADDF** F32, F12, F22

➔ 12: **ADDF** F33, F13, F23

- 13: ST M[Rc], F30
- 14: ST M[Rc+8], F31
- 15: ST M[Rc+16], F32
- 16: ST M[Rc+24], F33

- 17: ADD Ra, #32
- 18: ADD Rb, #32
- 19: ADD Rc, #32
- 20: SUB Rcont, #4
- 21: BNE bu

[--]

□ Ciclos / iteración =

□ CPF(ciclos/flop) =

□ TPF(ns/flop) =

□ R(MFLOPS) =

DESENROLLANDO $D = \infty$

○ LD F1n, M[Ra+n*8]

○ LD F2n, M[Rb+n*8]

[---]

→ **ADDF** F3n, F1n, F2n

[---]

○ ST M[Rc+n*8], F3n

$n = 0 \dots \text{max}-1$

¡ IRREALIZABLE ! → 3 x max registros

□ CPF(ciclos/flop) =

□ TPF(ns/flop) =

□ R(MFLOPS) =

PROCESADOR SUPERSEGMENTADO

enteros:

↑ @ dst salto

F1	F2	D1	D2	ALU	M1	M2	ER
----	----	----	----	-----	----	----	----

 $T_c = 0.25 \text{ ns} \rightarrow 4 \text{ GHz}$

coma flotante:

F1	F2	D1	D2	U1	U2	U3	U4	ER
----	----	----	----	----	----	----	----	----

bu:

F1	F2	D1	D2	A	M1	M2	E
----	----	----	----	---	----	----	---

2:

F1	F2	D1	D2	A	M1	M2	E
----	----	----	----	---	----	----	---

3:

F1	F2	D1	D2	A	-	-	E
----	----	----	----	---	---	---	---

4:

F1	F2	D1	D2	A	-	-	E
----	----	----	----	---	---	---	---

5:

F1	F2	D1	D2	U1	U2	U3	U4	E
----	----	----	----	----	----	----	----	---

6:

F1	F2	D1	D2	A	-	-	E
----	----	----	----	---	---	---	---

7:

F1	F2	D1	D2	A	-	-	E
----	----	----	----	---	---	---	---

8:

F1	F2	D1	D2	A	M1	M2	
----	----	----	----	---	----	----	--

9:

F1	F2	D1					
----	----	----	--	--	--	--	--

1:

		F1	F2	D1	D2	A	...
--	--	----	----	----	----	---	-----

☐ CPF =☐ TPF =☐ R =

bu: LD F1, M[Ra]

2: LD F2, M[Rb]

3: ADD Ra, #8

4: ADD Rb, #8

5: ADDF F3, F1, F2

6: ADD Rc, #8

7: DEC Rcont

8: ST M[Rc], F3 ; inicializamos

9: BNE bu ; Rc con 8 menos!

[--]

PROCESADOR SUPERESCALAR

- ❑ Cambiamos el control

Mantenemos rutas de datos (\approx coste)

$$T_c = 0.5 \text{ ns} \rightarrow 2 \text{ GHz}$$

- ❑ Suponemos grado 2:
se buscan dos instrucciones por ciclo
- ❑ Pero sólo se ejecutan en paralelo
si utilizan las dos rutas de datos: entera y coma flotante
- ❑ Hay que planificar para formar parejas de instrucciones
sin dependencias entera-coma flotante
- ❑ El efecto es “esconder” los ciclos de ejecución de ADDF
- ❑ Ciclos/FLOP

	d = 1	d = 4	d = ∞
Ciclos/ FLOP	9	4.5	3

RESUMEN

- TPF obtenidos en los diferentes casos en **ns**
(entre paréntesis, la velocidad **R** en **MFLOPS**)

	d = 1	d = 4	d = ∞	Tc ns
Segmentado $R_{\max} = 2$ GFLOPS	5 (200)	2.75 (363.6)	2 (500)	0.5
Supersegmentado $R_{\max} = 4$ GFLOPS	2.75 (363.6)	1.4375 (695)	1 (1000)	0.25
Superescalar $R_{\max} = 2$ GFLOPS	4.5 (222)	2.25 (444)	1.5 (666)	0.5

bu: LD **F11**, M[Ra]

LD F12, M[Ra+8]

LD F1N, M[Ra+(N-1)*8]

LD **F21**, M[Rb]

LD F22, M[Rb+8]

LD F2N, M[Rb+(N-1)*8]

ADDF F31, F11, F21

ADDF F32, F12, F22

ADDF F3N, F1N, F2N

ST M[Rc], **F31**

ST M[Rc + (N-1)*8], F3N

ADD Ra, #N*8

ADD Rb, #N*8

ADD Rc, #N*8

SUB Rcont, #N

BNE *bu*

LV **V1**, M[Ra]

LV **V2**, M[Rb]

ADDV V3, V1, V2

SV **V3**, M[Rc]

ADD Ra, #N*8

ADD Rb, #N*8

ADD Rc, #N*8



max/N iteraciones

cuerpo-N
desenrollado

❑ 1 instrucción vectorial corresponde a N instrucciones del bucle N-desenrollado

- ✓ Instrucciones por Iteración =
- ✓ Ancho de Banda con Memoria =
- ✓ Tiempo de ejecución =

2.1 ARQUITECTURA Y ORGANIZACIÓN

- ❑ Primeros supercomputadores tipo CISC (M-M)
 - ✓ instrucción compleja: **ADDV @md, @mf1, @mf2**
 - lee dos flujos desde memoria
 - calcula
 - escribe el flujo resultado hacia memoria
 - ✓ Organización segmentada
 - ⇓
 - Latencia de obtención de operandos muy grande, sobre todo para paso no secuencial.
Es difícil encadenar
 - CDC Star 100 ('72), TI ASC ('72)

- ❑ Después, y en la actualidad, filosofía RISC (R-R), o sea, desacoplar instrucciones de cálculo y de acceso a memoria:
 - ✓ Instr. vectoriales de acceso a memoria para carga/descarga de registros vectoriales: **LV**, **SV**
 - ✓ Instr. vectoriales de cálculo sobre los registros vectoriales: **ADDV**
 - ✓ VLR: Vector Length Register
n° de elementos a procesar
p.ej. para vectores de **64** elementos necesito **6 + 1** bits
 - ✓ VMR: Vector Mask Register
impide que algunas operaciones se realicen

❑ Ventajas de un repertorio vectorial:

Compacto

- ✓ Una instrucción pequeña codifica N operaciones

Expresivo:

La instrucción indica al hardware

- ✓ que las operaciones son independientes
- ✓ el número de operaciones
- ✓ el “patrón” de acceso a memoria:
en secuencia (stride = 1 elemento)
o
a saltos (stride > 1 elemento)

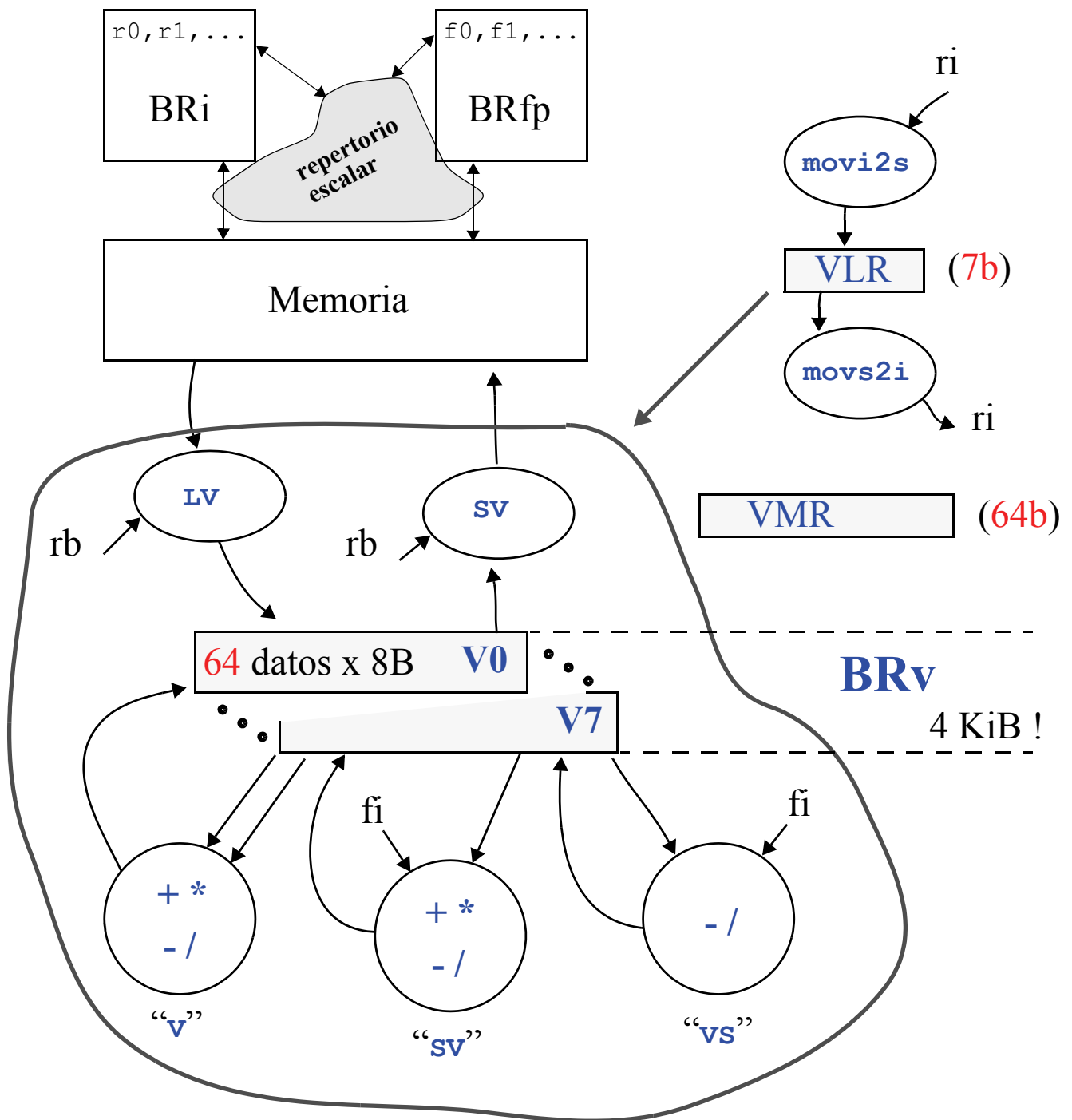
Escalable

- ✓ el mismo binario puede ejecutarse
en una o varias “pistas” segmentadas
(*parallel pipelines or lanes*)

Bajo consumo de energía

- ✓ Ahorramos energía en el acceso a la memoria de instrucciones y en la decodificación (1/N)
- ✓ Cuando se ejecutan las operaciones de una instrucción vectorial no hay que gastar energía en comprobar dependencias

□ Ejemplo: DLXV¹ Diagrama de flujo de la ALMa²



Simple precisión: `addv.F`

doble precisión: `addv.D`

1. [HePa12]: Appendix G

2. **ALMa**: Arquitectura de Lenguaje Máquina (en inglés *ISA, Instruction Set Architecture*)

EJEMPLOS

movi2s VLR, r2 ; VLR = r2_{6:0} ∈ {0:64}

lv V0, [r3] ; V0.i = mem [r3 + i*8]

sv [r5], V4 ; mem [r5 + i*8] = V4.i

addv V1, V2, V3 ; V1.i = V2.i + V3.i

addsv V1, f0, V3 ; V1.i = f0 + V3.i

subsv V1, f0, V3 ; V1.i = f0 - V3.i

subvs V1, V3, f5 ; V1.i = V3.i - f5

VLR	acción ^a	
0	NOP	i = 0:VLR-1
1	i = 0:0	
2	i = 0:1	
...		
64	i = 0:63	
65	i = 0:63	
...		
127		

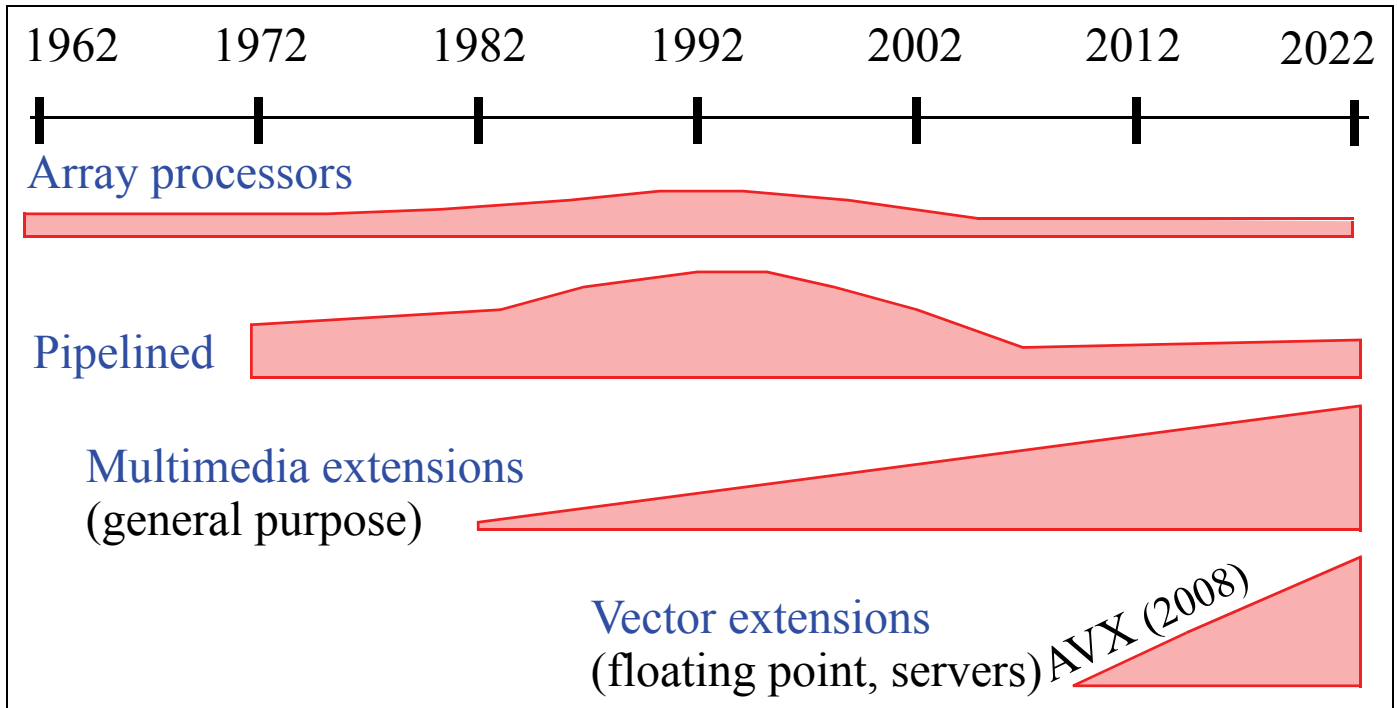
a. n° elementos a procesar = max (VLR mod 64, 64 x VLR div 64)

```
integer i, max  
  
parameter (max = multiplo de 64)  
  
real*8 C(max), A(max), B(max)  
  
do i= 1,max  
  
    C(i) = A(i) + B(i)  
  
enddo
```

Suponemos registros enteros inicializados:

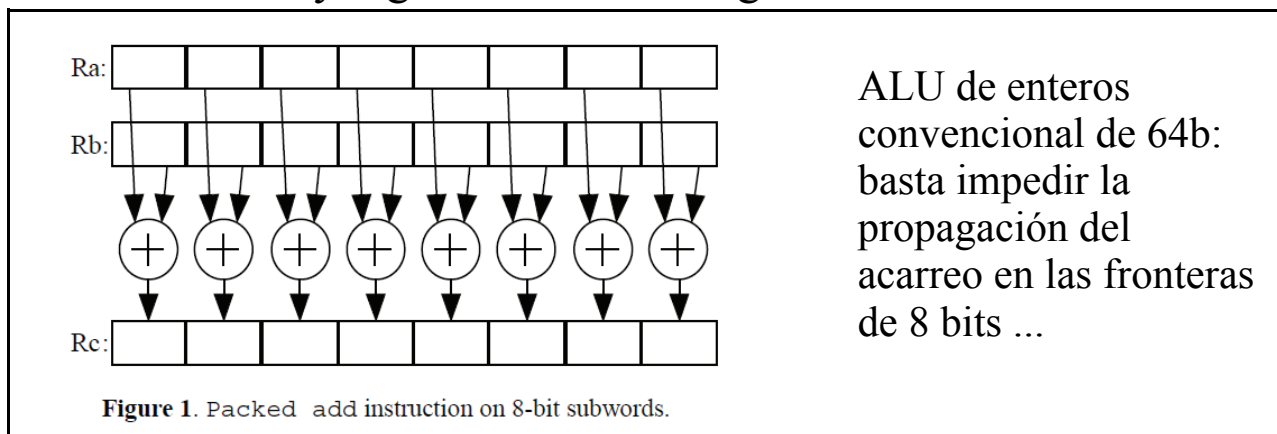
$R_A = \&A(0)$, $R_B = \&B(0)$, $R_C = \&C(0)$, $R_{\text{cont}} = \text{max}$, $R_{64} = 64$

SIMD: Single Instruction Multiple Data



❑ Extensiones Multimedia: Intel *MMX* (1982), luego Intel *SSE*

- ✓ Vectores cortos
- ✓ No hay registro VLR: la longitud vectorial va en el CO



- ✓ Freescale/Apple/IBM *AltiVec*, SPARC *VIS*, ARM *Neon*

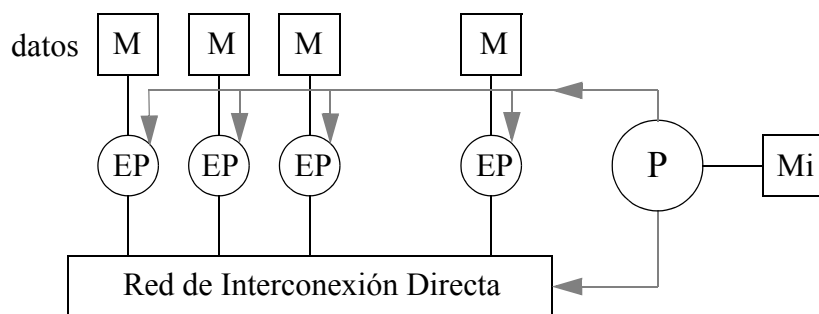
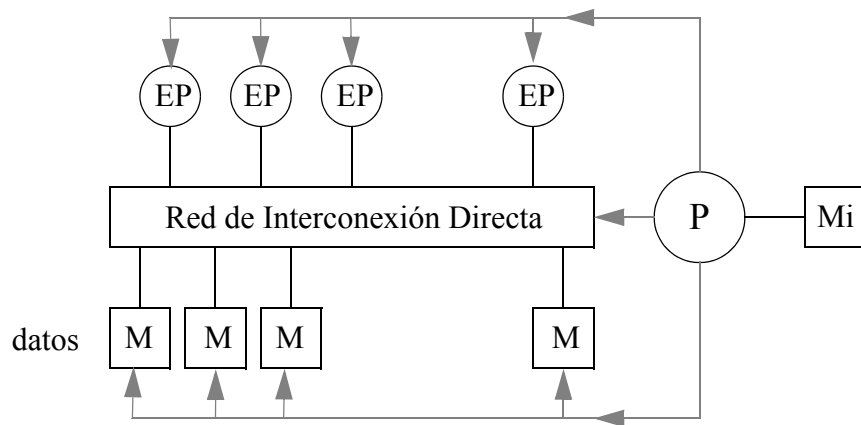
❑ Extensiones Vectoriales: Intel *AVX* (4 real*8, 8 real*4, ...)

- ✓ Soporte hw: ALUs replicadas y BR vectorial (BRV)
- ✓ BRV de AVX → 8 registros de 256 bits (YMM0-7)

❑ Procesadores Matriciales (Array Processors)

Históricamente fueron los primeros:

- ✓ Solomon computers
(Westinghouse Electric Corporation, 1960-62)
→ ILLIAC IV (8x8) en 1968-74
- ✓ Soporte hw: **Nodos (mem+regs+ALU) + Red Interconexión Directa entre nodos**
- ✓ STARAN (74), BSP (82)
Connection Machines:
 - CM1 - CM2 - CM5 (1985 - 87 - 93)
- ✓ Actualidad: procesamiento de imagen y de señal



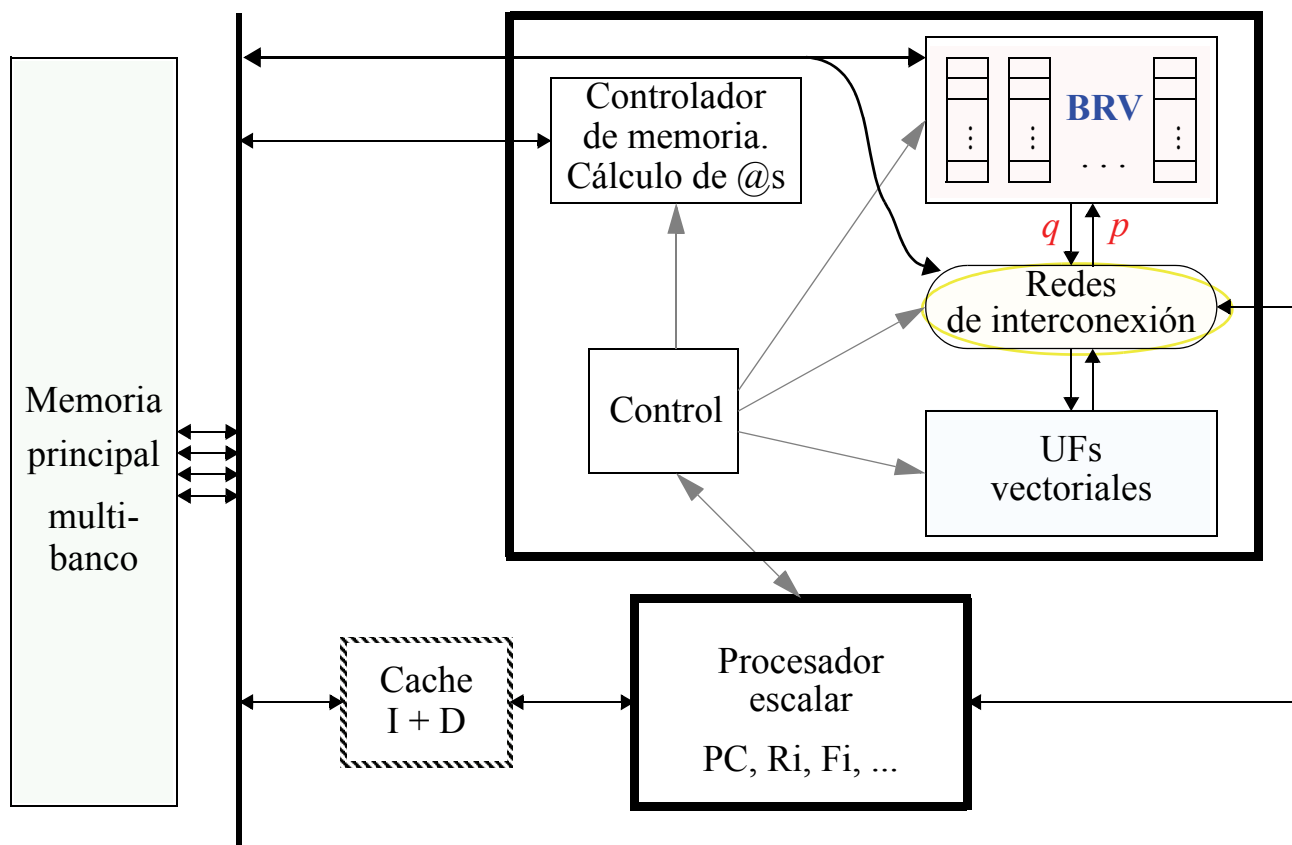
❑ Procesadores Vectoriales Segmentados (pipelined)

Memoria-memoria

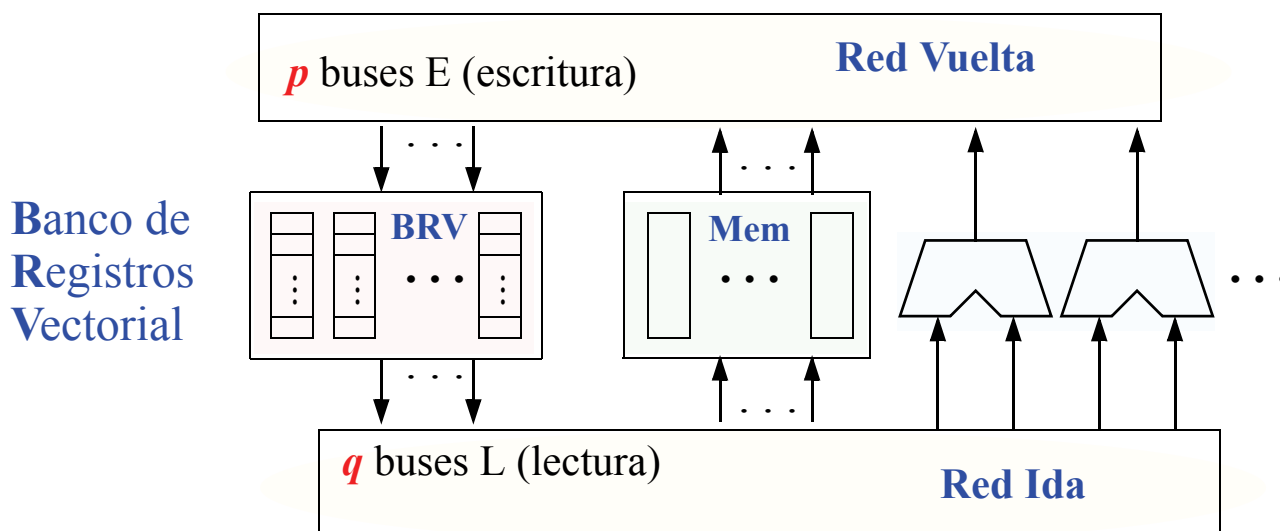
- ✓ **ADDV** @md, @mf1, @mf2
- ✓ CDC Star 100 (72), TI ASC (72)

Memoria-registros (Seymour Cray, 1975)

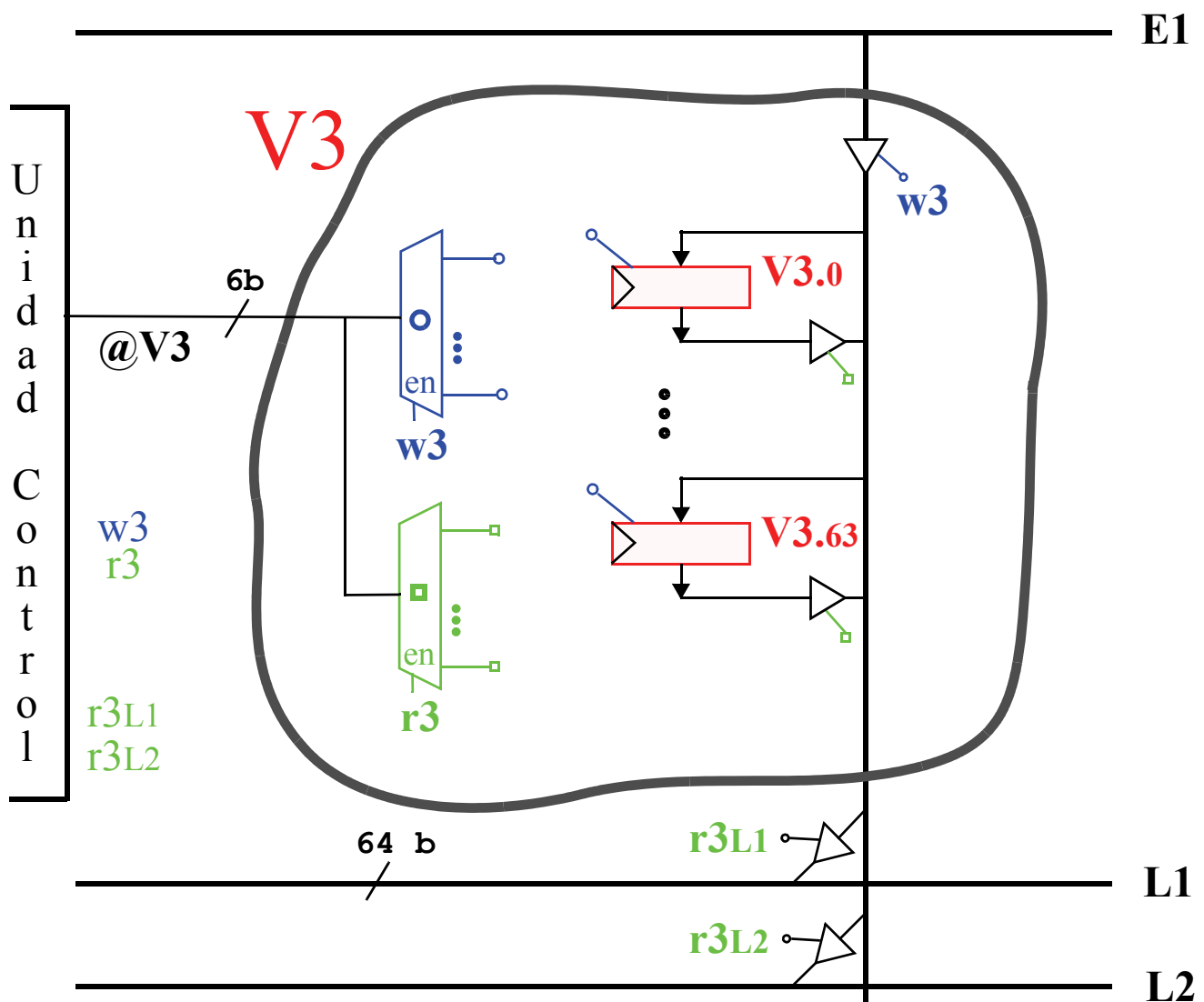
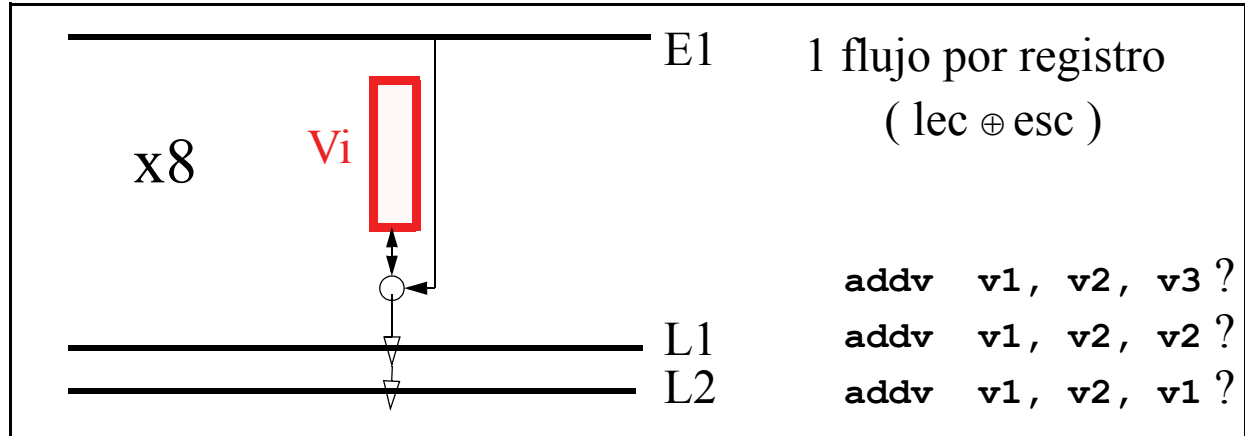
- ✓ Vectores de tamaño grande, hasta 4K elementos
- ✓ **Soporte hw:** Banco Registros Vectorial +
Memoria Multibanco +
ALUs + segmentación
- ✓ Cray 1 (76), Cray2, Cray X-MP y Cray Y-MP
CDC Eta y Cyber
IBM 3090 VF
Fujitsu VP200
Hitachi S-810
Convex
Alliant
NEC SX/2 - ... - SX/9 - SX/ACE (85 - 07 -13)



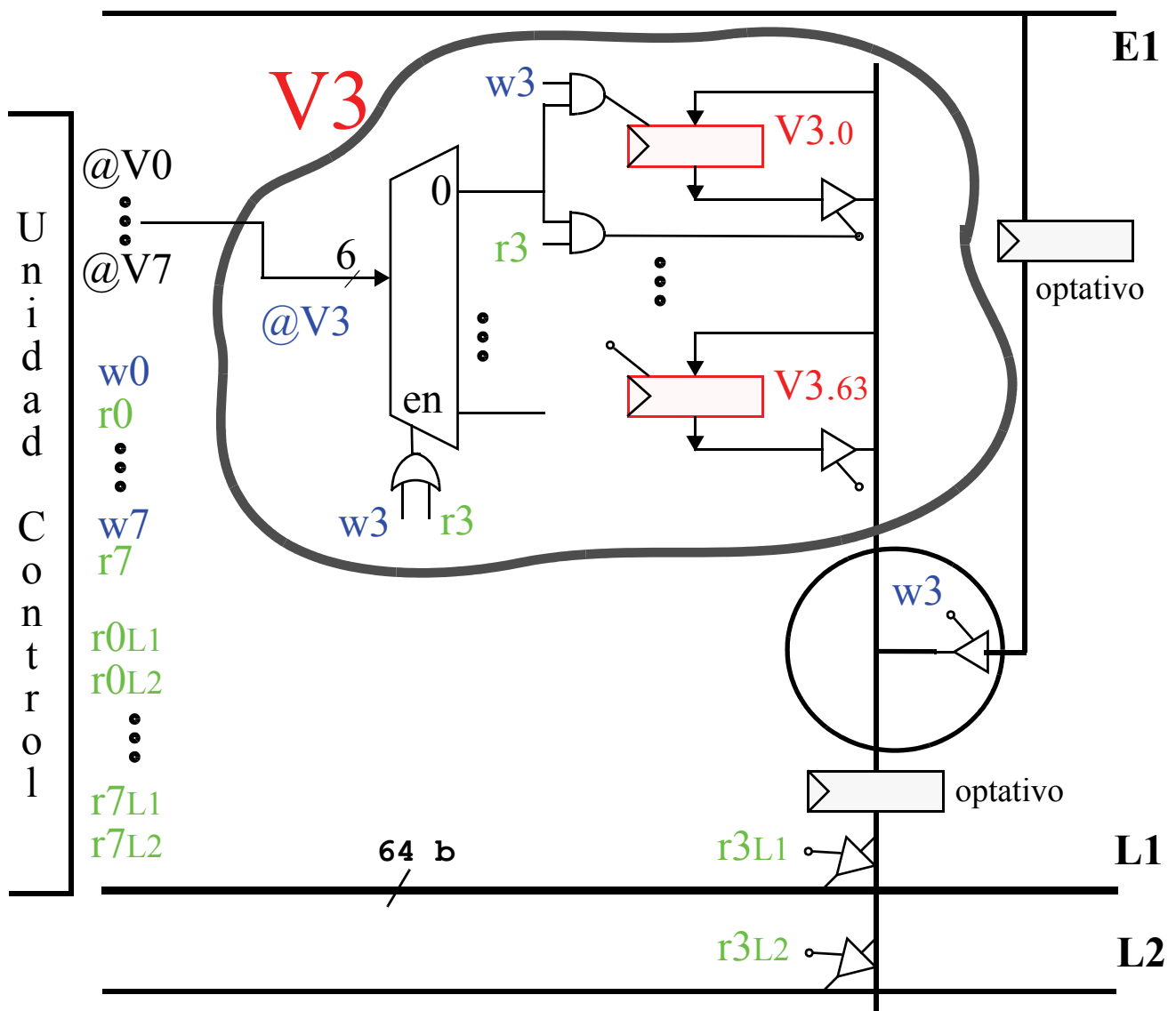
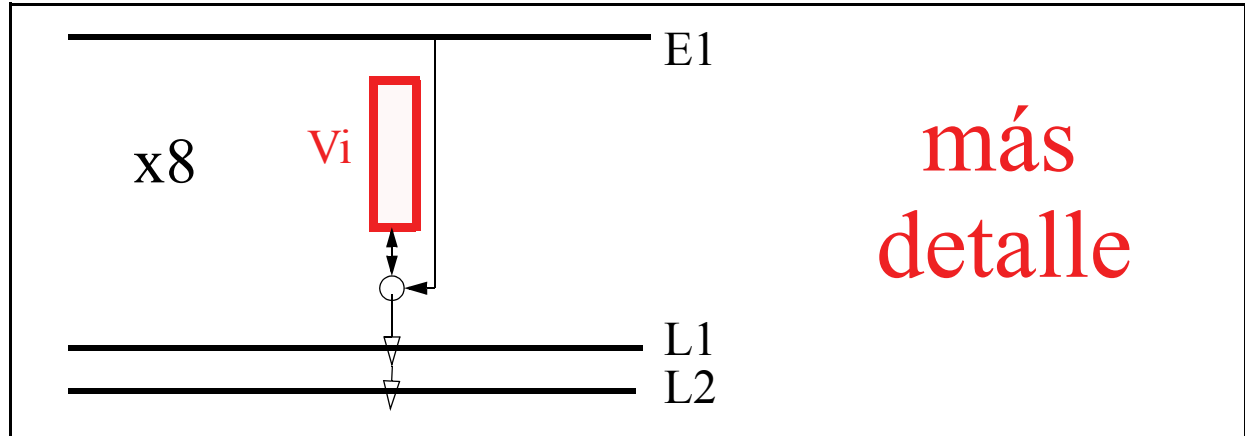
otro punto de vista:



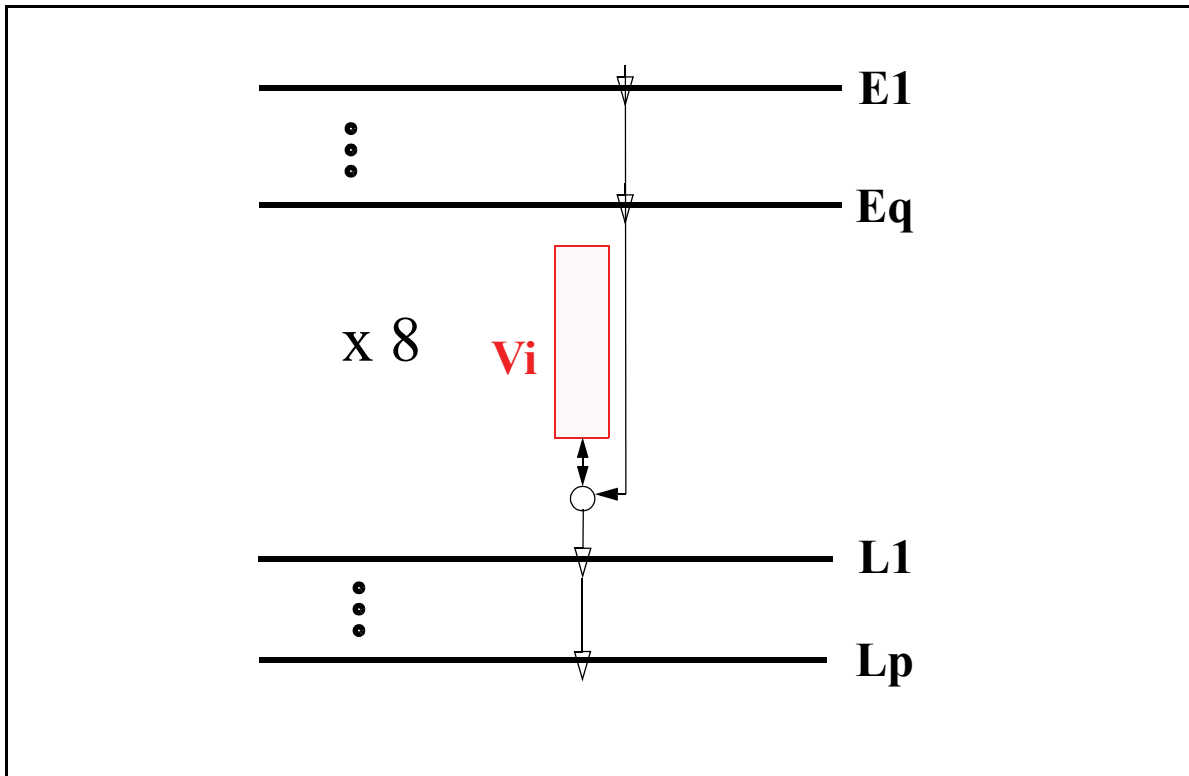
- ❑ BRV con 2 buses L, 1 bus E y Reg Vec de 1 puerto (1L⊕1E)



- ❑ BRV con 2 buses L, 1 bus E y Reg Vec de 1 puerto (1L⊕1E)

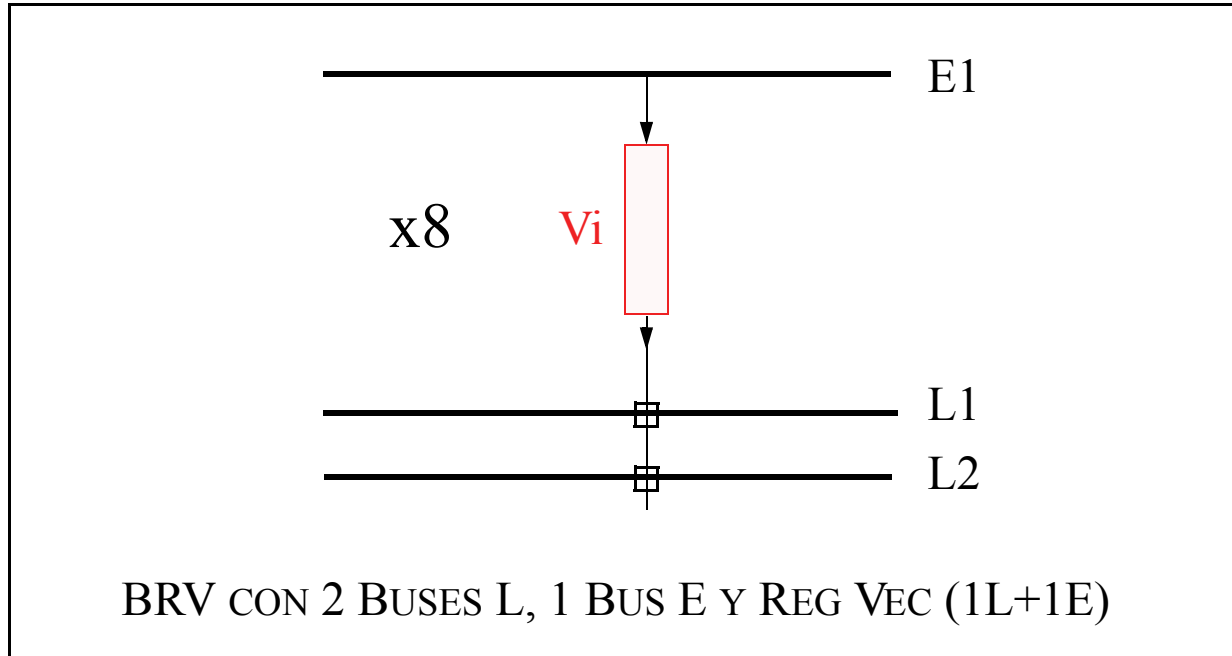


- ❑ BRV con p buses L, q buses E y **Reg Vec de 1 puerto**



- ❑ 1 flujo por registro
 - ✓ o bien lectura desde un bus L ($L1..Lp$)
 - ✓ o bien escritura desde un bus E ($E1..Eq$)
- ❑ Globalmente max. 8 flujos
- ❑ Ejemplo: 5 buses L y 3 buses E permiten a la vez
 - ✓ dos operaciones de dos operandos, y
 - ✓ una operación de un operando

- Mejora: **2 puertos por V_i** : lectura y escritura (1L+1E)



- Ya es posible `addv v1, v5, v1`

- Solapar dependencias?

`addv v1, v2, v3`
`subv v4, v5, v1`

- Solapar antidependencias?

`addv v1, v2, v3`
`subv v3, v4, v5`

SI, se llama

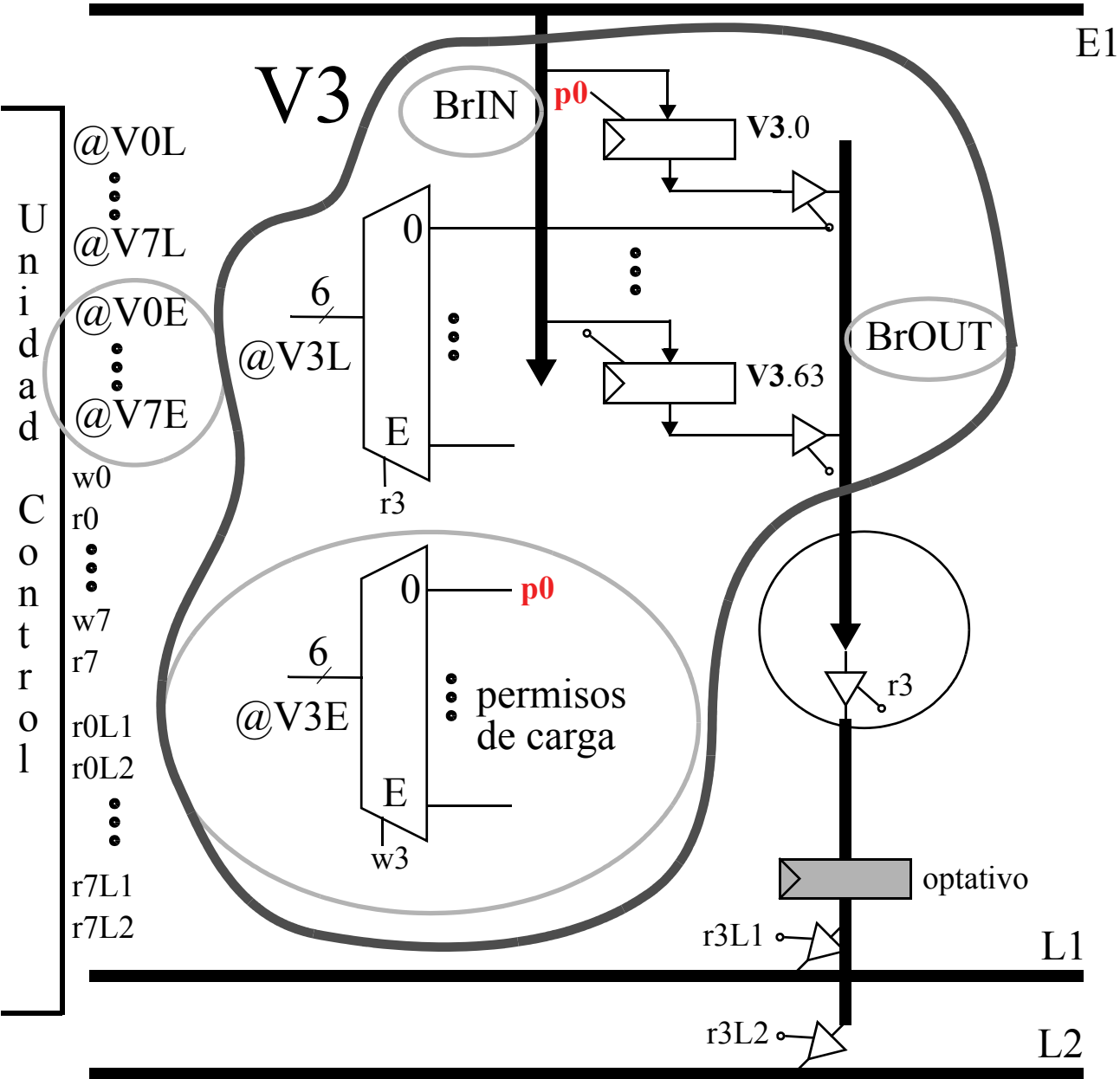
encadenamiento general

y necesita un control elaborado (lo veremos)

También, y usa un control similar

- Podemos

- ✓ añadir puertos de lectura en cada V_i
- ✓ tener un BRV con p buses L y q buses E

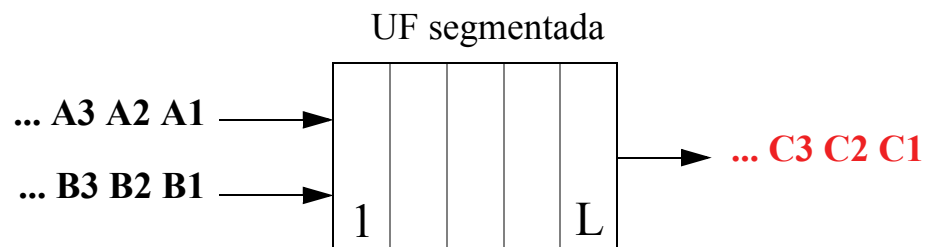


❑ Descripción temporal:

- ✓ Latencia (ciclos): tiempo desde que entran los operandos hasta que sale el resultado.
- ✓ Lat. iniciación (LI) = lat. finalización (ciclos/operación)¹:
nº min. ciclos entre entradas consecutivas de operandos

❑ Sumadores, Multiplicadores

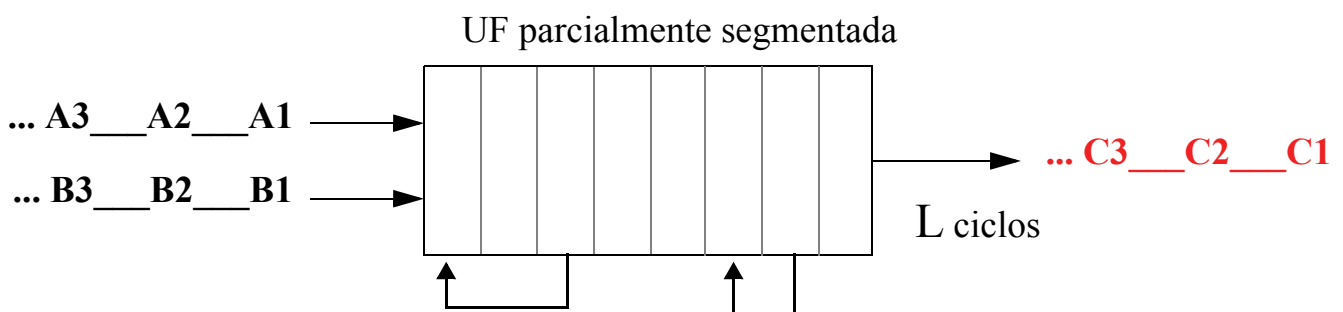
$L = \text{nº de etapas, } LI = 1 \text{ ciclo/op}$



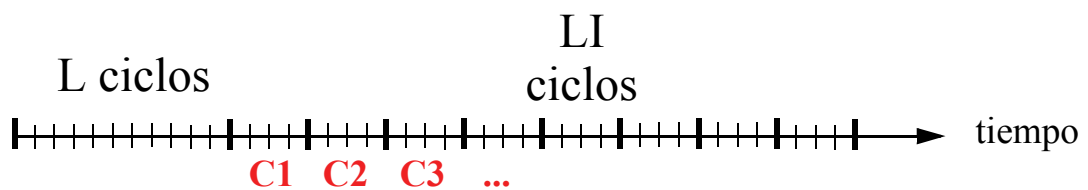
❑ División, Raíces, Exponenciación, Trigonómicas

$L > \text{nº etapas, por reutilización}$

$LI > 1$



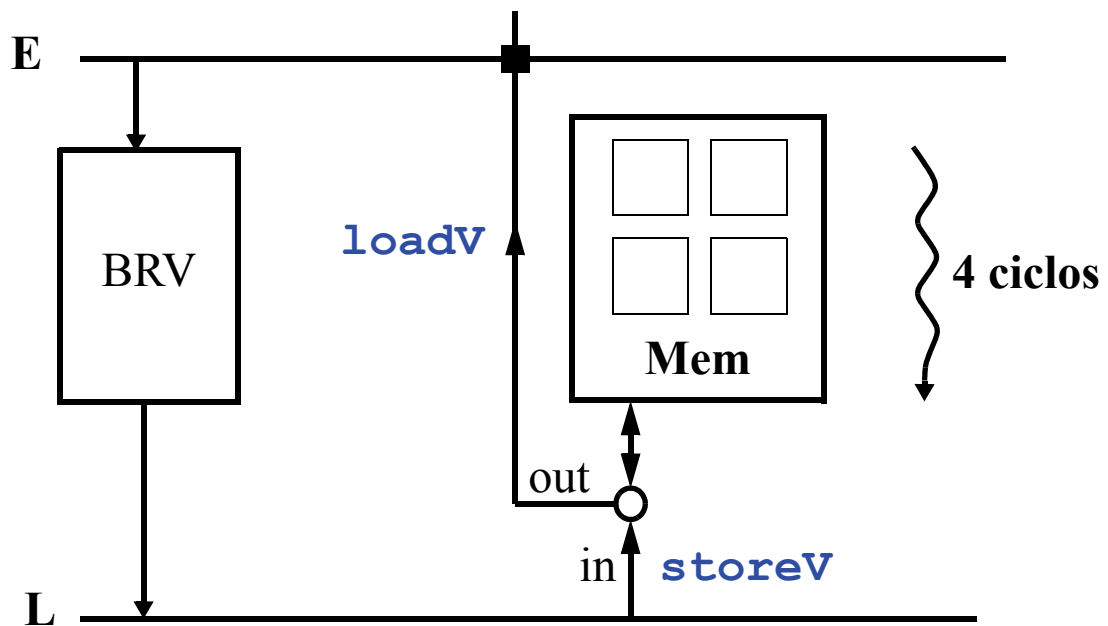
❑ Modelo de ejecución:



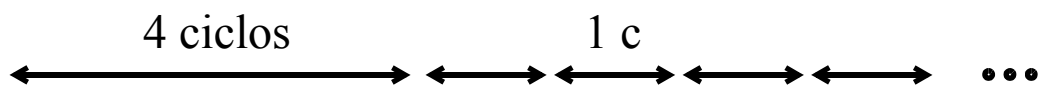
1. La inversa se llama tasa iniciación = tasa finalización (operaciones/ciclo)

❑ Ejemplo de memoria:

- 4 GiB con **1 puerto** compartido de lec/escr
- 4 bancos¹ de 4 ciclos de latencia



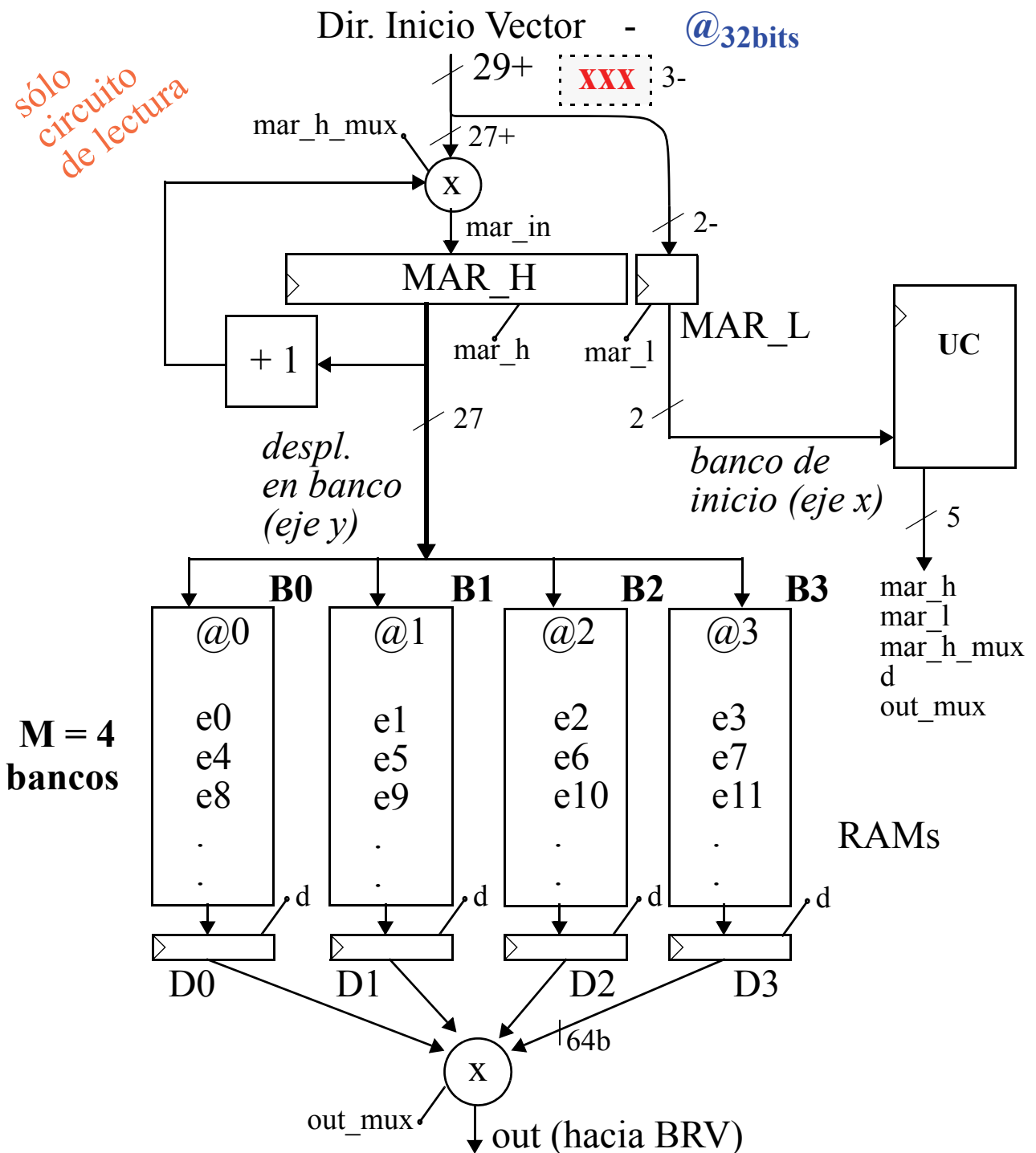
❑ Queremos este comportamiento en lectura:



❑ ¿CÓMO?

- Multibanco: n° bancos = latencia memoria
- Entrelazado “por palabras”:
 “palabras consecutivas en bancos consecutivos”,
 asumimos números de 8 bytes **alineados**
- $L = 4$ ciclos $\rightarrow M = 4$ bancos de 1GiB
- Acceso **síncrono** o acceso **desfasado**

1. También podemos decir “módulos”



❑ @32bits → n° banco = ($@_{32\text{bits}}/8$) mod 4 → {0, 1, 2, 3}

❑ Parte del control:

- Cargar **MAR_H** cada cuatro ciclos (ojo primer acceso!)
- Cargar todos los **Di** cada cuatro ciclos
- Multiplexar **Di** en **out** cada ciclo

Ej1. @ = 32 → fila 1, columna 0 (e4, e5, e6, ...)

MAR	mar_in	1											
	MAR_H												
B0	D0												
B1	D1												
B2	D2												
B3	D3												
		1	2	3	4	5	6	7	8	9	10	11	ciclos
OUT													

Ej2. @ = 48 → fila 1, columna 2 (e6, e7, e8, ...)

MAR	mar_in	1											
	MAR_H												
B0	D0												
B1	D1												
B2	D2												
B3	D3												
		1	2	3	4	5	6	7	8	9	10	11	ciclos
OUT													

Ej1. @ = 32 → fila 1, columna 0 (e4, e5, e6, ...)

MAR	mar_in	1	2	2	2	2	3	3	3	3	4	4
	MAR_H		1				2				3	
B0	D0						e4				e8	
B1	D1						e5				e9	
B2	D2						e6				e10	
B3	D3						e7				e11	
		1	2	3	4	5	6	7	8	9	10	11
OUT			L = 4				e4	e5	e6	e7	e8	e9

ciclos

Ej2. @ = 48 → fila 1, columna 2 (e6, e7, e8, ...)

MAR	mar_in	1	2	2	2	2	3	3	3	3	4	4
	MAR_H		1				2				3	
B0	D0						e4				e8	
B1	D1						e5				e9	
B2	D2						e6				e10	
B3	D3						e7				e11	
		1	2	3	4	5	6	7	8	9	10	11
OUT			L = 4				-	-	e6	e7	e8	e9

ciclos

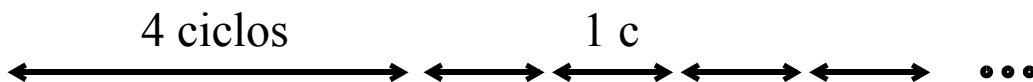
❑ Problemas con acceso **Síncrono**:

- ✓ lectura: latencia variable o irregular
- ✓ escritura
 - un registro MDRin para cada banco
 - es lenta: (1 1 1 1 4) (1 1 1 1 4) ...

❑ Una solución: acceso **desfasado**
los bancos son más independientes

- ✓ un registro MAR para cada banco
- ✓ control más complejo
- ✓ resultado:

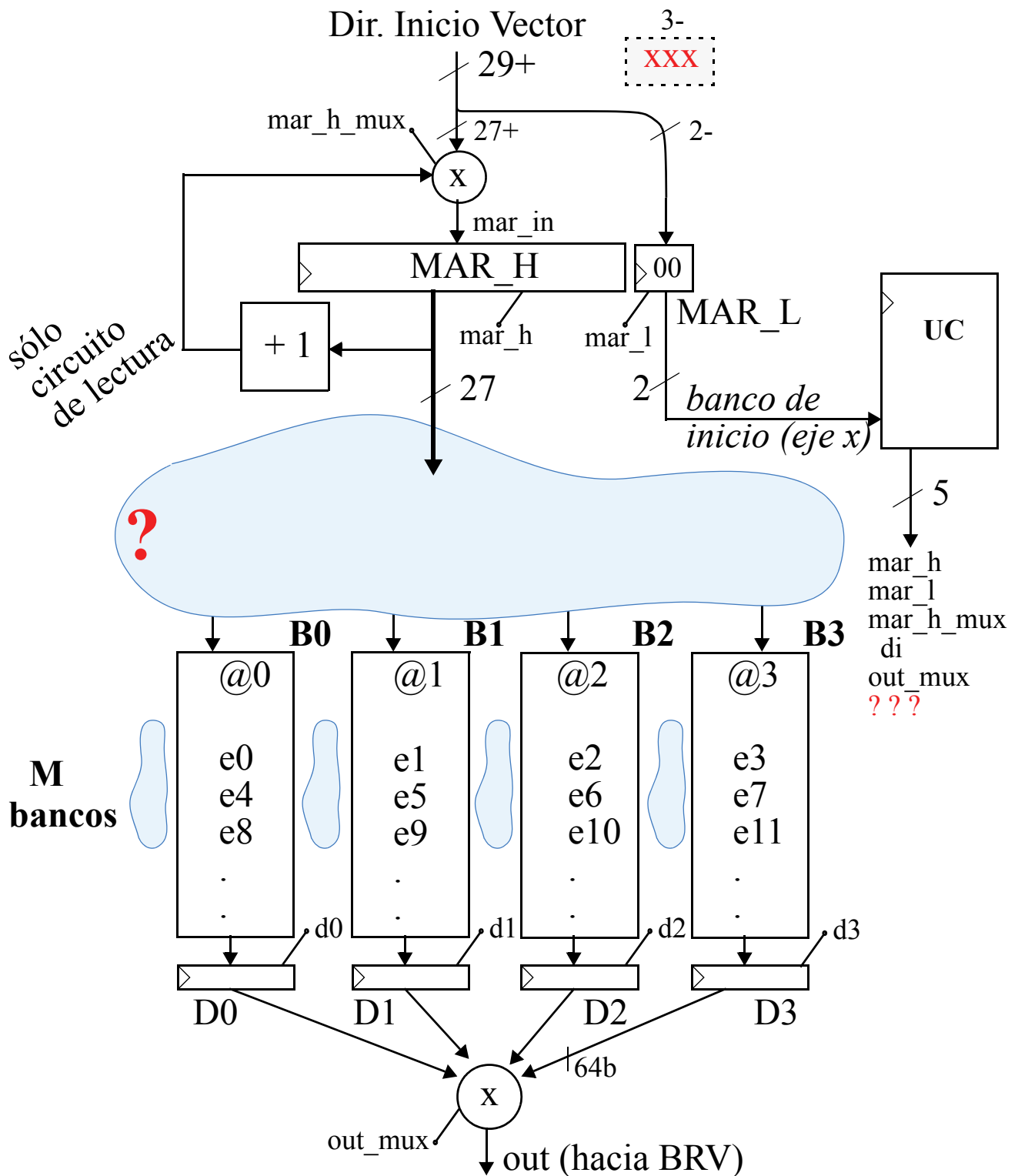
LECTURA



ESCRITURA



Veamos cómo se implementa -->



Control ?:

mar_h, mar_l, mar_h_mux, d_i,
outmux, ???

Ej1. @ = 32 → fila 1, columna 0 (e4, e5, e6, ...)

MAR	mar_in	1				2				3			
	MAR_H		1				2				3		
MAR_0													
MAR_1													
MAR_2													
MAR_3													
Actividad B0													
Actividad B1													
Actividad B2													
Actividad B3													
		1	2	3	4	5	6	7	8	9	10	11	ciclos
OUT													

Ej2. @ = 48 → fila 1, columna 2 (e6, e7, e8, ...)

MAR	mar_in	1				2				3			
	MAR_H		1				2				3		
MAR_0													
MAR_1													
MAR_2													
MAR_3													
Actividad B0													
Actividad B1													
Actividad B2													
Actividad B3													
		1	2	3	4	5	6	7	8	9	10	11	ciclos
OUT													

Ej1. @ = 32 → fila 1, columna 0 (e4, e5, e6, ...)

MAR	mar_in	1	1	1	1	2				3				
	MAR_H		1	1	1	1	2				3			
MAR_0				1				2				3		
MAR_1				1				2						
MAR_2					1				2					
MAR_3						1				2				
Actividad B0			e4				e8							
Actividad B1				e5				e9						
Actividad B2					e6				e10					
Actividad B3						e7				e11				
		1	2	3	4	5	6	7	8	9	10	11	ciclos	
OUT			L = 1 + 4				e4	e5	e6	e7	e8	e9		

□ Sincrono

Ej2. @ = 48 → fila 1, columna 2 (e6, e7, e8, ...)

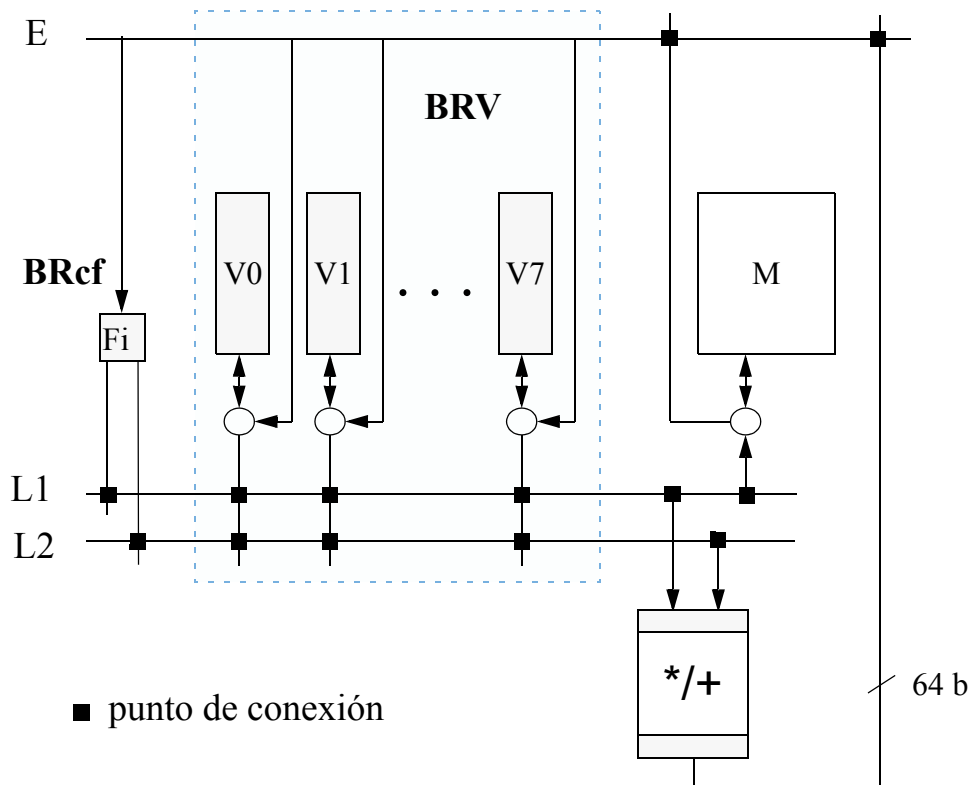
MAR	mar_in	1				2				3			
	MAR_H		1				2				3		
Actividad B0					e4				e8				
Actividad B1					e5				e9				
Actividad B2					e6				e10				
Actividad B3					e7				e11				
		1	2	3	4	5	6	7	8	9	10	11	ciclos
OUT		L = 4 +	e6	e7	e8	e9	e10	

□ Desfasado

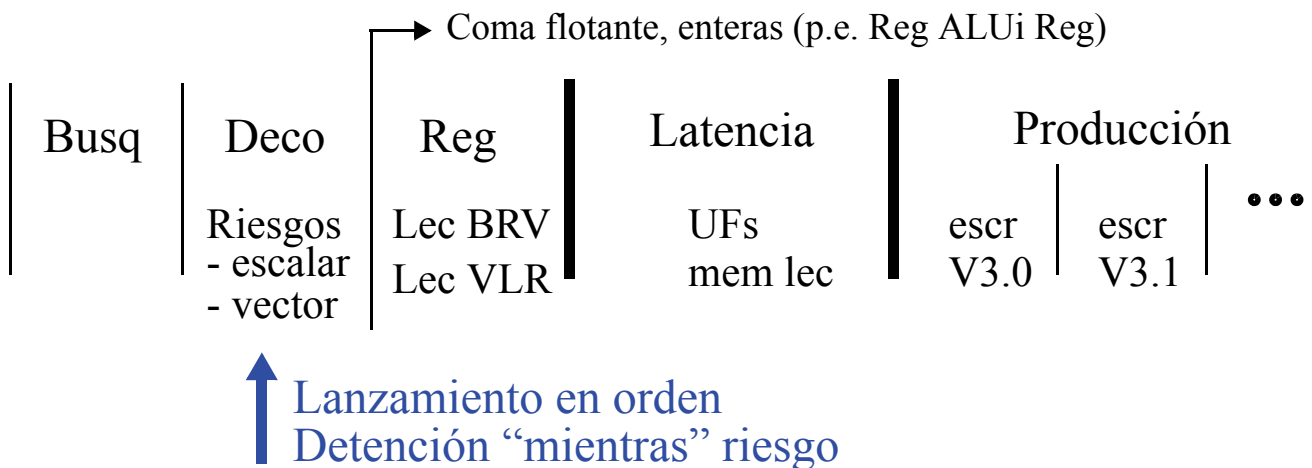
Ej2. @ = 48 → fila 1, columna 2 (e6, e7, e8, ...)

MAR	mar_in	1		2				3				
	MAR_H		1		2				3			
MAR_0						2				3		
MAR_1							2					
MAR_2				1				2				
MAR_3					1				2			
Actividad B0						e8				e12		
Actividad B1						e9					e13	
Actividad B2				e6				e10				
Actividad B3				e7				e11				
		1	2	3	4	5	6	7	8	9	10	11
												<i>ciclos</i>
OUT			L = 1 + 4				e6	e7	e8	e9	e10	

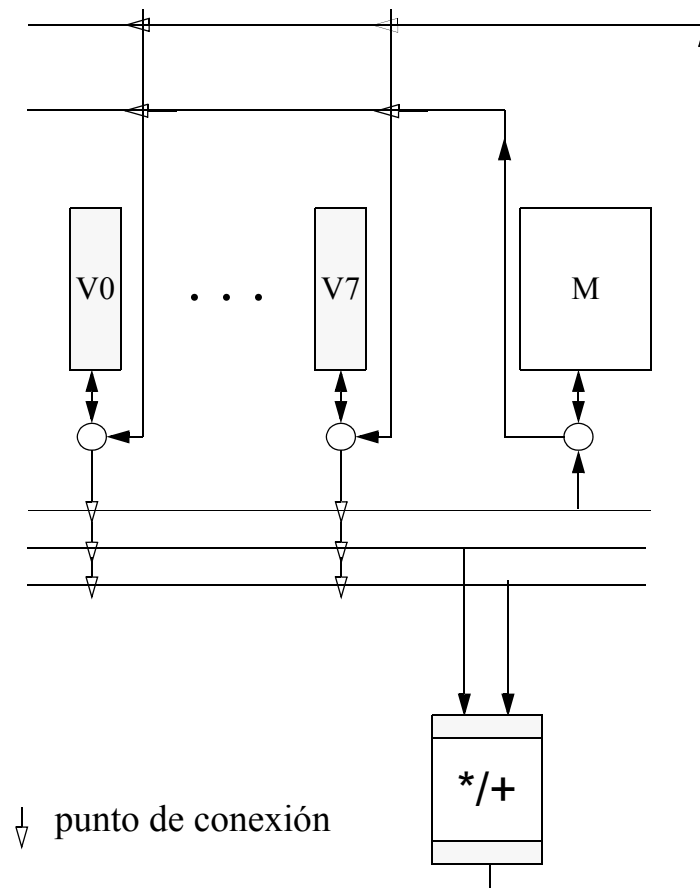
- ❑ BRV con $2L + 1E$ (y reg. vec. mínimos)
 - ✓ 2 buses L: $BRV \rightarrow \{M \oplus ALU\}$
 - ✓ 1 bus E: $\{M \oplus ALU\} \rightarrow BRV$
- ❑ Memoria de 1 puerto $\{L \oplus E\}$. Sólo 1 UF



- ❑ Segmentación. Ejemplo.



- ❑ 3 buses lectura: $BRV \rightarrow \{M, ALU\}$
- ❑ 2 buses escritura: $\{M, ALU\} \rightarrow BRV$



- ❑ Instrucciones independientes y sin riesgos estructurales

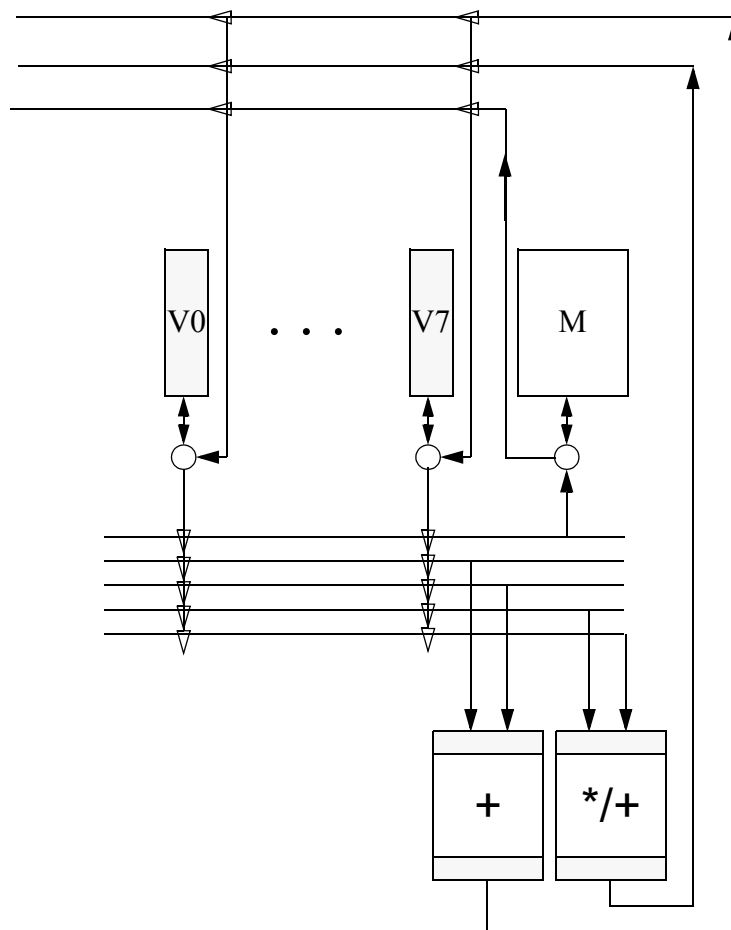
B	D ¹	R	L		W ^{VLR}			
	B	D ¹	R	L		W ^{VLR}		

- ❑ Instrucciones dependientes o con riesgos estructurales

B	D ¹	R	L			W ^{VLR}								
	B	D ¹	D ²	D ³	D ⁴	D ⁵	D ⁿ	R	L	

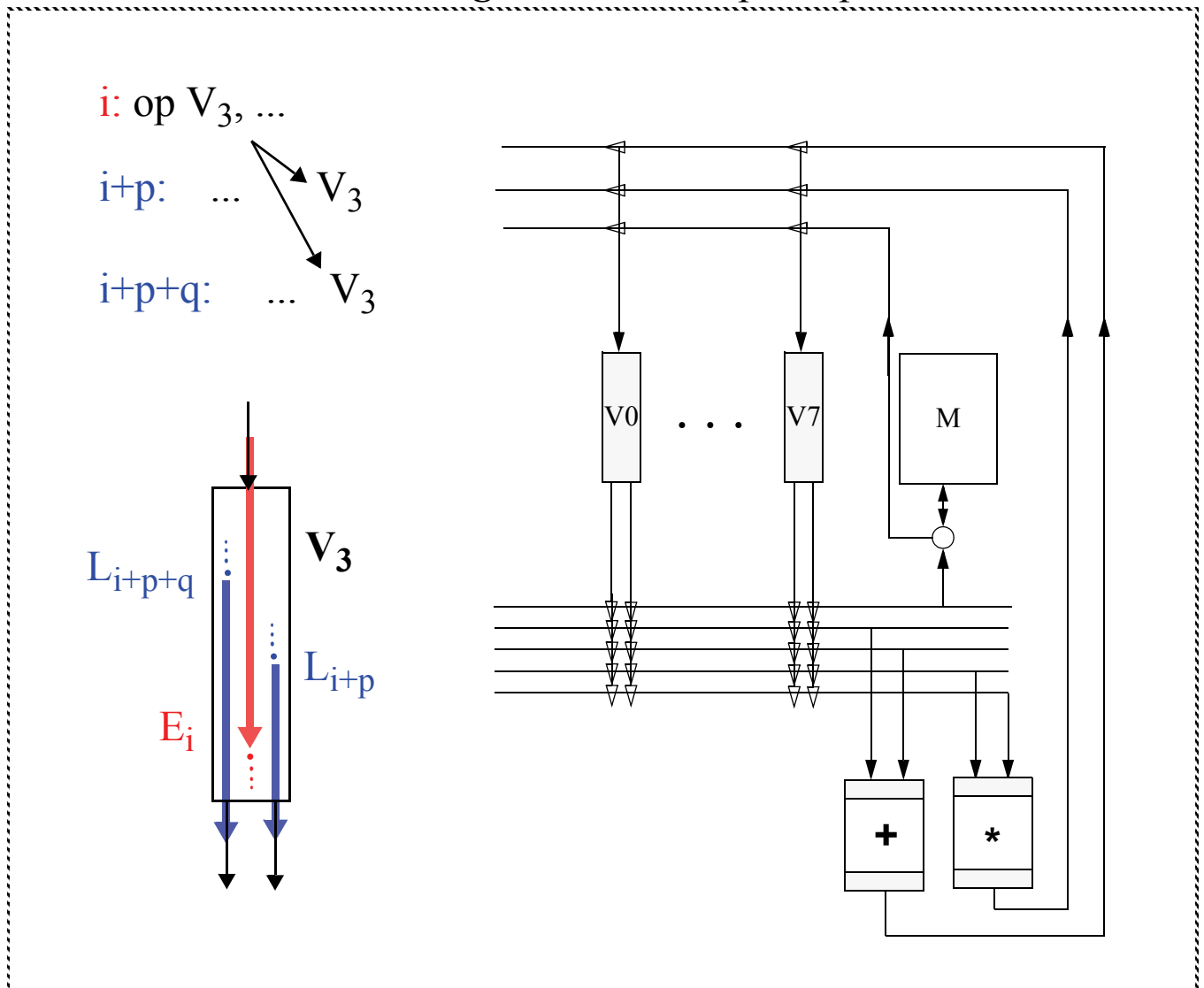
✓ Regla: último D = último W

- ❑ 2 UF's
 - ✓ ALU para sumar
 - ✓ ALU para sumar/multiplicar
- ❑ 3 buses de Escritura y 5 de Lectura



- ❑ Dos ALUs dedicadas
- ❑ **Dos** puertos de lectura y **uno** de escritura por reg. vectorial. La Unidad de Control permite dos flujos de lectura y uno de escritura en cada registro vectorial

Encadenamiento general → solapar *dependencias* datos



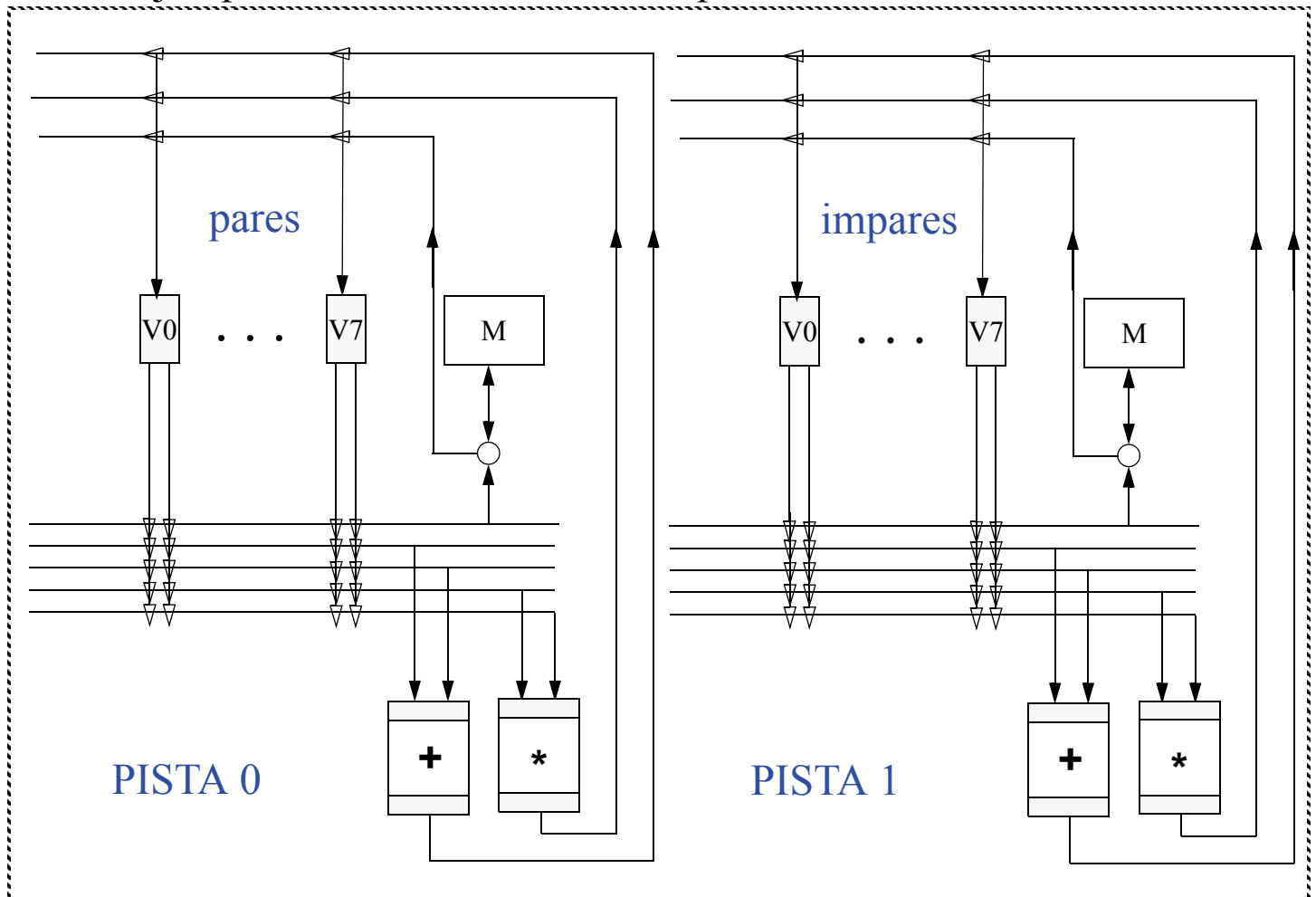
- ❑ Las *antidependencias* tampoco bloquean !!
- ❑ Temporización riesgos de datos productor-consumidor:

[illegible]

Víctor Viñals, Jesús Alastruey

- ❑ Replicar los caminos de proceso y memoria.

Ejemplo: el Procesador 4 con 2 *parallel lanes*



- ❑ El Banco de Registros Vectorial (BRV) y la memoria se *particionan*: pares en una partición e impares en la otra.
- ❑ Cada una de las dos ALUs (FP add, FP mul) se *duplican*.
- ❑ En el primer ciclo de ejecución entran en paralelo a las dos pistas los elementos (0, 1) de los registros vectoriales fuente. En el segundo ciclo entran los elementos (2, 3), y así sucesivamente ...
- ❑ Comparando con el Procesador 4 de una pista, el número de ciclos por elemento, C_e , se divide por el número de pistas.

código
ejemplo

lv	V0,Ra
mulv	V1,V2,V3
addv	V4,V5,V6

GFLOPS ?

- El rendimiento depende de n , el tamaño del vector
supongamos $n \leq \text{MVL}^1$

$$✓ \quad C_n (n \leq \text{MVL}) = C_{\text{fijos}} + n \cdot C_e \text{ ciclos}$$

C_{fijos} = latencia UFs + penalizaciones

penal. = riesgos estructurales y dependencias

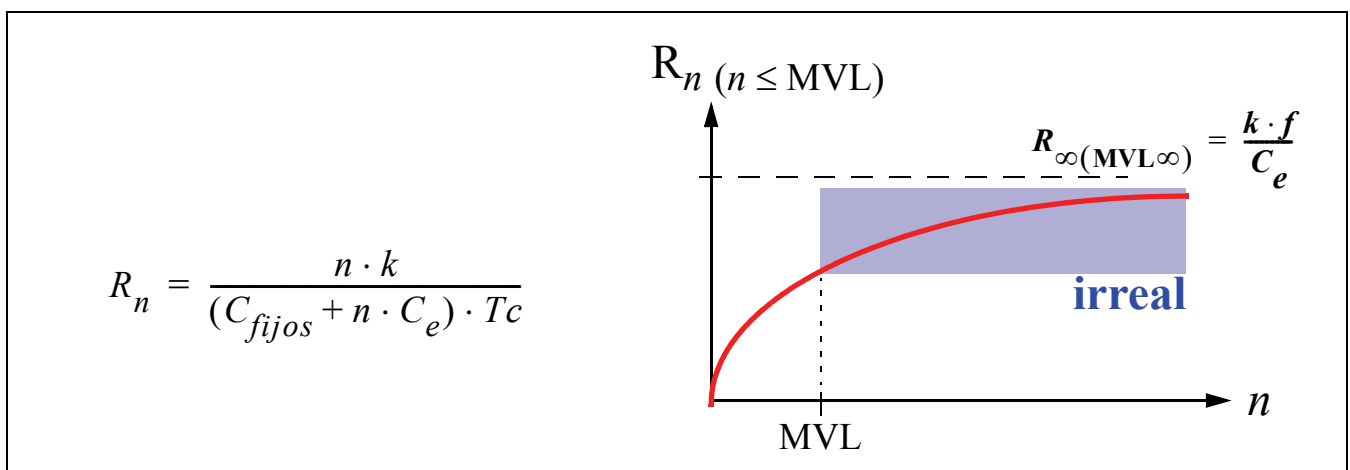
C_e = ciclos por elemento

$$✓ \quad T_n (n \leq \text{MVL}) = C_n \cdot T_c$$

- R = velocidad en FLOPS; se aplican k FLOPs a cada elemento

$$✓ \quad R_n (n \leq \text{MVL}) = n \cdot k / T_n = n \cdot k \cdot F / C_n$$

R en *GFLOPS* \rightarrow T_c en *ns* o F en *GHz*

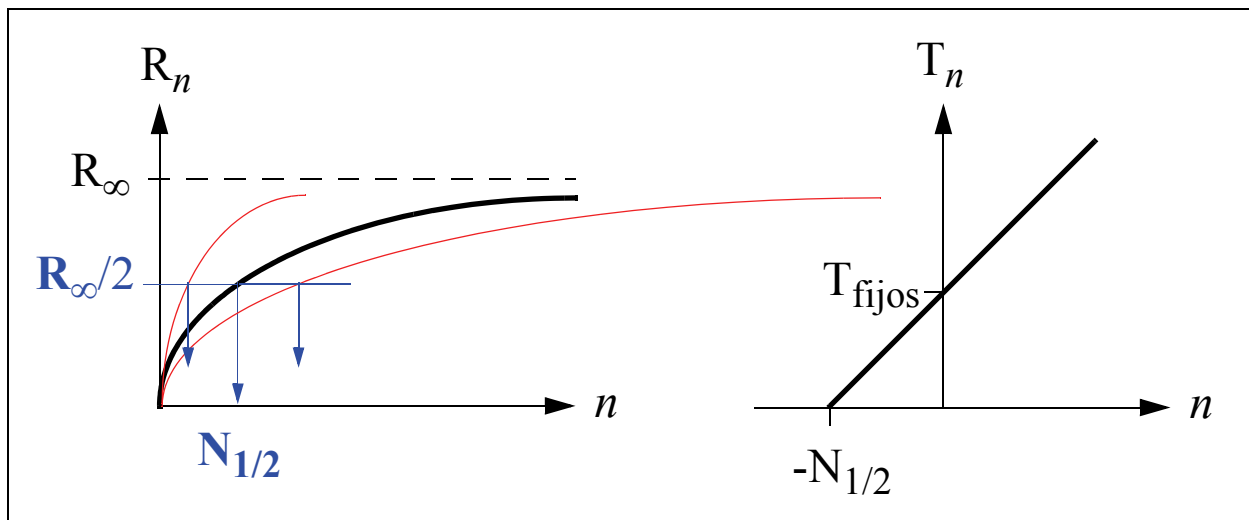


1. MVL: Maximum Vector Length, número de elementos de un Registro Vectorial.

□ $N_{1/2}$

longitud de vector que consigue
MITAD del rendimiento máximo

$$R_{N_{1/2}} = R_{\infty}(MVL_{\infty})/2 \rightarrow N_{1/2} = C_{\text{fijos}}/C_e = T_{\text{fijos}}/T_e$$



- ✓ Da idea de la *sobrecarga* vectorial (latencias y penalizaciones)
- ✓ es independiente de T_c
sólo depende de la arquitectura y del algoritmo

□ N_v

Longitud mínima de vector que consigue
igual velocidad
en las versiones escalar y vectorial

$$R_{N_v} = R_{\text{esc}} \rightarrow N_v = C_{\text{fijos}} / (C_{\text{esc}} - C_e) = T_{\text{fijos}} / (T_{\text{esc}} - T_e)$$

- ✓ mide la *sobrecarga* vectorial y la velocidad relativa entre los procesadores vectorial y escalar

- ❑ Supongamos los siguientes parámetros para los cinco procesadores anteriores:

- ✓ $T_c = 1 \text{ ns}$
 - ✓ $\text{Lat } \{ +, *, \text{mem} \} = 6, 7, 12 \text{ ciclos}$
 - ✓ $\text{LI} = 1 \text{ ciclo}$
 - ✓ Segmentación página 39:
 - aritméticas, loads: $B, D^n, R, \text{Lat}, W^{\text{VLR}}$
 - stores: $B, D^n, R, W^{\text{VLR}}, \text{Lat}$
- $n = 1$ si no hay riesgo estructural ni dependencia;
 $n > 1$ en caso contrario
 $\text{VLR} = \text{valor del registro } \text{VLR} \in \{0..64\}$

- ❑ Vamos a calcular

$$C_n \ (n \leq \text{MVL})$$

$$R_{64} \ (\text{MVL} = 64)$$

$$R_{\infty} \ (\text{MVL} \infty)$$

$$R_{\text{pico}}$$

C_n se calcula con reg. vect de n elem.

R_{64} con $n = 64$

R_{∞} con $n = \infty \rightarrow R_{\infty} (MVL \infty)$

Proc i
(Rpico, GFLOPS)

Código

C_n ciclos

R_{64} GFLOPS

R_{∞} GFLOPS

Cód 1

Cód 2

Cód 3

Cód 4

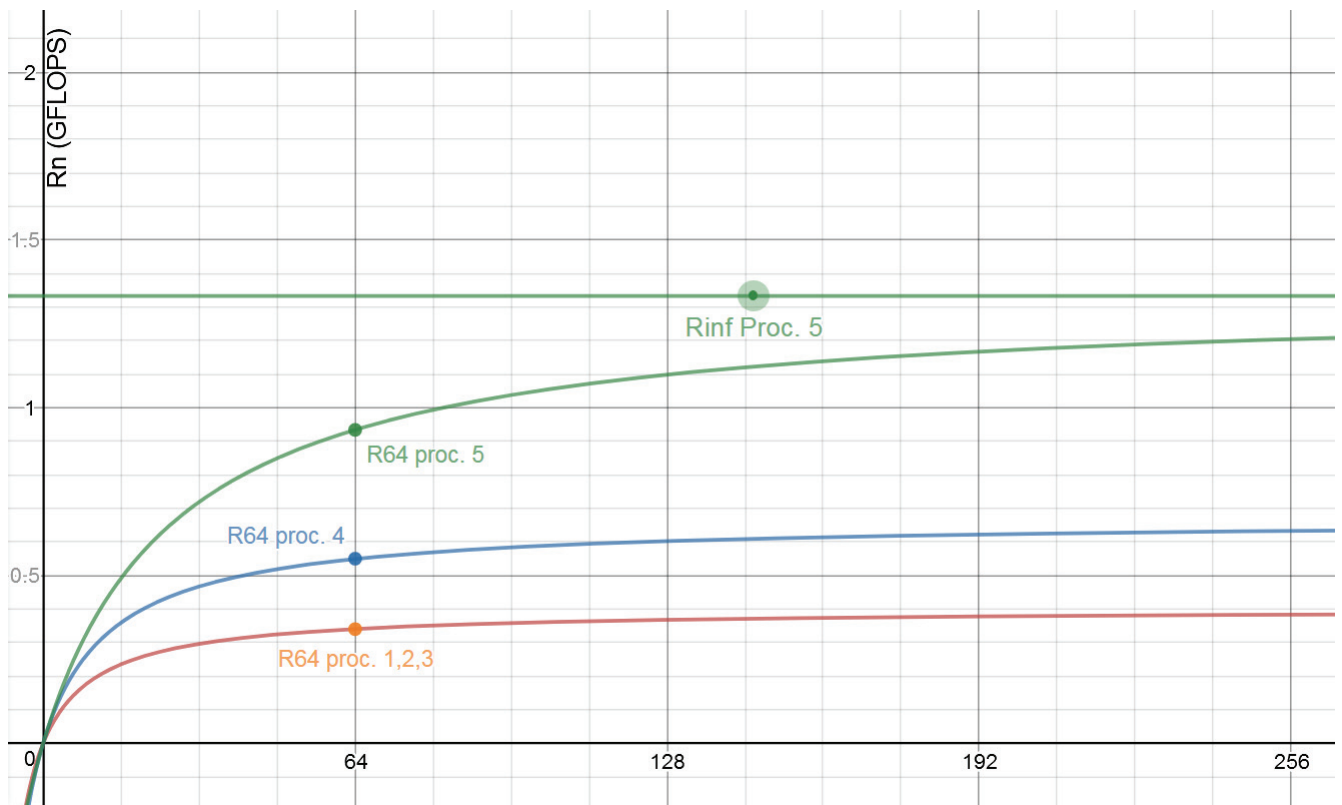
```
lv  V0, Ra
mulv V1, V2,V3
addv V4, V5,V6
```

```
mulv V1, V2,V3
addv V4, V5,V2
```

```
lv  V0, Ra
mulsv V1, F0,V0
addv V2, V2,V1
```

```
lv  V0, Ry
lv  V1, Rx
mulsv V2, Fa,V1
addv V3, V0,V2
sv  Ry, V3
```

	Proc1 (1)	Proc2 (1)	Proc3 (2)	Proc4 (2)	Proc5 (4)
Cód 1	$30 + 3n$ 0.58 0.66	$18 + 2n$ 0.88 1	$15 + n$ 1.62 2	=	?
Cód 2	$17 + 2n$ 0.88 1	=	=	$10 + n$ 1.73 2	?
Cód 3	no	no	no	$30 + n$ 1.36 2	?
Cód 4	$56 + 5n$ 0.34 0.4	=	=	$41 + 3n$ 0.55 0.66	?



(página en blanco)

2.6 ZV: Organización de un procesador vectorial segmentado con ALMa tipo DLXV

Un procesador vectorial de una arquitectura tipo DLXV se segmenta en 8 etapas para las instrucciones vectoriales básicas: `aluV`, `LV`, `SV`. El trabajo en cada etapa para cada instrucción se muestra en la tabla:

		<code>aluV</code>	<code>LV</code>	<code>SV</code>
B	búsqueda instrucción			
D1	deco 1/ @dst salto			
D2	deco 2			
L	¿riesgos en recursos? ¿riesgos en registros? lanzamiento y reserva recursos	lectura reg vectoriales	<i>nada</i>	lectura reg vectoriales
X1	red ida	paso crossbar 1	paso crossbar 1 generación @	
alu/ mem	operación o acceso a memoria	op	mem_rd	mem_wr liberar recursos en último ciclo
X2	red vuelta	paso crossbar 2		
W	escritura	escritura registro vectorial destino liberar recurso en último ciclo		

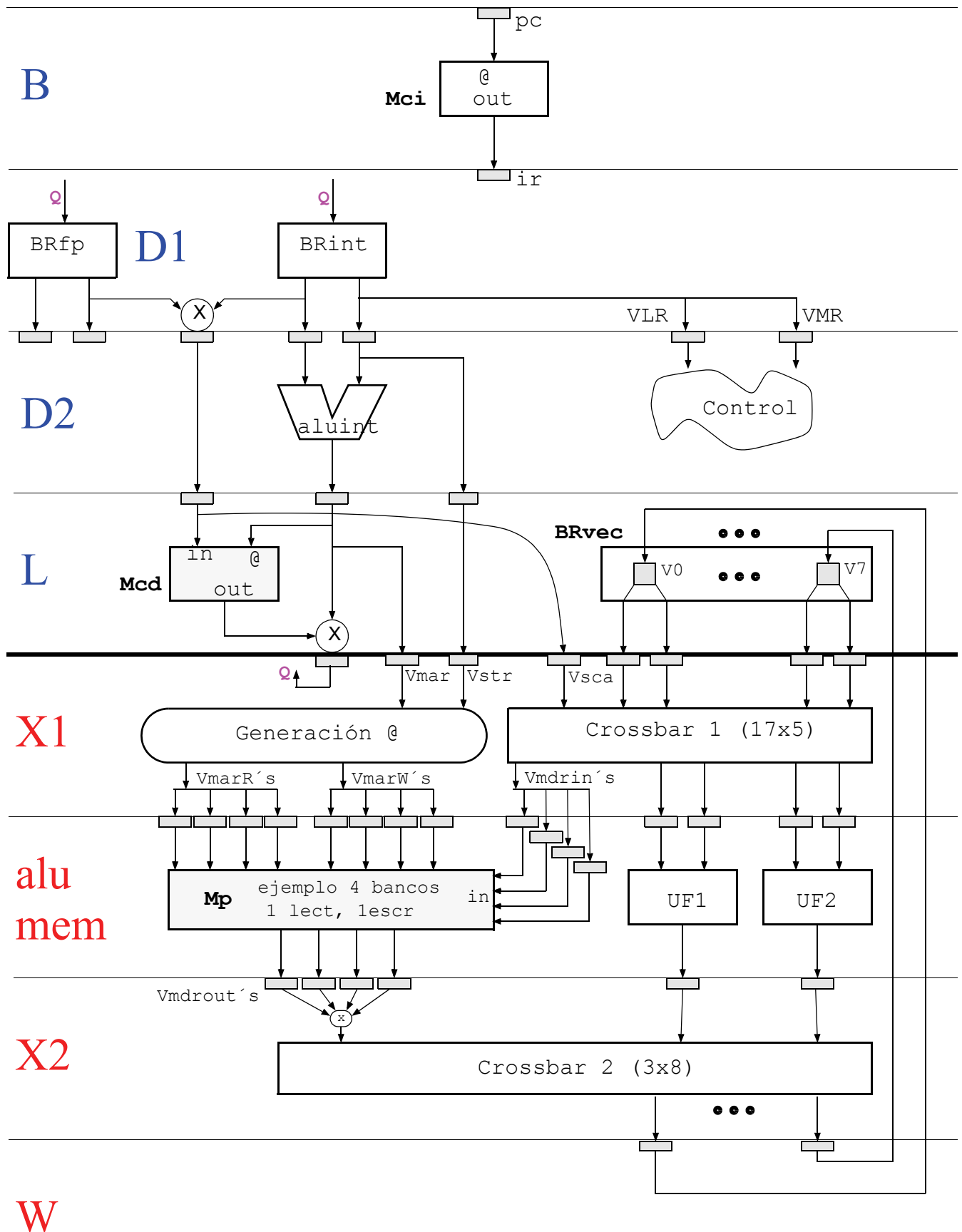
La ruta de datos se muestra en la siguiente hoja. La parte izquierda corresponde a un trozo del procesador escalar. Su operación y su relación con la parte vectorial puede deducirse observando la ruta de datos.

Hay dos unidades de cálculo multioperación (UF1 y UF2) y una memoria *multibanco de dos puertos*, uno de escritura y otro de lectura. O sea, cada banco de memoria soporta una lectura y una escritura simultáneas, con la misma latencia que una lectura o una escritura por separado. La latencia de la primera lectura no depende del banco de inicio porque el *acceso es desfasado*.

No hay control hardware de colisión entre lectura y escritura en memoria. Se solapan los flujos y el compilador ya se preocupará de *no* generar código vectorial si las *dependencias en memoria* pueden dar lugar a riesgos (RAW, WAR o WAW).

Cada registro vectorial tiene dos puertos de lectura y uno de escritura. La unidad de control permite ejecutar de forma solapada y segura a instrucciones dependientes en cualquier caso:

- Dependencia prod-cons. Encadenamiento general.
- Antidependencia. Flujo de escritura por detrás del de lectura.



En una instrucción vectorial de tamaño mínimo (VLR = 1) se ocupa cada etapa un ciclo, salvo alu/mem, que se ocupan un número de ciclos igual a la latencia de las unidades funcionales o de memoria. Para VLR>1 todas las etapas por debajo de L pueden ser multiciclo. Los siguientes ejemplos muestran la ocupación de las etapas para **VLR=6**:

Ejemplo 1: ADDV, Lat_sum = 2, Tasa iniciación = 1

B	x														
D1		x													
D2			x												
L				x	x	x	x	x	x						
X1					x	x	x	x	x	x					
+						x	x	x	x	x	x	x			
X2						Lat_sum		x	x	x	x	x	x		
W									x	x	x	x	x	x	x

resumen:

B	D1	D2	L	X1	+ 2c	X2	W 6c
----------	-----------	-----------	----------	-----------	-------------	-----------	-------------

:

Ejemplo 2: LV, Lmem = 4, sin conflictos de banco

B	x														
D1		x													
D2			x												
L				x											
X1					x	x	x	x	x	x					
mem						x	x	x	x	x	x	x	x	x	
X2						Lat_mem				x	x	x	x	x	
W											x	x	x	x	x

resumen:

B	D1	D2	L	X1	mem 4c	X2	W 6c
----------	-----------	-----------	----------	-----------	---------------	-----------	-------------

:

Ejemplo 3: SV, Lmem = 4, sin conflictos de banco

B	x														
D1		x													
D2			x												
L				x	x	x	x	x	x		Lat_mem				
X1					x	x	x	x	x	x					
mem						x	x	x	x	x	x	x	x	x	x
X2															
W															

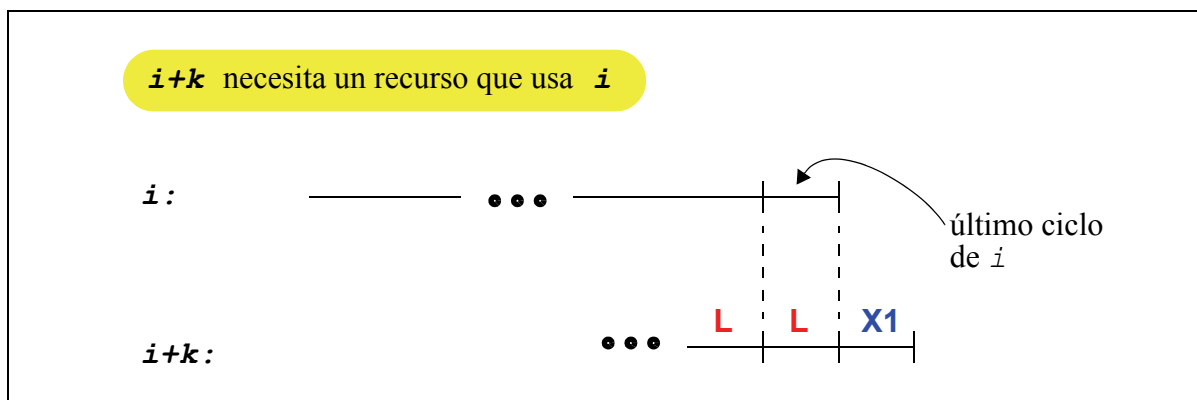
resumen:

B	D1	D2	L	X1 6c	mem 4c
----------	-----------	-----------	----------	--------------	---------------

Este procesador lanza en orden hasta 1 instrucción vectorial o escalar por ciclo. Instrucciones escalares y vectoriales independientes pueden solaparse. Pueden ejecutarse a la vez hasta cuatro instrucciones vectoriales (1 LV, 1 SV y 2 aluV)

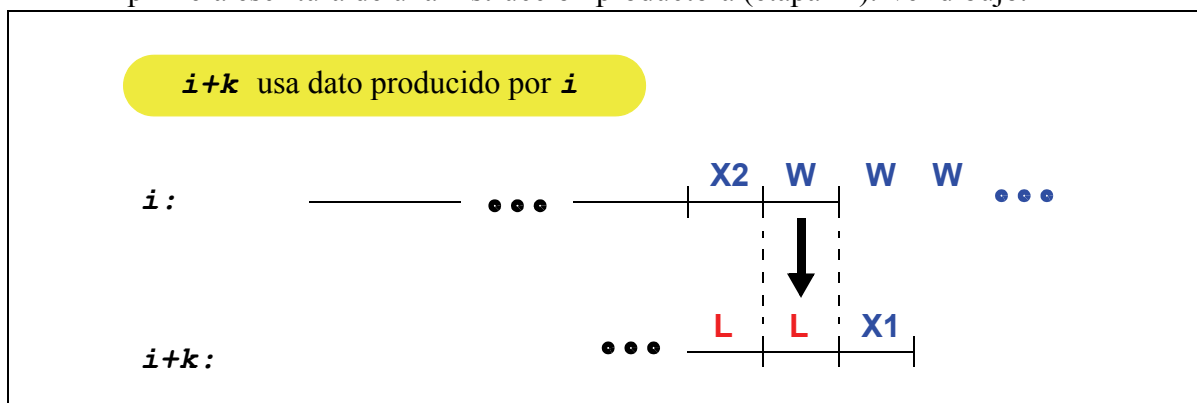
La etapa **L** se encarga de controlar los riesgos en las instrucciones vectoriales por dependencias de datos (en registros o memoria) y por falta de recursos:

- **RECURSOS** (UFs, puertos de memoria, puertos de registros, ...). Se reservan en la etapa **L** y se liberan en el *último* ciclo de ejecución. En caso de bloqueo por recurso, el último ciclo de detención coincide con el ciclo de liberación. Ver dibujo.



- **DEPENDENCIAS prod/cons en REGISTROS**. Existe encadenamiento **general**. No existen restricciones, incluso Loads y Stores encadenan.

Puede solaparse la primera lectura de una instrucción consumidora (etapa **L**) con la primera escritura de una instrucción productora (etapa **W**). Ver dibujo.



El procesador ZV de referencia tiene las siguientes características:

MVL	64	tamaño registros vectoriales
+	Lsum = 2c, TI = 1c	Latencia, Tasa de Iniciación
*	Lmul = 3c, TI = 1c	
mem	8 bancos, Lmem = 4c	sistema sobrado
Tc (f)	1ns (1000 Mhz = 1GHz)	
Coherencia	entre Mcd y Mp asegurada por hardware (ya veremos)	

3. DOS ASPECTOS DE PROGRAMACIÓN: *VECTOR LENGTH Y VECTOR STRIDE*

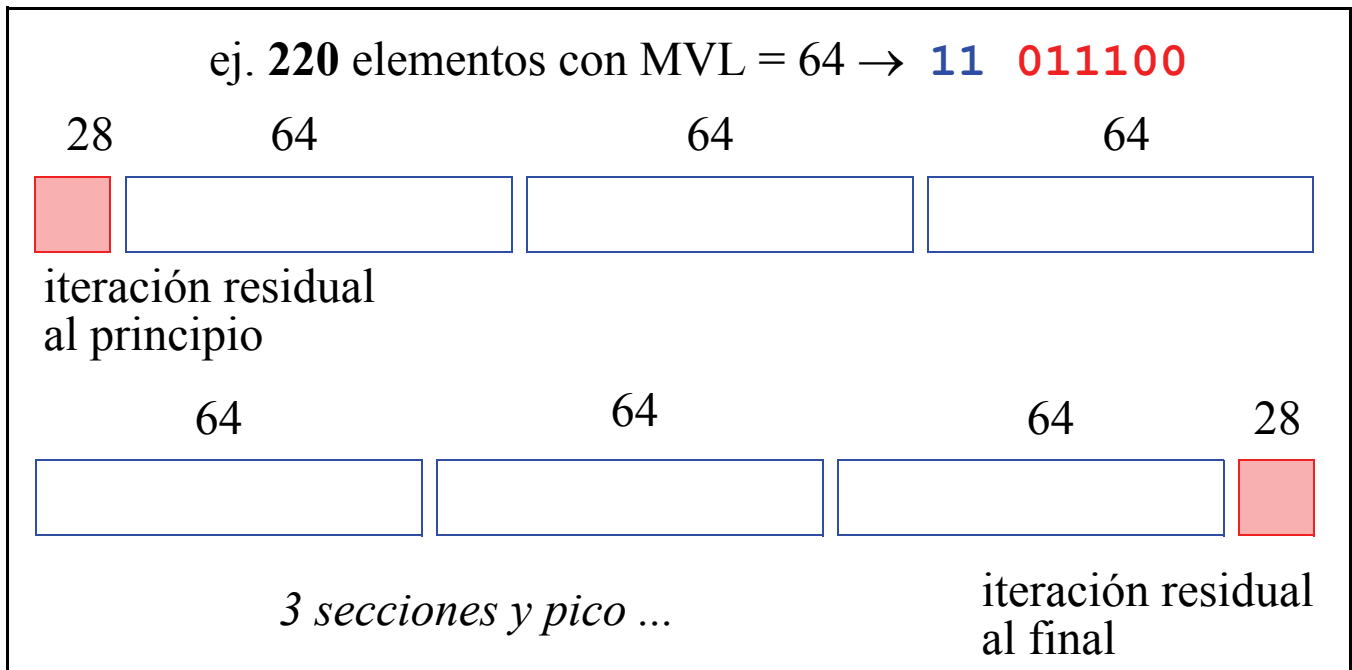
3.1 VECTOR LENGTH

Vectores de n elementos, pero
Registros vectoriales de MVL elementos:

SOLUCIONES:

- ❑ $n \leq 64$
se programa VLR con n (instr. DLXV movi2s)
- ❑ $n > 64$

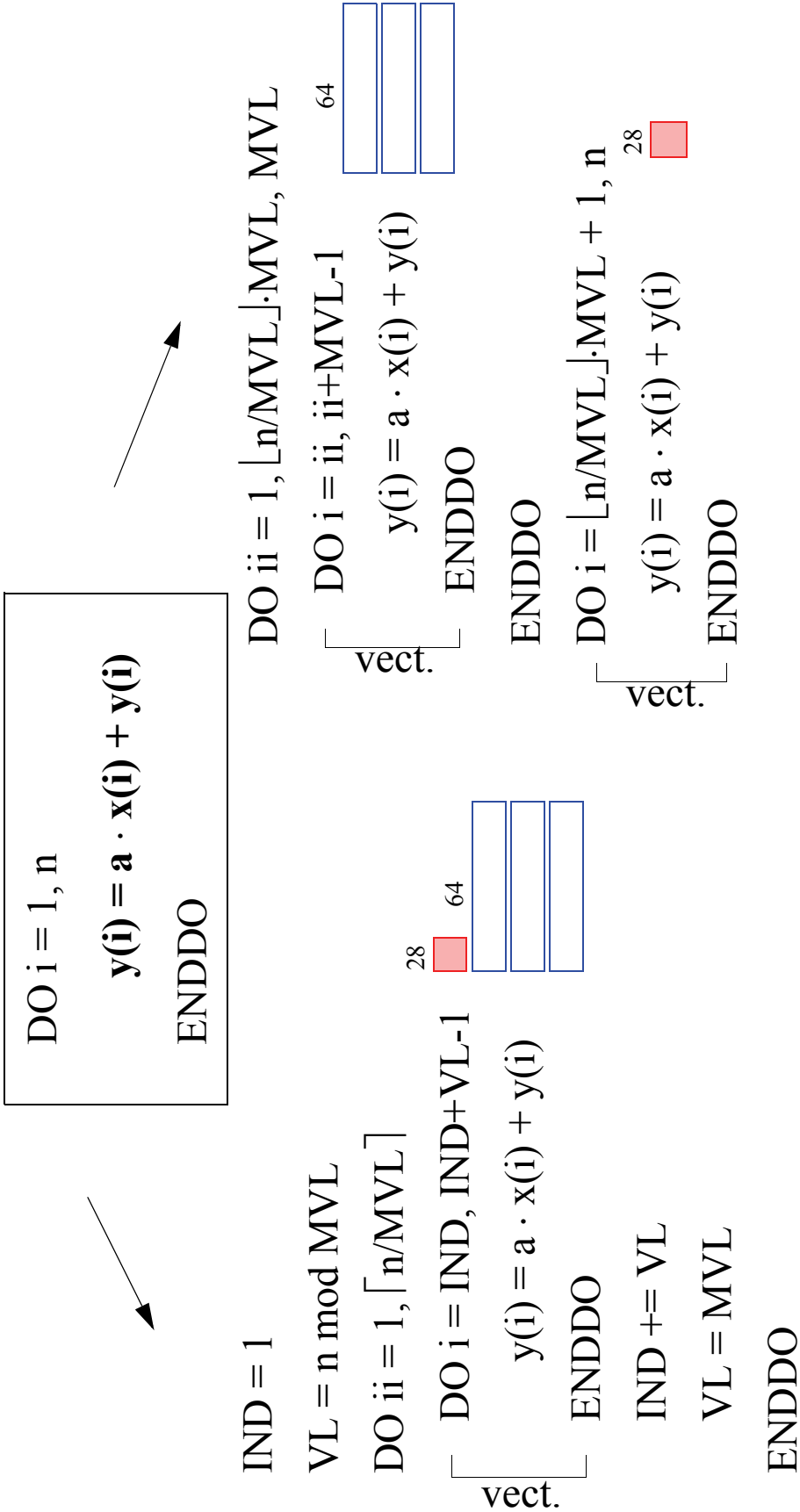
Técnica de Seccionado o *Strip-Mining*



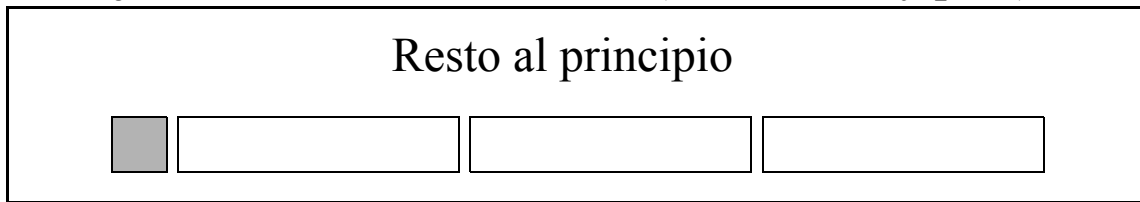
- ❑ Módulo y división entera por MVL (potencia de dos):
 - Operaciones con máscara, desplazamientos
 - Soporte especial: p.e. repertorio de Convex¹

1. Paqui Quintana, R. Espasa y M. Valero. "An ISA comparison between Superscalar and Vector Processors". En Proc. of the Int. Conf. on Vector and Parallel Processing (VECPAR98). Lecture Notes in Computer Science, issue 1573, pp. 548-560, 1999. Springer Verlag publisher.

Código AXPY (SAXPY o DAXPY: simple o doble prec.)



Bucle ejecutado en modo vectorial (3 secciones y pico)



Código AXPY seccionado con MVL = 64

```
1  and  Rmod, Rn, R3F
2  srl  Rent, Rn, #6
3  movi2s VLR, Rmod
```

$R_x = \&X[0];$

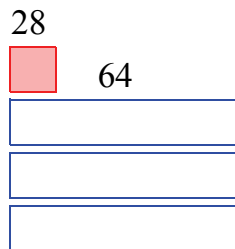
$Fa = a$

$R_y = \&Y[0];$

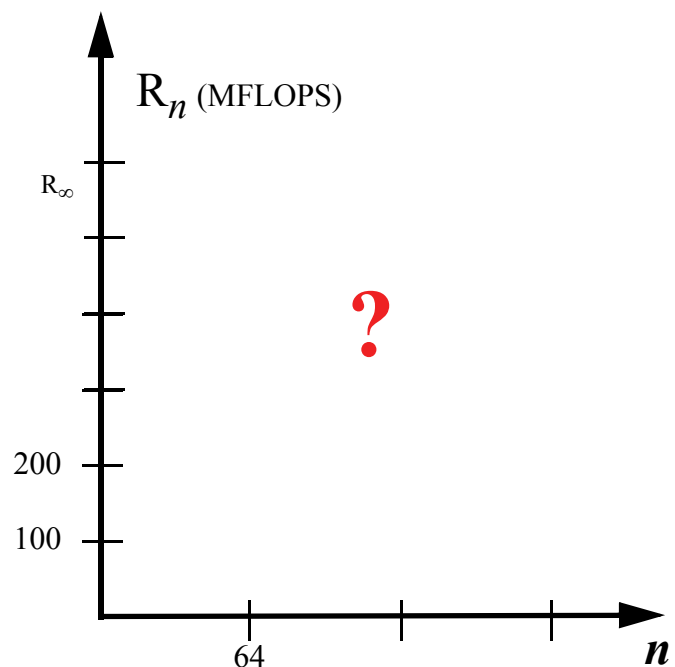
$R_n = n;$

$R_{3F} = 0x3F$

```
buc: LV    V0, Rx
      5 MULSV V1, Fa, V0
      6 LV    V2, Ry
      7 ADDV  V3, V2, V1
      8 SV    Ry, V3
```



```
9  movs2i  Rn, VLR
10  bz      out, Rent
11  nop
12  sll     Rsize, Rn, #3
13  add     Rx, Rx, Rsize
14  add     Ry, Ry, Rsize
15  movi2s  VLR, #64
16  jmp     buc
out: dec  Rent
```



Rendimiento de ZV
en función de n ?

Modelo de ejecución

- ❑ C_{base} ciclos del prólogo y del epílogo:
 - instr. escalares {1:3}
 - tiempo desde último bloque n hasta final
- ❑ C_{fij} costes fijos (incluso con $VLR = 1$)
 - Latencias + Penalizaciones
- ❑ C_{bucle} ciclos de control de bucle
 - instr. escalares {9:17}
 - en paralelo con las vectoriales anteriores
 - ✓ ciclos
 - desde que se lanza la última instr. vect. de una sección hasta que se lanza la primera de la sección siguiente.
 - ✓ a veces, $C_{bucle} = 0$,
 - porque la primera instr. vect. de la sección siguiente está detenida por recursos o dependencias.

$$C_n = C_{base} + \lceil n/MVL \rceil \cdot (C_{fij} + C_{bucle}) + n \cdot C_e$$

- ❑ A partir de C_n , podemos derivar
 - $R_n, R_{64}, R_{\infty}, \dots$
 - $N_{1/2}, N_v$

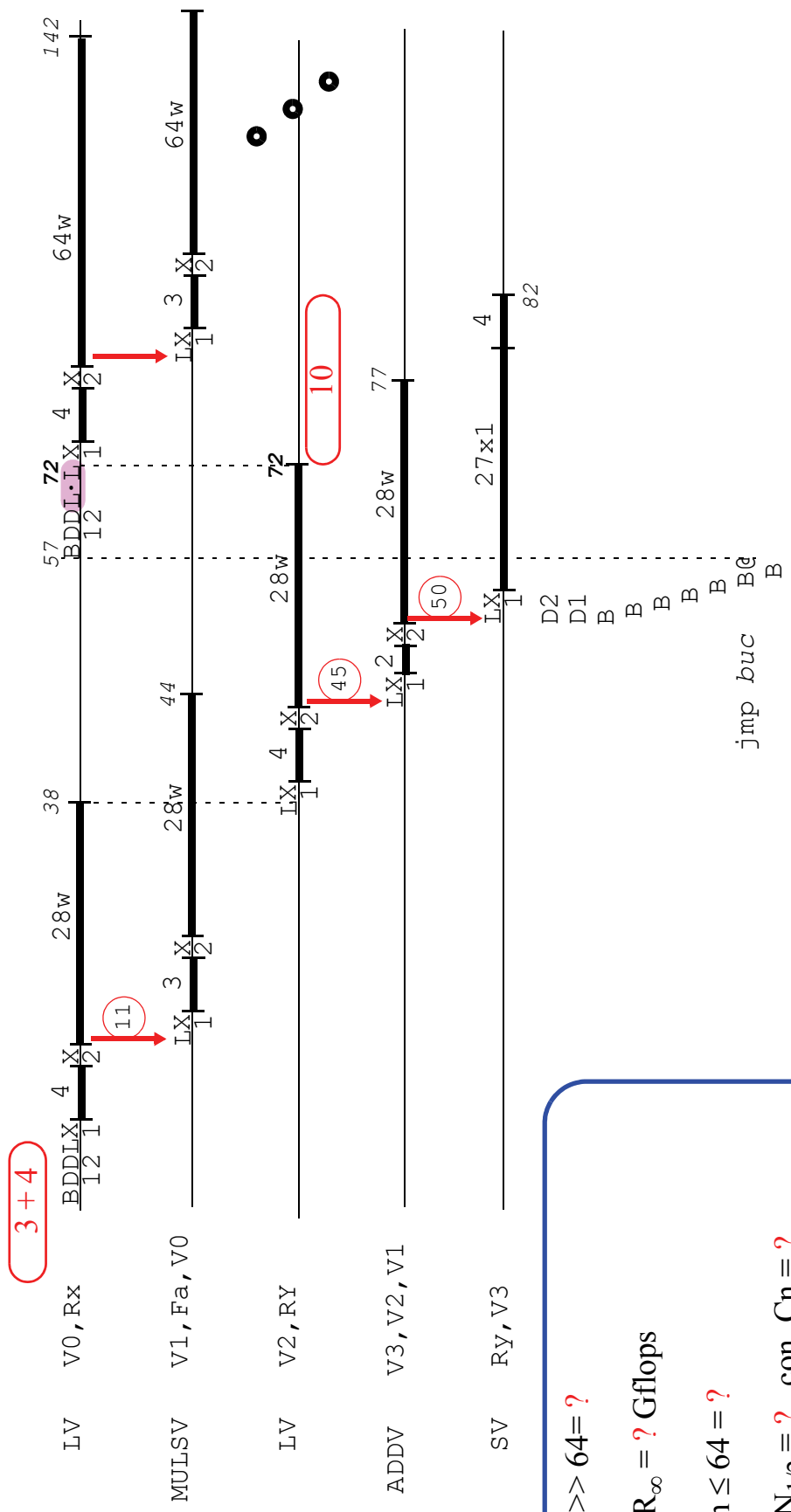
$\left[\begin{matrix} n \\ \end{matrix} \right]$	$Rn(GFLOPS) = \frac{FLOPn}{Cn} \times F(GHz)$
---	---

Ejercicio

obtener fórmulas para $N_{1/2}$ y N_v

Hipótesis: son menores que MVL

Ejecución detallada de AXPY de 220 elementos en ZV



Cn, n>> 64=?
R_∞ = ? Gflops
Cn , n ≤ 64 = ?
N_{1/2} = ? con Cn = ?
Cn, ∀n = ?

The diagram illustrates a 142-cycle pipeline with the following stages and instructions:

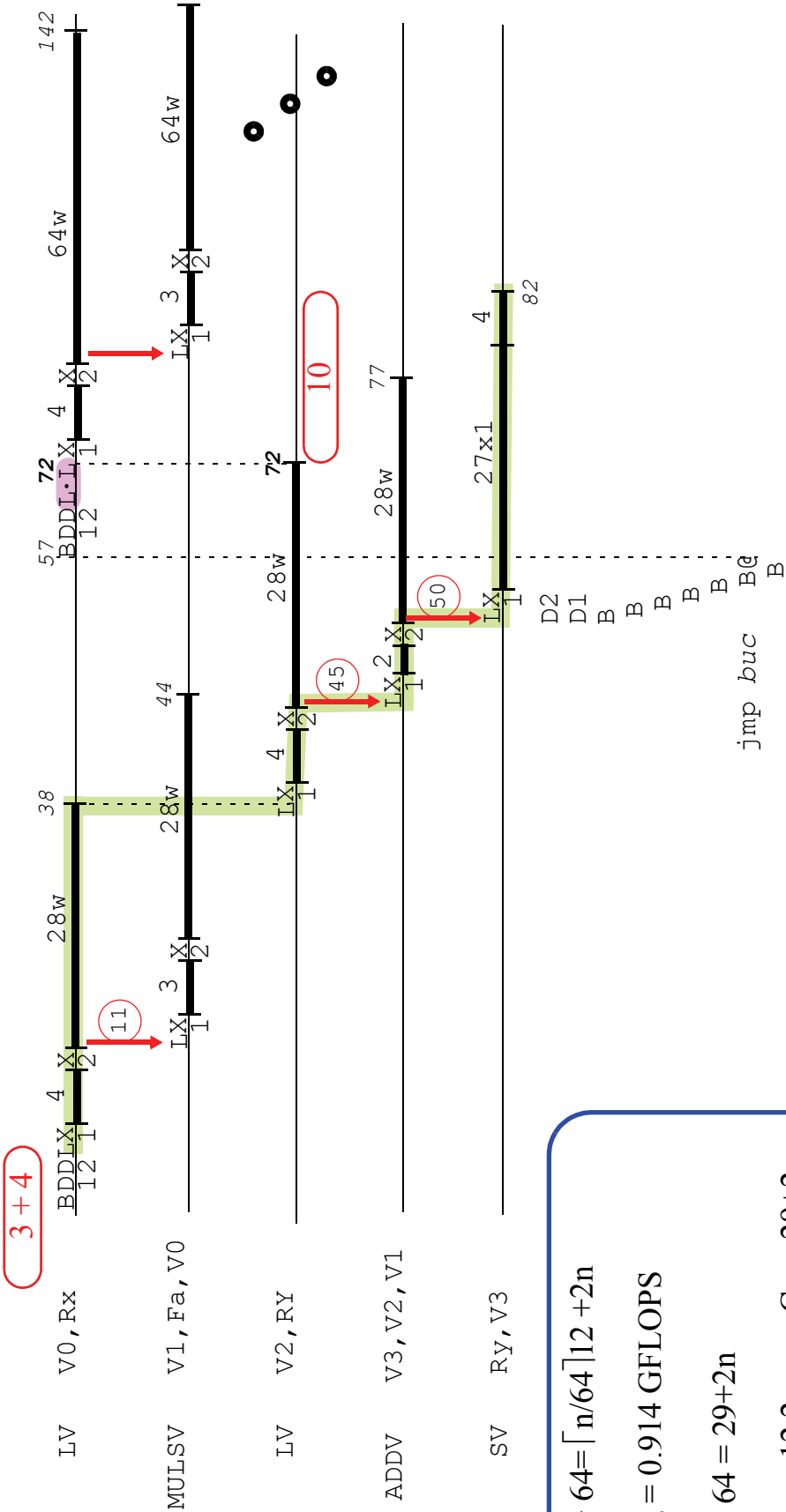
- V0, Rx**: Instruction LX (12 cycles), BDDI (12 cycles), LX (4 cycles). Total: 28w.
- MULSV V1, Fa, V0**: Instruction LX (12 cycles), BDDI (12 cycles), LX (4 cycles). Total: 28w.
- V2, RY**: Instruction LX (12 cycles), BDDI (12 cycles), LX (4 cycles). Total: 28w.
- ADDV V3, V2, V1**: Instruction LX (12 cycles), BDDI (12 cycles), LX (4 cycles). Total: 28w.
- SV Ry, V3**: Instruction LX (12 cycles), BDDI (12 cycles), LX (4 cycles). Total: 28w.

Summary box:

- $64 = \lceil n/64 \rceil 12 + 2n$
- $_{\circ} = 0.914 \text{ GFLOPS}$
- $\leq 64 = ?$

$$C_n, n > 64 =$$

Ejecución detallada de AXPY de 220 elementos en ZV



$$Cn, n \gg 64 = \lceil n/64 \rceil 12 + 2n$$

$$R_{\infty} = 0.914 \text{ GFLOPS}$$

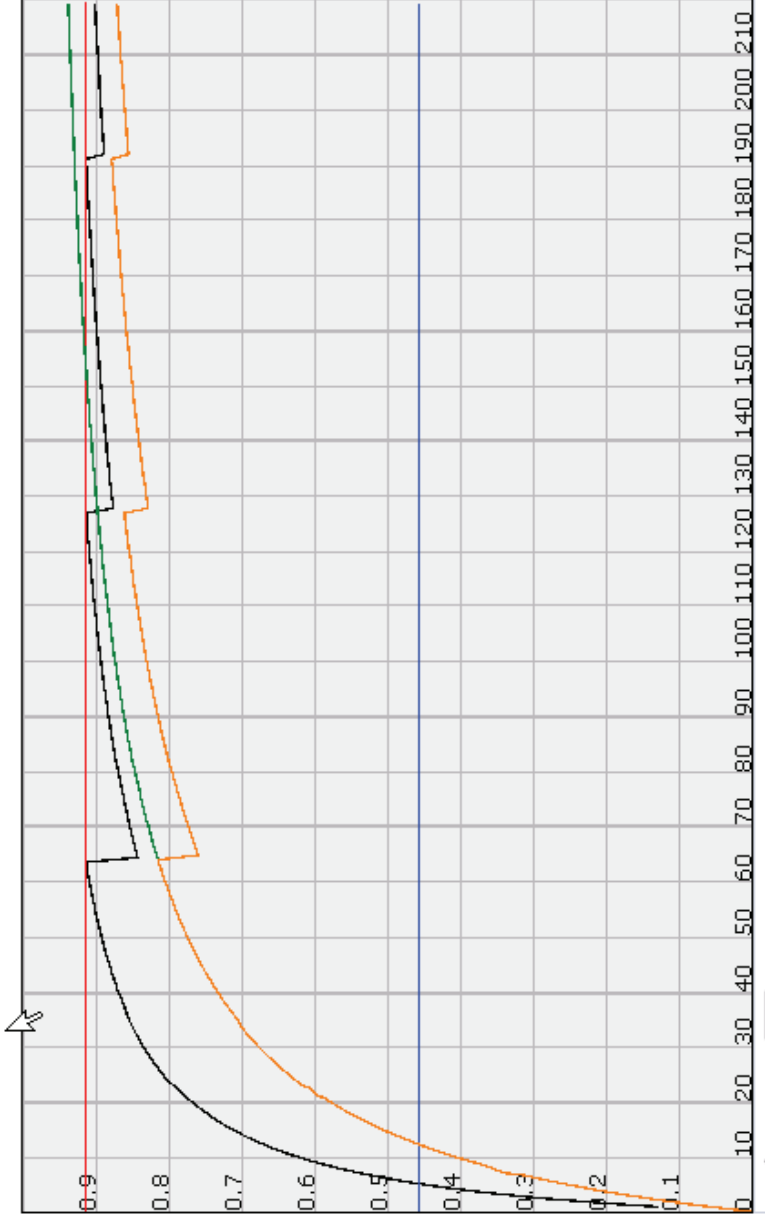
$$Cn, n \leq 64 = 29 + 2n$$

$$N_{1/2} = 12.2 \quad \text{con} \quad Cn = 29 + 2n$$

$$Cn, \forall n = 7 + \lceil n/64 \rceil 12 + 2n + 10$$

$$Cn, n > 64 =$$

$$Cn, n \leq 64 =$$



2-D plotting

Aqui puede graficar funciones, funciones en coordenadas polares y curvas parametricas.

Función

$y(x) = 2 * x / (\text{ceil}(x / 64) * 12 + 2 * x)$

Función

$y(x) = 0.914$

Función

$y(x) = 0.914 / 2$

Función

$y(x) = 2 * x / (29 + 2 * x)$

Función

$y(x) = 2 * x / (17 + \text{ceil}(x / 64) * 12 + 2 * x)$

permalink:

http://fooplot.com/index.php?type0=0&type1=0&type2=0&type3=0&type4=0&y0=2*x/%28ceil%28x/64%29*12%2B2*x%29&y1=0.914&y2=0.914/2&y3=2*x/%2829%20%2B%202*x%29&y4=2*x/%2817%2Bceil%28x/64%29*12%2B2*x%29&r0=&r1=&r2=&r3=&r4=&px0=&px1=&px2=&px3=&px4=&py0=&py1=&py2=&py3=&py4=&smin0=0&smin1=0&smin2=0&smin3=0&smin4=0&smax0=2pi&smax1=2pi&smax2=2pi&smax3=2pi&smax4=2pi&thetamin0=0&thetamin1=0&thetamin2=0&thetamin3=0&thetamin4=0&thetamax0=2pi&thetamax1=2pi&thetamax2=2pi&thetamax3=2pi&thetamax4=2pi&ipw=1&ixmin=0&ixmax=220&iymin=0&iymax=1&igx=10&igy=0.1&igl=1&igs=0&iia=1&ila=1&xmin=0&xmax=220&ymin=0&ymax=1

- ❑ Procesado de elementos no contiguos
a separación constante (*stride* = paso)

```
for (i=0; i<n; i++)
```

```
  for (j=0; j<n; j++)
```

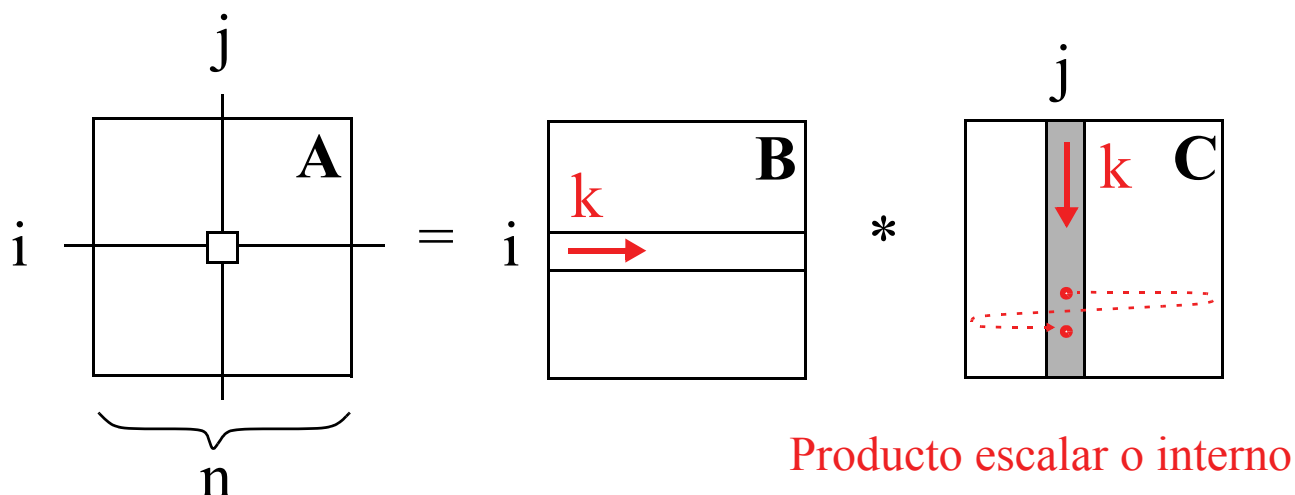
```
    { A[i][j] = 0.0;
```

$$\mathbf{A} = \mathbf{B} * \mathbf{C}$$

```
      for (k=0; k<n; k++)
```

```
        A[i][j] = A[i][j] + B[i][k]*C[k][j];
```

```
    }
```



- ❑ $B[i][k] * C[k][j]$ es vectorizable en K si cargamos
 $C[0:n-1][j]$ en un reg. vectorial (si $n > 64$ secc.)

- ❑ $s = n * 8$ bytes = separación en un registro **R**

```
LVWS V1, (R1, R2)  →  for (i=0; i<VLR; i++)
                           V1.i = mem (R1 + i*R2);

SVWS (R1, R2), V1
```

DLXV

□ Vectorización de la multiplicación de matrices en forma *ijk*:

LV V0, R_B

LVWS V1, (R_C, R_{n*8})

MULV V2, V0, V1

RED_SUM R_{tmp}, V2 ; $R_{tmp} = R_{tmp} + \sum_{i=0}^{VLR} V2.i$

st R_A, R_{tmp}

4. CONFLICTOS EN EL ACCESO A BANCOS DE MEMORIA

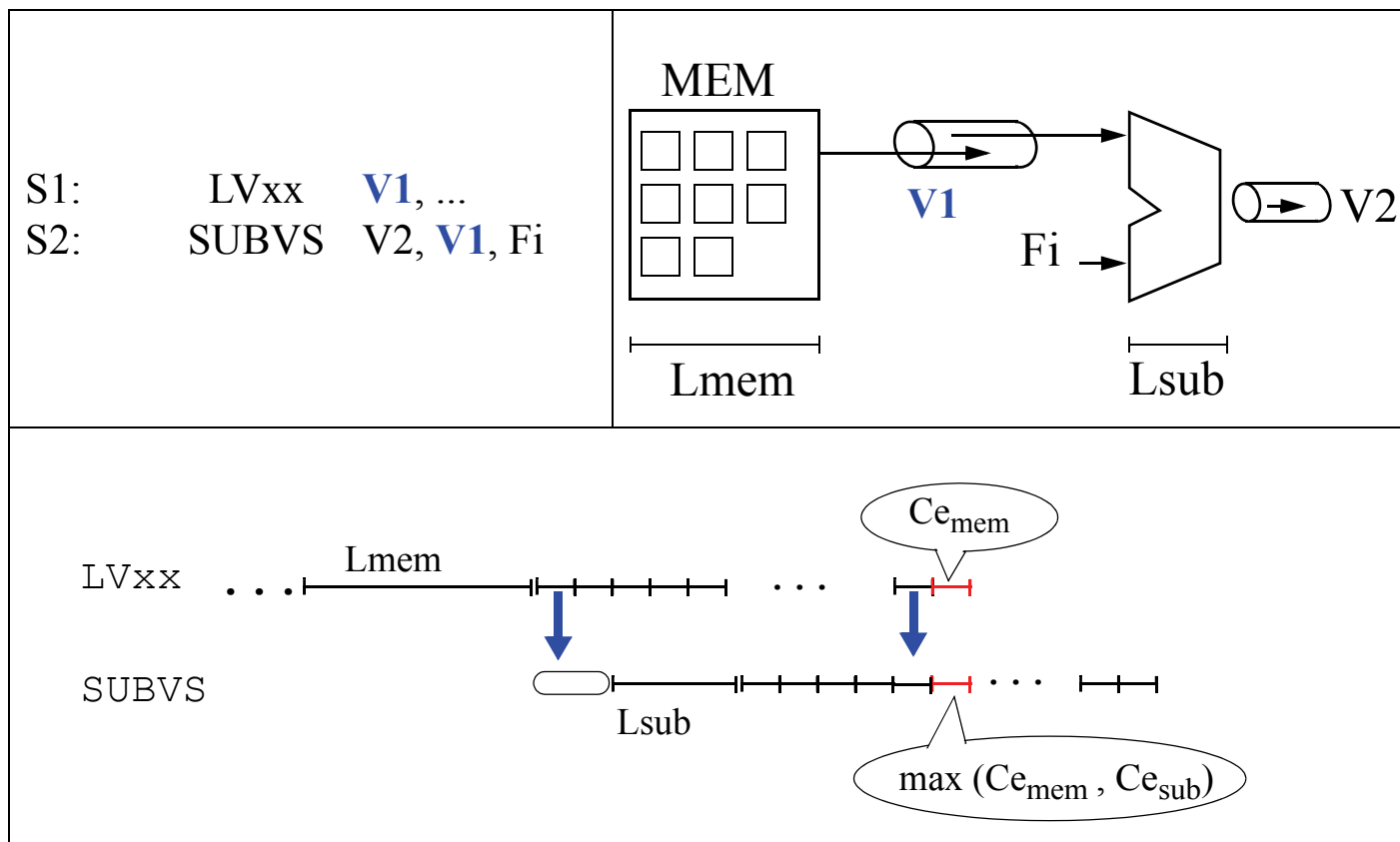
4.1 INTRODUCCIÓN

4.2 SISTEMAS AJUSTADOS

4.3 SISTEMAS SOBRADOS

4.1 INTRODUCCIÓN

Memoria = carga y descarga de registros vectoriales



LV V1, R@ $\rightarrow C_{\text{mem}} = 1 \text{ ciclo/dato}$

LV**WS** V1, (R@,Rs) $\rightarrow C_{\text{mem}} = ?$

ALMACENAMIENTO

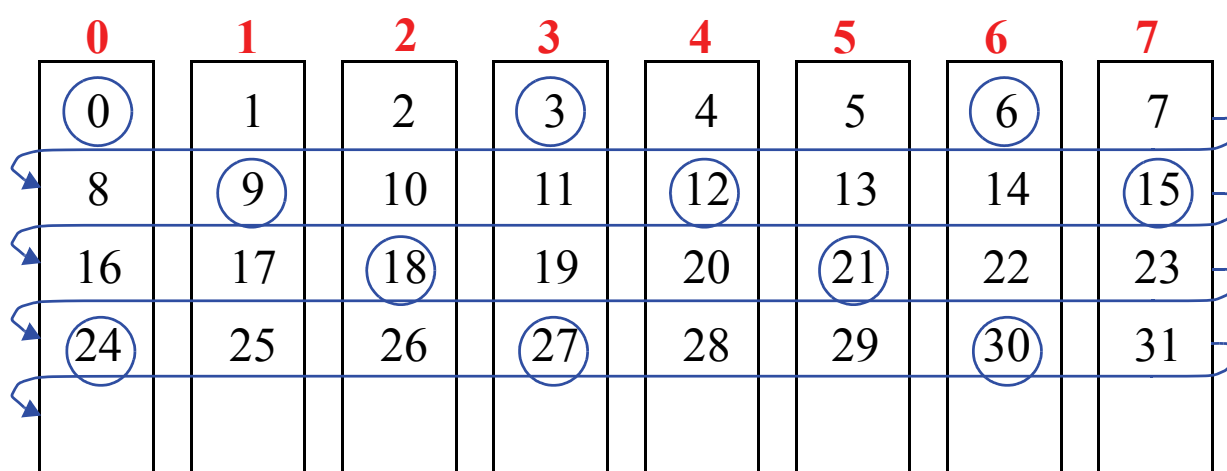
Palabras consecutivos en bancos consecutivos¹

“bancos entrelazados por palabras ó entrelazado de menos peso”

Suponemos:

ancho banco = ancho BUS memoria/registros vect. = tamaño palabra

Ejemplo: **M = 8 bancos**, **S = 3**



PROPIEDAD FUNDAMENTAL

Una secuencia de accesos con separación S bancos visita un subconjunto de P bancos del total de M

$$M = 8$$

S	1	2	3	4	5	8
P	8	4				

$$P = \frac{M}{\text{mcd}(S, M)}$$

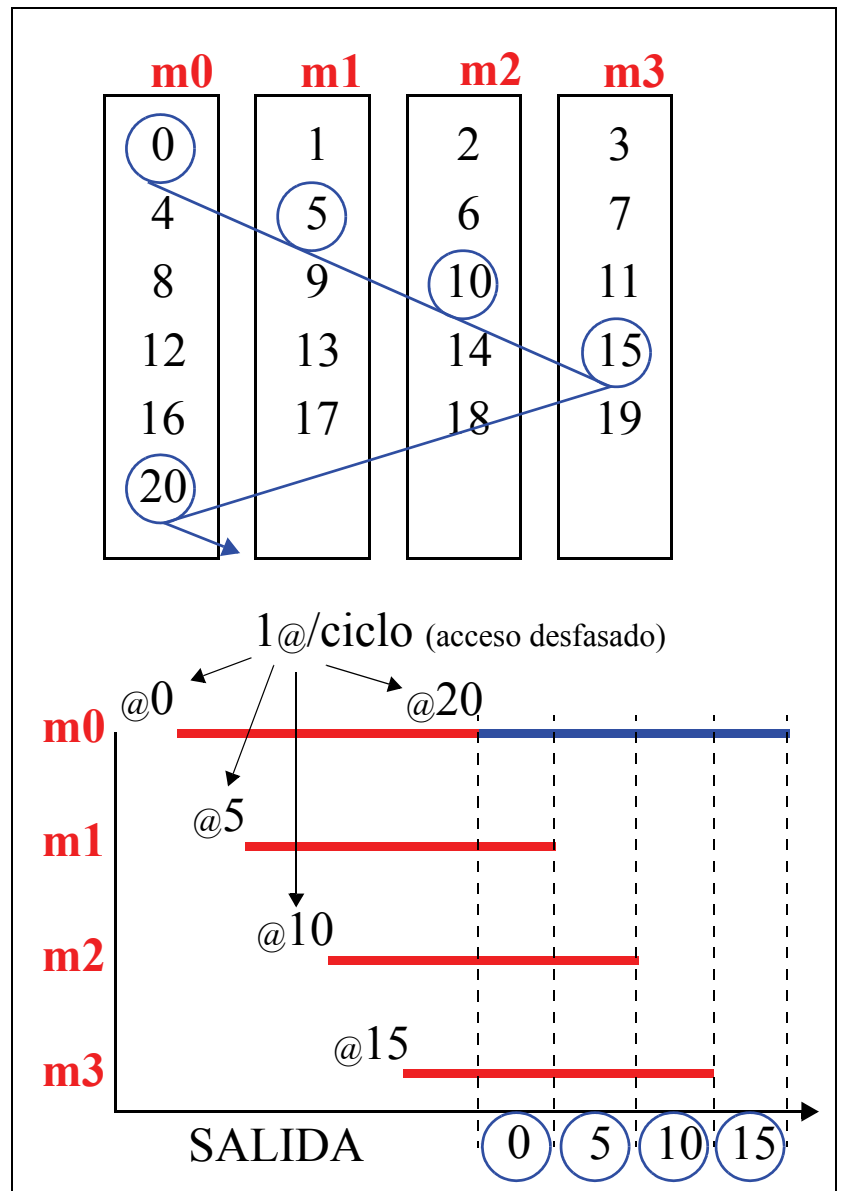
- 1) En otros sistemas, p.e. en un multiprocesador, las palabras que forman un bloque de cache pueden residir consecutivas en el mismo banco de memoria principal. Se dice entonces que los bancos están entrelazados por bloques (no por palabras).

4.2 SISTEMAS AJUSTADOS: $M = L$

Ejemplo: $L = 4$ ciclos y $M = 4$ bancos

□ Acceso sin conflicto $S = 5$ módulos

- ✓ $P = \frac{4}{\gcd(5, 4)} = 4$
- ✓ Cada banco diferente visitado aporta 1 dato / L ciclos
- ✓ Visitamos P bancos
Flujo = BW =
 P datos / L ciclos
- ✓ $Ce_{mem} = Flujo^{-1} =$
 L/P ciclos/dato =
1 cpe



□ Acceso con conflicto $S = 4$ módulos

- ✓ únicamente visitamos el banco 0: $P = \frac{4}{\gcd(4, 4)} = 1$
- ✓ $Ce_{mem} = L/P = 4$ cpe

4.2 SISTEMAS AJUSTADOS: $M = L$

$$Ce = \frac{L}{P} = \frac{L \cdot mcd(M,S)}{M} = mcd(M,S)$$

□ *Strides* IMPARES \rightarrow sin conflicto $\rightarrow \mathbf{Ce_{mem} = 1}$

$$S = \sigma \cdot 2^0, \text{ con } \sigma \text{ impar}$$

DEMO, suponiendo n° de bancos es potencia de 2, $M = 2^m$
$Ce = mcd(M,S) = mcd(2^m, \sigma \cdot 2^0) = 2^0 = 1$

□ *Strides* PARES \rightarrow con conflicto $\rightarrow \mathbf{Ce_{mem} > 1}$

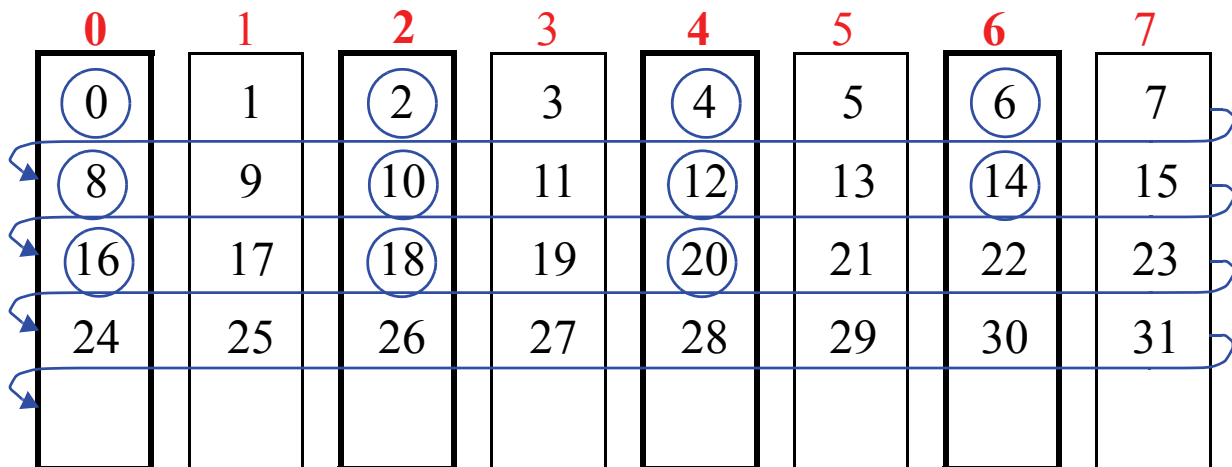
$$S = \sigma \cdot 2^k, \text{ con } \sigma \text{ impar y } k \geq 1$$

DEMO ($M = 2^m$)
Conflicto \rightarrow $Ce_{mem} > 1$
$Ce = mcd(M,S) = mcd(\sigma \cdot 2^k, 2^m) = 2^{\min(k,m)}$

4.3 SISTEMAS SOBRADOS: $M > L$

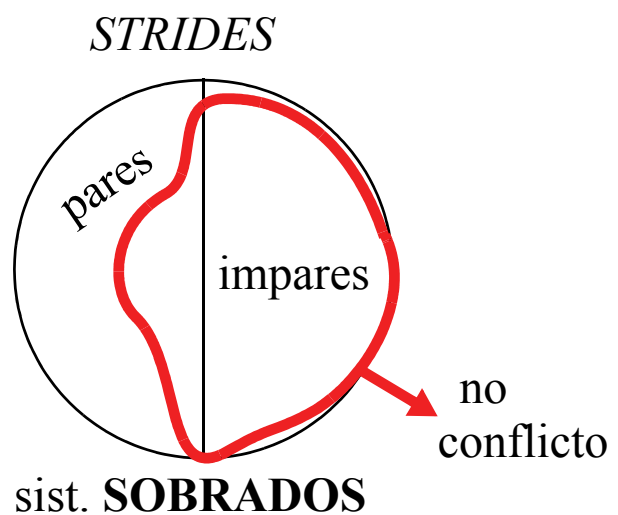
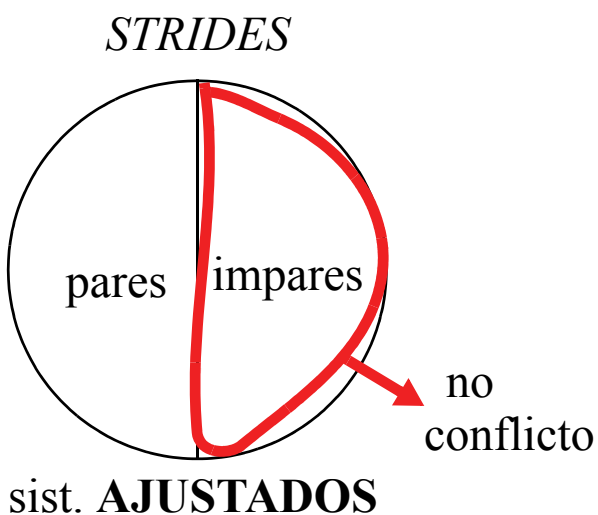
Ejemplo: $M = 8, L = 4$

□ Acceso sin conflicto $S = 2$



$$\text{Flujo} = \text{BW} = \frac{\text{Pelementos}}{\text{Lciclos}} \Big|_{P=L=4} = 1 \text{ elemento/ciclo}$$

permitimos más *strides* libres de conflicto



4.3 SISTEMAS SOBRADOS: $M > L$

Suponemos $L = 2^l$ (irreal)

□ *Strides* sin conflicto: IMPARES y algunos PARES

$$S = \sigma \cdot 2^k, \text{ con } \sigma \text{ impar y } 0 \leq k \leq m-l$$

DEMO ($M = 2^m > L = 2^l$)

Sin conflicto $\rightarrow P \geq L$

$$P = \frac{2^m}{\text{mcd}(S, 2^m)} \geq 2^l \rightarrow 2^{m-l} \geq \text{mcd}(S, 2^m) = 2^{\min(k, m)}$$

$$\rightarrow m-l \geq \min(k, m) \rightarrow m-l \geq k$$

No hay conflicto con PARES si tienen “pocas” potencias de dos.

□ *Strides* con conflicto: parte de los PARES

$$S = \sigma \cdot 2^k, \text{ con } \sigma \text{ impar y } k > m-l$$

DEMO ($M = 2^m$)

Conflicto $\rightarrow C_{\text{mem}} > 1$

$$\text{mcd}(2^m, S) > 1$$

$$C_{\text{mem}} = \frac{L}{P} = \frac{L \cdot \text{mcd}(S, M)}{M} = \frac{2^l \cdot 2^{\min(k, m)}}{2^m} = 2^{l + \min(k, m) - m}$$

Ejemplo: $C_e(\text{sistema } L, M, S) = C_e(8, 16, 30) =$

5. ARQUITECTURA DLXV

Repertorio completo de instrucciones

DLXV - VMIPS: CÓDIGOS DE OPERACIÓN

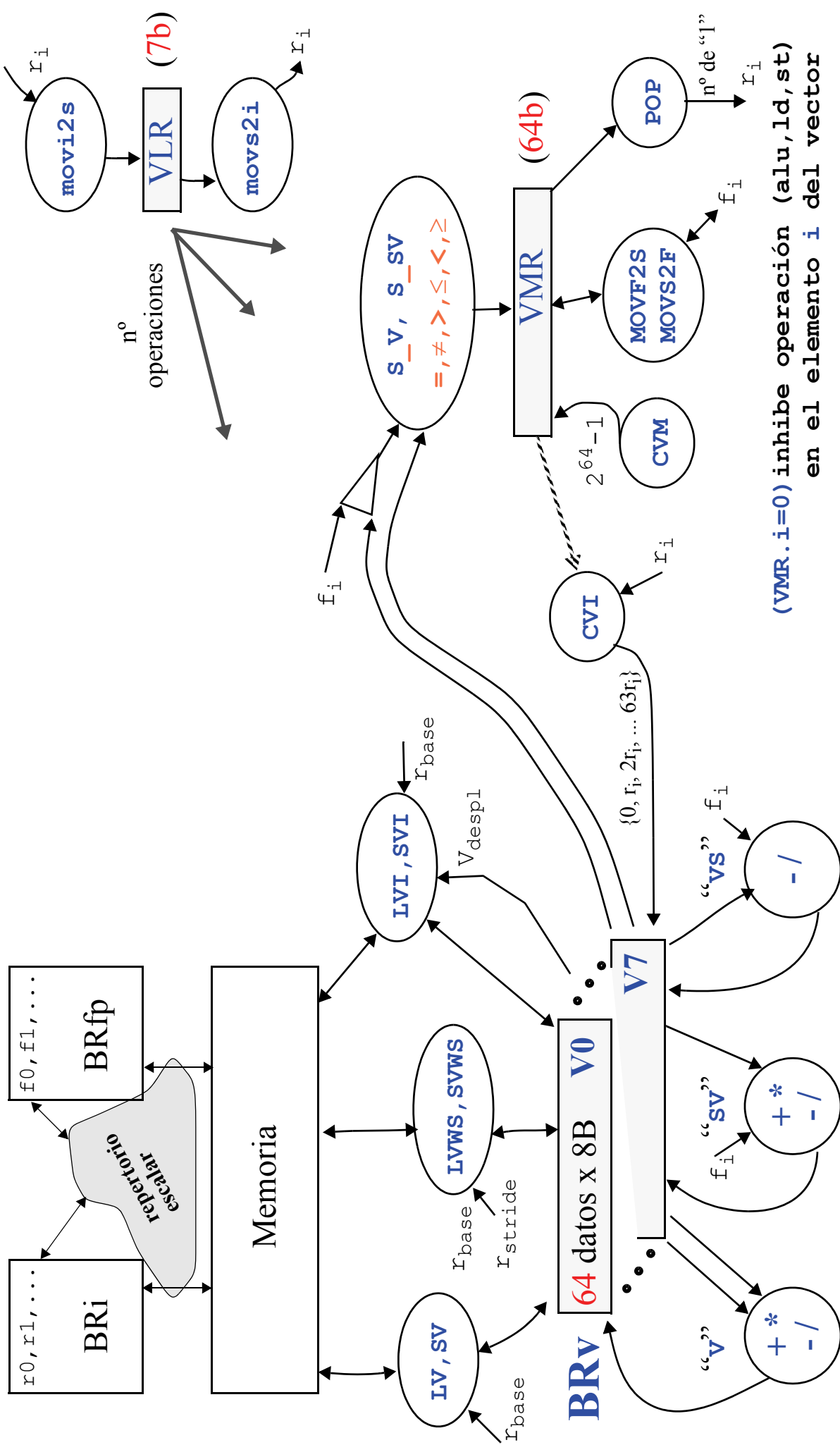
F-8 ■ Appendix F *Vector Processors* (4th edition)

Instruction	Operands	Function
ADDV.D	V1, V2, V3	Add elements of V2 and V3, then put each result in V1.
ADDVS.D	V1, V2, F0	Add F0 to each element of V2, then put each result in V1.
SUBV.D	V1, V2, V3	Subtract elements of V3 from V2, then put each result in V1.
SUBVS.D	V1, V2, F0	Subtract F0 from elements of V2, then put each result in V1.
SUBSV.D	V1, F0, V2	Subtract elements of V2 from F0, then put each result in V1.
MULV.D	V1, V2, V3	Multiply elements of V2 and V3, then put each result in V1.
MULVS.D	V1, V2, F0	Multiply each element of V2 by F0, then put each result in V1.
DIVV.D	V1, V2, V3	Divide elements of V2 by V3, then put each result in V1.
DIVVS.D	V1, V2, F0	Divide elements of V2 by F0, then put each result in V1.
DIVSV.D	V1, F0, V2	Divide F0 by elements of V2, then put each result in V1.
LV	V1, R1	Load vector register V1 from memory starting at address R1.
SV	R1, V1	Store vector register V1 into memory starting at address R1.
LVWS	V1, (R1, R2)	Load V1 from address at R1 with stride in R2, i.e., $R1+i \times R2$.
SVWS	(R1, R2), V1	Store V1 from address at R1 with stride in R2, i.e., $R1+i \times R2$.
LVI	V1, (R1+V2)	Load V1 with vector whose elements are at $R1+V2(i)$, i.e., V2 is an index.
SVI	(R1+V2), V1	Store V1 to vector whose elements are at $R1+V2(i)$, i.e., V2 is an index.
CVI	V1, R1	Create an index vector by storing the values $0, 1 \times R1, 2 \times R1, \dots, 63 \times R1$ into V1.
S--V.D	V1, V2	Compare the elements (EQ, NE, GT, LT, GE, LE) in V1 and V2. If condition is true, put a 1 in the corresponding bit vector; otherwise put 0. Put resulting bit vector in vector-mask register (VM). The instruction S--VS.D performs the same compare but using a scalar value as one operand.
S--VS.D	V1, F0	
POP	R1, VM	Count the 1s in the vector-mask register and store count in R1.
CVM		Set the vector-mask register to all 1s.
MTC1	VLR, R1	Move contents of R1 to the vector-length register.
MFC1	R1, VLR	Move the contents of the vector-length register to R1.
MVTM	VM, F0	Move contents of F0 to the vector-mask register.
MVFM	F0, VM	Move contents of vector-mask register to F0.

Figure F.3 The VMIPS vector instructions. Only the double-precision FP operations are shown. In addition to the vector registers, there are two special registers, VLR (discussed in Section F.3) and VM (discussed in Section F.4). These special registers are assumed to live in the MIPS coprocessor 1 space along with the FPU registers. The operations with stride are explained in Section F.3, and the uses of the index creation and indexed load-store operations are explained in Section F.4.

- [HePa12] J. HENNESSY and D. PATTERSON, Computer Architecture: a quantitative approach. 5th Edition, Morgan Kaufmann, 2012.
 - Chapter 4 and Appendix G: Vector Processors in More Depth

DLXV completo: diagrama de flujo de la ALMa



6. COMPILACIÓN

6.1 Introducción.

Fases en el *back-end* del compilador:
extracción automática de paralelismo
vectorial

6.2 Transformaciones previas que simplifican el análisis de dependencias

6.3 Análisis y grafo de dependencias. Tests aproximados

6.4 Optimizaciones independientes de la arquitectura: renombrar, expansión escalar, copia de vectores

6.5 Vectorización

- Procedimiento básico
- Vectorización parcial vs. total:
distribución e intercambio de bucles
- Reducción

6.1 FASES EN EL *BACK-END* DEL COMPILADOR

Front-end: análisis léxico y sintáctico

Programa escalar representado en lenguaje intermedio

⇓

Transformaciones
que simplifican el
análisis de las
dependencias

- Propagación y evaluación de expres. constantes
- Extracción de invariantes
- Normalización de bucles

⇓

Grafo y análisis
de dependencias

- Construcción del grafo
- Tipos y distancia de las dependencias
- Análisis de dependencias. Test MCD

⇓

Optimizaciones
independientes de
la arquitectura
vect/par

- Eliminar antidependencias y dep. de salida
 - renombrar
 - expansión escalar
 - copia de vectores

⇓

Vectorización

- Vectorización total y parcial
 - distribución e intercambio de bucles
- Reducción, punteros en C

⇓

Generación de
código.
Optimización

- Asignación de registros vectoriales
- Selección de instrucciones vectoriales
- Seccionado

⇓

**Programa
vectorial
optimizado**

$$\text{EJEC}(\text{progr. escalar}) = \text{EJEC}(\text{prog. vectorial}) + f \text{ estática vect } \uparrow\uparrow$$

6.2 TRANSFORMACIONES QUE SIMPLIFICAN EL ANÁLISIS DE LAS DEPENDENCIAS

- Propagación y evaluación de expresiones constantes + Extracción de invariantes

<pre>DO i = 1,N PI = 3.14 PD = 2*PI A(i) = PD * R(i) **2 ENDDO</pre>	⇒	<pre>PI = 3.14 PD = 6.28 DO i = 1,N A(i) = 6.28 * R(i) **2 ENDDO</pre>
--	---	--

Otro ejemplo de invariantes

<pre>DO i = 1,N DO j = 1,M A(i,j) = B(j) * C(i) ENNDDO ENDDO</pre>	⇒	<p style="color: red; font-size: 2em;">?</p>
--	---	--

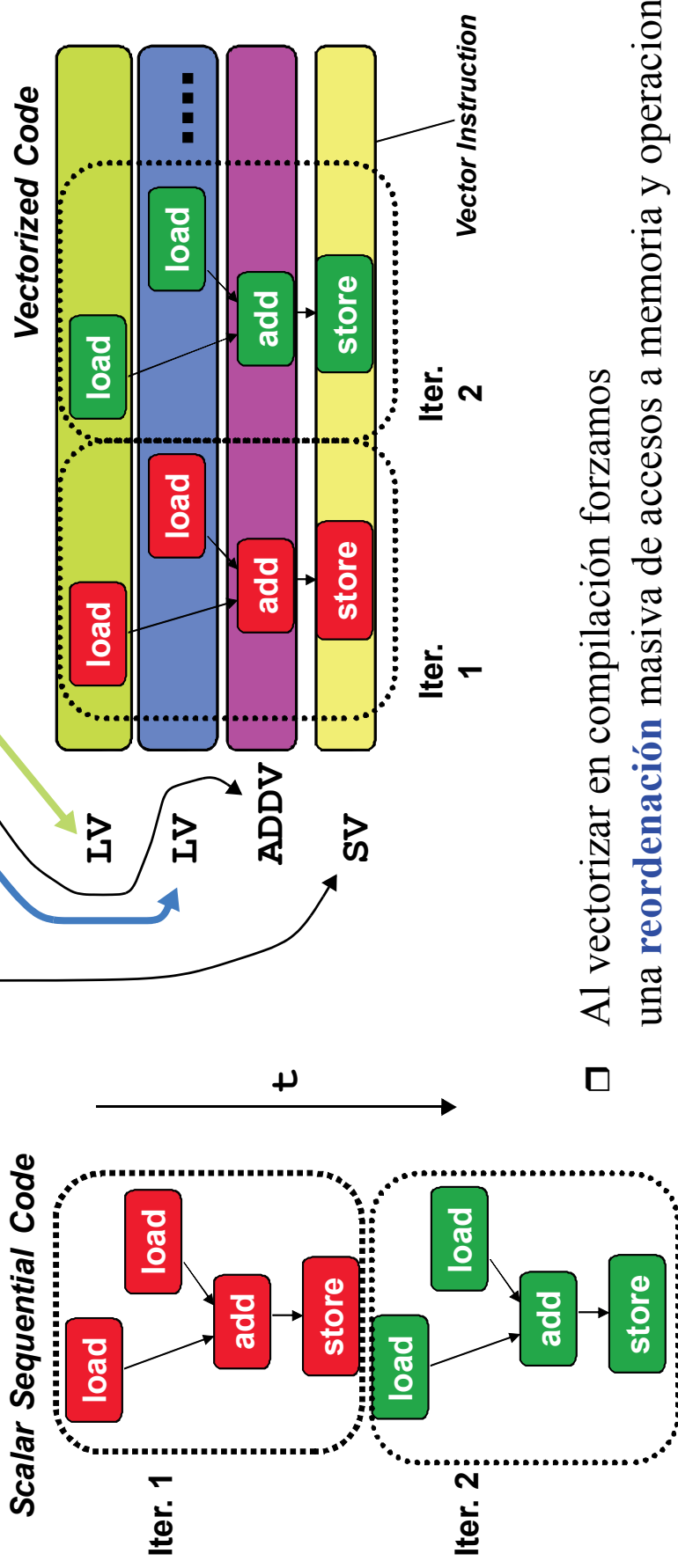
- Normalización de bucles:
 - normalizar el paso de la variable de control
 - eliminar variables de inducción, o sea, indexar vectores únicamente con las variables de control del bucle

<pre>DO i = 1,100 aux = i DO j = 1,300,3 aux = aux + 2 A(j) = A(j-1) + B(aux) ENNDDO ENDDO</pre>	⇒	<pre>DO i = 1,100 aux = i DO nj = 1,100 aux = aux + 2 A(3*nj-2) = A(3*nj-3) + B(i+2*nj) ENNDDO ENDDO</pre>
--	---	--

$3*nj-2, 3*nj-3, 2*nj+i \Rightarrow$ combinación lineal de las variables de control de bucle

6.3 ANÁLISIS Y GRAFO DE DEPENDENCIAS

```
for (i=0; i < N; i++) C[i] = A[i] + B[i];
```



□ Al vectorizar en compilación forzamos

una **reordenación** masiva de accesos a memoria y operaciones

□ Ejemplo de reordenación: $C[2] <_p B[3]$, pero $B[3] <_m C[2]$!

□ Es necesario analizar las dependencias con cuidado:

que pasa si el **store rojo** escribe donde lee el **load verde** ?

$$\Rightarrow \text{c}[i] = \text{c}[i-1] + \dots$$

DEPENDENCIAS EN MEMORIA ENTRE SENTENCIAS¹

R: -----

S y **R** son sentencias de asignación

...

S: -----

S depende de **R** si:

- ☐ **R** antes que **S** en orden de programa ($R <_p S$)
- ☐ **R** y **S** referencian a la misma posición de memoria
- ☐ una referencia al menos es una escritura

- Las dependencias establecen relaciones de orden parcial que cualquier ejecución “legal” debe respetar.

Tres tipos		
R escribe y S lee Dependencia verdadera <i>flow dependence, RAW hazard</i>	$R \delta S$	$R \longrightarrow S$
R lee y S escribe Antidependencia <i>antidependence, WAR hazard</i>	$R \delta^- S$	$R \dashrightarrow S$
R y S escriben Dependencia de salida <i>output dependence, WAW hazard</i>	$R \delta^o S$	$R \ominus \longrightarrow S$
Si R y S están en un bucle, R_i y S_i indican la ejecución de la iteración i -ésima		

1) La investigación en vectorización automática se inicia en los años 70, y el siguiente trabajo de la Universidad de Illinois se considera uno de los mas importantes para su herramienta básica, el análisis de dependencias.

Kuck, Kuhn, Padua, Leasure, and Wolfe. "Dependence graphs and compiler optimizations". In Proc. of the 8th ACM Symposium on Principles of Programming Languages, pp. 207-218. 1981.

- Ejemplo 1 -

$$\left[\begin{array}{l} i = 1, 64 \\ S1: A(i) = B(i) * C \\ S2: B(i) = B(i-1) + A(i) \end{array} \right.$$

por inspección: $S1_i \delta S2_i$ posición $A(i)$
 $S1_i \delta^- S2_i$ posición $B(i)$

Pero ...

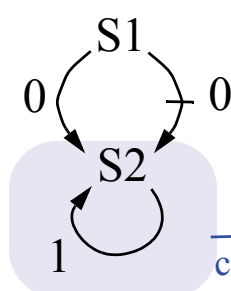
$S1_1$	$A(1) = B(1) * C$	¡ Dependencias entre iteraciones !
$S2_1$	$B(1) = B(0) + A(1)$	
$S1_2$	$A(2) = B(2) * C$	$S2_i \delta S2_{i+1}$ en pos $B(i)$
$S2_2$	$B(2) = B(1) + A(2)$	
	\vdots	<i>loop carried dependence, LCD</i> o sea, dependencias generadas por el bucle
	\vdots	

λ = distancia de la dependencia: número de iteraciones que separa el uso de las mismas posiciones de memoria

$$R_i \delta S_j \rightarrow \lambda = j - i \geq 0$$

$$\mathbf{R} \xrightarrow{\lambda} \mathbf{S}$$

□ Grafo de dependencias: grafo dirigido, anotado con distancias:



• Vectorizable en parte

$$A(1:64) = B(1:64) * C$$

$$i = 1, 64$$

$$S2: B(i) = B(i-1) + A(i)$$

ciclo unitario de dependencia

- Ejemplo 2 -

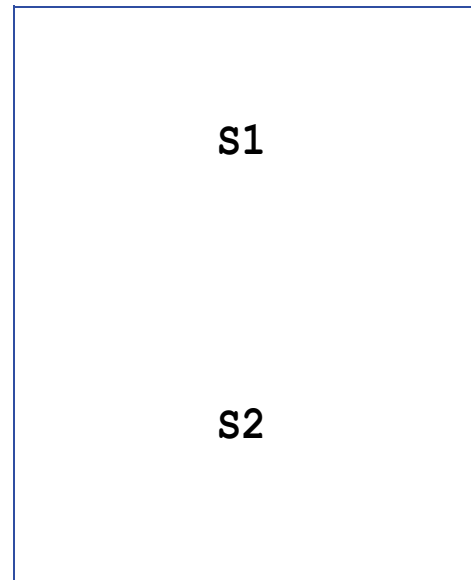
```

i= 1, 64
[
  S1: A(i+1)= A(i) + B(i)
  S2: B(i+1)= B(i) - A(i)

```

ejecución

S1 ₁
S2 ₁
S1 ₂
S2 ₂
⋮



❑ O sea: NO es vectorizable !

```

A(2:65) = A(1:64) + B(1:64)
B(2:65) = B(1:64) - A(1:64)

```

```

LV B(1:64)
LV A(1:64)
ADDV +
SV A(2:65) =
A(1:64)
B(1:64)
-
B(2:65) =

```

- Ejemplo 3: AXPY -

```

i= 1, 64
[  S: Y(i) = a*X(i) + Y(i)

```

**S**

... o simplemente

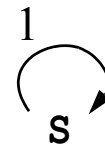
Autociclos de antidep. a distancia 0, pueden obviarse: *Vectorizable*

- Ejemplo 4 -

```

i= 1, 64
[  S: fac(i) = fac(i-1) * i

```



Autociclos de dependencias verdaderas distancia > 0, *No Vectorizable*

- Ejemplo 5 -

```

i= 1, 64
[  S: A(i) = A(i+3) + A(i)

```

secuencia de cálculos:

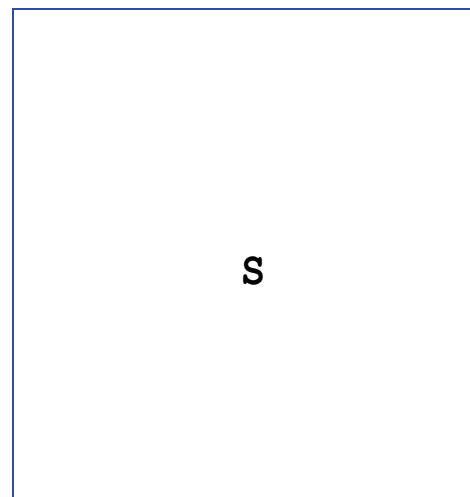
 S_1

 S_2

 S_3

 S_4

 S_5



Vectorizable: $A(1:64) = A(4:67) + A(1:64)^a$

- a. Esta notación corresponde a un lenguaje real y en uso: FORTRAN 90, el sucesor directo de FORTRAN 77

□ Vectorizable: $A(1:64) = A(4:67) + A(1:64)^1$

- Ejemplo 6: escalado -

$\left[\begin{array}{l} i = 1, 64 \\ S: A(i) = k * A(i) \end{array} \right.$	S VECT, PAR
---	---------------------------

- Ejemplo 7: progresión geométrica-

$\left[\begin{array}{l} i = 1, 64 \\ S: A(i+1) = k * A(i) \end{array} \right.$	S no VECT, no PAR
---	---------------------------------

- Ejemplo 8: escalado con desplazamiento a la izquierda -

$\left[\begin{array}{l} i = 1, 64 \\ S: A(i) = k * A(i+1) \end{array} \right.$	S si VECT, no PAR
---	---------------------------------

- Ejemplo 9 -

$\left[\begin{array}{l} i = 1, 64 \\ S1: A(i) = B(i+2) \\ S2: B(i) = A(i-1) \end{array} \right.$	S1 S2 si VECT, no PAR
---	---

- Escribir el código vectorial

1) Esta notación corresponde a un lenguaje real y en uso: FORTRAN 90, el sucesor directo de FORTRAN 77

- Ejemplo 10 -

$$\left[\begin{array}{l} i = 1, 64 \\ S1: A(i+2) = B(i) \\ S2: B(i+2) = A(i+1) \end{array} \right. \quad \begin{array}{l} S1 \\ S2 \end{array}$$

no VECT, no PAR

ANÁLISIS DE DEPENDENCIAS

Caso bastante general:

- conjunto de bucles anidados normalizados
- código en el bucle mas interno y sin sentencias condicionales
- espacio de iteraciones conocido en compilación

real*8 A(dim₁, dim₂, ...dim_m)

$$\left[\begin{array}{l} i_1 = 1, \max_1 \\ \dots \\ \dots \\ \left[\begin{array}{l} i_{k-1} = 1, \max_{k-1} \\ \dots \\ \left[\begin{array}{l} i_k = 1, \max_k \\ S1: A(f_1(i_1, \dots, i_k), f_2(i_1, \dots, i_k), \dots, f_m(i_1, \dots, i_k)) = \dots \\ S2: \dots = A(g_1(i_1, \dots, i_k), g_2(i_1, \dots, i_k), \dots, g_m(i_1, \dots, i_k)) \end{array} \right. \end{array} \right. \end{array} \right.$$

En cada iteración del código S1 S2, las variables de control forman una k-tupla (i_1, \dots, i_k) diferente. La unión de todas ellas se llama *Espacio de Iteraciones* (EspIt). EspIt es el producto cartesiano de los k conjuntos de números naturales que recorre cada variable de control.

Puede establecerse un *orden de programa total* entre k-tuplas (\leq_p)

Las funciones **f** y **g** son *funciones de indexación*. Calculan un índice a partir de los valores (i_1, \dots, i_k) de las variables de control

PREGUNTAS ...

$\exists P, Q \in \text{EspIt} \mid \text{escr en } S1_P \delta \text{ lect en } S2_Q \text{ para } P \leq_p Q; \quad ?^a$
 $\exists T, Z \in \text{EspIt} \mid \text{lect en } S2_T \delta^- \text{ escr en } S1_Z \text{ para } T <_p Z; \quad ?^b$

- a. P.e. para $P = (i_1, i_2, i_3) = (2, 3, 4)$, en S1 se escribe el elemento A(1,2) y para $Q = (4, 1, 1)$, en S2 se lee *el mismo* elemento.
 b. P.e. para $T = (3, 3, 3)$, en S2 se lee el elemento A(9,1) y para $Z = (4, 1, 1)$, en S1 se escribe *el mismo* elemento.

Caso sencillo:

- $K = 1$ variable de control
- $m = 1$ dimensión

$i = 1, \max$
 $S1: A(f(i)) = \dots$
 $S2: \dots = A(g(i)) \dots$

\exists dependencia prod/cons a distancia λ sii^a

$S1_P$
 $\downarrow \lambda$
 $S2_Q$
 $\exists p, q \in \text{EspIt} \mid f(p) = g(q) \text{ para } p \leq q$
 $\lambda = q - p$

- a. P.e. para $p = 6$ en S1 se escribe el elemento A(19) y para $q = 6$, en S2 se lee *el mismo* elemento.

\exists antidependencia a distancia λ sii^a

$S1_Z$
 $\uparrow \lambda$
 $S2_t$
 $\exists t, z \in \text{EspIt} \mid g(t) = f(z) \text{ para } t < z$
 $\lambda = z - t$

- a. P.e. para $t = 10$, en S2 se lee el elemento A(31) y para $z = 12$, en S1 se escribe *el mismo* elemento.

- Ejemplo 11 - f y g son lineales

$$\begin{array}{l}
 i = 1, 100 \\
 \left[\begin{array}{l}
 S1: A(2i+7) = B(i) - 3 \\
 S2: C(i) = A(3i+1)
 \end{array} \right.
 \end{array}
 \quad
 \begin{array}{l}
 S1 \\
 S2
 \end{array}$$

secuencia de cálculos:

VECT ?, PAR ?

i	escri +2 lect +3
1	9
2	
3	
4	
5	
6	
7	
8	
9	
10	
11	
12	

Eqn. dependencias prod/cons: S1 δ S2

$$\begin{array}{l}
 \exists p, q \in [1..100] \mid f(p) = g(q) \text{ para } p \leq q \\
 2p+7 = 3q+1 \quad (\lambda = q - p)
 \end{array}$$

Despejamos q y enumeramos p; **q = 2p/3+2**

p	1	2	3	4	5	6	7	8	9	10 ...
q	no	no	4	no	no	6	no	no	8	no
$\lambda \geq 0$			1			0			-1	

Eqn. antidependencias: S2 δ^- S1

$$\begin{array}{l}
 \exists t, z \in [1..100] \mid g(t) = f(z) \text{ para } t < z \\
 3t+1 = 2z+7 \quad (\lambda = z - t)
 \end{array}$$

Despejamos z y enumeramos t; **z = 3t/2-3**

t	1	2	3	4	5	6	7	8	9	10 ... 68
z	no	0	no	3	no	6	no	9	no	12 ... 99
$\lambda > 0$		-2		-1		0		1		2 ... 31

❑ Este tipo de ecuación se llama ecuación Diofántica

- Ejemplo 12 - f y g son lineales

i= 1,100

S1: C(i)= A(3i+1)

S2: A(2i+7)= B(i)-3

S1

S2

VECT ?, PAR ?

i	<div>lect +3</div> <div>escri +2</div>	<div>Eqn. dependencias prod/cons: S2 δ S1</div> <div></div> <div>Despejamos q y enumeramos p;</div> <table><tr><td>p</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td><td>8</td><td>9</td><td>10 ...</td></tr><tr><td>q</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></tr><tr><td>λ</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></tr></table>	p	1	2	3	4	5	6	7	8	9	10 ...	q											λ										
p	1	2	3	4	5	6	7	8	9	10 ...																									
q																																			
λ																																			
1																																			
2																																			
3																																			
4																																			
5																																			
6																																			
7		<div>Eqn. antidependencias: S1 δ^- S2</div> <div></div> <div>Despejamos z y enumeramos t;</div> <table><tr><td>t</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td><td>8</td><td>9</td><td>10 ...</td></tr><tr><td>z</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></tr><tr><td>λ</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></tr></table>	t	1	2	3	4	5	6	7	8	9	10 ...	z											λ										
t	1	2	3	4	5	6	7	8	9	10 ...																									
z																																			
λ																																			
8																																			
9																																			
10																																			
11																																			
12																																			

TESTS APROXIMADOS

Resolver las ecuaciones directamente es costoso en tiempo

Los límites inferior y superior a veces no son constantes

❑ Tests aproximados, fáciles de calcular:

- ✓ **condición necesaria** para que existan ciclos de dependencias
 - se cumple ? → no vectorizamos
 - muy conservador ...
- ✓ **condición suficiente** para la independencia
 - se cumple ? → vectorizamos
 - muy exigente

❑ Test MCD para funciones de indexación lineales (**cond. necesaria**)

$$f \rightarrow \mathbf{a \cdot i + b} \quad g \rightarrow \mathbf{c \cdot i + d}$$

Asumimos posible ciclo de dependencias si **mcd(a, c)** divide **(d-b)**

- ejemplos 11 y 12 -

$$\left. \begin{array}{l} a = 3, b = 1 \\ c = 2, d = 7 \end{array} \right\} \text{mcd}(3, 2) = 1 \text{ divide a } (6) \rightarrow \text{SI}$$

... pues asumimos dependencias cíclicas
y no vectorizamos

Ejercicios:

- ✓ demostración test MCD
- ✓ buscar un ejemplo vectorizable que no pase el test MCD

EJERCICIO E5

En bucles con un solo nivel de anidación pueden aparecer varios tipos de dependencias.

$$\left[\begin{array}{l} i \in \text{rango } N \\ R: \quad \text{var}(r) \\ S: \quad \text{var}(s) \end{array} \right. \quad \begin{array}{l} R \\ S \end{array}$$

Los índices r y s

son funciones lineales de i : $r(i)$, $s(i)$

VECT ?, PAR ?

En la tabla siguiente se resumen todas las posibilidades.

$\text{var}(r)$		L	ind	ind	ind
		E	ind	ind	ind
			L	E	
			$\text{var}(s)$		

Por ejemplo:

$$\left[\begin{array}{l} i \in \text{rango } N \\ R: \text{var}(\mathbf{r}) = \dots \quad \mathbf{r} \text{ y } \mathbf{s} \text{ son funciones lineales de } i \\ S: \text{var}(\mathbf{s}) = \dots \quad \lambda \in N \text{ es la distancia de la dependencia} \end{array} \right.$$

Podemos formular las tres ecuaciones de dependencias:

- ✓ $\mathbf{R} \delta^0 \mathbf{S} \text{ sii } \exists \lambda \mid r(i) = s(i + \lambda); \lambda \geq 0; i, i + \lambda \in \text{rango}$
- ✓ $\mathbf{S} \delta^0 \mathbf{R} \text{ sii } \exists \lambda \mid s(i) = r(i + \lambda); \lambda \geq 1; i, i + \lambda \in \text{rango}$
- ✓ $\mathbf{R} \text{ ind } \mathbf{S} \text{ caso contrario}$
ejemplos de este caso:

```
DO i= 2, 100
R:  A(2i-1)=
S:  A(6i+4)=
ENDDO
```

```
DO i= 1, 100
R:  A(2i)=
S:  A(2i+1)=
ENDDO
```

Ejercicio. Se pide lo siguiente:

Repetir el análisis anterior para el caso

“**var(r)** es Lectura y **var(s)** es Escritura”

- ✓ las tres ecuaciones de dependencias
- ✓ un ejemplo de cada caso, con funciones de indexación diferentes de las empleadas en los ejemplos anteriores.

6.4 OPTIMIZACIONES NO LIGADAS A LA ARQUITECTURA



Grafo y análisis
de dependencias



Optimizaciones
independientes de
la arquitectura
vect/par



Vectorización



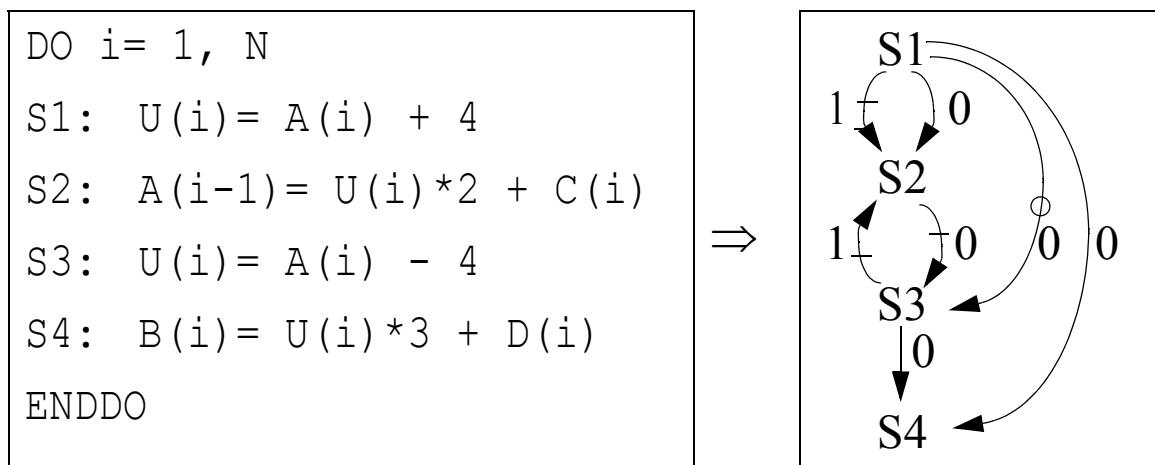
- **Eliminar antidependencias y dep. de salida**
 - **renombrar**
 - **expansión escalar**
 - **copia de vectores**

Renombrar

Problema: reutilizar vectores o escalares para ahorrar memoria o variables

Transformación:

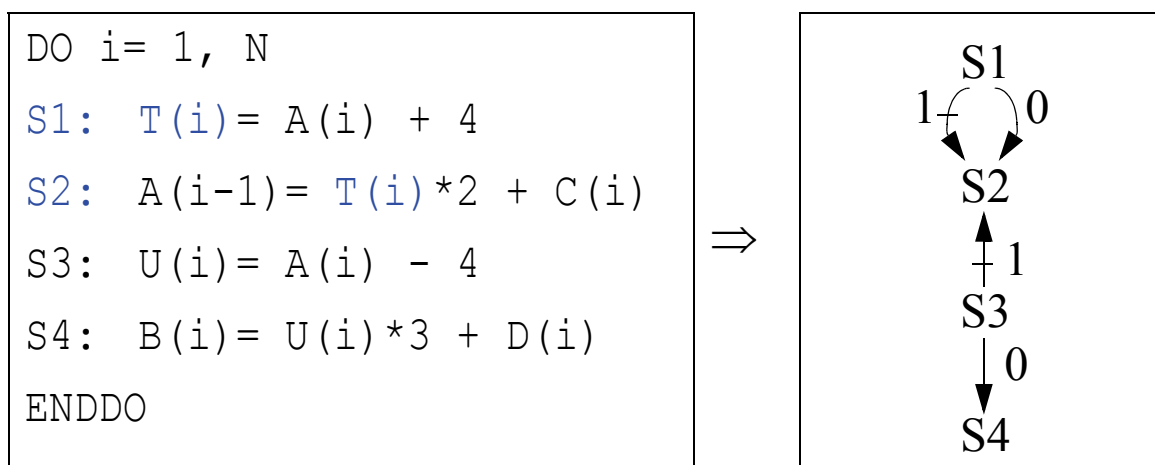
- ✓ nombre diferente a cada asignación (parte izquierda)
- ✓ propagar nombres nuevos a las sentencias posteriores (partes derechas)



- Ejemplo 13 - no vect

Transformación:

- $U(i) \rightarrow T(i)$ en S1
- propagar $T(i)$ en S2



casi vect

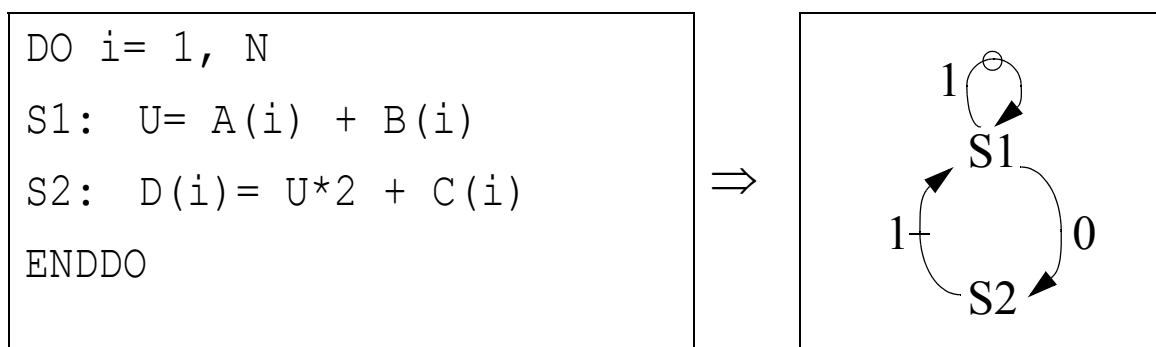
- ✓ ¿Aumentan las necesidades de almacenamiento ?

Expansión Escalar

Problema: variables escalares que se utilizan en iteraciones sucesivas para almacenar valores diferentes

Transformación:

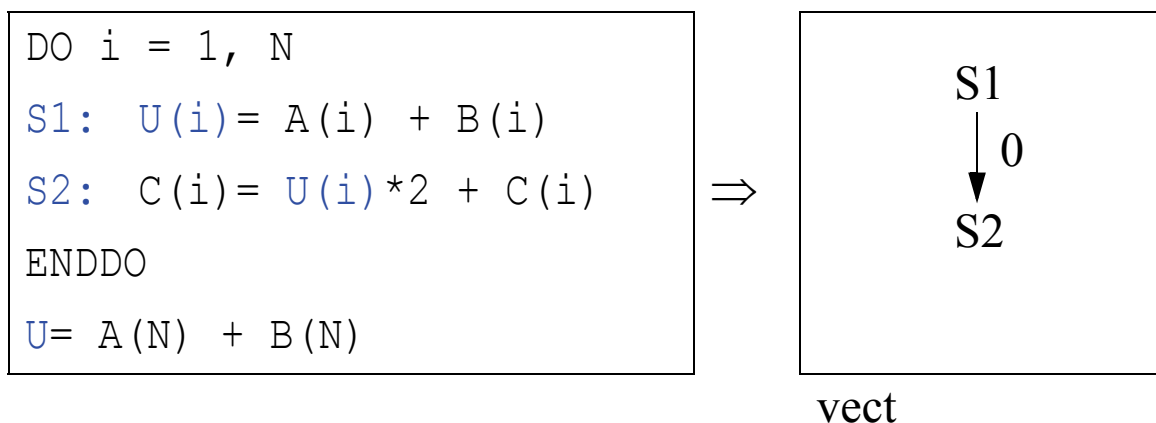
- ✓ promocionar el escalar a vector (parte izquierda)
- ✓ propagar nombres nuevos a las sentencias posteriores (partes derechas)



- Ejemplo 14 - no vect

Transformación:

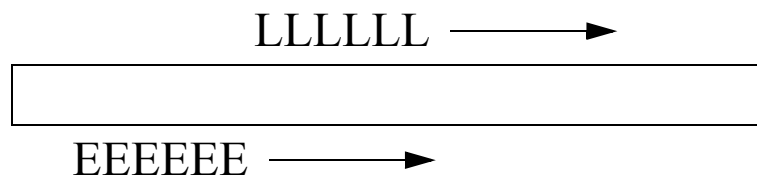
- $U \rightarrow U(i)$ en S1
- propagar $U(i)$ en S2



- ✓ ¿ Aumentan las necesidades de almacenamiento ?

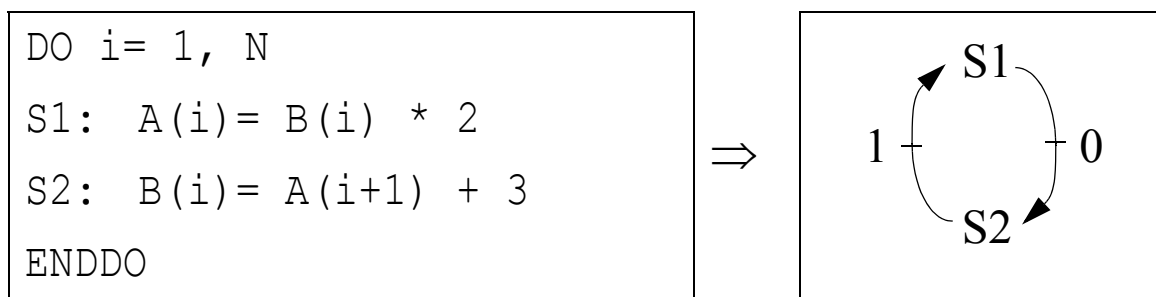
Copia de vectores¹

Problema: pre-uso y post-definición de los elementos de un vector



Transformación:

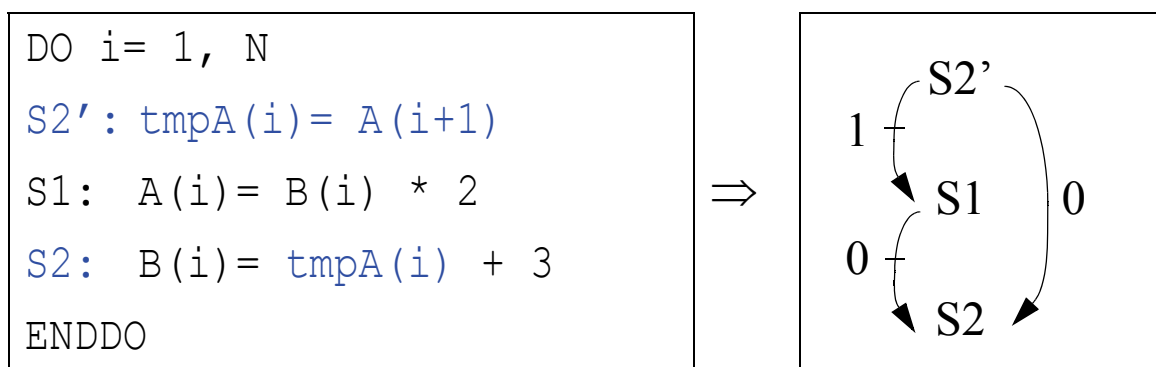
- ✓ crear una copia del vector de solo-lectura
- ✓ propagar nombre copia a las sentencias posteriores (partes derechas)



- **Ejemplo 15** - no vect

Transformación:

- $A \rightarrow \text{tmpA}$ en S2'



vect

- ✓ ¿ Aumentan las necesidades de almacenamiento ?

1) También llamado *node splitting*, particionado de nodos

6.5 VECTORIZACIÓN

Optimizaciones
independientes de
la arquitectura
vect/par



Vectorización

- Vectorización total y parcial
 - distribución e intercambio de bucles
- Reducción
- Punteros en C



Generación de
código.
Optimización

Vectorización total

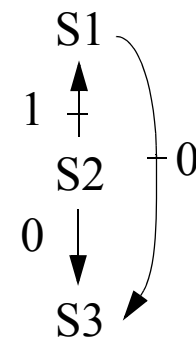
No existen ciclos de dependencias

Transformación:

- ✓ cambiar el orden de las sentencias para que todos los arcos vayan hacia abajo¹
- ✓ generar código vectorial sentencia a sentencia, empezando por cualquiera que sea independiente

```
DO i= 1, N
S1:  A(i)= D(i) - C(i)
S2:  B(i)= A(i+1) + 1
S3:  C(i)= B(i) * 2
ENDDO
```

⇒

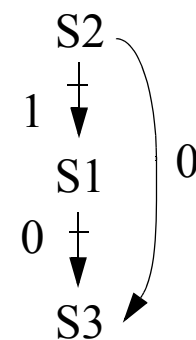


- Ejemplo 16 -

- algoritmo de “flotación”: peso proporcional al nº de arcos entrantes

```
S2 (1:N) es B(1:N)= A(2:N+1)+1
S1 (1:N) es ...
S3 (1:N)
```

⇒



- ✓ **con varios niveles de anidación**, una sentencia puede vectorizarse con respecto a un bucle si no está en un ciclo de dependencias generadas por ese bucle → veremos ejemplos

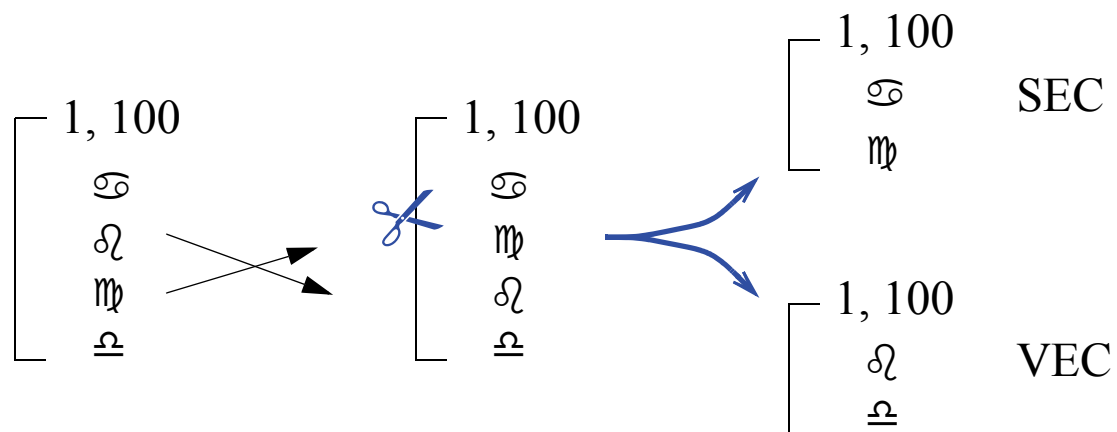
1) Comrobad que este nuevo orden de jecución secuencial es correcto

Vectorización parcial, distribución de bucles¹

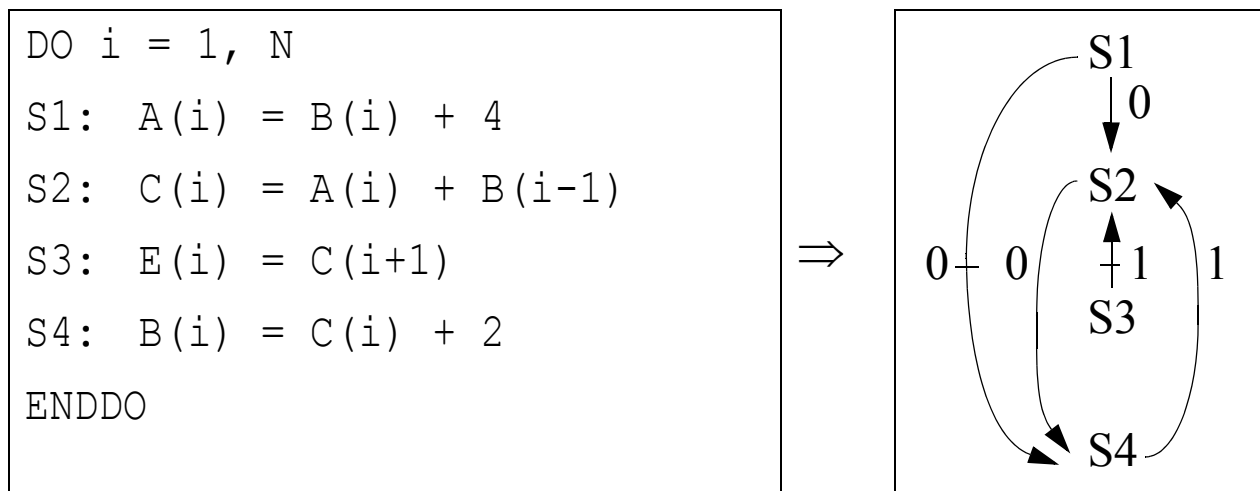
Problema 1 de 2: existen ciclos de dependencias

Transformación: extraer las sentencias sin ciclos de dependencias

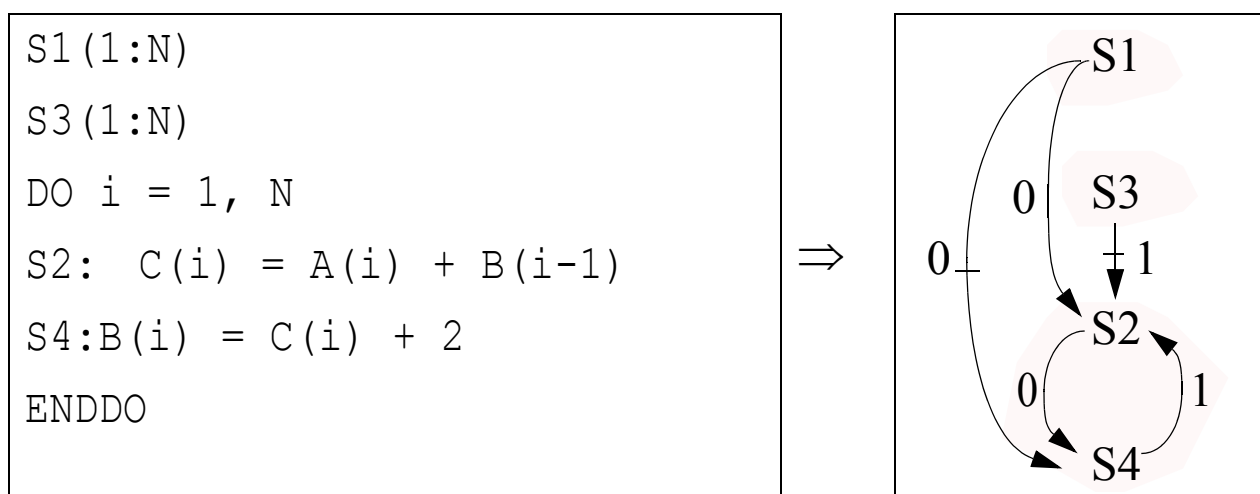
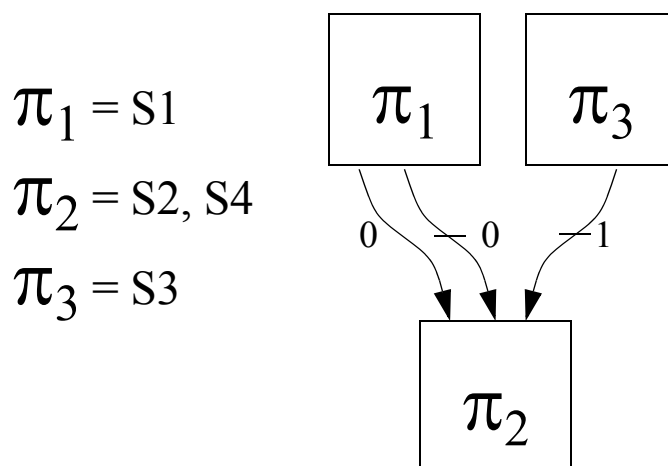
- Particionar el grafo en subgrafos nodo-disjuntos relacionados entre sí de forma acíclica (π bloques)
- reordenar con las dependencias hacia abajo
- Distribuir/cortar en varios bucles



1) También llamado *loop fision*, corte de bucles



- Ejemplo 17 -

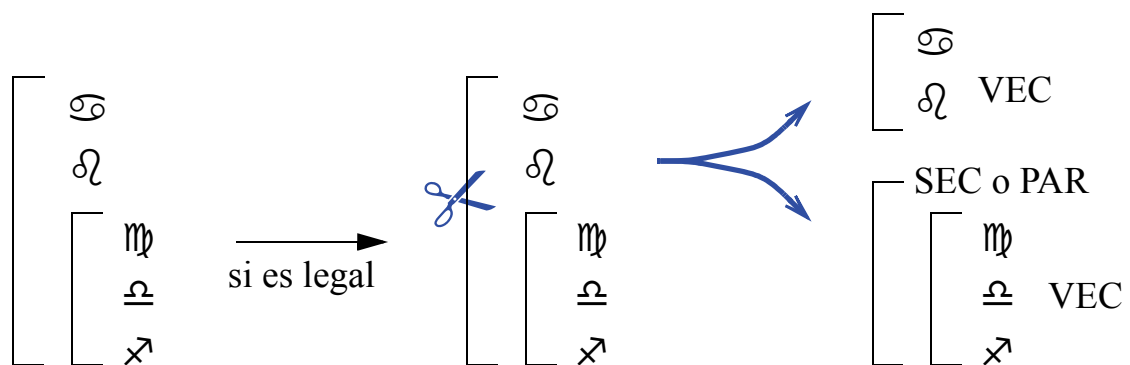


Vectorización parcial, distribución de bucles¹

Problema 2 de 2: dos niveles de anidación de bucles no perfectos

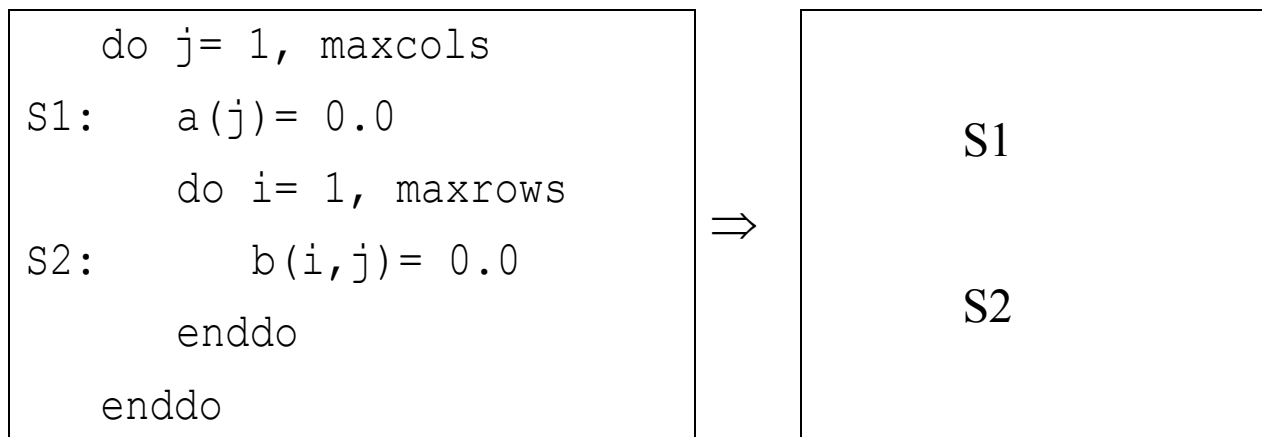
Transformación:

- si es legal, distribuir el bucle mas externo
- **vectorizar el bucle mas interno**, si es posible, aplicando lo anterior
- **paralelizar el bucle mas externo**, si es posible



no tocamos los espacios de iteración

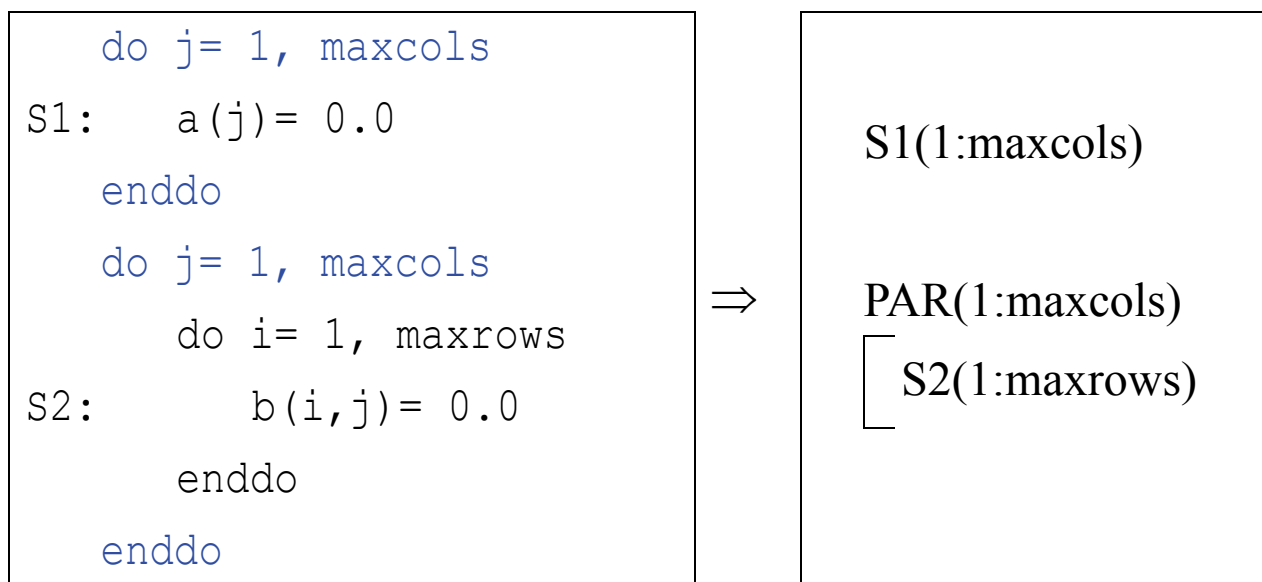
1) También llamado *loop fision*, corte de bucles



- Ejemplo 18 -

Transformación:

- distribuir el bucle exterior



Vectorización parcial/total: intercambio de bucles

Problema:

- ✓ en bucles anidados el recorrido del espacio de iteraciones impide:
 1. una fracción vectorial elevada
 2. un acceso eficiente a memoria
 3. vectorización

Transformación:

- ✓ intercambiar el orden de anidación de los bucles, si es legal

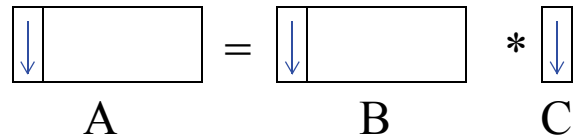
```
do j= 1, 64
```

```
  do i=1, 8
```

```
    A(i,j) = B(i,j) * C(i)
```

```
  enddo
```

```
enddo
```

 \Rightarrow


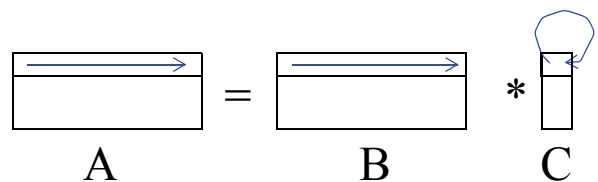
- Ejemplo 19 -

1. Intercambiar para aumentar la fracción vectorial

```
do i= 1, 8
```

```
  A(i,1:64) = B(i,1:64) * C(i)
```

```
enddo
```

 \Rightarrow


```
integer maxlen
parameter (maxlen = 64)
real A(maxlen, maxlen)

do i= 1, maxlen          ! intercambiar para mejorar el
  do j= 1, maxlen        ! acceso a la memoria multibancoa
    A(i,j)= 0.0
  enddo
enddo

do i= 1, 8               ! no intercambiamos para
  do j= 1, maxlen        ! preservar vectores largos
    A(i,j)= float(i+j)
  enddo
enddo
```

- Ejemplo 20 -

2. Intercambiar para mejorar el acceso a memoria

-
- a. En FORTRAN las matrices se almacenan por columnas: $a_{\text{fila columna}} \rightarrow a_{11}, a_{21}, a_{31}, \dots, a_{12}, a_{22}, a_{32}, \dots$
En C, se almacenan por filas.

```
real A(64,64)
do j= 2, 63
  do i= 1, 64
S: A(i,j)= (A(i,j-1) + 2*A(i,j) + A(i,j+1))*0.25
  enddo
enddo
```

- Ejemplo 21 -

3. No es necesario intercambiar para poder vectorizar

En este ejemplo el bucle j genera ciclos de dependencias, pero el bucle i no, y por tanto la sentencia es vectorizable en i

```
do j= 2, 63
  A(1:64,j)= (A(1:64,j-1) + 2*A(1:64,j) + A(1:64,j+1))*0.25
enddo
```

¿ Como lo hemos sabido ?


```
do j= 1, 64
  do i= 1, 64
S: A(i+1,j)= A(i,j) * B(i)
  enddo
enddo
```

El bucle i tiene
una recurrencia a distancia 1
que impide vectorizar

- Ejemplo 22 -

3. Intercambiar para poder vectorizar

En este ejemplo

1. el intercambio es legal, y
2. posibilita la vectorización del bucle interno

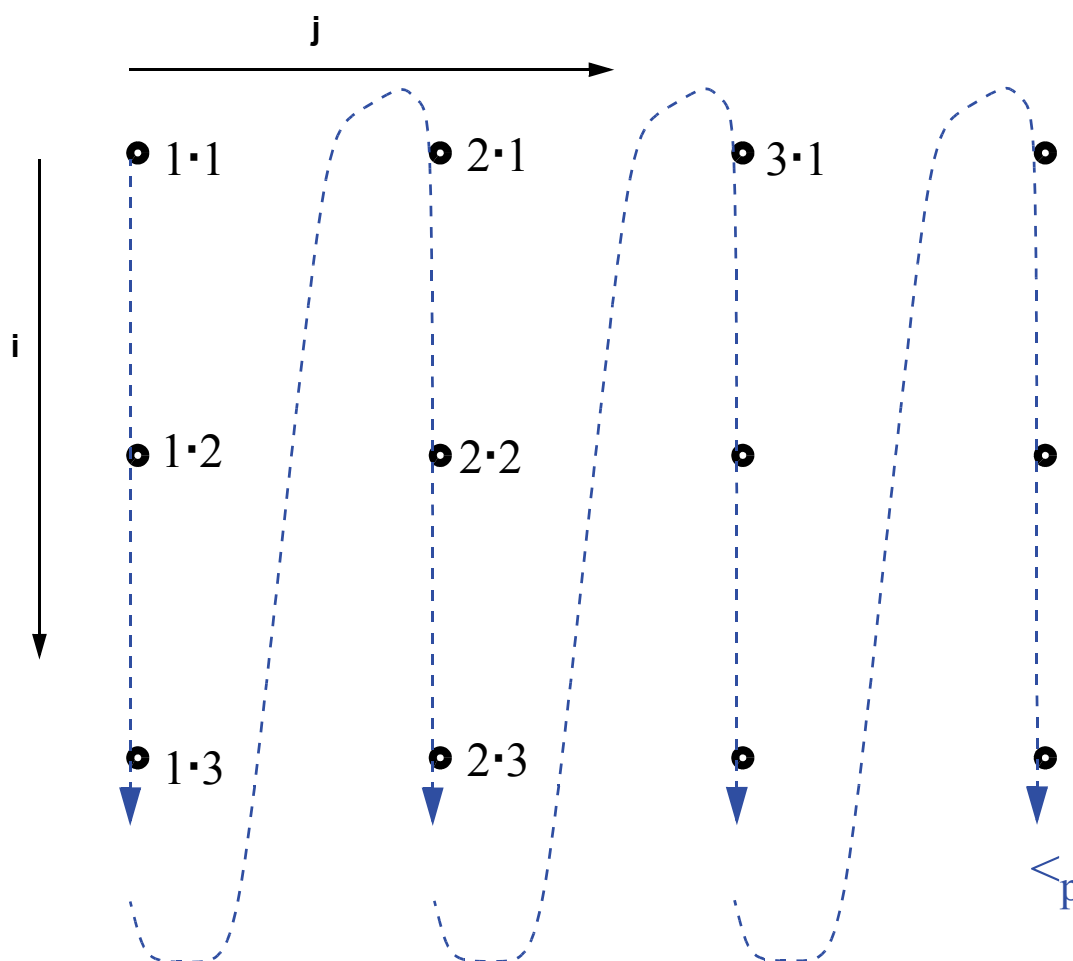
```
do i= 1, 64
  A(i+1,1:64)= A(i,1:64) * B(i)
enddo
```

¿ Como lo hemos sabido ?

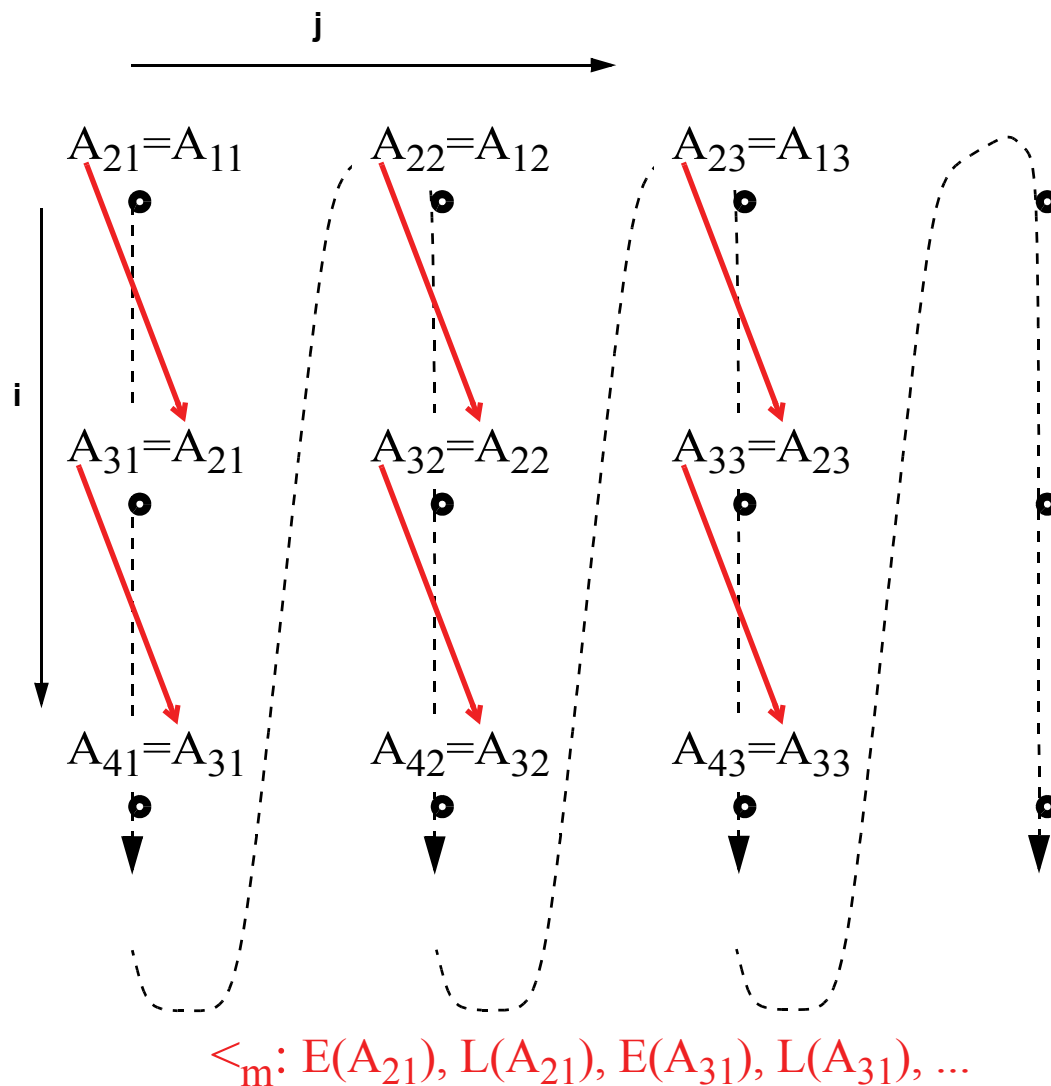
DIAGRAMA DE DEPENDENCIAS CON ANIDACIONES

```
do j= 1, 64
  do i= 1, 64
S: A(i+1,j)= A(i,j) * B(i)
  enddo
enddo
```

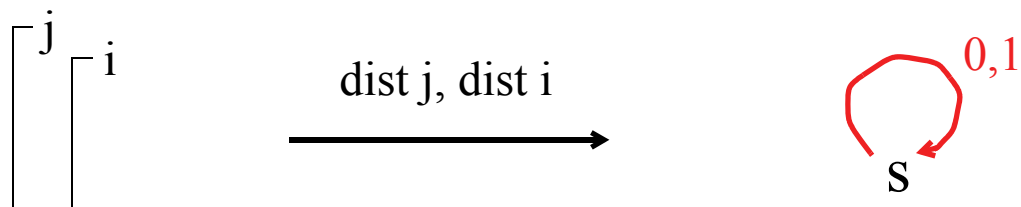
Esplter $\rightarrow \{ j=1:64 \times i=1:64 \} = \{ 1 \cdot 1, 1 \cdot 2, 1 \cdot 3, \dots 2 \cdot 1, 2 \cdot 2, 2 \cdot 3, \dots 3 \cdot 1, \dots \}$



- ❑ Al recorrer EspIt en el orden de programa, camino a trazos, los accesos a memoria determinan las dependencias. Cualquier otro orden que las respete, *es legal*

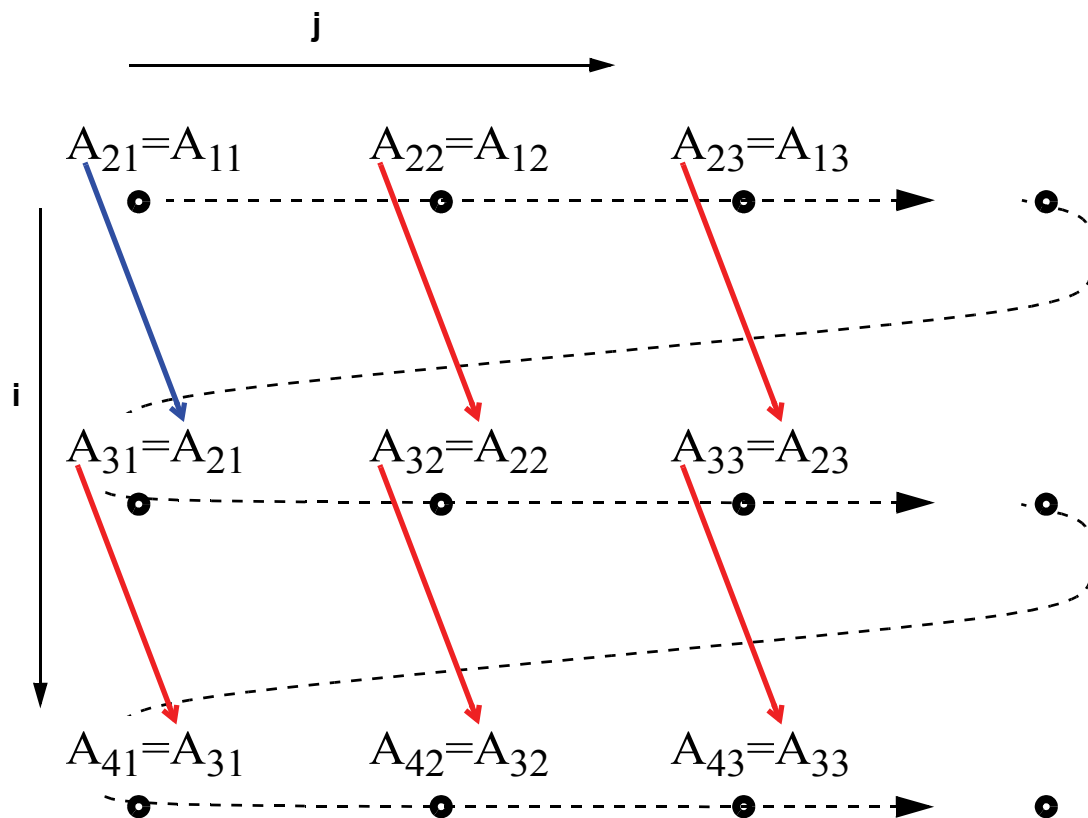


¿ Como representar las dependencias ?



Es una recurrencia en (i), y por tanto no es vectorizable !!

- ❑ Intercambiar los bucles es recorrer EspIt en otro orden:



\langle_m : E(A₂₁), E(A₂₂), E(A₂₃), ... L(A₂₁), E(A₃₁), L(A₂₂), E(A₃₂)...

- ❑ El nuevo recorrido, respeta las dependencias ?
- ✓ SI → el intercambio es legal
 - ✓ NO → no es legal
- ❑ Ahora el bucle interno (j) ejecuta los cálculos en otro orden.
En este nuevo orden NO hay dependencias entre cálculos consecutivos
Es posible vectorizar en (j) !!

```
do i= 1, 64
  A(i+1,1:64)= A(i,1:64) * B(i)
enddo
```

RESUMEN DEL PROCEDIMIENTO MANUAL

- 1) desplegar EspIt,
anotando referencias a memoria y dependencias
- 2) grafo de dependencias
- 3) ¿se puede vectorizar y paralelizar ?
- 4) ¿es legal el intercambio?
en caso afirmativo, repetir (3)
- 5) escoger la opción mejor

```
do i= 2, 64
  do j= 2, 32
S1:   A(i,j)= B(i,j)
S2:   B(i-1, j-1)= A(i-1,j)
  enddo
enddo
```

- Ejemplo 23 -

Estudiar las posibilidades de codificación y valorar su rendimiento

PASO (3)

$\left[\begin{array}{l} i \text{ PAR} \\ j \text{ SEC} \end{array} \right]$	sí/no
	<input type="checkbox"/>

$\left[\begin{array}{l} j \text{ PAR} \\ i \text{ SEC} \end{array} \right]$	sí/no
	<input type="checkbox"/>

$\left[\begin{array}{l} i \text{ SEC} \\ j \text{ VEC} \end{array} \right]$	sí/no
	<input type="checkbox"/>

$\left[\begin{array}{l} j \text{ SEC} \\ i \text{ VEC} \end{array} \right]$	sí/no
	<input type="checkbox"/>

$\left[\begin{array}{l} i \text{ PAR} \\ j \text{ VEC} \end{array} \right]$	sí/no
	<input type="checkbox"/>

$\left[\begin{array}{l} j \text{ PAR} \\ i \text{ VEC} \end{array} \right]$	sí/no
	<input type="checkbox"/>

```
void calc_array (int** a)
{
    for (i=1; i<N; ++i)
        for (j=1; j<N; ++j)
            a[i][j] = a[i-1][j] + a[i][j-1]
}
```

- Ejemplo 24 -

[Leer el texto y reflexionar](#)

“This code shows a nested loop operating on a 2D array with cross-iteration dependencies over both loops, making it appear serial. From the dependence graph it can be shown that iterations can be grouped into independent sets, allowing parallel execution if *loop skewing* and *interchange* are used. Techniques relying on dependence testing would overlook this parallelism. Furthermore, the 2D array in (a) is represented as an array of pointers to arrays, thwarting a parallelizing compiler’s attempt to statically analyze this section of code.”

Garcia, S.; Donghwan Jeon; Louie, C.; Taylor, M.B.; , "The Kremlin Oracle for Sequential Code Parallelization," IEEE Micro, vol.32, no.4, pp. 42-53, July-Aug. 2012
IEEE Micro Special Issue on Parallelization of Sequential Code

Reducción

En este contexto *reducir* es una operación matemática que quita una dimensión a una estructura de datos:

matriz \rightarrow vector, **vector** \rightarrow **escalar**, ...

Es una operación frecuente en cálculo científico.

Los repertorios vectoriales tienen instrucciones para operadores de reducción *asociativos*:

- ✓ máximo, mínimo, OR, AND, suma, producto
- ✓ ojo con el problema asociatividad / precisión finita

Transformación:

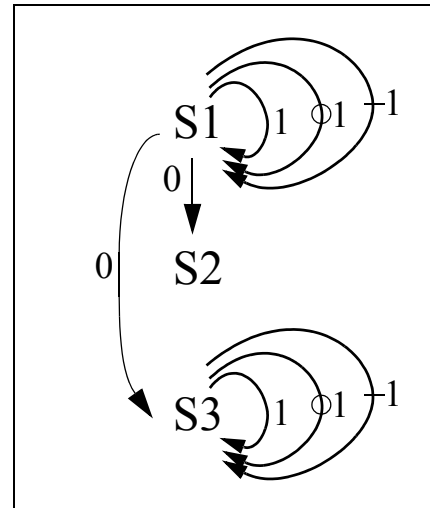
- ✓ detectar los autociclos de dependencias debidos al reuso de la variable escalar de reducción
- ✓ sustituir el bucle de reducción por la instrucción apropiada

```

s=0
p=1
DO i= 1, 128
S1:  A(i)= B(i) + C(i)
S2:  s = s + A(i)
S3:  p = p * A(i)
ENDDO

```

⇒



- Ejemplo 25 -

- vectorizar usando las instrucciones de lenguaje máquina apropiadas,
y en este caso seccionar

```

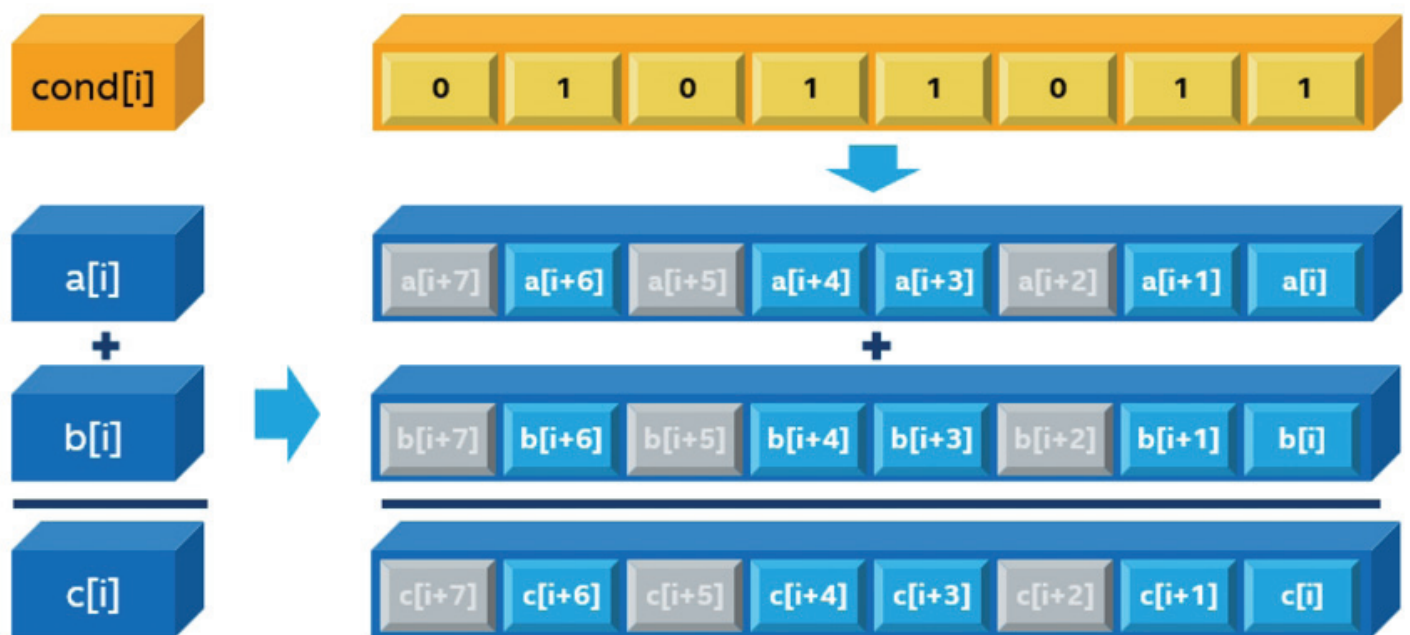
s= 0
p= 1
do i=1,2
  S1: A(64(i-1)+1:64i)= B(64(i-1)+1:64i) + C(64(i-1)+1:64i)
  S2: s = s + SUM(64(i-1)+1:64i)
  S3: s = p * PROD(64(i-1)+1:64i)
enddo

```


Bucles con condicionales: operaciones con máscara

Códigos con muchos saltos, por ejemplo, un bucle con una sentencia condicional en su cuerpo, pueden ser vectorizados usando máscaras vectoriales que bloquean las operaciones para las cuales la condición no es cierta

```
for (i = 0; i < MAX; i++)  
    if (cond(i))  
        c[i] = a[i] + b[i]
```



Punteros en C¹

The following loop may not get vectorized because of a potential aliasing problem between pointers a, b and c

```
void add (float *a, float *b, float *c) {  
    for (int i=0; i<SIZE; i++)  
        c[i] += a[i-1] + b[i];  
}
```

If the **restrict** keyword is added to the parameters, the compiler will trust you, that you will not access the memory in question with any other pointer and vectorize the code properly

```
void add (float * __restrict a,  
          float * __restrict b, float * __restrict c) {  
    for (int i=0; i<SIZE; i++)  
        c[i] += a[i-1] + b[i];  
}
```

The downside of using `restrict` is that not all compilers support this keyword, so your source code may lose portability.

- ❑ `restrict` es un cualificador complicado, que puede ayudar al compilador a generar ejecutables mas rápidos en contextos no relacionados con vectorizar. Para mas información:

How to Use the restrict Qualifier in C, Douglas Walls, Sun ONE Tools Group, Sun Microsystems, July 2003 (revised March 2006)

http://web.archive.org/web/20120225055041/http://developers.sun.com/solaris/articles/cc_restrict.html

1) http://software.intel.com/sites/products/documentation/studio/composer/en-us/2011Update/compiler_c/optaps/common/optaps_vec_use.htm

7. LEY DE AMDHAL

- ❑ Modelo matemático sencillo que calcula la mejora en rendimiento al ejecutar parte de las operaciones de un programa en menos tiempo:
 - ♦ un modo lento (escalar)
 - ♦ un modo rápido (vectorial)
- ❑ Sea N el número de operaciones en coma flotante de un programa y f_v el porcentaje que puede realizarse en modo vectorial
 - ♦ $TPF_{(s,v)}$: tiempo por FLOP en modo (escalar, vectorial)
 - ♦ $G = TPF_s / TPF_v$

$$\begin{aligned}
 \text{Speed_up} = T_{\text{escalar}} / T_{\text{vectorial}} &= \frac{N \cdot TPF_s}{N \cdot f_v \cdot TPF_v + N \cdot (1 - f_v) \cdot TPF_s} = \\
 &= \frac{1}{f_v \cdot \frac{TPF_v}{TPF_s} + (1 - f_v)} = \frac{1}{(1 - f_v) + f_v \cdot \frac{1}{G}} = \frac{1}{f_s + f_v / G}
 \end{aligned}$$

♦ $\text{Speed_up}(f_v \rightarrow 1) \rightarrow G$

♦ $G \rightarrow \infty \Rightarrow \text{Speed_up} \rightarrow 1/f_s$



$f_s \uparrow \uparrow \Rightarrow \text{Speed_up} \downarrow \downarrow$

"El modo lento va a limitar, aunque el modo rápido sea fantástico"



- ❑ Valor habitual de f_v entre 0,4 y 0,75 para programas compilados

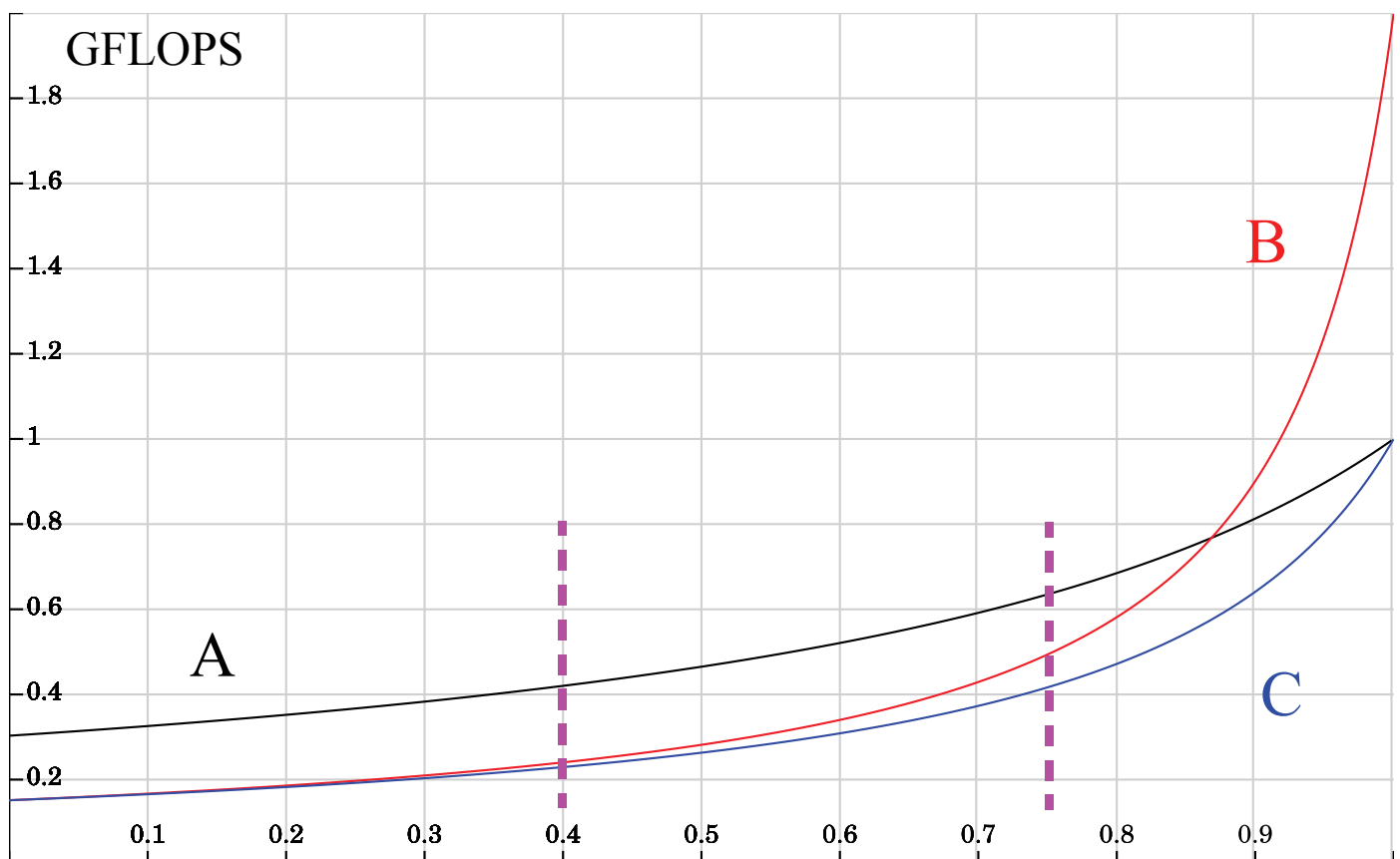


Un buen computador vectorial debe tener un buen procesador escalar

- ❑ Caso de estudio:

- Computador A: $TPF_s = 3,3ns$; $TPF_v = 1ns$ (G= 3,3)
- **Computador B**: $TPF_s = 6,6ns$; $TPF_v = 0,5 ns$ (G= 13,2)
- **Computador C**: $TPF_s = 6,6ns$; $TPF_v = 1 ns$ (G= 6,6)

$$R_v = R_s \cdot Sup = R_s \cdot \frac{1}{f_s + f_v / G}$$



Permalink to this graph:

<http://fooplots.com/plot/2hqbw5aoc7>

