Memoria Cache en Ripes

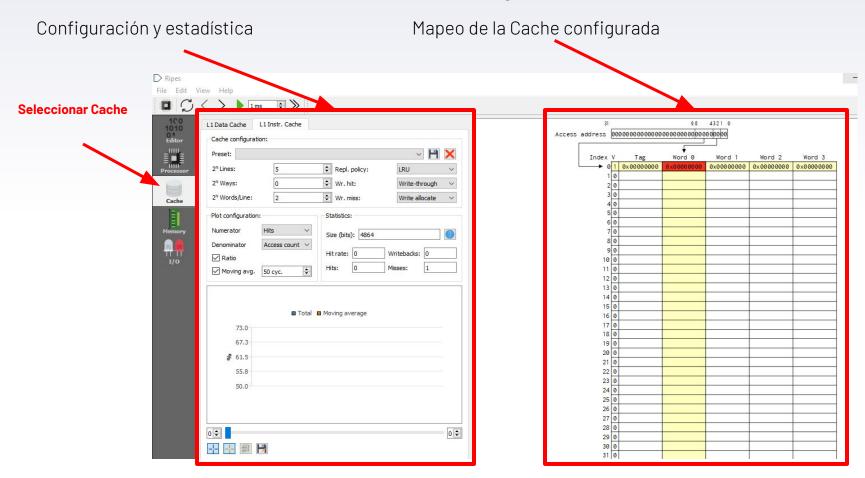
Memoria Cache en Ripes

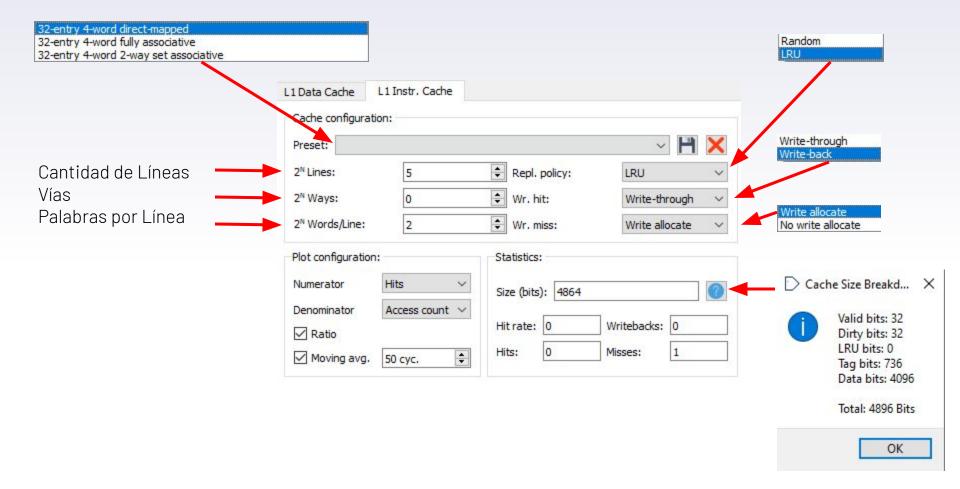
En la clase de Memoria Cache hemos aprendido

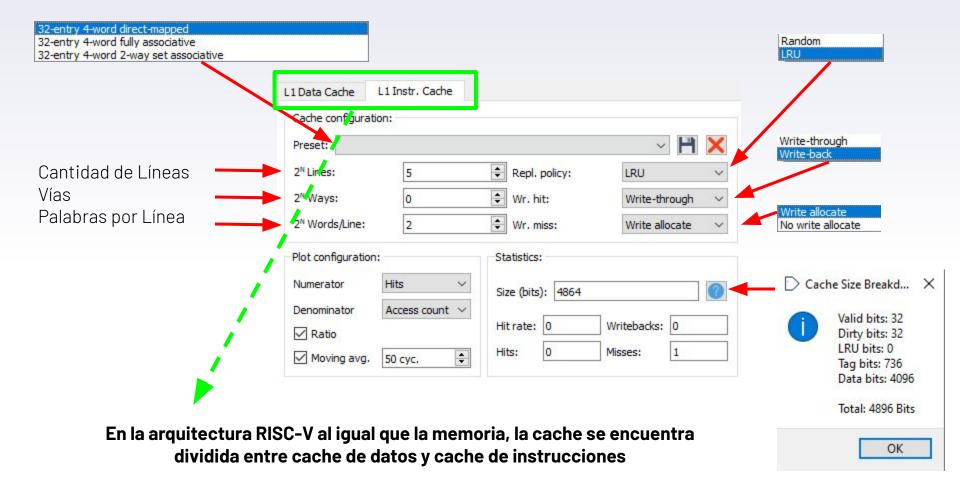
- Su principio de funcionamiento
- Métodos de asignación
- Ventajas y desventajas de cada método

En esta clase mediante algunos ejemplos veremos cómo aprovechar Ripes para comprender cómo esto se ve reflejado en la ejecución de un programa.

Al seleccionar cache aparecerá una pantalla como la siguiente

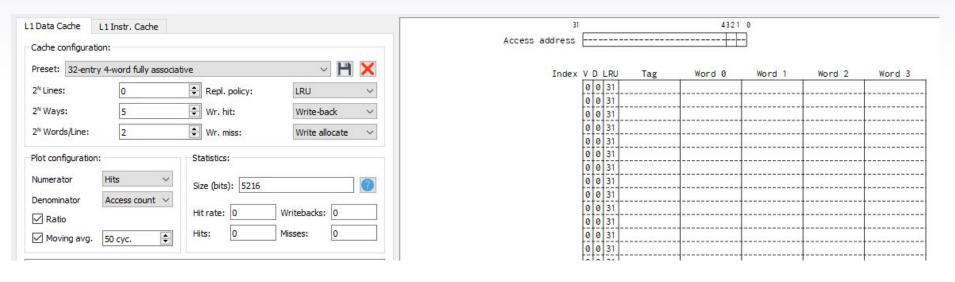




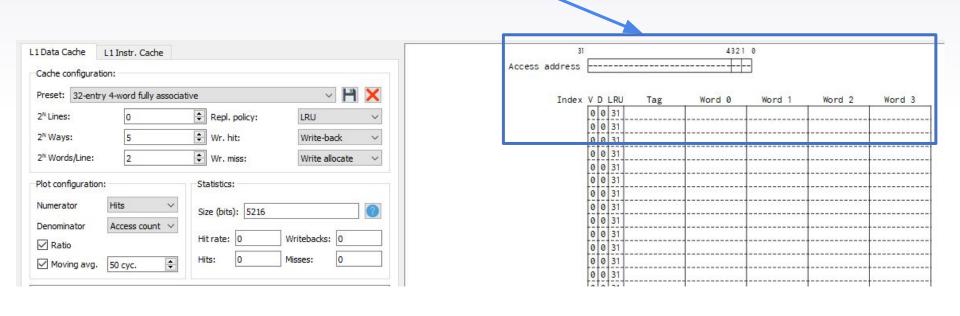


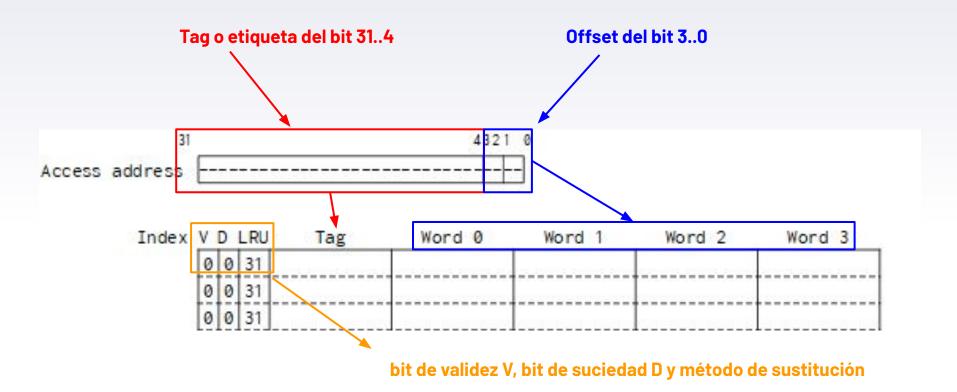
Asignación Asociativa

Veamos cómo queda configurada en Ripes una memoria cache asociativa de 32 líneas y 4 palabra por línea que es una de las opciones que esta pre-seteada



Entendamos el detalle





Tener en cuenta que los bit 0 y 1 del **Offset** corresponden a cada byte dentro de la palabra **Word**

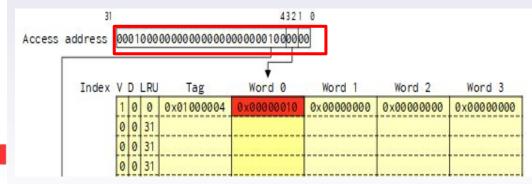
Analicemos un ejemplo

```
.data
 2 vector1: .word 0x0,0x10,0x20,0x30,0x40,0x50,0x60,0x70,
 3 vector2: .word 0x80,0x90,0xa0,0xb0,0xc0,0xd0,0xe0,0xf0
 4 cont: .word 16
5 .text
 6 lui x7,0x10000
 7 la x13, cont
8 lw x5, 0, x13
9 loop:
       lw x10, 0, x7
       addi x10, x10,1
       sw x10, 0, x7
       addi x5, x5, -1
       addi x7, x7, +4
       bne x0, x5, loop
16 fin:
       beg x0, x0, fin
```

Un vector de 16 posiciones (o dos de 8 contiguos) al cual le vamos a sumar 1 (uno) a cada posición.

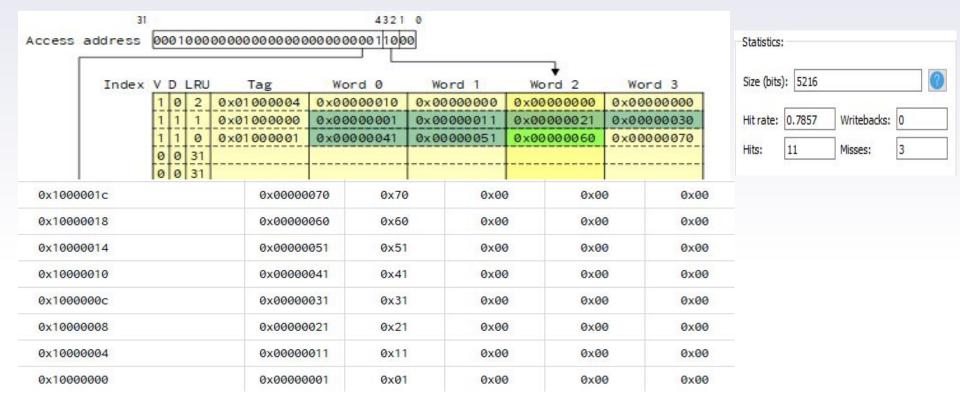
Vamos analizar la memoria cache de datos, tener en cuenta que esta comienza en 0x10000000

```
1 .data
2 vector1: .word 0x0,0x10,0x20,0x30,0x40,0x50,0x60,0x70,
3 vector2: .word 0x80,0x90,0xa0,0xb0,0xc0,0xd0,0xe0,0xf0
4 cont: .word 16
5 .text
6 lui x7,0x10000
7 la x13,cont
8 lw x5, 0,x13
9 loop:
```



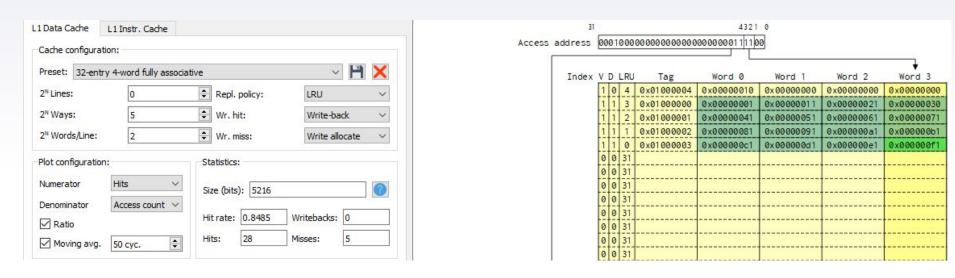
Vemos que el primer fallo se produce al cargar x5 con el valor contenido por cont que se encuentra en la posición de memoria 0x10000040, TAG 0x0100004, Offset 0x0. Vemos que el bit de validez es **V=1** indicando que esa línea de cache corresponde al entorno actual

Address	Word	Byte 0	Byte 1	Byte 2	Byte 3
0x10000048	X	Х	X	X	Х
0x10000044	X	X	X	X	X
0x10000040	0x00000010	0x10	0x00	0x00	0x00
0x1000003c	0x000000f0	0xf0	0x00	0x00	0x00
0x10000038	0x000000e0	0xe0	0x00	0x00	0x00



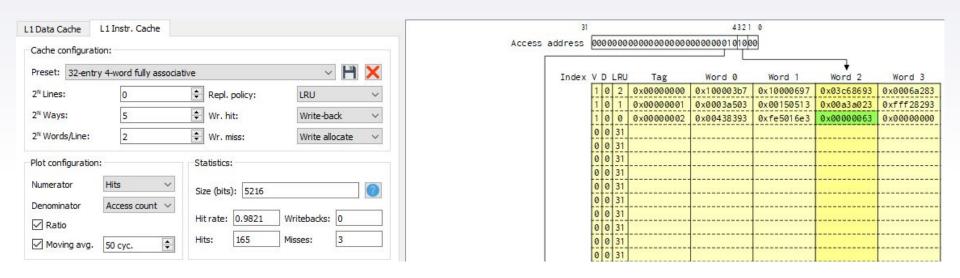
Se puede apreciar que a medida que el programa se ejecuta, los bit de validez y suciedad se modifican como así también el bit de LRU (Least Recently Used) indicando cual es la última línea a la que se accedió. En el cuadro de estadísticas se ve el Hit Rate y la cantidad de Hits y Misses al momento de la ejecución

Al finalizar la ejecución del programa podremos evaluar la eficiencia de la configuración de la cache ayudándonos del cuadro de estadísticas



Podemos cambiar la cantidad de líneas y palabras por línea para analizar las diferencias en su funcionamiento y eficiencia

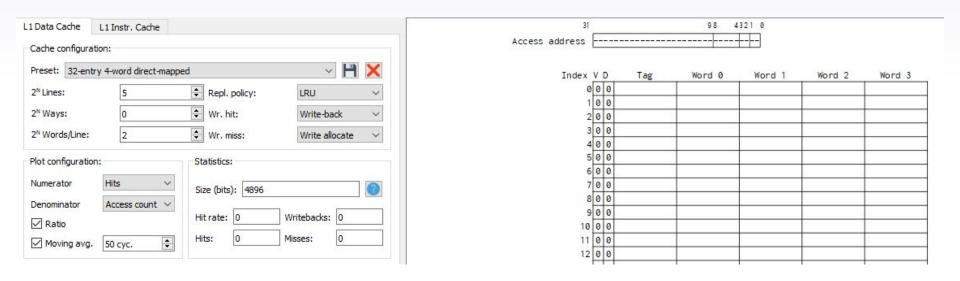
De la misma forma se puede analizar el comportamiento de la cache de instrucciones, y al finalizar la ejecución del programa podremos evaluar cómo se comportó la misma



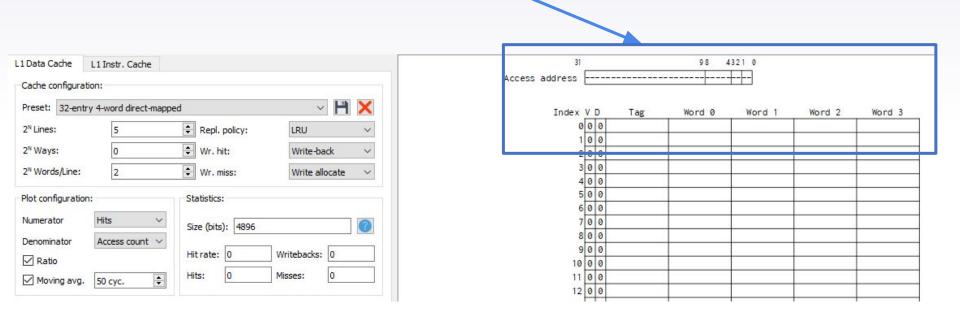
También podemos apreciar que el bit de suciedad **D** siempre se mantuvo en cero al no actualizarse ningún valor en la memoria de instrucciones

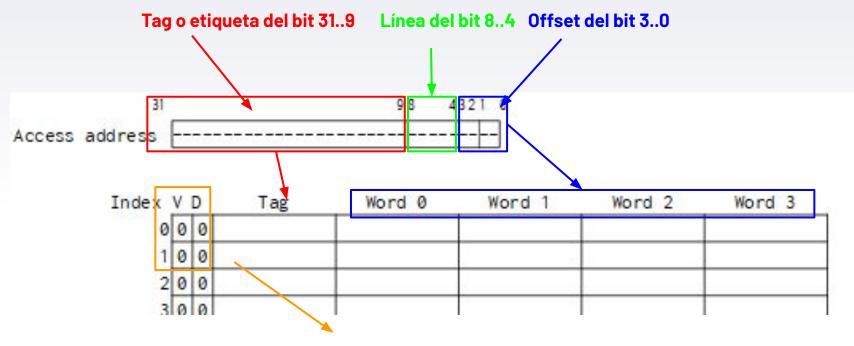
Asignación Directa

Veamos cómo queda configurada en Ripes una memoria cache con asignación directa de 32 líneas y 4 palabra por línea otra de las opciones pre-seteada



Entendamos el detalle





Bit de validez V bit de suciedad D Observar que en asignación directa no posee un método de reemplazo explícito

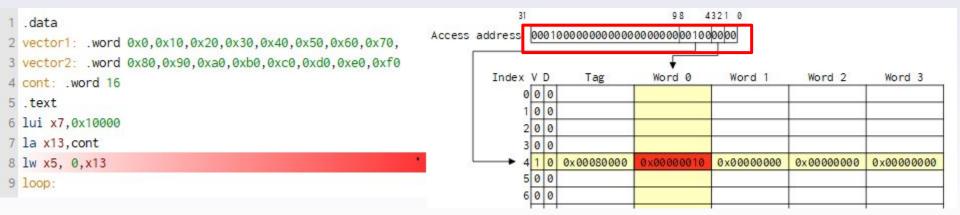
En asignación directa recordemos que se reduce el tamaño de la etiqueta y además tenemos el valor de la Línea. Tener en cuenta que los bit 0 y 1 del **Offset** corresponden a cada byte dentro de la palabra **Word**

Analicemos el mismo ejemplo

```
.data
 2 vector1: .word 0x0,0x10,0x20,0x30,0x40,0x50,0x60,0x70,
 3 vector2: .word 0x80,0x90,0xa0,0xb0,0xc0,0xd0,0xe0,0xf0
 4 cont: .word 16
5 .text
 6 lui x7,0x10000
 7 la x13, cont
8 lw x5, 0, x13
9 loop:
       lw x10, 0, x7
       addi x10, x10,1
       sw x10, 0, x7
       addi x5, x5, -1
       addi x7, x7, +4
       bne x0, x5, loop
16 fin:
       beg x0, x0, fin
```

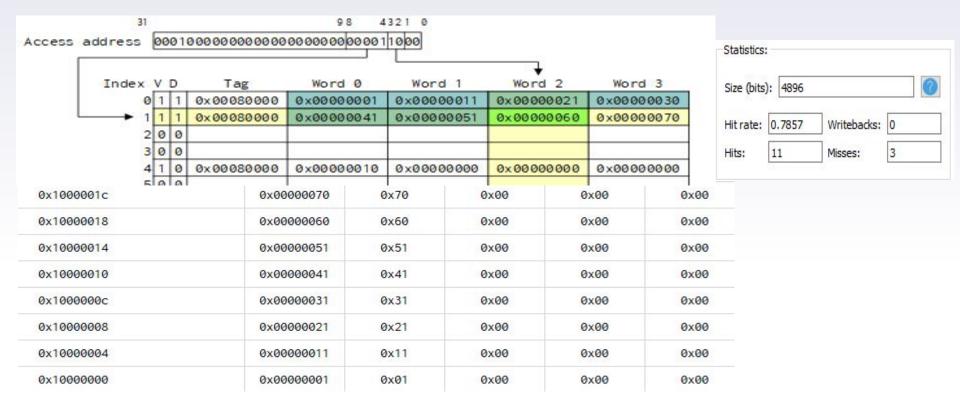
Un vector de 16 posiciones (o dos de 8 contiguos) al cual le vamos a sumar 1 (uno) a cada posición.

Vamos analizar la memoria cache de datos, tener en cuenta que esta comienza en 0x10000000



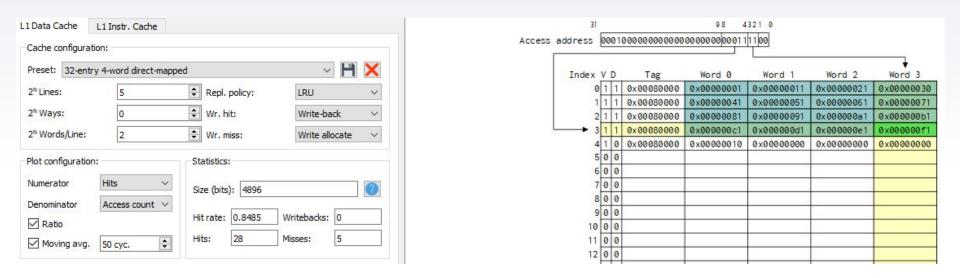
Vemos que el primer fallo, al igual que antes, se produce al cargar x5 con el valor que apunta cont en la posición de memoria 0x10000040, TAG 0x00080000, Línea 0x04, Offset 0x0. Observemos que el bit de validez es **V=1** indica que esa línea de cache corresponde al entorno actual.

Address	Word	Byte 0	Byte 1	Byte 2	Byte 3
0x10000048	X	X	X	X	X
0x10000044	X	X	X	X	×
0x10000040	0x00000010	0x10	0x00	0x00	0x00
0x1000003c	0x000000f0	0xf0	0x00	0x00	0x00
0x10000038	0x000000e0	0xe0	0x00	0x00	0x00



Se puede apreciar que a medida que el programa se ejecuta los bit de validez y suciedad se modifican. En el cuadro de estadísticas se ve el Hit Rate y la cantidad de Hits y Misses al momento de la ejecución

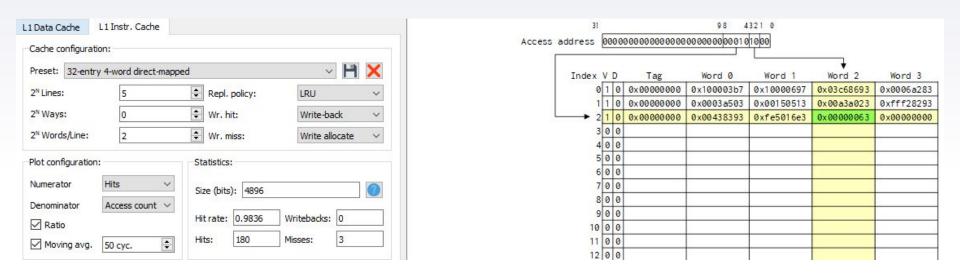
Al finalizar la ejecución del programa podremos evaluar la eficiencia de la configuración de la cache ayudándonos del cuadro de estadísticas



Podemos analizar la diferencia entre un método y otro comparando los resultados, pudiendo apreciar que no existen diferencias en este programa.

Se puede cambiar la cantidad de líneas y palabra por línea para analizar las diferencias en su funcionamiento y eficiencia

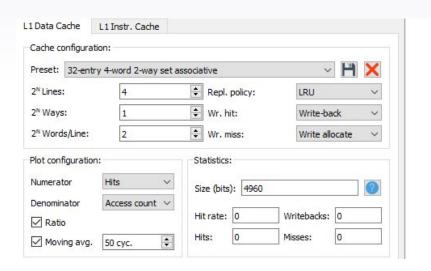
De la misma forma se puede analizar el comportamiento de la cache de instrucciones, y al finalizar la ejecución del programa podremos evaluar cómo se comportó la misma

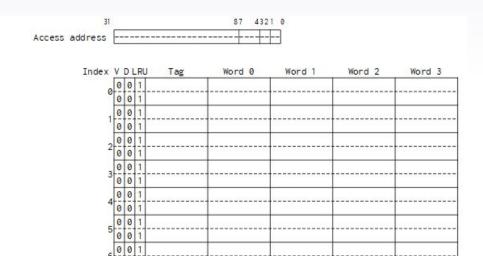


Podemos apreciar que el bit de suciedad **D** siempre se mantuvo en cero al no actualizarse ningún valor en la memoria de instrucciones

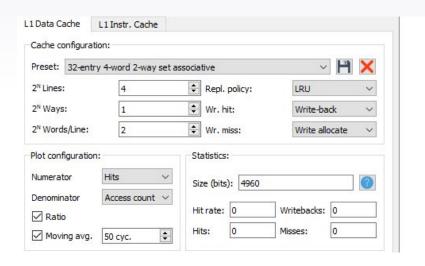
Asignación Asociativa por Conjuntos

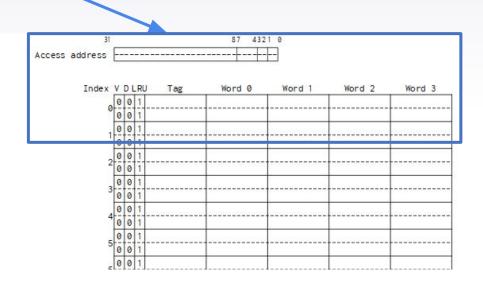
Veamos cómo queda configurada en Ripes una memoria cache con asignación asociativa por conjuntos de 2 vías y 4 palabra por línea, otra de las opciones pre-seteadas

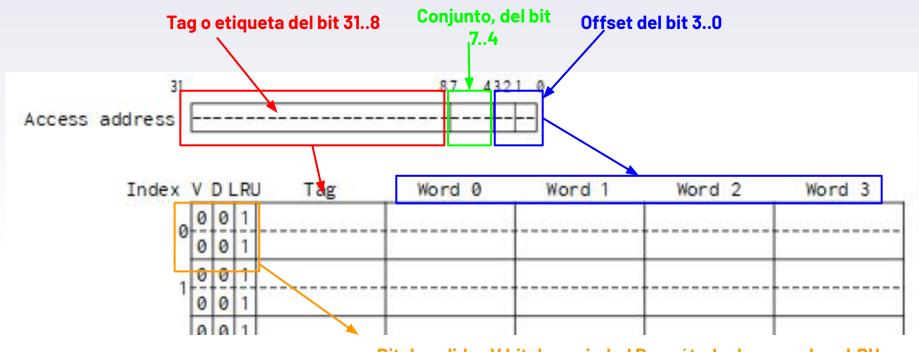




Entendamos el detalle







Bit de validez V bit de suciedad D y método de reemplazo LRU. Observar que aquí sí existe método de reemplazo explícito.

En asignación asociativa por conjunto recordemos que se incrementa el tamaño de la etiqueta dependiendo de la cantidad de vías y por ende el tamaño de la Línea se reduce.

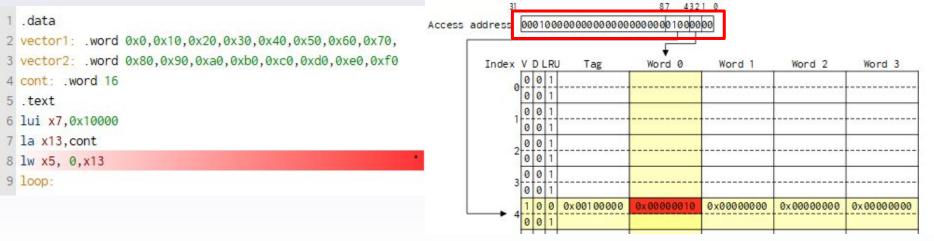
Tener en cuenta que los bit 0 y 1 del **Offset** corresponden a cada byte dentro de la palabra **Word**

Analicemos el mismo ejemplo

```
.data
 2 vector1: .word 0x0,0x10,0x20,0x30,0x40,0x50,0x60,0x70,
 3 vector2: .word 0x80,0x90,0xa0,0xb0,0xc0,0xd0,0xe0,0xf0
 4 cont: .word 16
5 .text
 6 lui x7,0x10000
 7 la x13, cont
8 lw x5, 0, x13
9 loop:
       lw x10, 0, x7
       addi x10, x10,1
       sw x10, 0, x7
       addi x5, x5, -1
       addi x7, x7, +4
       bne x0, x5, loop
16 fin:
       beg x0, x0, fin
```

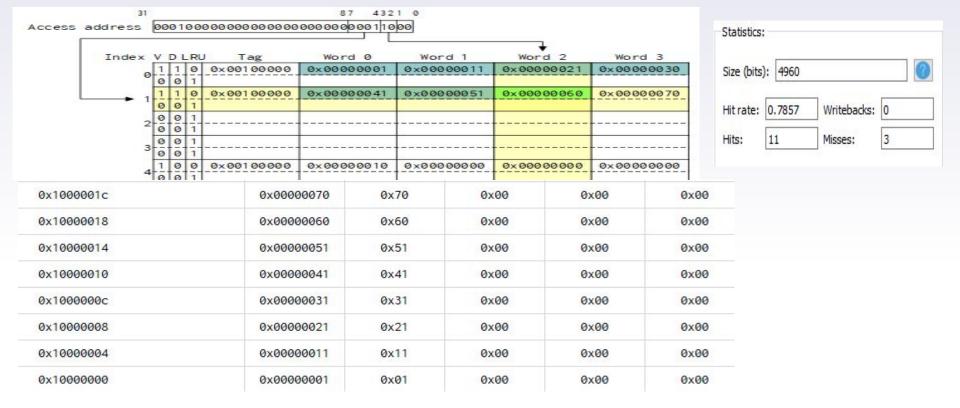
Un vector de 16 posiciones (o dos de 8 contiguos) al cual le vamos a sumar 1 (uno) a cada posición.

Vamos analizar la memoria cache de datos, tener en cuenta que esta comienza en 0x10000000



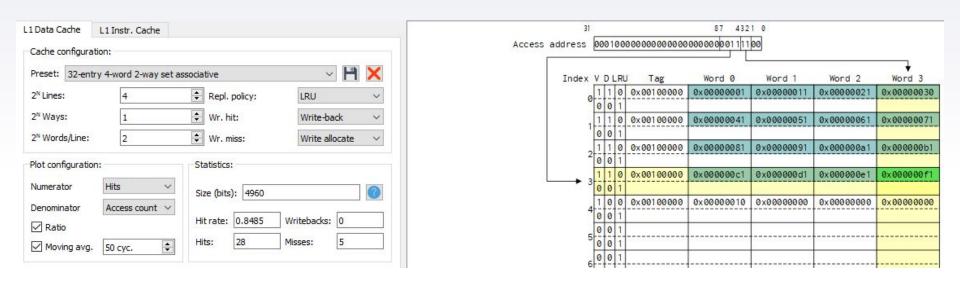
Vemos que el primer fallo, al igual que antes, se produce al cargar x5 con el valor del contador en la posición de memoria 0x10000040, TAG 0x00100000, Línea 0x4, Offset 0x0 vemos que el bit de validez es **V=1** indicando que esa línea de cache corresponde al entorno actual

Address	Word	Byte 0	Byte 1	Byte 2	Byte 3
0x10000048	X	Х	X	X	X
0x10000044	X	X	X	X	X
0x10000040	0x00000010	0x10	0x00	0x00	0x00
0x1000003c	0x000000f0	0xf0	0x00	0x00	0x00
0x10000038	0x000000e0	0xe0	0x00	0x00	0x00



Se puede apreciar que a medida que el programa se ejecuta los bit de validez y suciedad se modifican como así también el bit LRU. En el cuadro de estadísticas se ve el Hit Rate y la cantidad de Hits y Misses al momento de la ejecución.

Al finalizar la ejecución del programa podremos evaluar la eficiencia de la configuración de la cache ayudándonos del cuadro de estadísticas

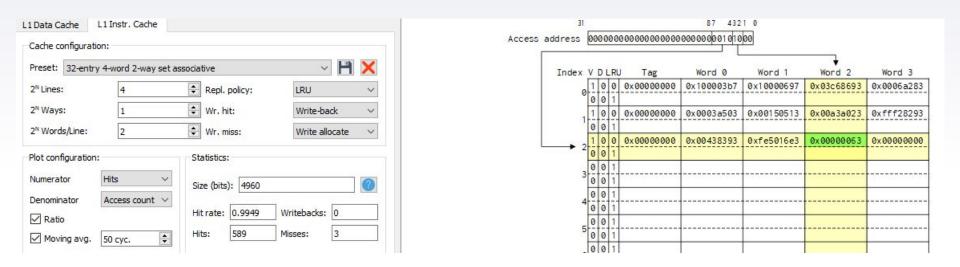


Podemos analizar las diferencias entre los distintos métodos de asignación comparando los resultados, pudiéndose apreciar que no existen diferencias en el ejemplo dado.

Se pueden cambiar la cantidad de líneas y palabras por línea y vías para analizar las diferencias en su funcionamiento y eficiencia

En este caso podemos apreciar la diferencia existente respecto a la asignación directa de la forma que se ubican los datos en la cache.

De la misma forma se puede analizar el comportamiento de la cache de instrucciones, y al finalizar la ejecución del programa podremos evaluar cómo se comportó la misma



Podemos observar que el bit de suciedad **D** siempre se mantuvo en cero al no actualizarse ningún valor en la memoria de instrucciones, como así también la disposición de los datos.

¿ Cuál es el mejor método de asignación?

A esta altura ya podrán inferir que no existe un método mejor que otro.

Solo podemos afirmar que cada uno tiene sus ventajas y/o desventajas.

Dependerá no solo del tamaño de la cache y su método de asignación sino también de la optimización en la utilización de la memoria. Es decir donde se almacenan los datos y cómo se acceden a ellos tiene un impacto en la eficiencia del uso de la cache.

Veamos qué sucede si....

Tenemos dos vectores de 32 words cada uno

V1 - ubicado a partir de la Dirección 0x10000000

V2 - ubicado a partir de la Dirección 0x20000000

Y queremos intercambiar los valores de V1 en V2 y los de V2 en V1

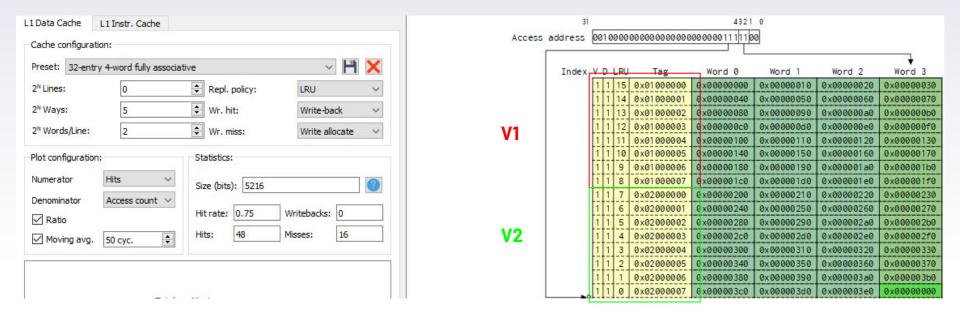
Configurando la memoria cache en cada uno de los formatos pre-seteados en Ripes

Generamos los vectores de 32 posiciones uno apuntado por x7 y otro por x8

```
lui x7,0x10000
2 addi x10, x10, 32
  add x6, x6, x0
4 carga1:
       sw x6, 0, x7
       addi x7, x7, 4
       addi x6,x6,0x10
      addi x10, x10, -1
       bne x0, x10, carga1
10 lui x8,0x20000
11 addi x10,x10,32
12 carga2:
         sw x6, 0, x8
       addi x8,x8,4
       addi x6,x6,0x10
       addi x10, x10, -1
       bne x0,x10,carga2
```

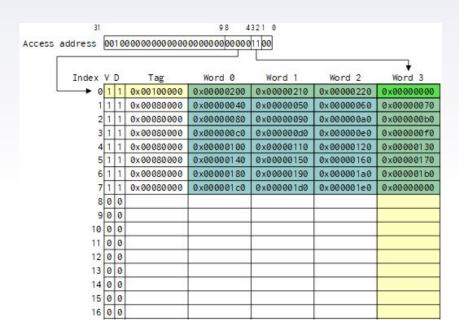
y realizamos el intercambio de sus valores

```
19 addi x7,x7,-128
20 addi x8,x8,-128
21 addi x10,x0,32
23 intercambio:
       lw x11,0,x7
       lw x12,0,x8
       sw x12,0,x7
       sw x11,0,x8
       addi x7,x7,4
       addi x8, x8, 4
       addi x10, x10, -1
       bne x0,x10,intercambio
32 fin:
33
       beg x0, x0, fin
```

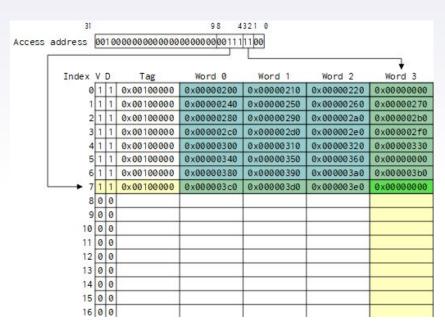


Vemos que al finalizar la carga, ambos vectores se encuentran cargados en la memoría cache si la asignación es asociativa

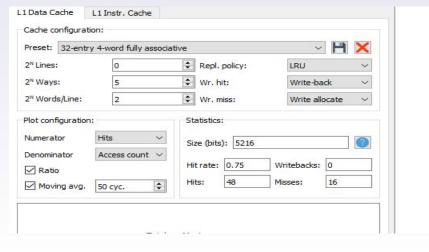
Que sucede en asignación directa

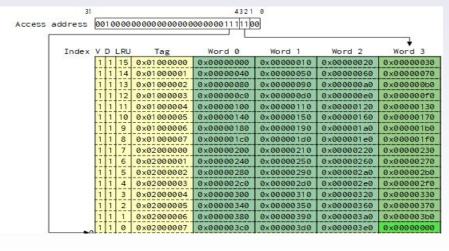


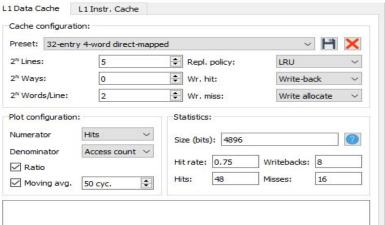
Vemos que al comenzar a cargar el segundo vector este corresponde a la misma línea por lo que se reemplaza el contenido anterior provocando un fallo

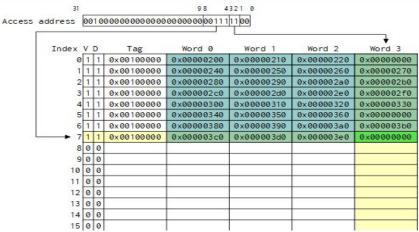


Al finalizar la carga solo los valores correspondientes al vector 2 permanecen en la memoria cache









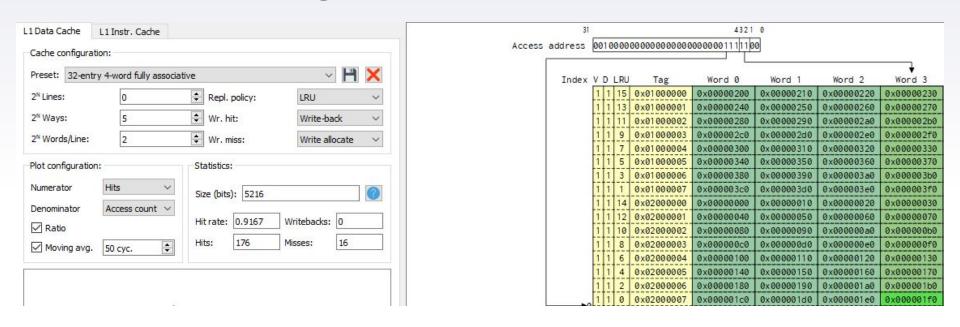
Podemos apreciar que la cantidad de fallos en esta etapa fue la misma en ambos casos, pero

¿Cuál es el mayor inconveniente a partir de ahora?

Veamos que sucede al intercambiar los valores de los vectores

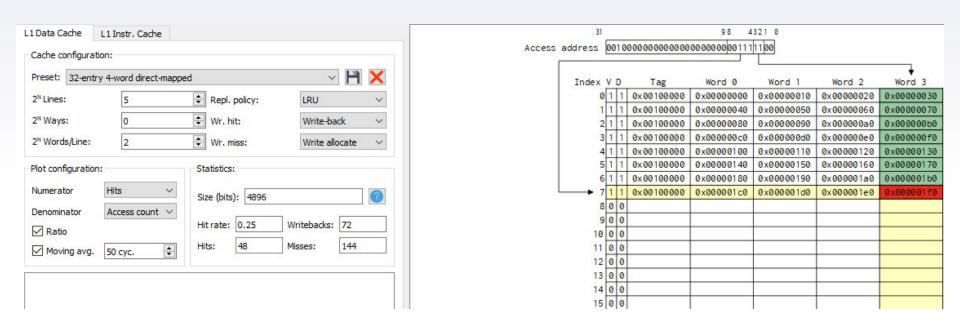
 ¿Cuál cree Ud que será el método de asignación que mejor se comporte en este caso?

Asignación Asociativa



Podemos observar una tasa de aciertos del 91,67% con 176 aciertos y 16 fallos

Asignación Directa



Podemos observar una tasa de aciertos del 25,00% con 48 aciertos y 144 fallos



,										67			A	sign	ació	ón Di	irec	ta					г						1805		
										Etiqu	eta -	TAG													Línea	9			Of	set	
31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	7	0	0
	0			- 8	3			()			()			()			()		0			0				0	

										As	igna	ciór	As	ocia	tiva	por	con	junt	o de	2 V	ias										
9	Si 103	Etiqueta - TAG																Con	junto			Off	fset								
31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	43 40	1			()			())				0			()				0				5	

													28	Addr	ess -	Dire	cción	rgi													
31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	2	2			()			())			()			()			- ()			0)	

3													Asi	gna	cion	Asc	ocial	iva													
												Eti	quet	a - Ta	4G													2000,000	Of	set	
31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
		2			()			- 1	D			- 1)			()			()				0				0	

													Α	sign	ació	in Di	irec	a											K		
										Etiqu	eta -	TAG		100								- 1		ı	Línea	•	1		Off	set	
31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	1 0 0								()		-	()		-		0	-	0		()			()				

										Asi	igna	ciór	As	ocia	tiva	por	con	junte	o de	2 V	ias			8			8	6			
07000	********		Etiqueta - TAG 9 28 27 26 25 24 23 22 21 20 19 18 17 16 15 14 13 12 11 10 9 8														Con	junto		2.00.000	Off	set									
31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	- 2	2			- ()			- 1	0			- 3)			()			()				0			()	

Podemos ver que en asignación Directa, ambas direcciones corresponden a la misma Línea dentro de la memoria Cache

Esto produce un fallo en cada reemplazo de los valores de los vectores

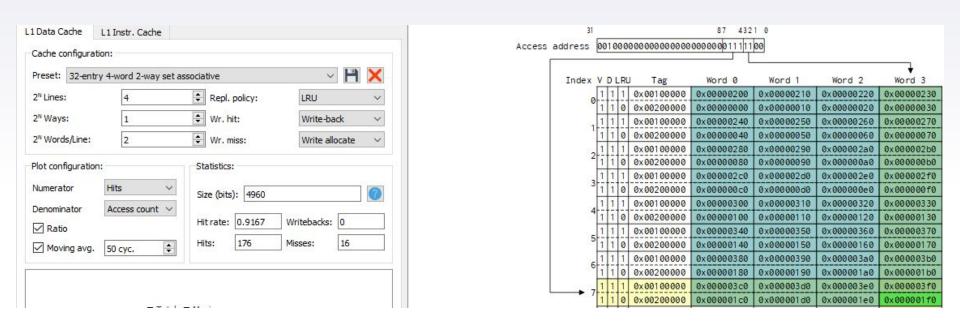
Con lo que acabamos de ver:

- ¿Podemos decir que un método es mejor que el otro?
- ¿depende solamente del tamaño de la memoria cache?
- ¿Cuál piensa usted que sería la solución si no podemos cambiar el programa?

Con lo que acabamos de ver:

- ¿Podemos decir que un método es mejor que el otro?
- ¿depende solamente del tamaño de la memoria cache?
- ¿Cuál piensa usted que sería la solución si no podemos cambiar el programa?
- ¿Podría solucionarse con asignación asociativa por conjunto de 2 vías?

Asignación Asociativa por conjuntos 2 vías



Podemos observar una tasa de aciertos del 91,67% con 176 aciertos y 16 fallos Vemos que para esta asignación también ambos vectores permanecen en memoria cache aunque ahora se ven intercalados en su posición

¿Qué pasa con la Cache de Instrucciones?

- ¿Cree Ud que tendremos el mismo inconveniente?
- ¿Es importante la manera de organizar nuestro algoritmo?

Veamos este fragmento de programa

```
addi x10, x0, 32
   intercambio:
                       #PC 0x000000048
                                                               137
       lw x11,0,x7
                                                               138 salto:
24
       lw x12,0,x8
                                                               139
                                                                       sw x11,0,x8
25
       sw x12,0,x7
26
                                                                       addi x7,x7,4
                                                               140
       jal salto #PC 0x000000054 salto a 0x00000026c
                                                               141
                                                                       addi x8,x8,4
27
       sw x11,0,x8
                                                                       addi x10, x10,-1
                                                               142
28
       addi x7,x7,4
                                                               143
                                                                       jal vuelve
                                                                                               #PC 0000027C salta a 00000068
29
       addi x8, x8, 4
                                                                       bne x0,x10,intercambio
                                                               144
30
       addi x10, x10, -1
                                                               145
31 vuelve:
                                                               146
       bne x0,x10,intercambio #PC 0x00000068
   addi x7, x7, -128
34 addi x8, x8, -128
                                                                                             ial x1 -532 <vuelve>
                                                                27c:
                                                                            dedff0ef
                                  jal x1 536 <salto>
     54:
                 218000ef
```

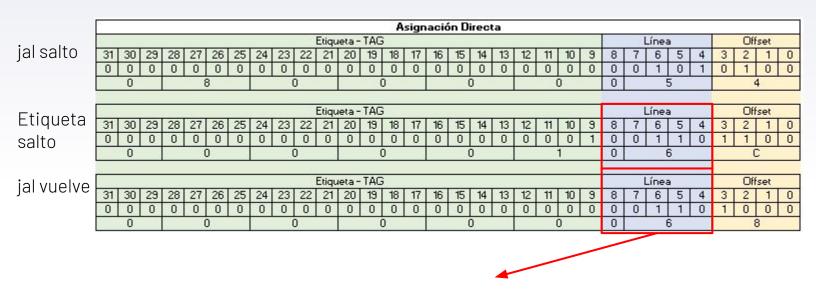
Analicemos que sucede con la memoria cache de instrucciones en el momento que se produce **jal salto** y **jal vuelve**

Asignación Asociativa

				Asignacion	Asociativa			
				Etiqueta - TAG				Offset
jal salto	31 30 29 28 3	27 26 25 24 3	23 22 21 20	19 18 17 16	15 14 13 12	11 10 9 8	7 6 5 4	3 2 1 0
,	0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0	0 1 0 1	0 1 0 0
	0	0	0	0	0	0	5	4
		.00				0	5	
	B			Etiqueta - TAG				Offset
Etiqueta	31 30 29 28 3	27 26 25 24 3	23 22 21 20	19 18 17 16	15 14 13 12	11 10 9 8	7 6 5 4	3 2 1 0
salto	0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0	0 0 1 0	0 1 1 0	1 1 0 0
Saitu	0	0	0	0	0	2	6	C
ialyulahya				Etiqueta - TAG				Offset
jal vuelve	31 30 29 28 3	27 26 25 24 3	23 22 21 20	19 18 17 16	15 14 13 12	11 10 9 8	7 6 5 4	3 2 1 0
	0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0	0 1 1 0	1 0 0 0
	0	0	0	0	0	0	6	8

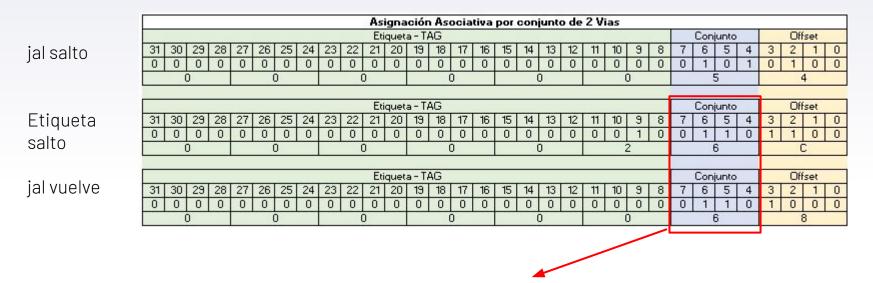
Podemos ver que dependiendo del tamaño de la memoria cache no tendríamos conflicto con este salto, en este ejemplo usamos la cache de 32 líneas 4 palabras por línea

Asignación Directa



Aquí podemos ver el conflicto que aparece dado que tanto la etiqueta salto con la instrucción **jal vuelve** pertenecen a la misma línea de cache al ser esto un salto repetitivo generará tantos fallos como iteraciones haya, en este ejemplo usamos la cache de 32 líneas 4 palabras por línea

Asignación Asociativa por conjuntos 2 vías



Podemos ver que a pesar de pertenecer al mismo conjunto al ser de dos vias no presentará inconvenientes como en asignación directa, en este ejemplo usamos la cache de 32 líneas 4 palabras por línea asociativa por conjunto en 2 vías

Hasta aquí hemos visto el comportamiento de la cache en Ripes con las configuraciones pre-establecidas en el simulador de la arquitectura Risc-V

Los invitamos a que realicen diferentes pruebas modificando las configuraciones de la memoria cache y analicen cómo se comporta ante diferentes situaciones

Esperamos que con lo que han aprendido se comprenda porque la memoria cache colabora a mejorar el funcionamiento de todo el sistema