

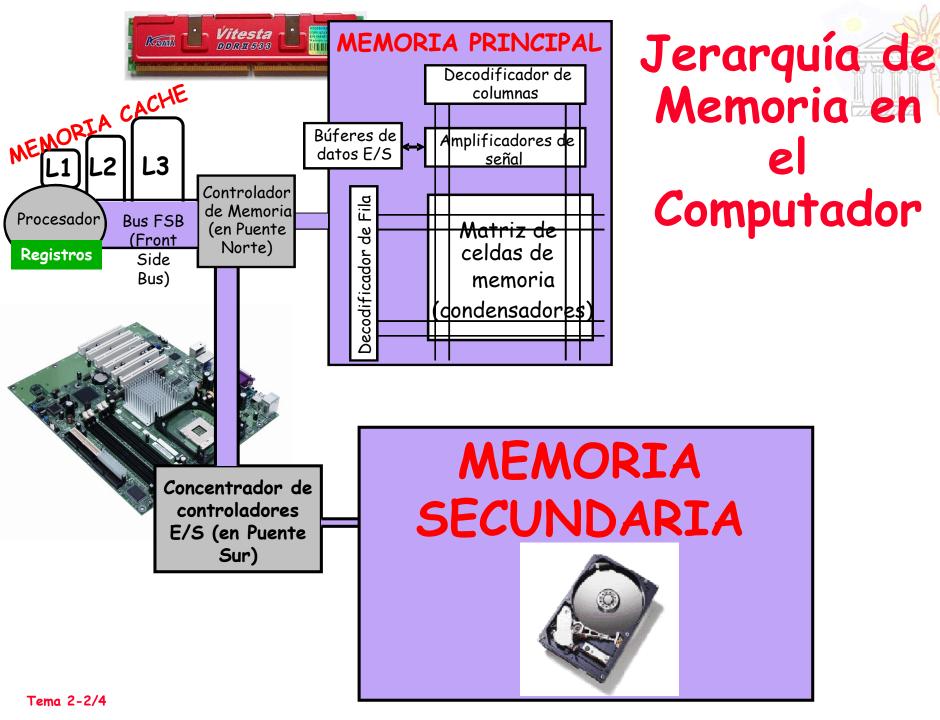




Arquitectura de Computadores Escuela de Ingeniería Informática Universidad de Las Palmas de Gran Canaria

#### Sumario

- · Jerarquía de memoria en un computador
- Tiempos típicos de acceso
- Direccionamiento
- Bloque de memoria
- Big/Little Endian
- Localidades
- · Circuitos: DIMM y controlador
- Prestaciones
- Influencia de las arquitecturas de memoria principal sobre el tiempo de penalización de la memoria cache
- Ejercicio



# Importantes parámetros de comportamiento de las tecnologías de jerarquía de memoria

Nivel	de	lα	jerarquía			
de memoria						

Capacidad de almacenamiento típica

Tiempo de acceso típico

Registros

Registros (caros)

64 B ... 32 KB

0,2 ns ... 100 ns



Cache

8 KB ... 32 MB

5,0 ns ... 200 ns



Principal

128 MB ... 512 GB

60 ns ... 500 ns

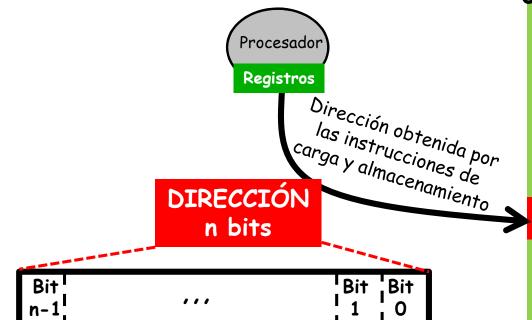


Secundaria (barata)

128GB ... 1TB

10ms...100ms

# Direccionamiento



\_ \_ \_ n = 32 bits

4 GB (gigabytes)

- - - - - - n = 64 bits

16 EB (exabytes)

Tema 2-2/6

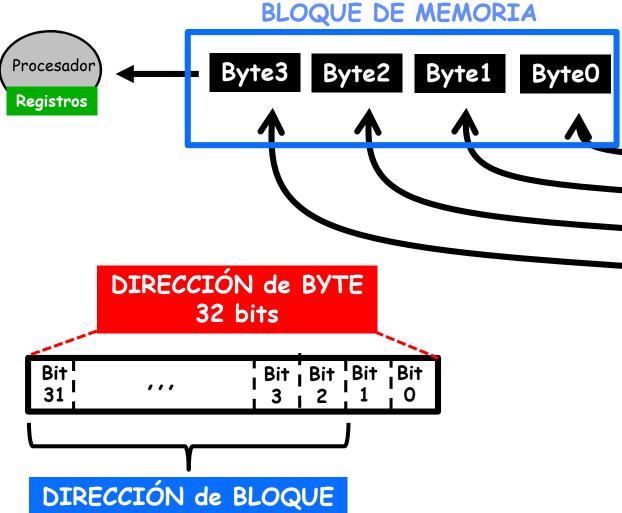
Espacio de direccionamiento del procesador mapeado en memoria

DIRECCIÓN

Byte

# Bloque de memoria

Espacio de direccionamiento del procesador mapeado en memoria



30 bits

DIRECCIÓN Byte0

DIRECCIÓN Byte1

DIRECCIÓN Byte2

DIRECCIÓN Byte3

# Big y Little Endian

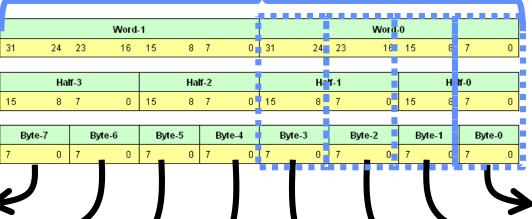
Espacio de direccionamiento del procesador

> BIG **ENDIAN**

Byte7 DIRECCIÓN Byte6 DIRECCIÓN Byte5 DIRECCIÓN Byte4

DIRECCIÓN

DATOS de un bloque de memoria: Palabra 1 Palabra O



Espacio de direccionamiento del procesador

**ENDIAN** 



#### **BIG ENDIAN:**

byte menos significativo almacenado en la dirección más ALTA

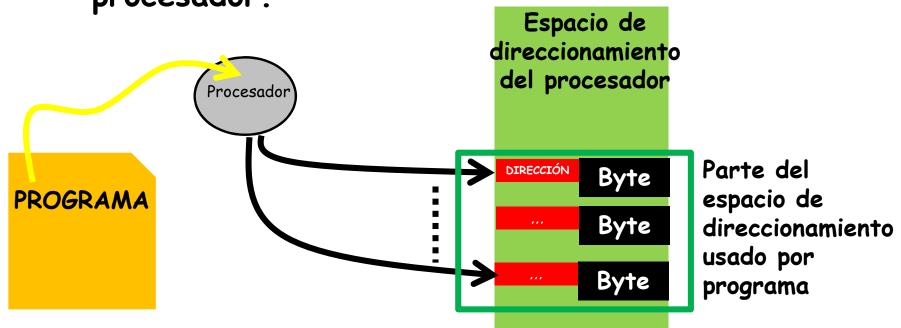
#### LITTLE ENDIAN:

byte menos significativo almacenado en la dirección más BAJA

#### Proximidad o Localidad Referencial

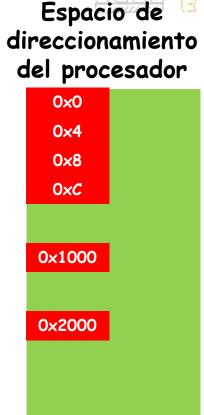


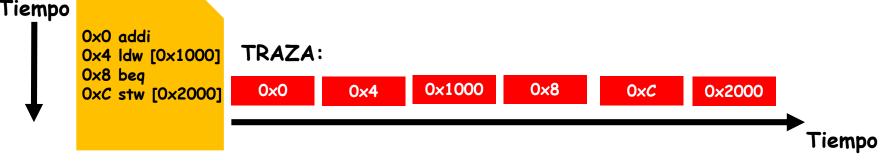
 Definición: característica de un determinado nivel de la jerarquía de memoria que establece que un programa en ejecución utiliza en cada momento sólo una pequeña parte del espacio de direccionamiento del procesador.



## Traza de un programa

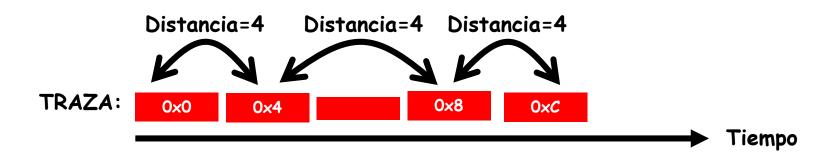
- Definición: lista ordenada en el tiempo de las direcciones de memoria que un programa utiliza.
- La traza ordenada temporalmente incluye tanto las direcciones asignadas a las instrucciones del programa como las direcciones asignadas a los datos del programa.





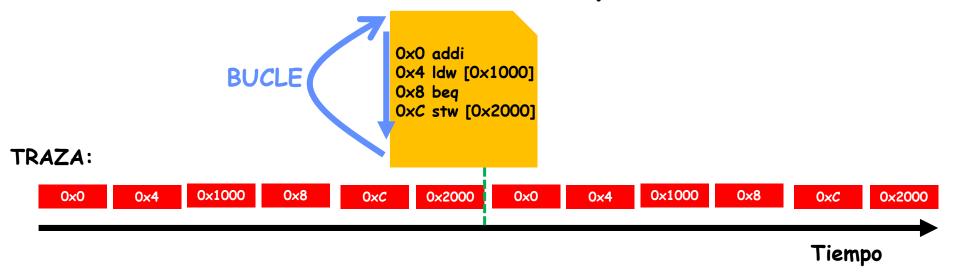
# Proximidad o localidad espacial

 <u>Definición</u>: Dadas dos direcciones del espacio de direccionamiento del procesador que han sido referenciadas recientemente en el tiempo, existe una alta probabilidad de que la distancia entre direcciones de memoria sea muy pequeña.

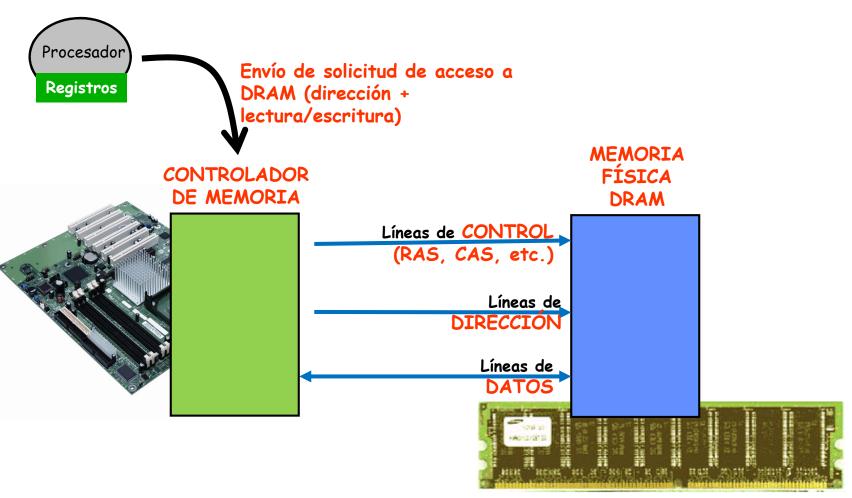


## Proximidad o localidad temporal

 <u>Definición</u>: Un programa en ejecución tiende a referenciar repetidamente las mismas direcciones del espacio de direccionamiento en un corto intervalo de tiempo.

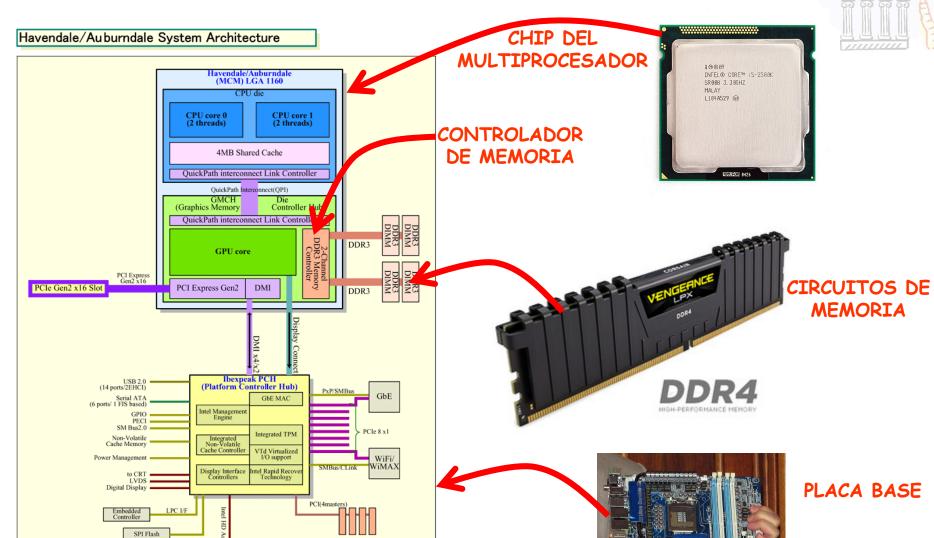


# Circuitos de la Memoria Principal

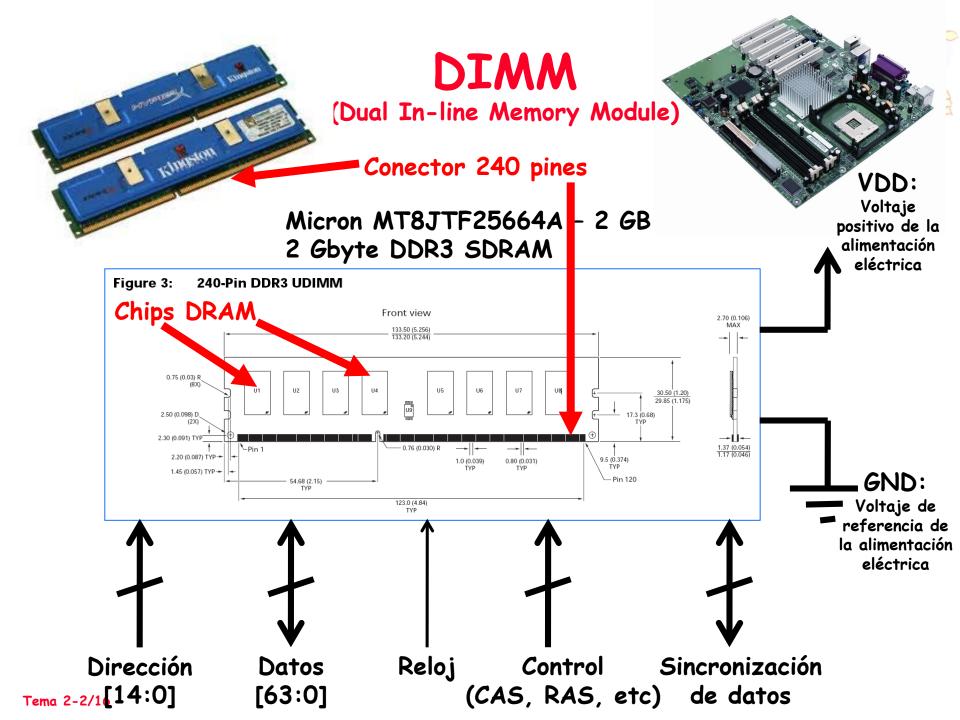


184 pin DDR SDRAM

#### Placas Base de ordenadores sobremesa



Copyright (c) 2007 Hiroshige Goto All rights reserved.





# Ejemplos de parámetros de prestaciones en memorias DRAM - DDR3

TIEMPOS DE SIN-CRONIZA-CIÓN: CAS  $(t_{CL})$ , RAS2CAS  $(t_{RCD})$ , PRECARGA  $(t_{RP})$ 

				RITMO		RITMO DE	$(t_{RCD}),$
NOMBRE TECNOLOGÍA	FRECUENCIA RELOJ	TIEMPO CICLO	FRECUEN- CIA E/S	TRANS- FEREN- CIA E/S	NOMBRE MÓDULO	TRANSFERENCIA de DATOS	PRECARGA (t <sub>RP</sub> )
DDR3-800	100 MHz	10 ns	400 MHz	800 MT/s	PC3-6400	6400 MB/s	5-5-5 6-6-6
DDR3-1066	133 MHz	7.5 ns	533 MHz	1066 MT/s	PC3-8500	8533 MB/s	6-6-6 7-7-7 8-8-8
DDR3-1333	166 MHz	6 ns	667 MHz	1333 MT/s	PC3-10600	10667 MB/s	7-7-7 8-8-8 9-9-9 10-10-10
DDR3-1600	200 MHz	5 ns	800 MHz	1600 MT/s	PC3-12800	12800 MB/s	8-8-8 9-9-9 10-10-10 11-11-11

Latencia DDR (CICLOS): t<sub>RP</sub> +

(Frecuencia E/S)

Precarga

Tema

t<sub>RCD</sub>

↓

Acceso Fila
(RAS2CAS)

t<sub>CL</sub>
↓
Acceso Columna
(CAS, latencia

# Ejemplos de parámetros de prestaciones en memorias DRAM-DDR4

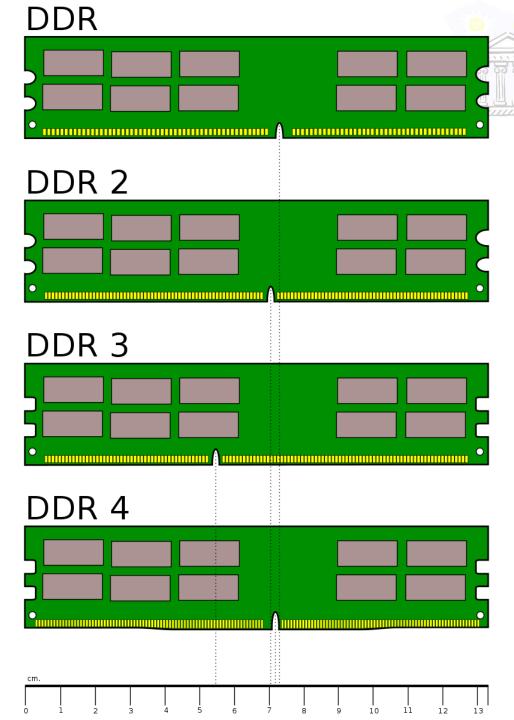


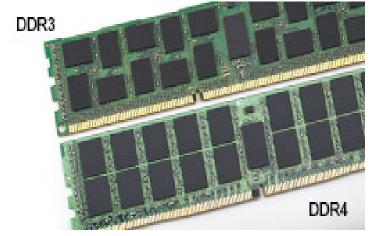
Módulos de memoria DDR4 8 GB

284 pines

	Standard name	Memory clock (MHz)	I/O bus clock (MHz)	Data rate (MT/s)	Module name	Peak trans- fer rate (MB/s)	Timings CL-tRCD-tRP	CAS latency (ns)
	DDR4-1600J* DDR4-1600K DDR4-1600L	200	800	1600	PC4-12800	12800	10-10-10 11-11-11 12-12-12	12.5 13.75 15
	DDR4-1866L* DDR4-1866M DDR4-1866N	233.33	933.33	1866.67	PC4-14900	14933.33	12-12-12 13-13-13 14-14-14	12.857 13.929 15
	DDR4-2133N* DDR4-2133P DDR4-2133R	266.67	1066.67	2133.33	PC4-17000	17066.67	14-14-14 15-15-15 16-16-16	13.125 14.063 15
	DDR4-2400P* DDR4-2400R DDR4-2400T DDR4-2400U	300	1200	2400	PC4-19200	19200	15-15-15 16-16-16 17-17-17 18-18-18	12.5 13.32 14.16 15
_	DDR4-2666T DDR4-2666U DDR4-2666V DDR4-2666W	333.33	1333.33	2666.67	PC4-21300	21333.33	17-17-17 18-18-18 19-19-19 20-20-20	12.75 13.50 14.25 15
	DDR4-2933V DDR4-2933W DDR4-2933Y DDR4-2933AA	366.67	1466.67	2933.33	PC4-23466	23466.67	19-19-19 20-20-20 21-21-21 22-22-22	12.96 13.64 14.32 15
	DDR4-3200W DDR4-3200AA DDR4-3200AC	400	1600	3200	PC4-25600	25600	20-20-20 22-22-22 24-24-24	12.5 13.75 15

# Incompatibilidad de las memorias DRAM-DDR







#### Modelos DDR4



x8 Mbps

**DDR4 SDRAM DIMMs** 

Chip	Module	Memory Clock	I/O Bus Clock	Transfer rate	Voltage
DDR4-1600	PC4-12800	200 MHz	800 MHz	1600 MT/s	1.2 V
DDR4-1866	PC4-14900	233 MHz	933 MHz	1866 MT/s	1.2 V
DDR4-2133	PC4-17000	266 MHz	1066 MHz	2133 MT/s	1.2 V
DDR4-2400	PC4-19200	300 MHz	1200 MHz	2400 MT/s	1.2 V
DDR4-2666	PC4-21300	333 MHz	1333 MHz	2666 MT/s	1.2 V
DDR4-3200	PC4-25600	400 MHz	1600 MHz	3200 MT/s	1.2 V

MBytes/segundo

**x4** 

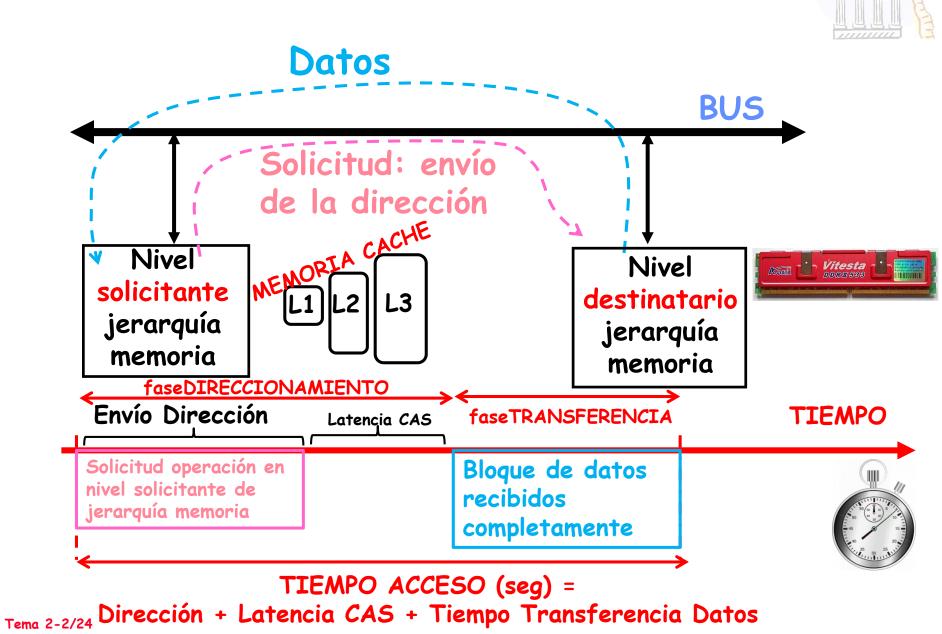
**x2** 





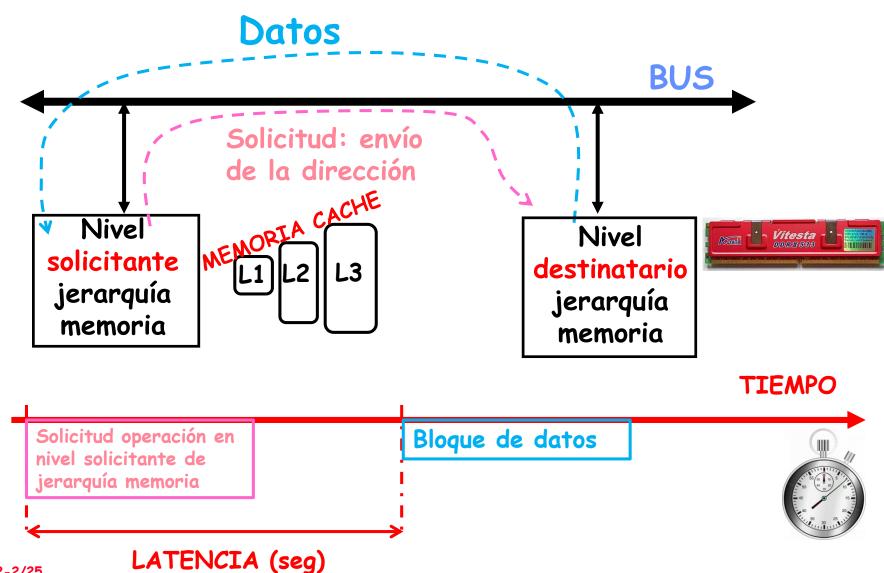
- Tiempo de acceso
- Latencia
- -Ritmo de Transferencia
- Tiempo de Ciclo

#### Prestaciones: TIEMPO DE ACCESO



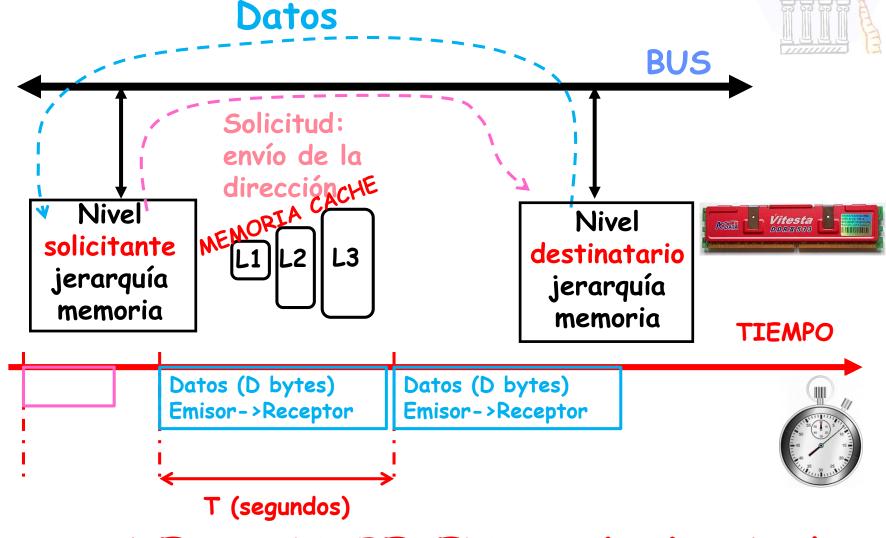
#### Prestaciones: LATENCIA





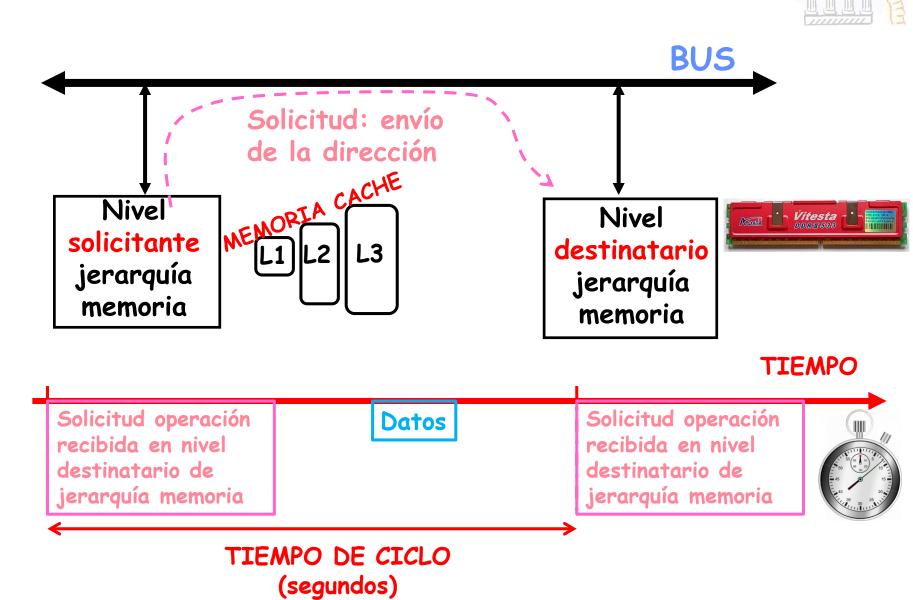
Tema 2-2/25

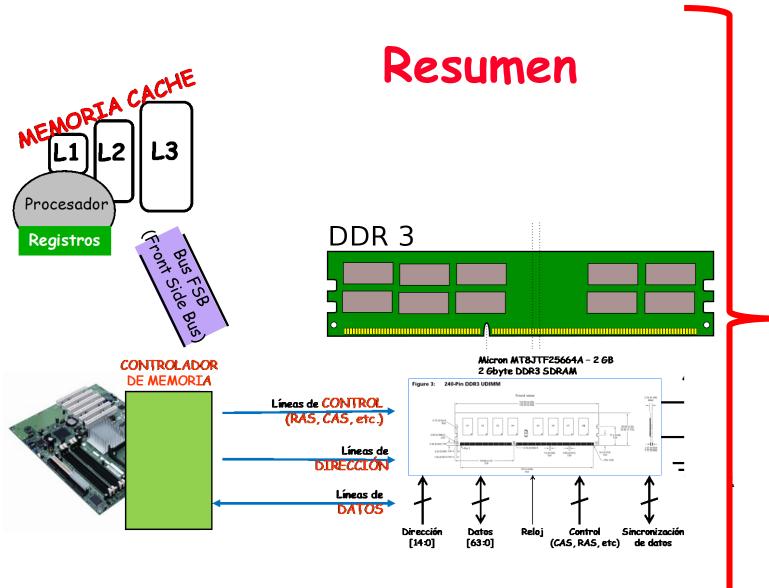
## Prestaciones: RITMO TRANSFERENCIA

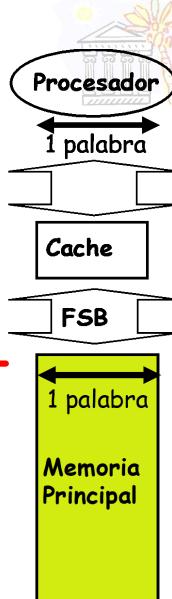


RITMO DE TRANSFERENCIA (RT) = D / T Tiempo de transferencia datos = D / RT

# Prestaciones: TIEMPO DE CICLO



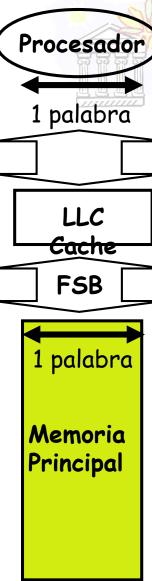




# Organización Simple

- Calcular <u>Penalización Fallo Cache LLC (Last Level Cache)</u> para un bloque de cache de 4 palabras (16 bytes).
- · Suponemos:
  - DIRECCION: 1 ciclo para enviar la dirección de 1 palabra.
  - OBTENCIÓN DATOS EN CIRCUITO DE MEMORIA: 15 ciclos para cada acceso a una posición de la DRAM.
  - ENVÍO DATOS: 1 ciclo para enviar una palabra de datos.
- Penalización de fallo = 4x1
   (direccionamiento) + 4x15 (obtención datos)
   + 4x1 (envío) = 68 ciclos
- · Cache Memoria Principal:

Ritmo Transferencia = 16 B / 68 clk =



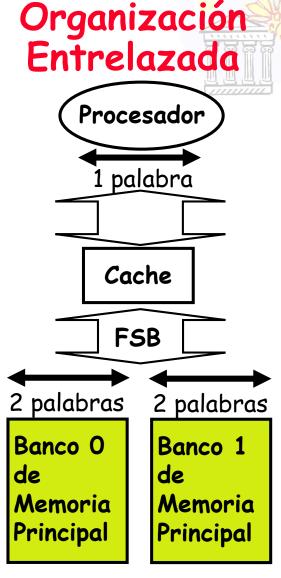
BANCO UNICO

CONEXIÓN SIMPLE

Tema 02/24 bytes/ciclo

- Banco: parte en la que se divide la memoria y que almacena direcciones entrelazadas
- Una única dirección para varios módulos de memoria
- · Direccionamiento:
  - Num Banco = direccion MOD num total bancos
  - Dir dentro banco= parte entera (direccion / num total bancos)
- · Aumenta el ritmo transferencia ensanchando la memoria pero no el bus
- Se sigue pagando el coste de transmitir cada palabra secuencialmente, pero se evita pagar más de una vez la latencia de acceso.
- Con 2 bancos de memoria, 2 palabras/banco Penalización fallo =  $1 + 1 \times 15 + 4 \times 1 = 20$ ciclos

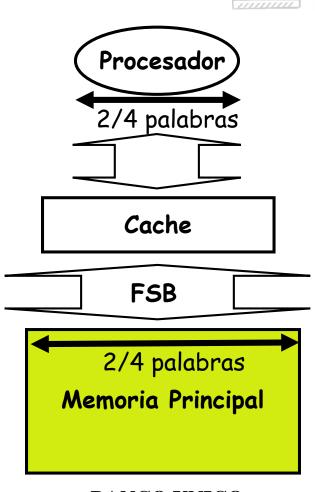
Ritmo Transferencia = 16 bytes / 20 ciclos = 0.80 bytes/ciclo (aumenta 233% vs. organización simple:  $233\% = 100 \times [0.8-0.24] / 0.24$ )



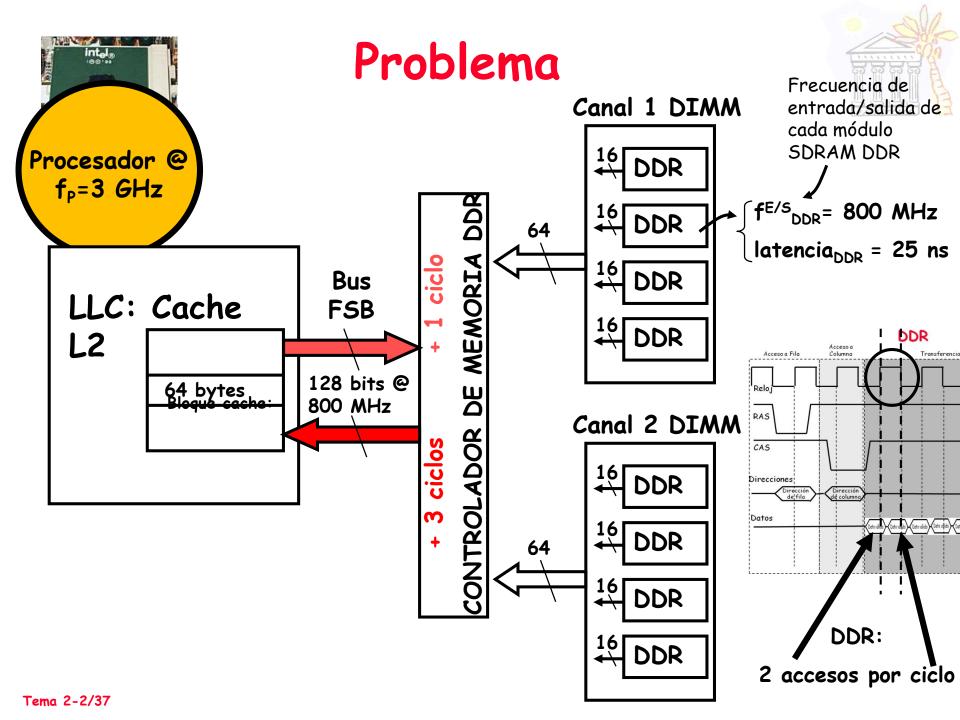
MULTIPLES BANCOS CONEXIÓN SIMPLE

## Organización Ancha

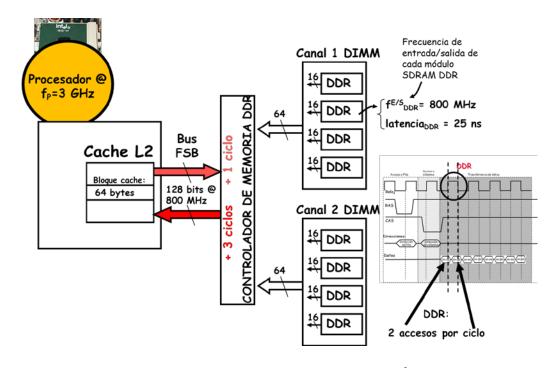
- Se aumenta el ancho de banda haciendo la memoria y el bus de varias palabras de ancho.
- Aumentar el ancho de la memoria y del bus disminuye el tiempo de transferencia
   A disminuye la penalización de fallo
- Si ancho de memoria es de 2 palabras
   Penalización de fallo = 2 + 2×15 + 2×1 = 34 ciclos
   Ritmo Transferencia = 16/34 = 0,47 bytes/ciclo
- Si ancho de memoria = 4 palabras
   Penalización de fallo = 1 + 1x15 + 1x1 = 17 ciclos
   Ritmo Transferencia = 16 bytes / 17 ciclos = 0.94 bytes/ciclo (aumenta 291% vs. Organización simple, 291% = 100 x [0.94-0.24] / 0.24)



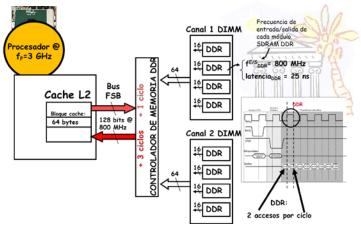
BANCO UNICO CONEXIÓN MULTIPLE







Se desea analizar la penalización del CPI de un programa debido a los accesos a la memoria principal utilizando un sistema computador como el que se muestra a continuación con dos canales independientes de acceso a memoria principal (Canales 1 y 2).

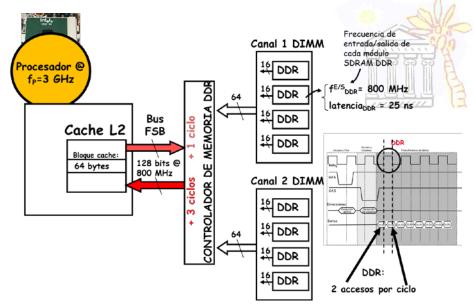


Por simulación se sabe que si el programa se ejecuta sin fallos de la cache L2, el CPI sería: CPIideal=0.3 ciclos del procesador ( $f_p$ = 3 GHz); y la frecuencia de fallos de la cache L2 por instrucción retirada es de FF=5%. Por otro lado, considerar que:

- Los accesos a la memoria principal se originan por fallos de cache L2 cuyo tamaño de bloque es de 64 bytes.
- El bus FSB tiene un ancho<sub>FSB</sub>= 128 bits (16 bytes), y una frecuencia de  $f_{FSB}$ = 800 MHz (periodo:  $T_{FSB}$  = 1.25 ns).
- El envío de la dirección del bloque de cache L2 que falla al controlador de memoria DDR tarda 4 ciclos de bus FSB (envío dirección desde L2 a controlador memoria:  $4 \times T_{FSB} = 5$  ns).
- El envío de la dirección del bloque de L2 desde el controlador a los módulos DIMM de memoria principal tarda 1 ciclo del bus FSB. Suponer que el bloque de L2 está distribuido homogéneamente de forma entrelazada sobre todos los chips DDR, y que el controlador de memoria sólo puede resolver un fallo de L2 en cada instante; es decir, que no acepta un nuevo acceso a memoria principal por fallo de L2 hasta que no haya resuelto el anterior fallo.
- La latencia de cada uno de los chips DDR que forman los módulos DIMM es de latencia<sub>DDR</sub>=25 ns (tiempo que transcurre desde que un chip DDR recibe la dirección hasta que proporciona -lectura- o almacena -escritura- los 16 bit (2 bytes) a través de su puerto de datos), y su frecuencia de entrada/salida es de f<sup>E/S</sup><sub>DDR</sub>= 800 MHz.
- El retraso de enviar los datos desde el controlador de memoria a la cache L2 es de 3 ciclos del bus FSB.

#### Calcular:

- (a) El ritmo de transferencia máximo BW<sub>DDR</sub> [Bytes/segundo] que exhibe cada chip de memoria DDR.
- (b) El ritmo de transferencia en cada acceso a memoria que es exhibido por el bus FSB: BW<sub>FSB</sub>.
- (c) Los ciclos de penalización por cada fallo de cache L2.
- (d) El CPI real que aparecería cuando se consideran los fallos de la cache L2.
- (e) Porcentaje del CPI debido a la penalización por transferencia del bloque de cache por el bus FSB.
- (f) Porcentaje del CPI debido a la penalización por la latencia de los módulos DDR.



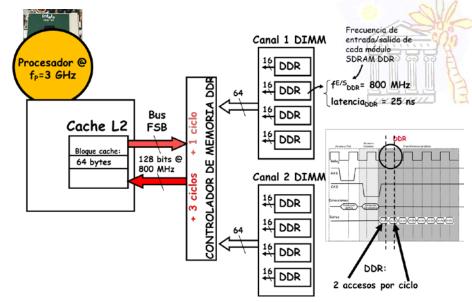
#### Calcular:

(a) El ritmo de transferencia máximo  $BW_{DDR}$  [Bytes/segundo] que exhibe cada chip (módulo) de memoria DDR.

Frecuencia de entrada/salida de cada módulo DRAM DDR

 $BW_{DDR} = f_{DDR}^{E/S} \times Bytes = 800MHz \times 2 \ Bytes \times 2 \ envios_{DDR} / ciclo_{DDR} = 3200MB / s = 3.2GB / s$ 

Ritmo de transferencia de cada módulo SDRAM DDR : 1600 MT/s (800 x 2) Ritmo de transferencia de cada módulo SDRAM DDR en bytes/segundo



#### Calcular:

(b) El ritmo de transferencia en cada acceso a memoria que es exhibido por el bus FSB:  $BW_{FSB}$ .

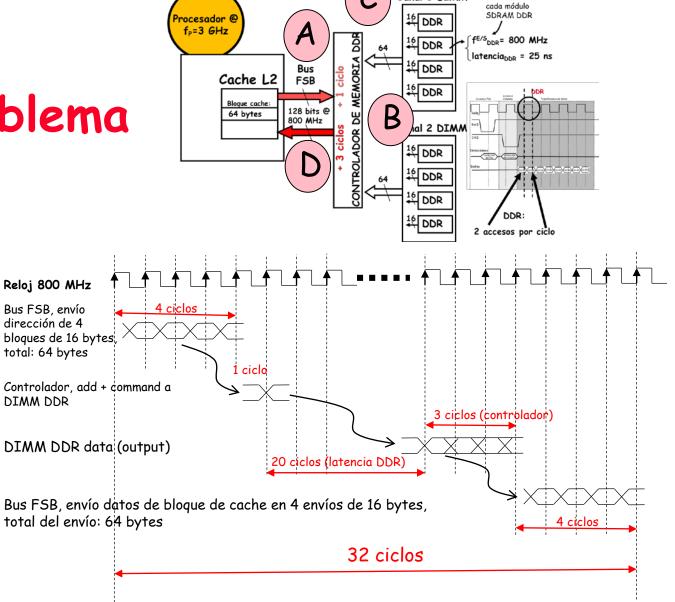
$$BW_{FSB} = f_{FSB} \times Bytes = 800MHz \times 16 \; Bytes = 12800MB/s = 12.8GB/s$$
16 bytes = 128 bits / 8



Reloj 800 MHz Bus FSB, envío

dirección de 4 bloques de 16 bytes total: 64 bytes

DIMM DDR

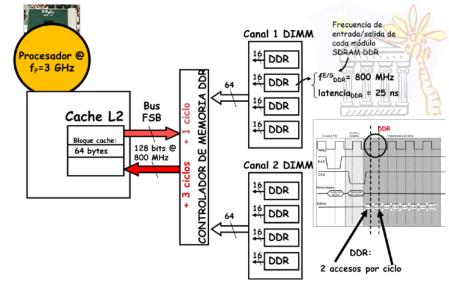


Frecuencia de

entrada/salida de

Canal 1 DIMM

B

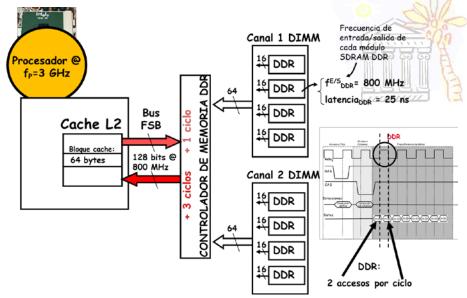


#### Calcular:

(c) Los ciclos de penalización por cada fallo de cache L2.

El ritmo de transferencia de todos los módulos DDR en paralelo (3.2GB/s × 8DDR = 25.6GB/s) es superior al del FSB (12.8GB/s), por lo que <u>el número de ciclos que se tarda en transferir un bloque de L2 viene determinado por el bus FSB</u>.

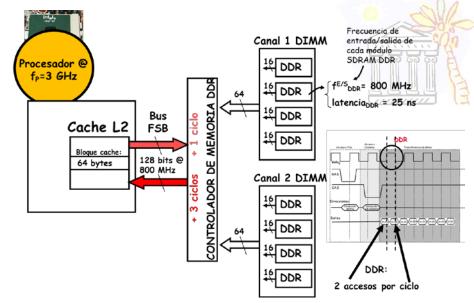
A	Envío de la dirección desde L2 al controlador de memoria	4 ciclos (64 bytes / 16 bytes)	
В	Envío de la dirección desde el controlador de memoria al módulo DIMM	1 ciclo $\frac{25ns}{1} = 20ciclos$	· @ 800 <i>MHz</i>
C	Latencia de los circuitos de memoria DDR	$\frac{1}{800MHz}$	
D	Envío del bloque desde DDR a L2 (pipeline con el controlador de memoria)	$\frac{64bytes/L2}{16bytes/ciclo_{FSB}} = 4ciclos @ 800MHz$	
Е	Retraso en enviar los primeros datos desde el controlador a la L2	3 ciclos	
	PENALIZACIÓN <sub>L2</sub> : P <sub>L2</sub> = A+B+C+D+E =	32 ciclos@800MHz = 120 ciclos @ 3 <i>G</i> Hz	



#### Calcular:

(d) El CPI real que aparecería cuando se consideran los fallos de la cache L2.

$$CPI = CPI_{ideal} + FF \times P_{L2} = \frac{1}{3} + 0.05 \times 120 = 6.333$$
 ciclos (del procesador, 3 GHz)

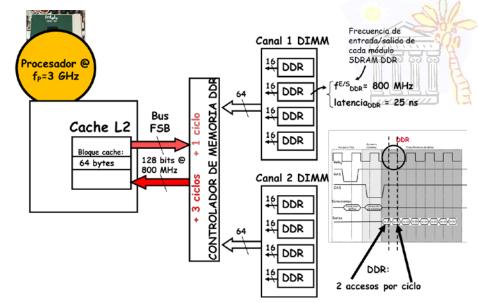


#### Calcular:

(e) Porcentaje del CPI debido a la penalización por transferencia del bloque de cache por el bus FSB.

De los 32 ciclos (@800 MHz) de la penalización de un fallo de L2, la originada por la transferencia del bloque de cache L2 es (ver apartado c): 4 ciclos en la frecuencia FSB o 4×(120/32) en ciclos del procesador (@3 GHz).

$$\frac{CPI_{FSB}}{CPI} = \frac{FF \times P_{L2}|_{FSB}}{CPI} = \frac{0.05 \times 4 \times \frac{120}{32}}{6.333} = 11.8\%$$



#### Calcular:

(f) Porcentaje del CPI debido a la penalización por la latencia de los módulos DDR.

La penalización originada por la latencia del módulo DDR es (ver apartado c): 20 ciclos en la frecuencia FSB o 20×120/32 en ciclos del procesador.

$$\frac{CPI_{LatenciaDDR}}{CPI} = \frac{FF \times P_{L2}|_{LatenciaDDR}}{CPI} = \frac{0.05 \times 20 \times \frac{120}{32}}{6.333} = 59.2\%$$