



UNIVERSIDAD  
POLITECNICA  
DE VALENCIA

# Configuración y Optimización de Sistemas de Cómputo

Master Universitario en Ingeniería Informática  
Depto. de Informática de Sistemas y Computadores (DISCA)  
Universidad Politécnica de Valencia

---

# Memoria

---

- Estructura memoria principal
- Evolución tecnológica
- Arquitectura
- Impacto en prestaciones
- Nuevas Tecnologías

# Memoria

---

- El objetivo de nuestro sistema informático es albergar la mayor cantidad de información posible en memoria (acceso rápido)
  - Accediendo a la información de forma rápida
  - Con un consumo/coste asumible
- Diferentes circuitos para implementar memorias
  - Proporcionan diferentes ventajas/inconvenientes
  - No existe la solución ideal (a día de hoy)

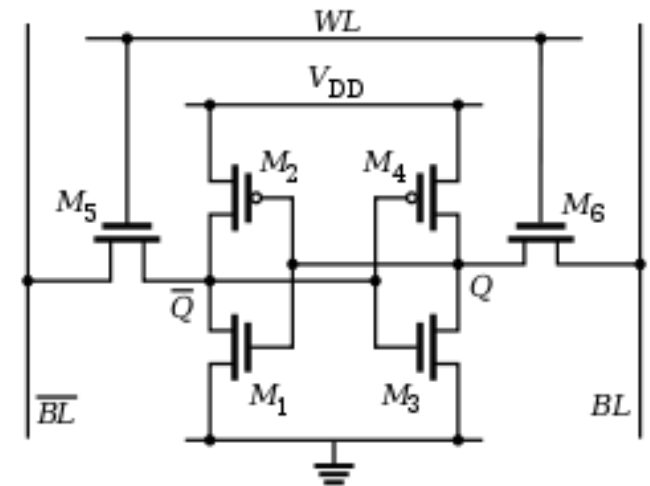
# Memoria SRAM

---

- Static Random Access Memory
  - Permite acceder a un dato partiendo de una dirección “Aleatoria”
  - Puede ser volátil o no volátil
    - Generalmente son volátiles
      - Es decir requieren ser alimentadas para mantener la información
    - NVRAM es el término que se suele utilizar para referirse a la SRAM no volátil

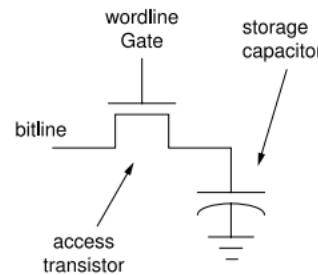
# Memoria SRAM

- Tecnología que se utiliza para implementar memoria en el chip
  - Tiempos de acceso rápidos
  - Elevado coste (área/bit)
  - El consumo es “reducido”
- La celda típica es la 6T
  - N-Transistores por bit
  - Más transistores aumentan fiabilidad

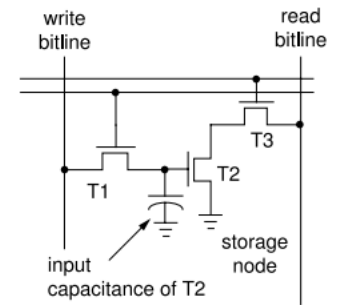


# Memoria DRAM

- Dynamic Random Access Memory
  - Permite una mayor densidad (vs SRAM)
    - Menor número de transistores por bit
  - Son más lentas (vs SRAM)
    - Aunque proporcionan buenos anchos de banda
  - Tienen elevado consumo (vs SRAM)
    - Hay que alimentar y refrescar



(a) 1T1C



(b) 3T1C

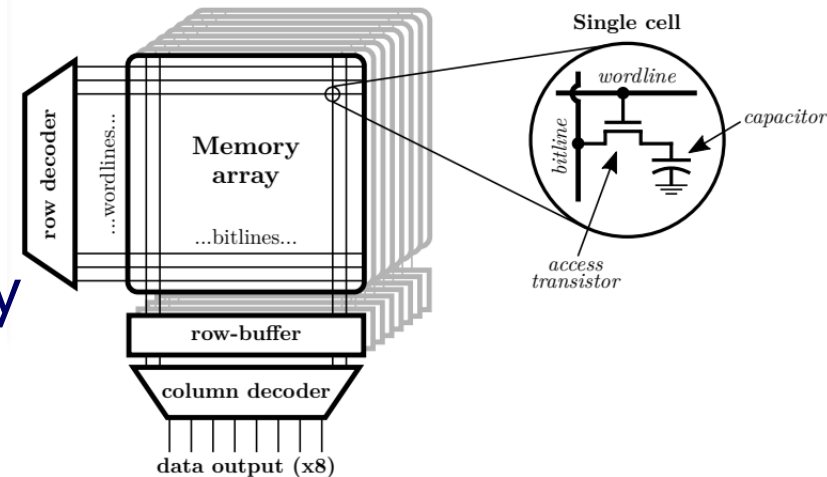
# Memoria DRAM

---

- Refresco
  - Los bits en una DRAM se almacenan en una condensador (capacidad)
  - La carga almacenada en el condensador se va vaciando debido a la corriente de fugas
  - Es necesario refrescar (leer y volver a escribir) los valores de forma periódica
    - Periodo de refresco (refresh period)

# Memoria DRAM

- Chips específicos de DRAM
  - Bajo coste por bit
  - Alto nivel de integración
- Celdas de bit → memory array
  - Leemos/escribimos filas
  - Varios arrays en paralelo
- Row buffer
  - Las lecturas son destructivas
  - Almacena los bits de una fila (página) y sirve las lecturas.
  - Vuelve a recargar el contenido de la fila





# DDR RAM

- DDR (Double Data Rate) RAM
  - La familia de las SDRAM más utilizada
  - Double data rate
    - Transmitimos datos en los flancos de subida y de bajada

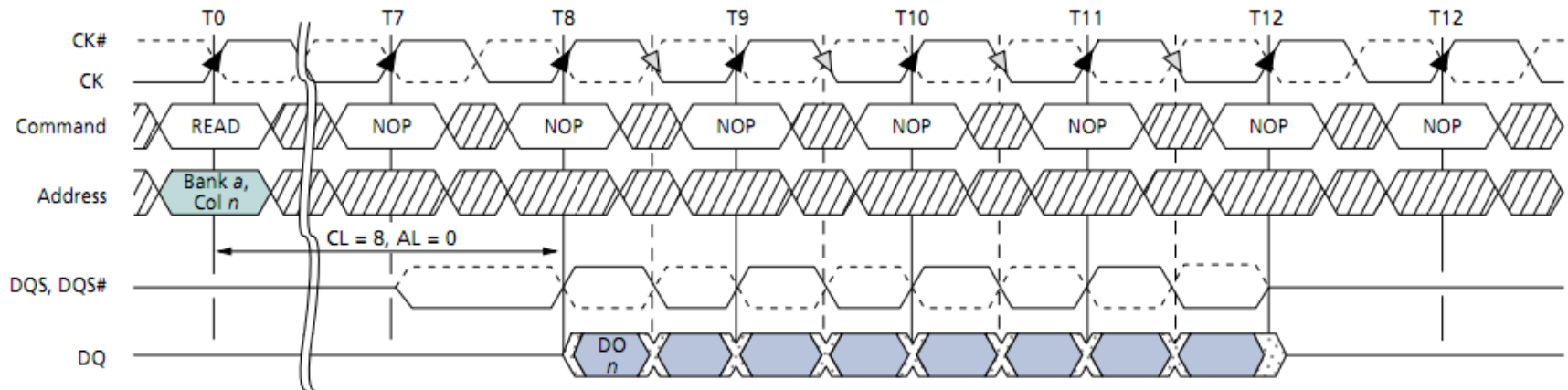
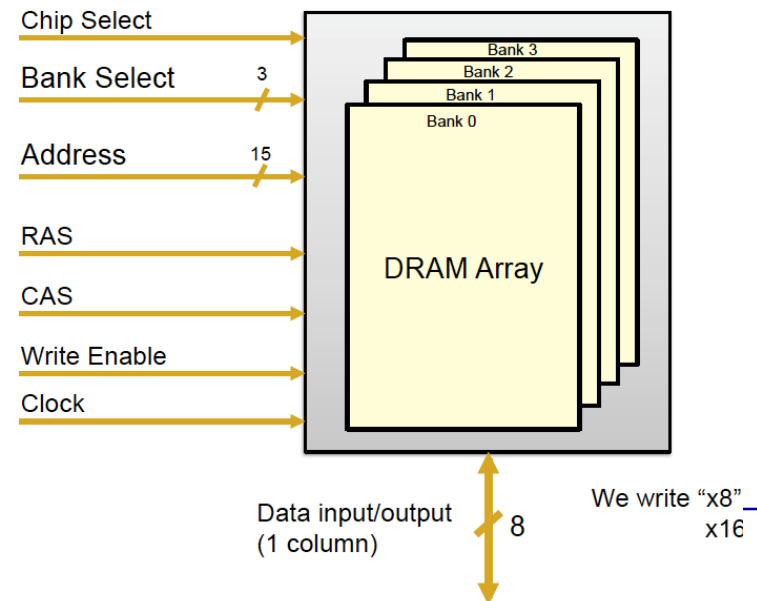
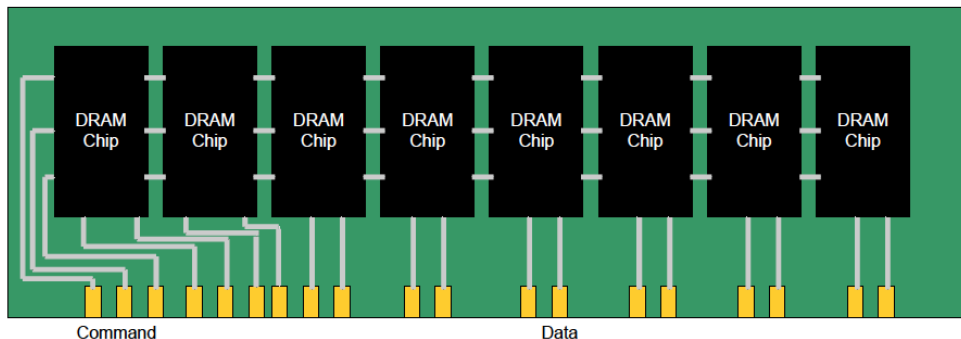


Diagrama SDRAM → en la DDR escribimos en los 2 flancos

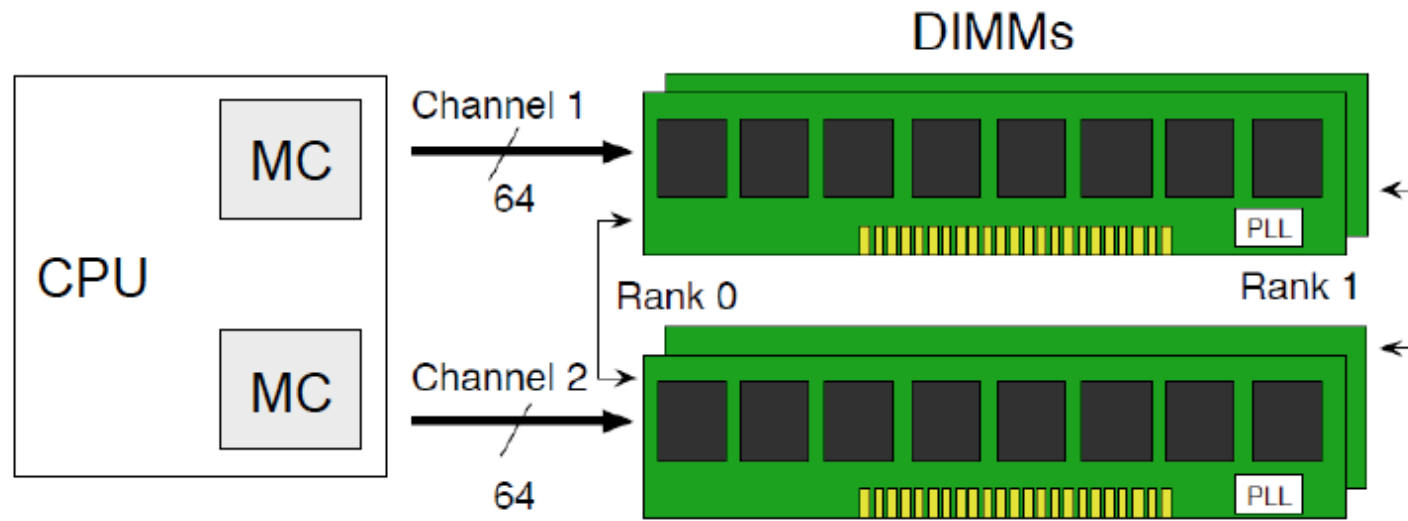
# DDR Arquitectura

- DIMMs
  - Compuestos por varios chips.
  - Ranks. Grupos de chips de 64 bit de ancho
    - El objetivo es maximizar el uso/ancho del bus
  - Banks. Cada chip tiene internamente varios *Banks*
    - Mejora el ancho de banda, permite el entrelazado de peticiones (interleaving)



# CPU → Memoria

- La CPU incluye uno o varios controladores de Memoria (MC)
  - Los MC envían comandos a la DDR
    - Según el tipo de operación
      - Lectura, Escritura, Refresco

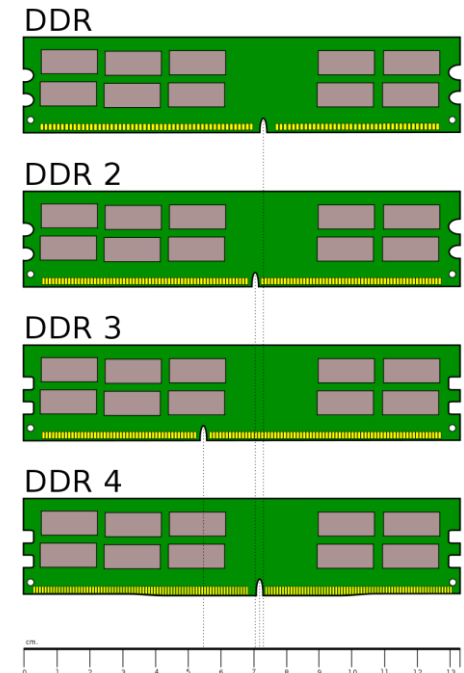


# DDR Standards

- JEDEC a estandarizado diferentes configuraciones de DDR
  - El objetivo es permitir interoperabilidad entre componentes

	DDR	DDR2	DDR3	DDR4	DDR5
Frec (MHz)	200	400	800	1600	3600*
Velocidad (MBps)	3200	6400	12800	25600	57600
Tensión (V)	2,5	1,8	1,5/1,35	1,2/1,05	1,1
Bits internos	2	4	8	8	8
Pines	184	240	240	288	288

\*valor más alto del rango



# DDR

- Prestanciones

- Latencia:

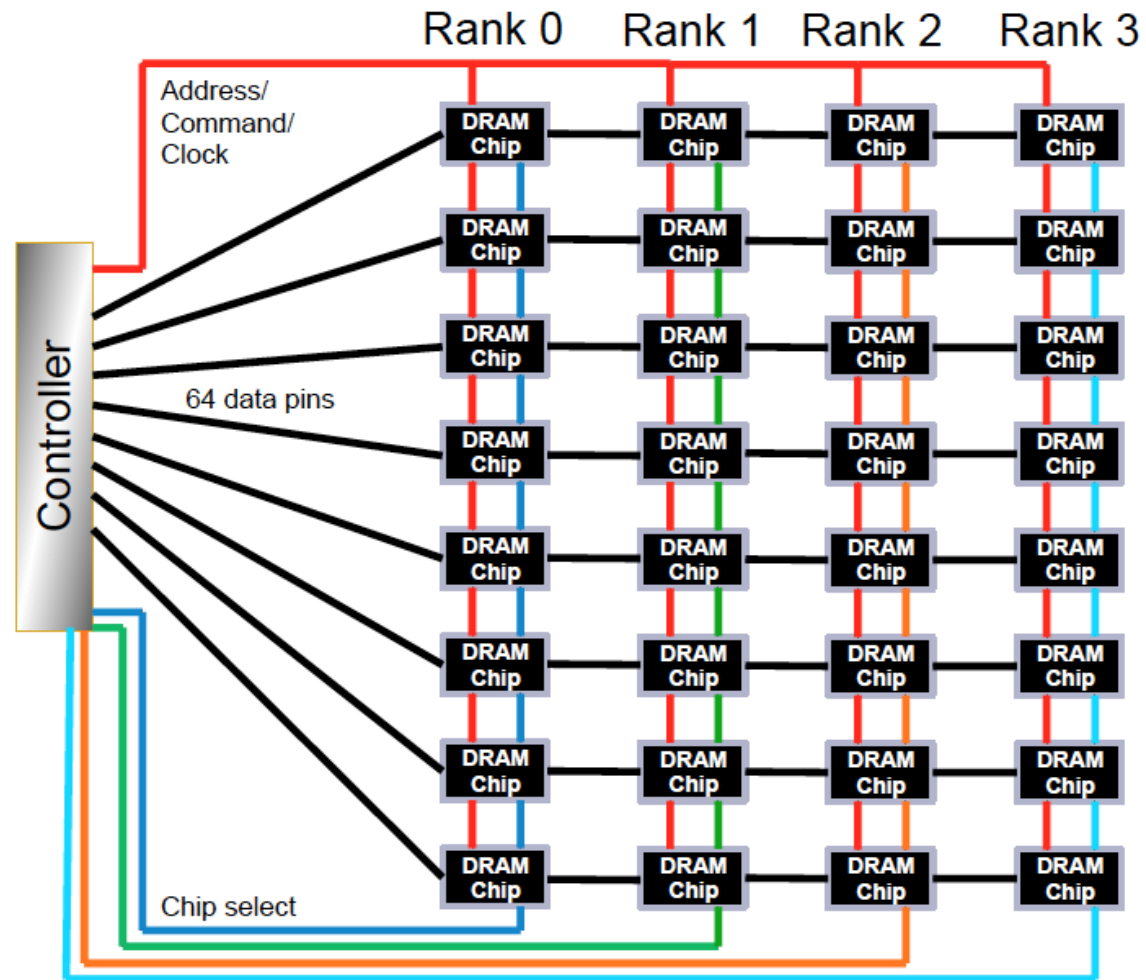
- De 20 a 40 ns

- Ancho de banda:

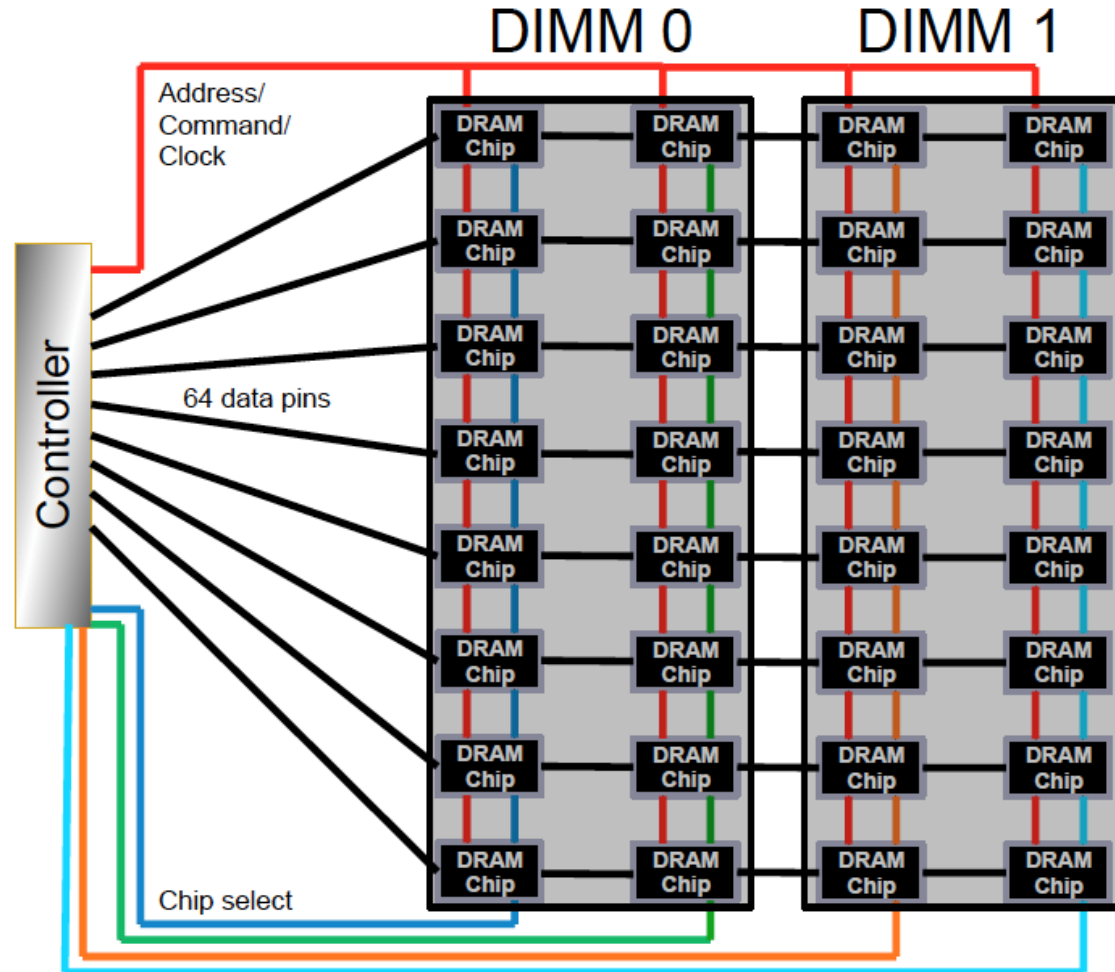
- Ancho de banda Pico (GB/s) = frecuencia bus x bytes/acceso x accesos/ciclo
    - Ejemplo: DDR3-1600 =  $800 \times 8 \times 2 = 12,8$  GB/s

\*valor más alto del rango

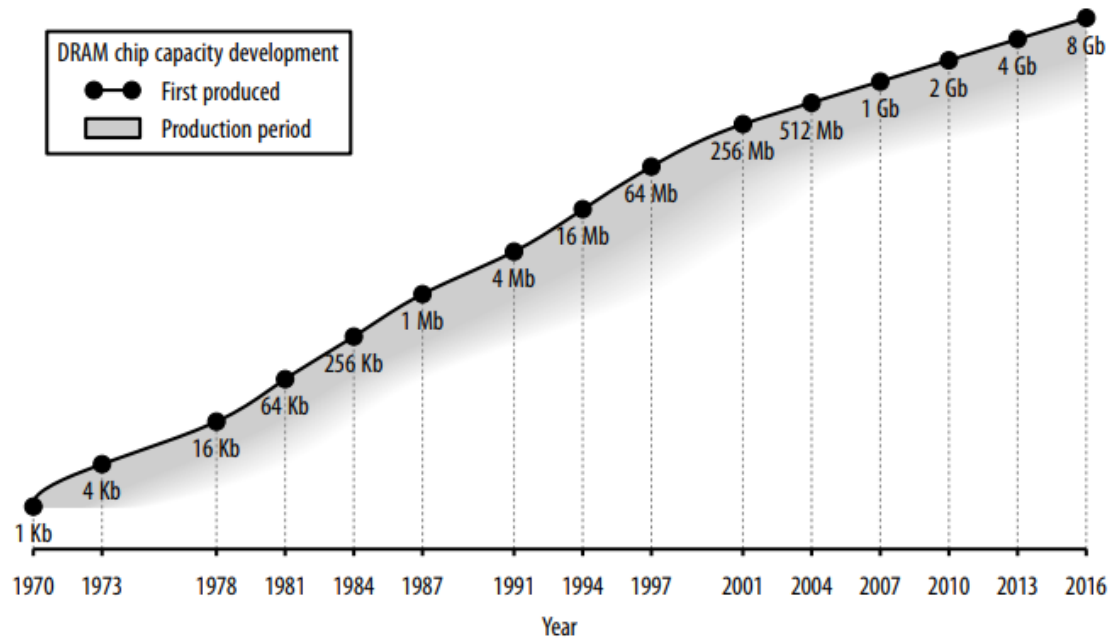
# DDR Ranks



# DDR DIMMs

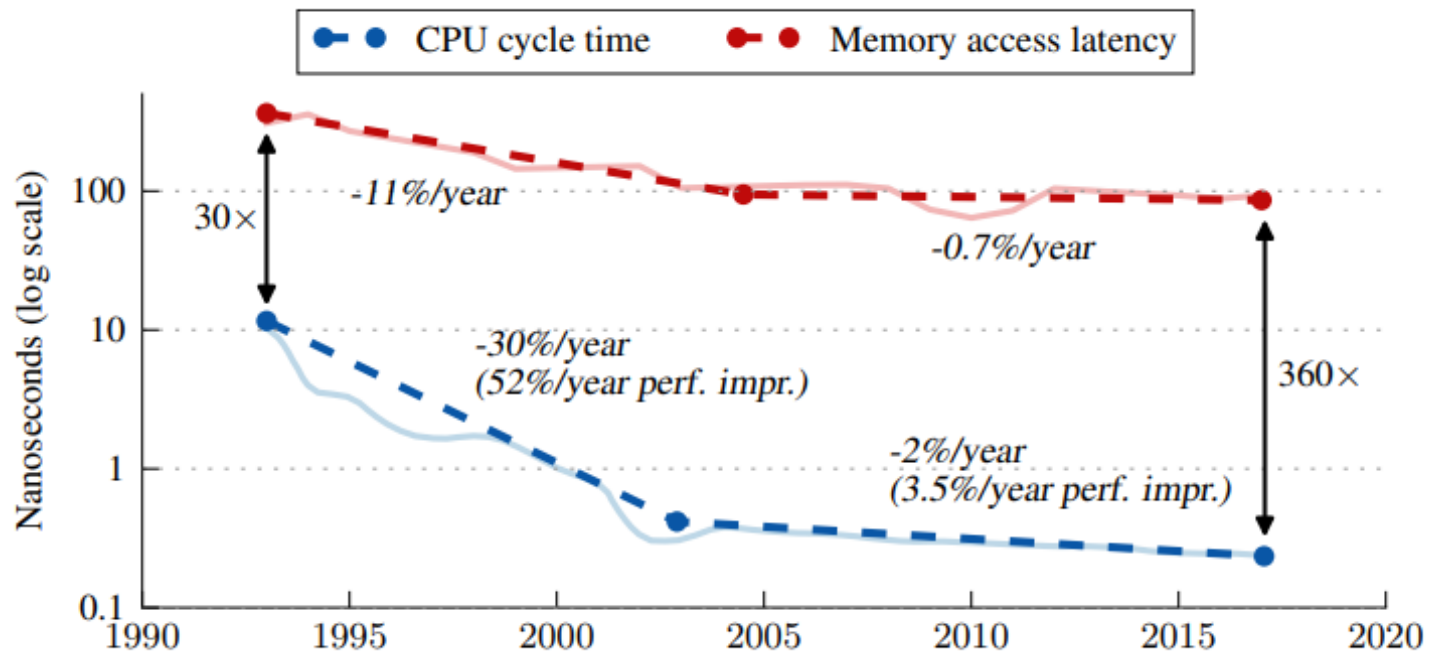


# Evolución Capacidad Memoria



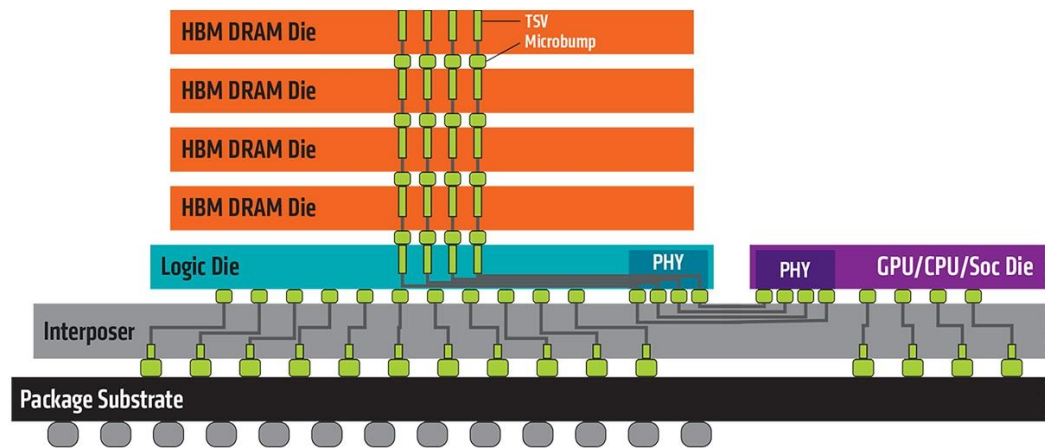


# Memory Wall



# High-Bandwidth Memory (HBM)

- HBM es una memoria “RAM” que apila de forma vertical chips de memoria
- Se conecta a través de un interposer
- HBM stacks no estan integrados en la CPU pero si muy cerca

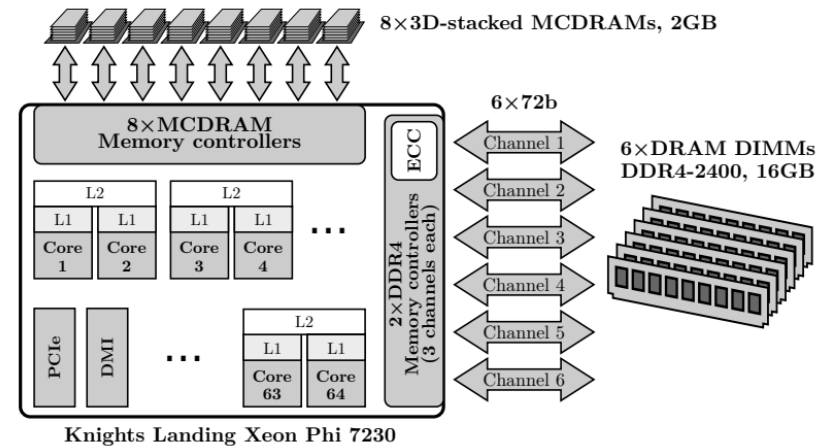


<https://www.amd.com/en/technologies/hbm>

# HBM - MCDRAM

- **Multi-Channel DRAM**

- Basado en Hybrid Memory Cube
- Incluida en el Xeon Phi Knights Landing
- 3D-stacked DRAM
- Mayor BW que DDR
- Latencia similar a DDR



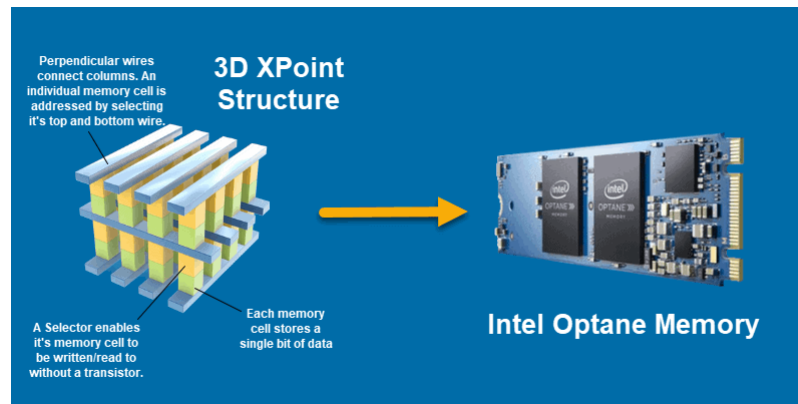
# Memoria No Volatil

---

- La mejora en tecnologías NVM permite el diseño de sistemas que integren NVRAM
  - La memoria no volátil es generalmente lenta
  - Muy bajo consumo
  - Gran capacidad de almacenamiento
- Diferentes tecnologías
  - ReRAM
  - PCM
  - ... no hay una solución ideal (todavía)

# Intel Optane

- Memoria de tipo NVRAM que utiliza tecnología 3D XPoint desarrollada por Intel y Micron
  - Las celdas más pequeñas permiten densidades cuatro veces mayor que la de DRAM
  - Se venden en formato NVDIMM para su integración en placas base



# Intel Optane

	DRAM	Intel Optane	Flash Memory (SSD)
<b>Speed</b>	Very Fast	Slower than DRAM, but much faster than flash memory	Slower than both DRAM and Intel Optane
<b>Cost</b>	Expensive	Costs less than DRAM but more than flash memory	Affordable
<b>Volatile / Non-Volatile</b>	Volatile	Non-Volatile	Non-Volatile
<b>Latency</b>	Low	Low	High
<b>Reliability</b>	High	Excellent read response times compared to flash-based drives	Low
<b>Endurance</b>	High	High	Low

# Elegir la Memoria

---

- Parametros a tener en cuenta
  - Ancho de Banda
  - Latencia
  - Capacidad
- Aplicaciones
  - Servidor de maquinas virtuales
    - Memoria vCPU = Memoria total / Número Threads
  - Tipo de carga
    - Sensible a latencia
    - Gran uso de memoria
    - Ancho de banda

# Referencias

---

- Radulović, M. Memory bandwidth and latency in HPC: system requirements and performance impact. Tesi doctoral, UPC, Departament d'Arquitectura de Computadors, 2019. Disponible en: <http://hdl.handle.net/2117/134613>
- Živanović, D. Memory systems for high-performance computing: the capacity and reliability implications. Tesi doctoral, UPC, Departament d'Arquitectura de Computadors, 2018. Disponible en: <http://hdl.handle.net/2117/121194>