# 了解内存

核心系统数据库组　余锋

http://yufeng.info
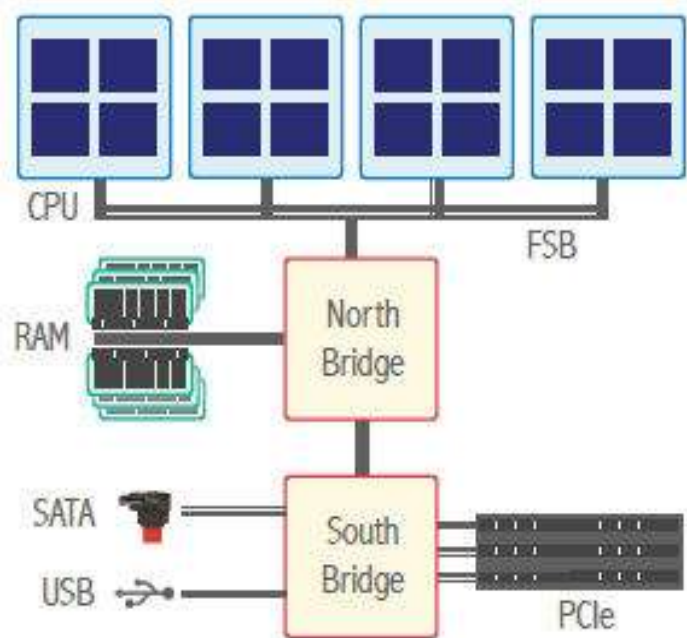
@淘宝褚霸

2012-03-17
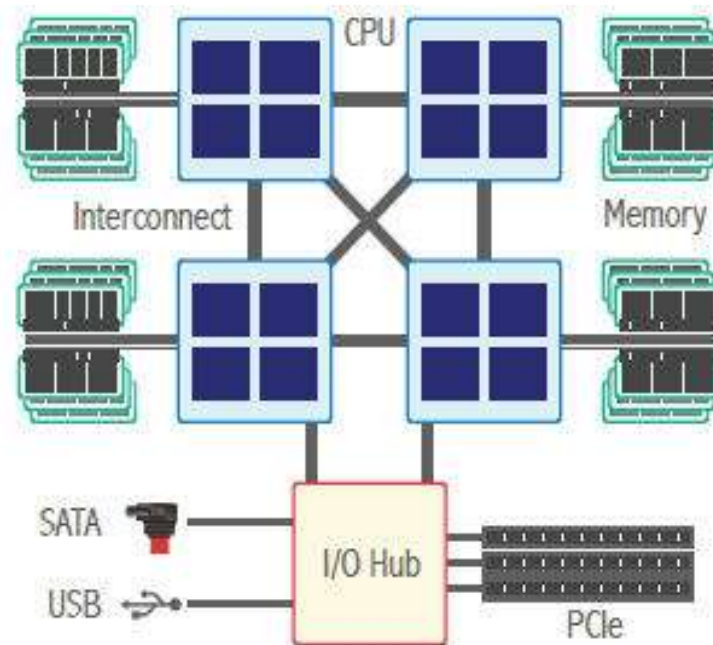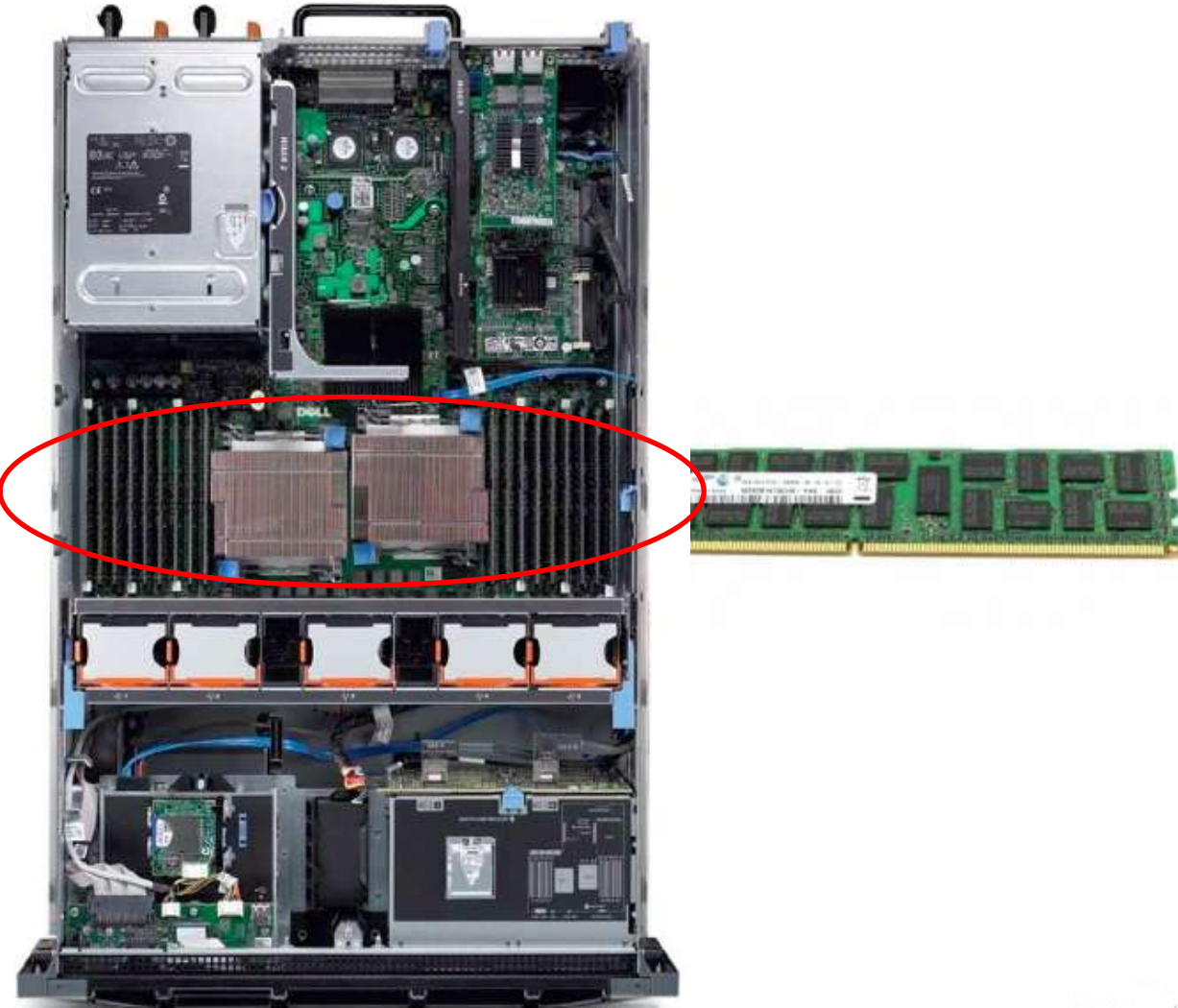
过去

Figure: PowerEdge R710 Main Components

**Memory:**

**31.5GB / 32GB 1333MHz DDR3 == 8 x 4GB - 4GB PC3-10600 Samsung DDR3-1333 ECC Registered CL9 2Rx8**

asset="02120761"

cas="9"

form="DIMM"

handle="76"

locator="DIMM_A1"
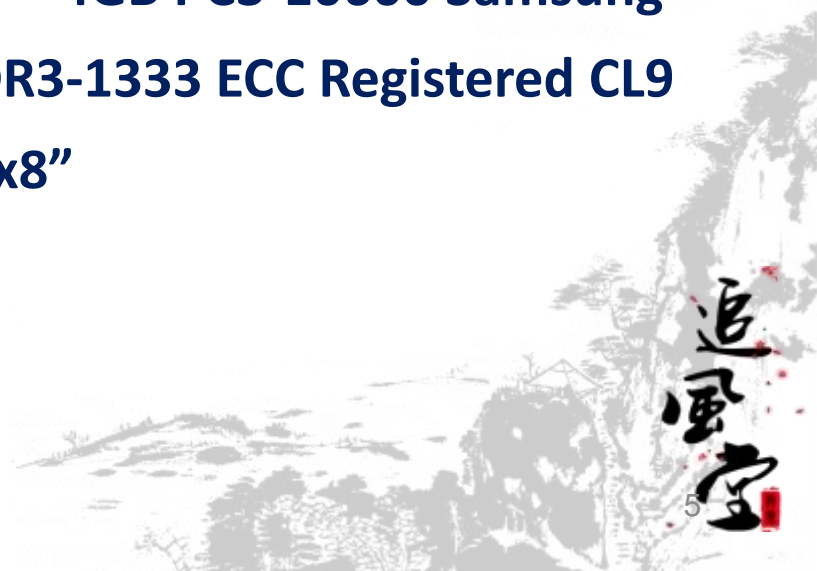
org="x8"

ranks="2"

serial="87BE9BB9"

size="4096 MB"

speed="1333MHz"

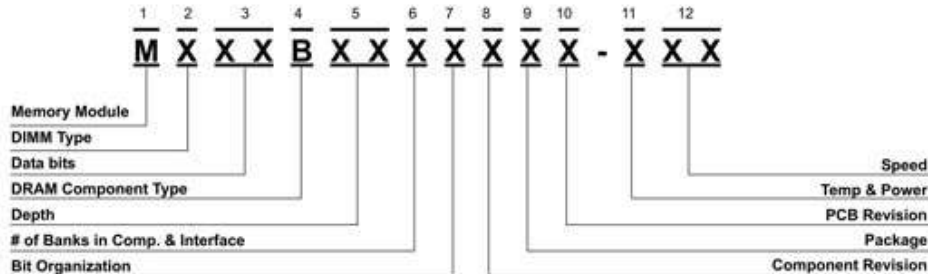type="DDR3"

width="72 bits"

part_number="M393B5273CH0-YH9"

pretty="4GB PC3-10600 Samsung DDR3-1333 ECC Registered CL9 2Rx8"

## 3. DDR3 SDRAM Module Ordering Information

```
     1  2  3  4  5  6  7  8  9 10  11 12
     M  X  X  B  X  X  X  X  X  X  - X X X
```

- Memory Module
- DIMM Type
- Data bits
- DRAM Component Type
- Depth
- # of Banks in Comp. & Interface
- Bit Organization

- Speed
- Temp & Power
- PCB Revision
- Package
- Component Revision

**1. Memory Module : M**

**2. DIMM Type**
- 3 : DIMM
- 4 : SODIMM

**3. Data Bits**
- 71 : x64 204pin Unbuffered SODIMM
- 78 : x64 240pin Unbuffered DIMM
- 91 : x72 240pin ECC unbuffered DIMM
- 92 : x72 240pin VLP Registered DIMM
- 93 : x72 240pin Registered DIMM

**4. DRAM Component Type**
- B : DDR3 SDRAM (1.5V VDD)

**5. Depth**

| | | | |
|---|---|---|---|
| 32 : 32M | | 33 : 32M (for 128Mb/512Mb) | |
| 64 : 64M | | 65 : 64M (for 128Mb/512Mb) | |
| 28 : 128M | | 29 : 128M (for 128Mb/512Mb) | |
| 56 : 256M | | 57 : 256M (for 512Mb/2Gb) | |
| 51 : 512M | | 52 : 512M (for 512Mb/2Gb) | |
| 1G : 1G | | 1K : 1G (for 2Gb) | |
| 2G : 2G | | 2K : 2G (for 2Gb) | |

**6. # of Banks in comp. & Interface**
- 7 : 8Banks & SSTL-1.5V

**7. Bit Organization**
- 0 : x4
- 3 : x8
- 4 : x16

**8. Component Revision**

| | | | |
|---|---|---|---|
| M : 1st Gen. | | A : 2nd Gen. | |
| B : 3rd Gen. | | C : 4th Gen. | |
| D : 5th Gen. | | E : 6th Gen. | |
| F : 7th Gen. | | G : 8th Gen. | |

**9. Package**
- Z : FBGA(Lead-free)
- H : FBGA(Lead-free & Halogen-free)
- J : FBGA(Lead-free, DDP)
- M : FBGA(Lead-free & Halogen-free, DDP)

**10. PCB Revision**

| | | | |
|---|---|---|---|
| 0 : None | | 1 : 1st Rev. | |
| 2 : 2nd Rev. | | 3 : 3rd Rev. | |
| 4 : 4th Rev. | | S : Reduced Layer | |

**11. Temp & Power**
- C : Commercial Temp.( 0°C ~ 85°C) & Normal Power
- Y : Commercial Temp.( 0°C ~ 85°C) & Low VDD(1.35V)

**12. Speed**
- F7 : DDR3-800    (400MHz @ CL=6, tRCD=6, tRP=6)
- F8 : DDR3-1066 (533MHz @ CL=7, tRCD=7, tRP=7)
- H9 : DDR3-1333 (667MHz @ CL=9, tRCD=9, tRP=9)
- K0 : DDR3-1600 (800MHz @ CL=11, tRCD=11, tRP=11)

NOTE: PC3-6400(DDR3-800),PC3-8500(DDR3-1066),
PC3-10600(DDR3-1333), PC3-12800(DDR3-1600)

### M393B5273CH0-YH9
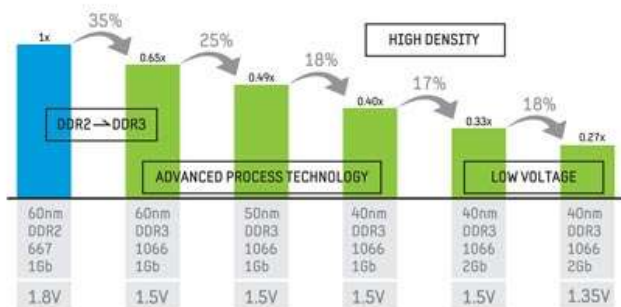DIMM、x72 240pin Registered、2Gb颗粒、x8位宽、第四代产品、无铅无汞FBGA封装、1.35V低电压、DDR3-1333

SAMSUNG DDR3 VS. DDR2
POWER SAVINGS FOR SERVER MODULES · All same density

FIGURE 3. SAMSUNG 48GB DDR3 POWER REDUCTION

Overall, a server built using Samsung's 40nm class, 1.35V, 2Gb Green DDR3 memory uses 73% less power than a 60nm, 1.8V, 1Gb DDR2 chip.
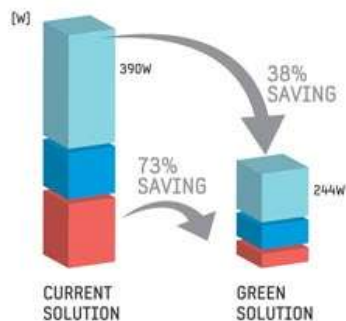In a 48GB server, this translates into a 38% power savings.
Higher density servers will show even greater savings.

FIGURE 4. 48GB MEMORY POWER SAVINGS COMPARISON
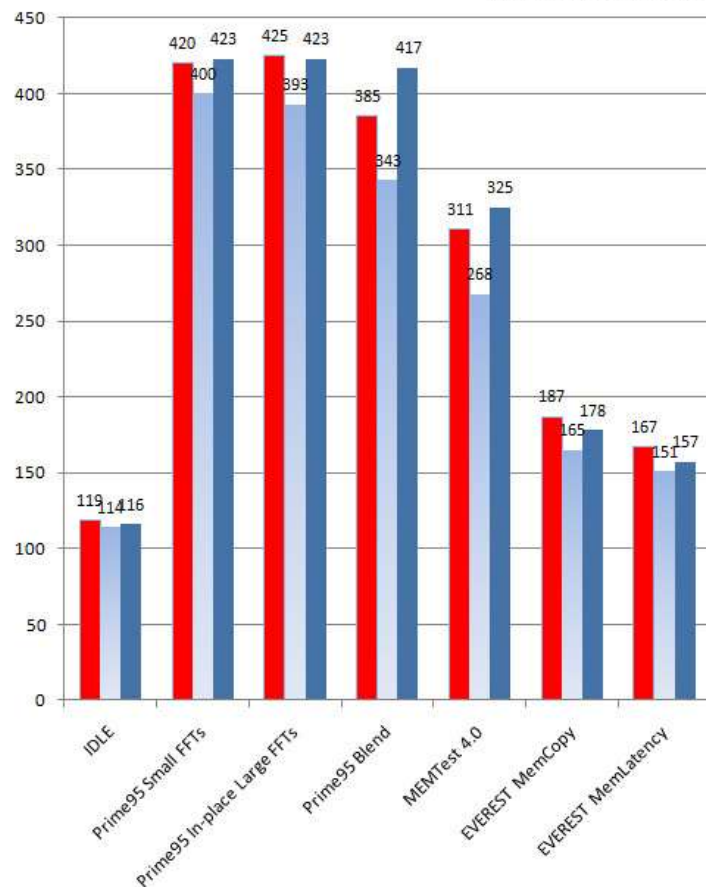
FIGURE 5.
48GB BASED SERVER SYSTEM-LEVEL POWER SAVINGS

*Source : Measured by SAMSUNG Lab

功耗测试
功耗

IT168评测中心
http://labs.it168.com

单位：W（数值越小越好）

■ 24GB(SAMSUNG 50nm 4GB x6)
■ 24GB(SAMSUNG 40nm 4GB x6)
■ 48GB(SAMSUNG 40nm 4GB x12)

| | |
|---|---|
| **L1 cache reference** | **0 . 5 n s** |
| **Branch mispredict** | **5 n s** |
| **L2 cache reference** | **7 ns** |
| **Mutex lock/unlock** | **25 ns** |
| **Main memory reference** | **100 ns** |
| **Compress 1K bytes with Zippy** | **3,000 ns** |

```
Memory latencies in nanoseconds — smaller is better
------------------------------------------------------------
----------------

Host OS Mhz L1 $ L2 $ Main mem Rand mem Guesses
------------------------------------------------------------
---

Dr4000 Linux 2.6.32- 2631 1.1590 5.7170  78.0 110.4
```

内存**带宽计算公式：带宽**=内存核心**频率x**内存**总线**位数**x**倍增系数。

**每个通道** (1333/8)*64*8 /8 = 10.6G Byte；

而我**们的CPU**是**3**个通道的，也就是**说这个CPU**的**总**的内存**带宽是** 10.6*3=31.8G

Nehalem-EP（Xeon 5500系列）双路服务器系统架构图

QPI (25.6 GB/s)

DDR3 memory 10.6GB/s
DDR3 memory 10.6GB/s
DDR3 memory 10.6GB/s

Nehalem-EP Xeon 5500

Nehalem-EP Xeon 5500

DDR3 memory 10.6GB/s
DDR3 memory 10.6GB/s
DDR3 memory 10.6GB/s

QPI (25.6 GB/s)

QPI (25.6 GB/s)

PCI Express* 2.0 I/O

Support for Multi-card configurations: 1x16, 2x16, 4x8 or other combination

up to 36 lanes

5500 5520 IOH

2 GB/s DMI

特性优势：

• 3条QPI总线各为25.6GB/s，CPU互联与I/O总带宽是Xeon 5400系统的3.6倍
• DDR3内存最大传输率是DDR2-800的1.67倍
• 6通道内存总带宽量高达64GB/s，是Xeon 5400系统的2.5倍
• CPU互联+I/O+内存总带宽是Xeon 5400系统的2.75倍
• DDR3内存总能耗比同容量DDR2低16.7%，远胜FB-DIMM
• CPU可进行调频以应对节能与高负载应用的需求

vtbwrun -c -A

# NUMA节点内存访问速度差异



```
root@dr4000 # numademo -t 1g memset
2 nodes available
memory with no policy memset                        Avg  6358.97  MB/s
local memory memset                                 Avg  6357.17  MB/s
memory interleaved on all nodes memset              Avg  4778.59  MB/s
memory on node 0 memset                             Avg  5836.55  MB/s
memory on node 1 memset                             Avg  3669.97  MB/s
memory interleaved on 0 1 memset                    Avg  4508.74  MB/s
        numa_foreign                    0            1709654
        interleave_hit              17395              17388
        local_node               66916103           72839750
        other_node                1727189              43003
root@dr4000 #
```
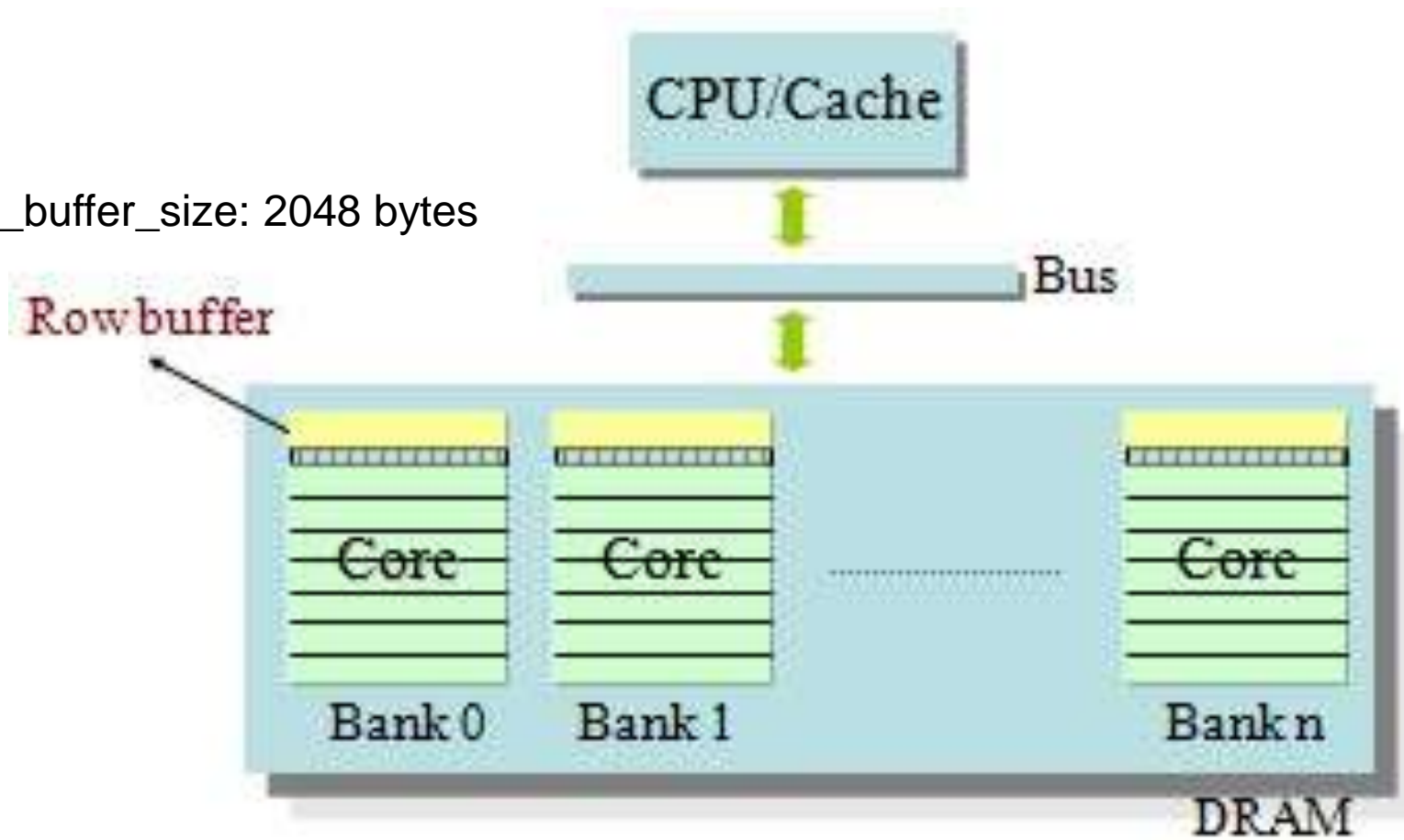
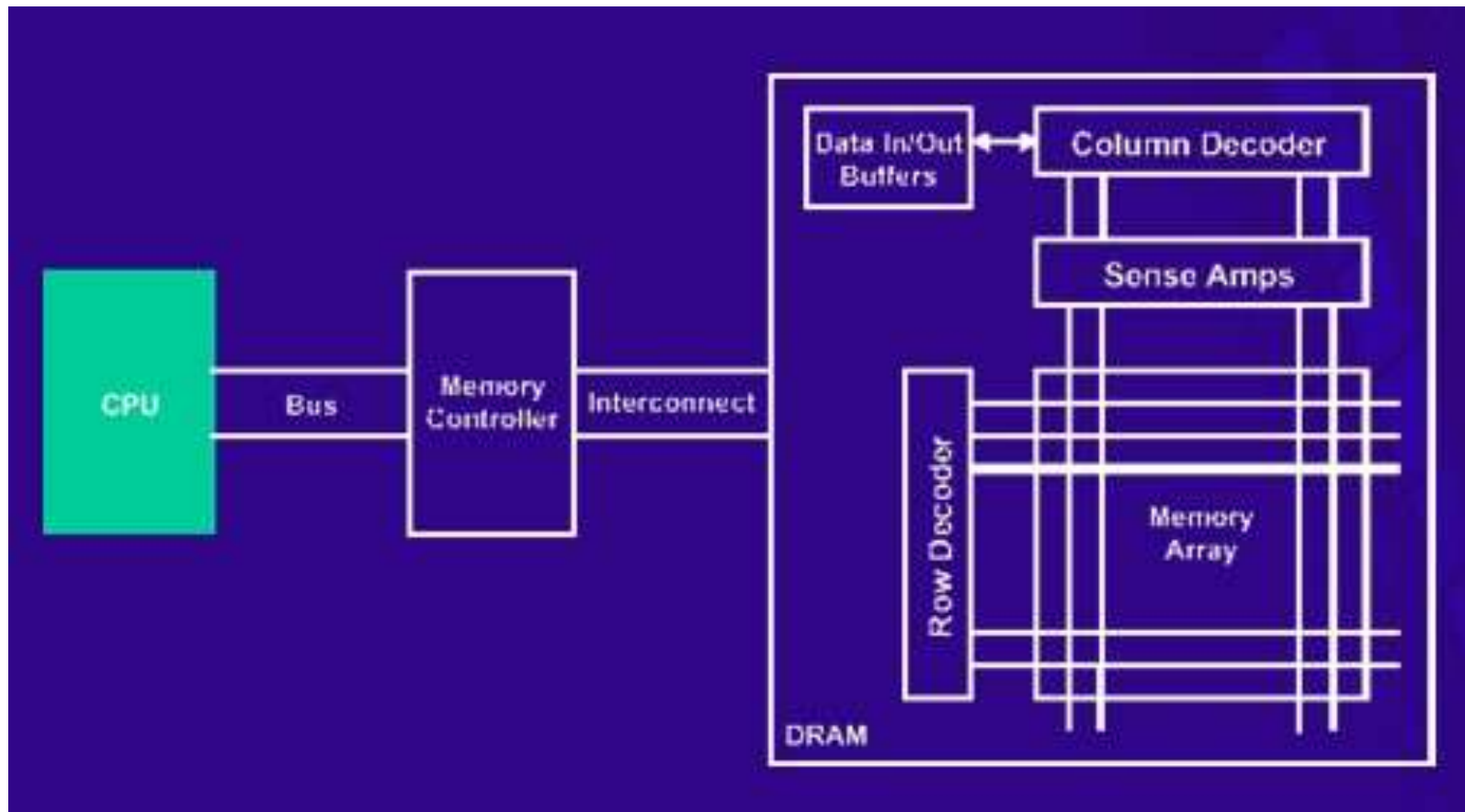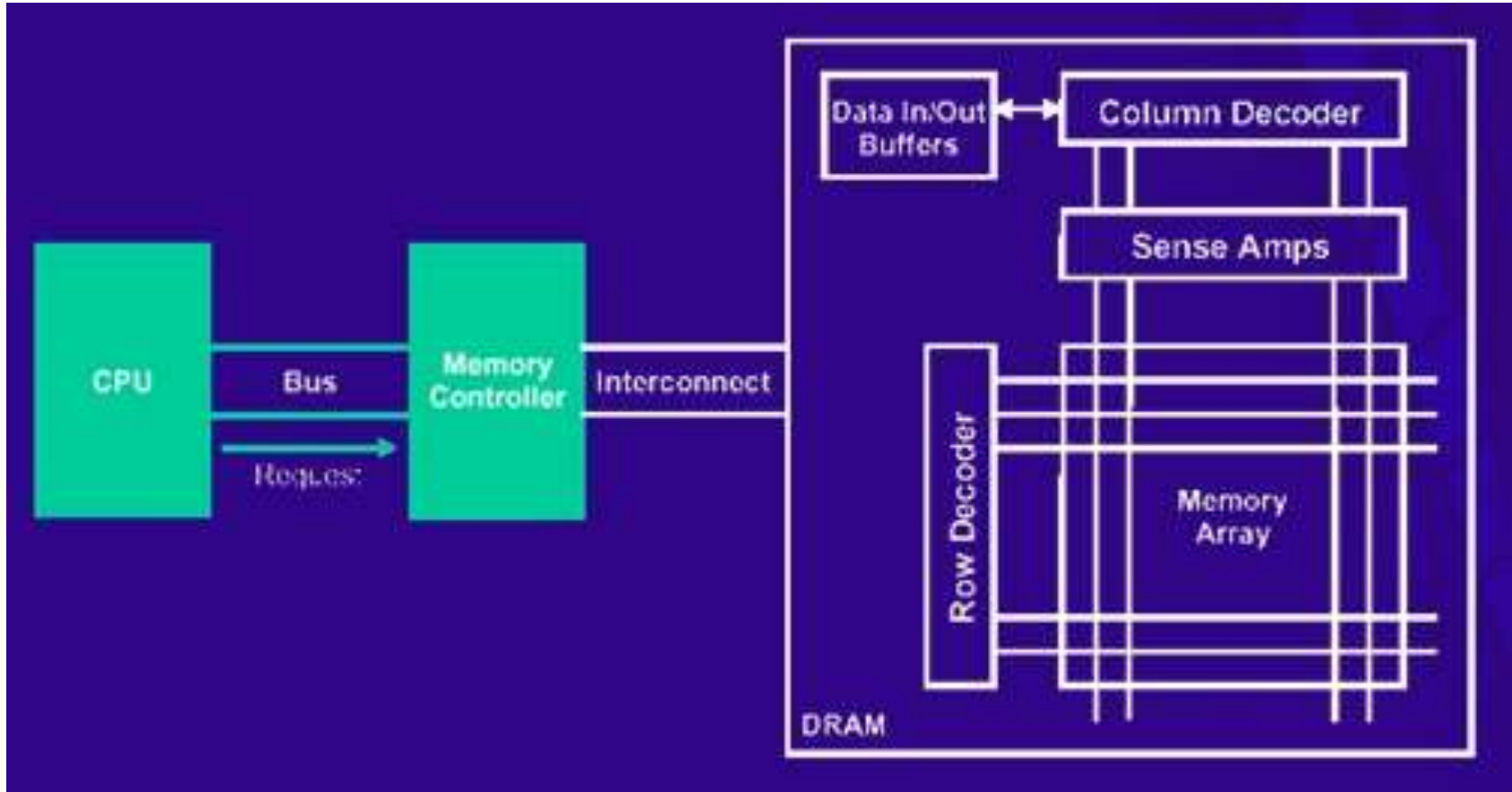Locality Exploitation in Row Buffer

# DRAM结构图

row_buffer_size: 2048 bytes

- **Precharge: charge a DRAM bank before arow access**

- **Row access: activate a row (page) of a DRAM bank**

- **Column access: select and return a block of data in an activated row**

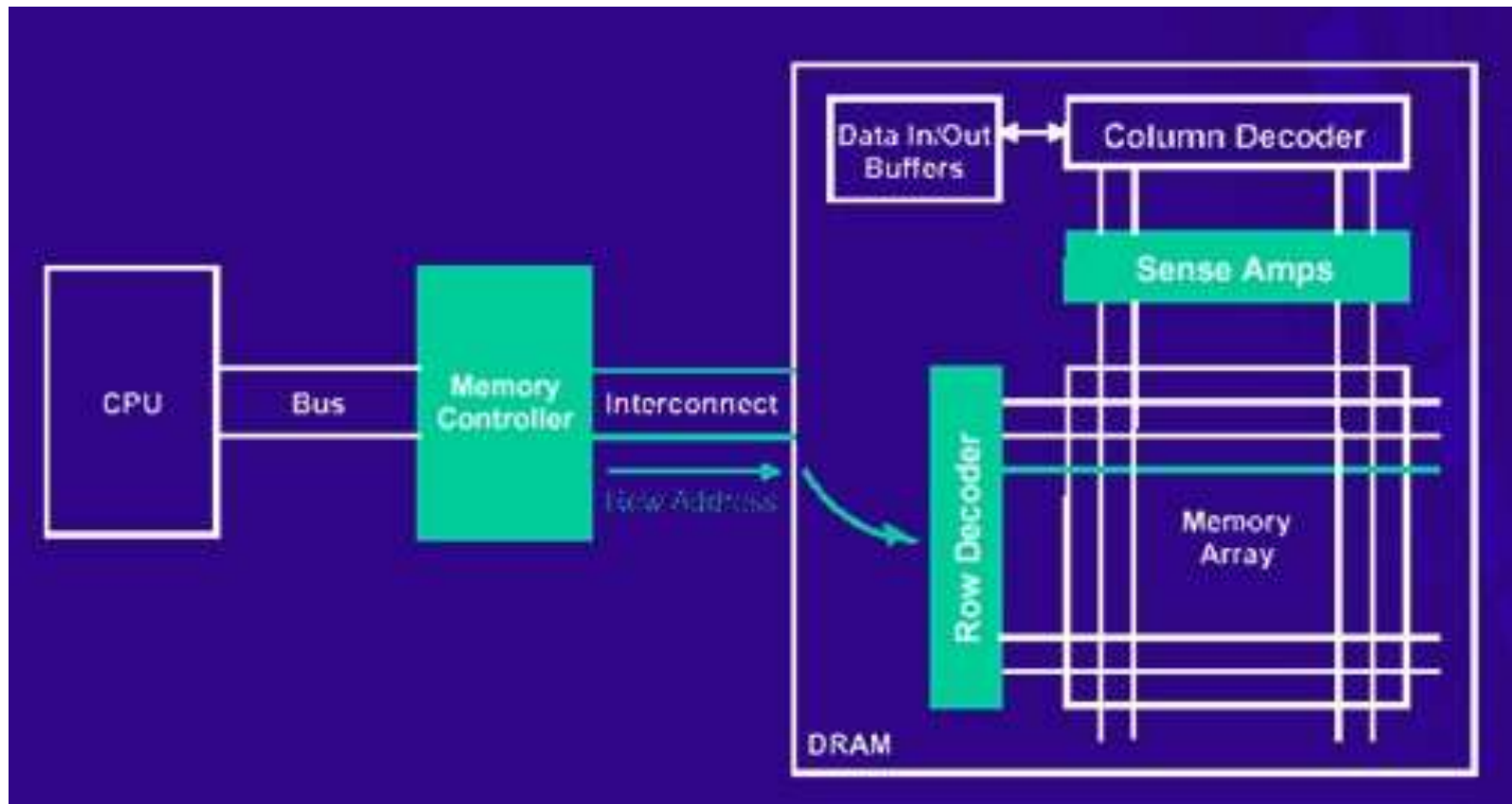- **Refresh: periodically read and write DRAM to keep data**

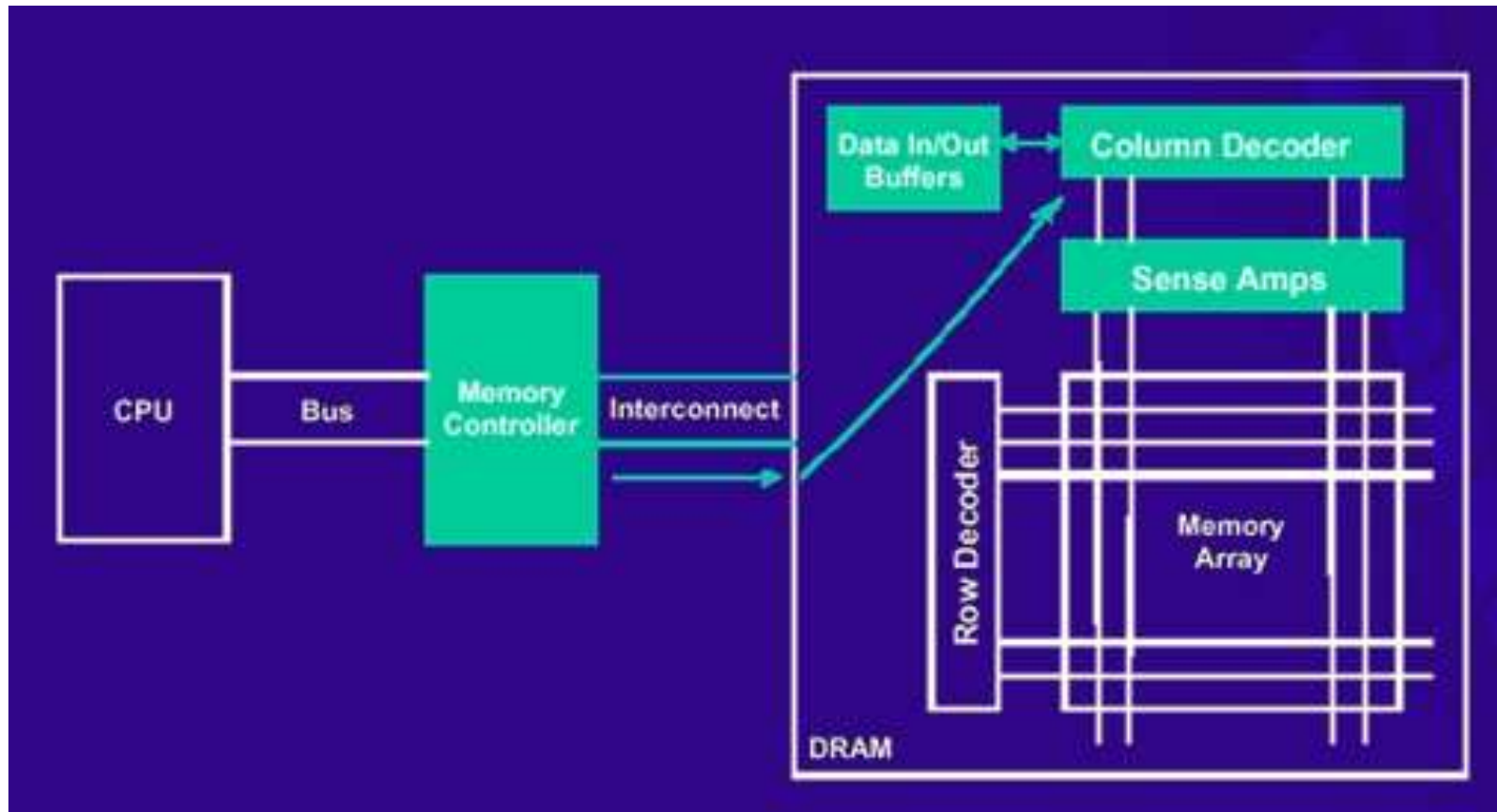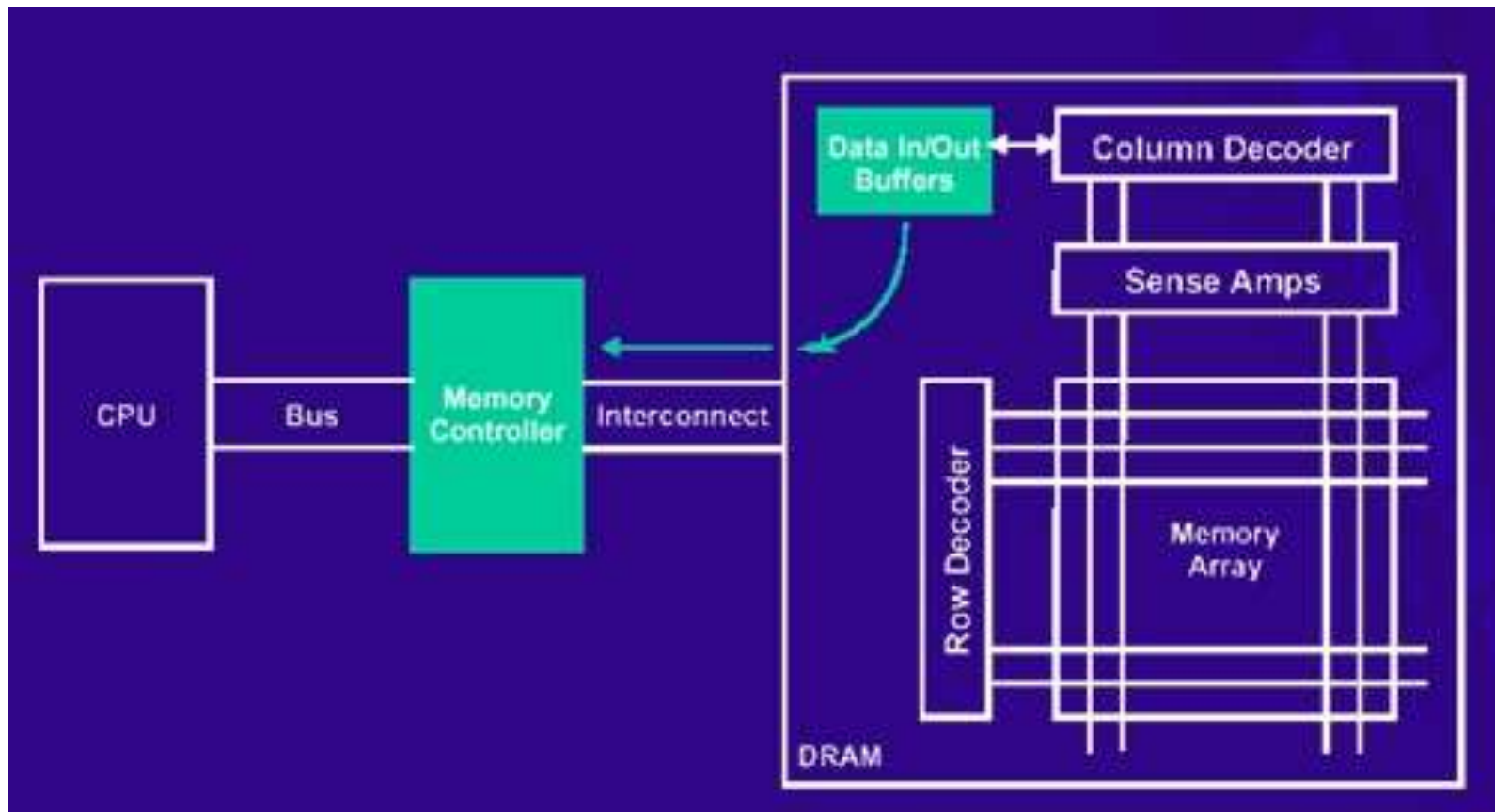## Nonuniform DRAM Access Latency

- Case 1: Row buffer hit (20+ ns)

  | col. access |

- Case 2: Row buffer miss (core is precharged, 40+ ns)

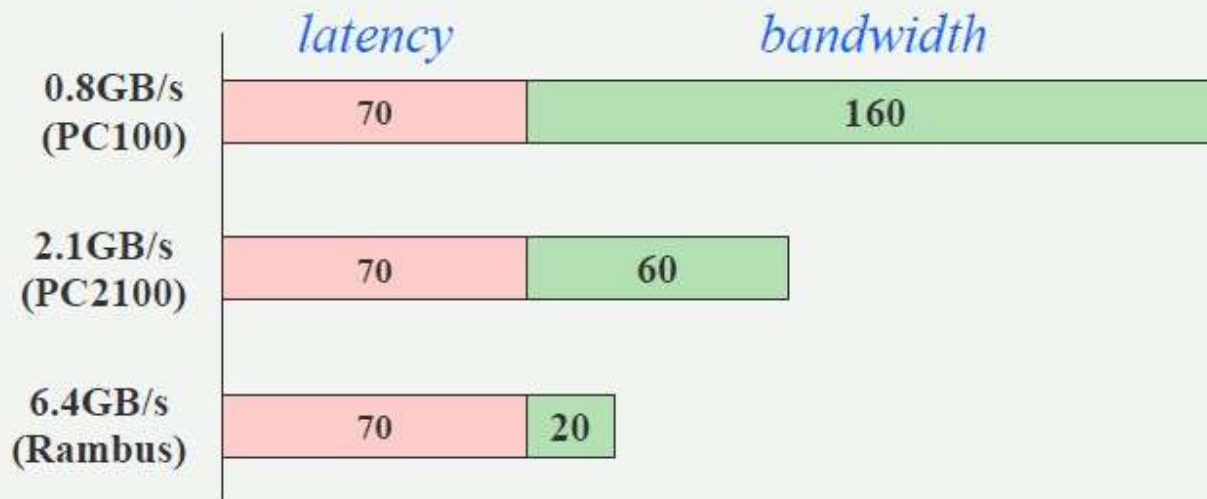  | row access | col. access |

- Case 3: Row buffer miss (not precharged, ≈ 70 ns)

  | precharge | row access | col. access |

# Amdahl's Law applies in DRAM

♦ Time (ns) to fetch a 128-byte cache block:

| | latency | bandwidth |
|---|---|---|
| 0.8GB/s (PC100) | 70 | 160 |
| 2.1GB/s (PC2100) | 70 | 60 |
| 6.4GB/s (Rambus) | 70 | 20 |

♦ As the bandwidth improves, DRAM latency will decide cache miss penalty.

- **http://en.wikipedia.org/wiki/Prefetch_buffer**
- **详解服务器内存带宽计算和使用情况测量:**
  **http://blog.yufeng.info/archives/1511**
- **DDR3 内存带宽如何计算:**
  **http://zhidao.baidu.com/question/107154668**
- **hwconfig查看硬件信息:**
  **http://blog.yufeng.info/archives/2086**
- **Exploiting Locality in DRAM, Xiaodong Zhang**

# 谢谢大家！