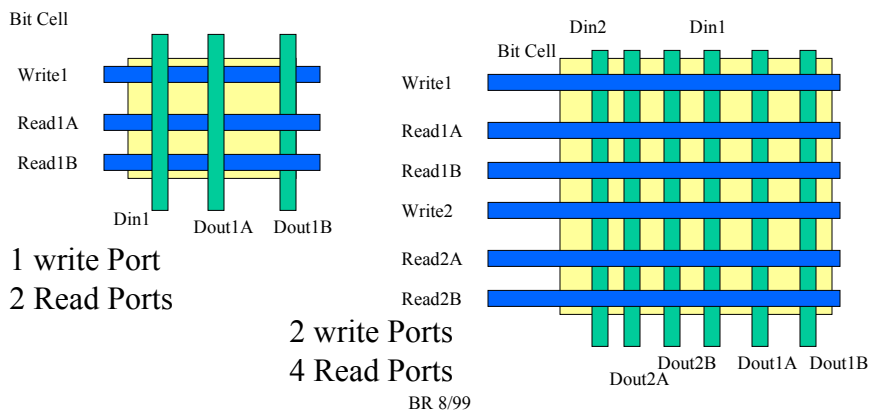# Register Files

- Register files on processors are simply small SRAMs
- Always multi-ported
  - Minimum of 2 read ports, 1 write port (fetch two operands, write one result)
- Superscalar, VLIW processors require register files with many ports
  - 4 instructions per clock requires 8 read ports, 4 write ports
- Because array size is small, only one Bitline needed, no sense amp.
  - Main concern on register file is SPEED.

BR 8/99

---

# Multi-Ported Register File Design has Limits

- Area of the register file grows approximately with the square of the number of ports
  - Typically routing limited, each new port requires adding new routing in both X and Y direction



1 write Port
2 Read Ports
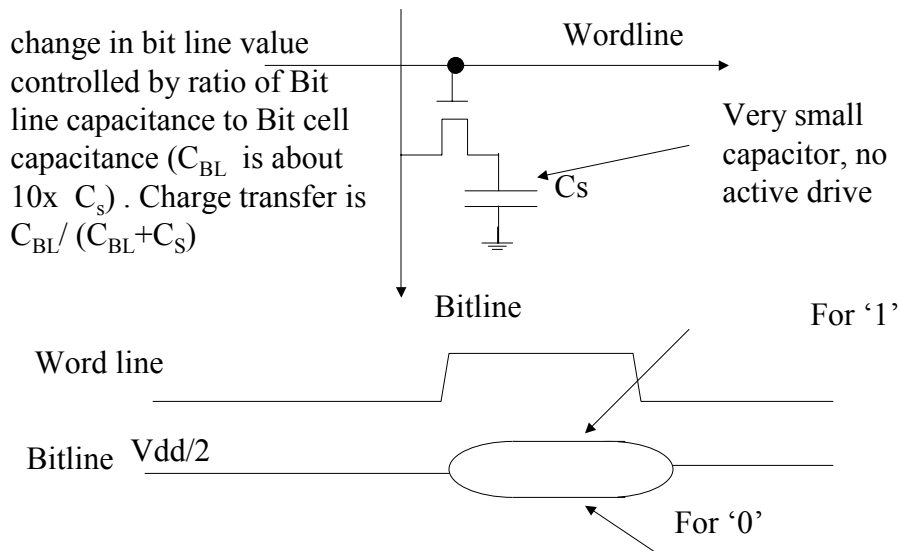
2 write Ports
4 Read Ports

BR 8/99

# Multiported Register Files (cont)

- Read Access time of a register file grows approximately linearly with the number of ports
  - Internal Bit Cell loading becomes larger
  - Larger area of register file causes longer wire delays
- What is reasonable today in terms of number of ports?
  - Changes with technology, 15-20 ports is currently about the maximum (read ports + write ports)
  - Will support 5-7 execution units simultaneous operand accesses from register file

---

# Dynamic Memory Cell

change in bit line value controlled by ratio of Bit line capacitance to Bit cell capacitance ($C_{BL}$ is about 10x $C_s$) . Charge transfer is $C_{BL}/ (C_{BL}+C_S)$

Wordline

Very small capacitor, no active drive
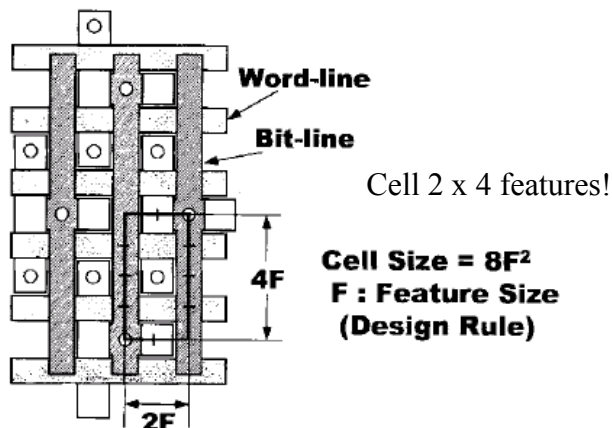
$Cs$

Bitline

Word line

Bitline   Vdd/2

For '1'

For '0'

# DRAM Comments

- Voltage swing on Bitline is small
  - Want Bitline capacitance as small as possible, Bit cell capacitance as large as possible to increase charge transfer
- Read is destructive – part of read cycle is used to restore level inside of bit cell capacitor
- Capacitor leaks, must be refreshed periodically
- Noise sources in DRAM are word line to bit line coupling, bit line to bit line coupling

---

DRAM Layout†



Word-line

Bit-line

Cell 2 x 4 features!

Cell Size = 8F²
F : Feature Size
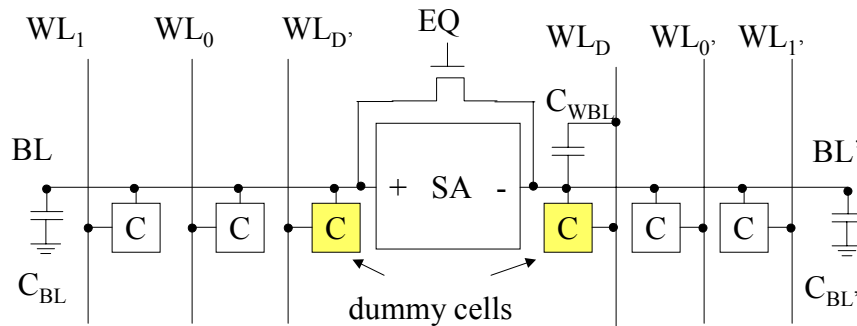(Design Rule)

4F

2F

Actual size for 256 Mb DRAM cell in 0.22 um  reported as 0.484um x 0.968 um   ( 2.2F  x  4.4F)

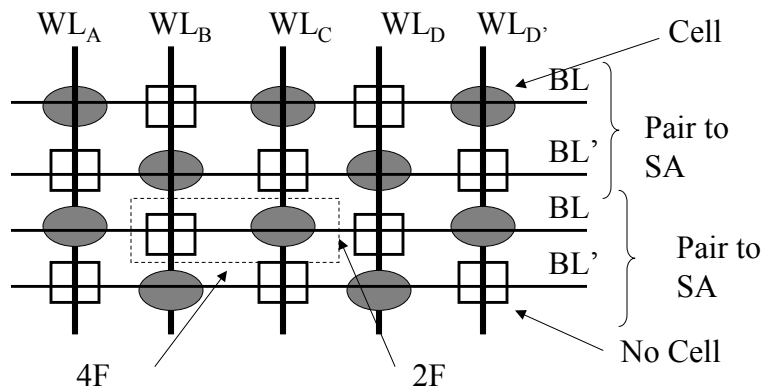†Okuda,  "A Four-level Storage 4GB DRAM", IJSCCS Vol 32, No11, Nov 1997

# Open Bitline Architecture

WL$_1$      WL$_0$      WL$_{D'}$      EQ      WL$_D$   WL$_{0'}$   WL$_{1'}$

C$_{WBL}$

BL                                    + SA -                                   BL'

C        C        C                              C        C        C

C$_{BL}$                                                           C$_{BL'}$

dummy cells

EQ raised, and L$_D$, L$_{D'}$ also raised to precharge bitlines and dummy bit cells to Vdd/2 . During read, if cell from left hand side is read (raise WL$_1$), then dummy word line on right cell is raised. Raising a word line couples noise into bitline, want same noise injected on both lines (only the same if both signals are matched, and bitlines are matched).

---

Folded Bitline Architecture  -- more noise supression

WL$_A$     WL$_B$     WL$_C$     WL$_D$    WL$_{D'}$          Cell

BL

Pair to
SA

BL'

BL

Pair to
SA

BL'

No Cell

4F                          2F

Cell is connected to BL or BL' on every other column.   WL$_D$, WL$_D$' are dummy bit lines.   Assume WL$_A$ is driven.  Then WL$_D$ is driven;  if WL$_B$ is driven, then WL$_{D'}$ is driven.
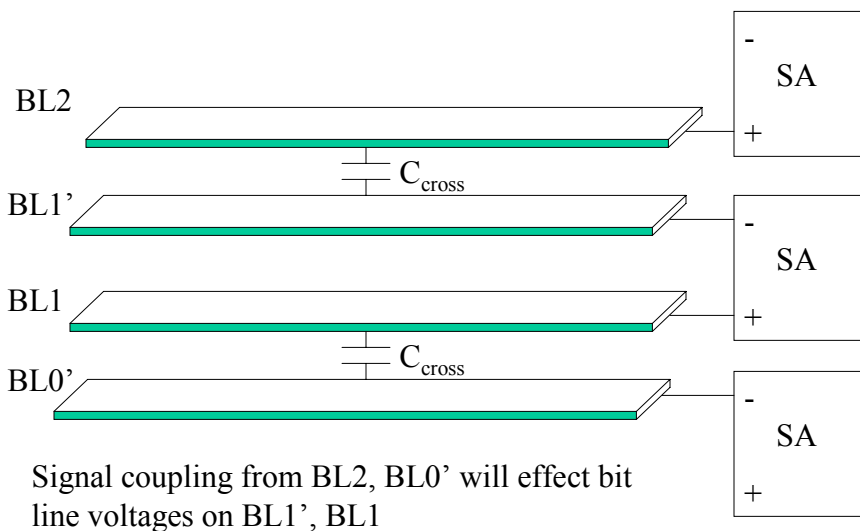
# Folded Bit lines – more noise suppression

- normal word lines ($WL_A$, $WL_B$, etc) and dummy word lines ($WL_D$, $WL_{D'}$) cross both bitlines
  - Even if signal characteristics of both word lines differ substantially, same coupling noise from both is coupled into both bitlines, which appears as common mode to SA, which rejects it.
  - bitlines more closely matched since they run side-by-side
- However, bitlines in folded architecture are somewhat longer than non-folded bitlines, so bit line capacitance is higher, reducing charge transfer from cell.
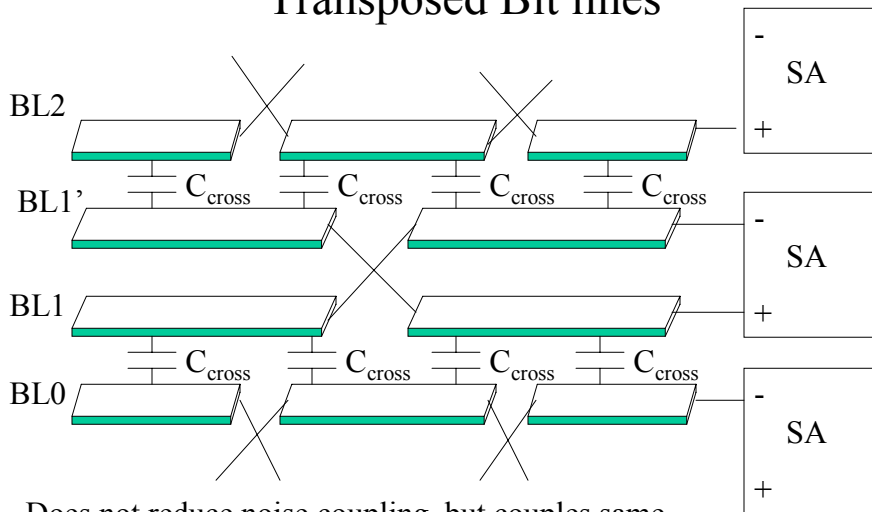
---

# Bit Line to Bit Line Coupling



Signal coupling from BL2, BL0' will effect bit line voltages on BL1', BL1

# Transposed Bit lines

BL2

BL1'

BL1

BL0

$C_{cross}$  $C_{cross}$  $C_{cross}$  $C_{cross}$

$C_{cross}$  $C_{cross}$  $C_{cross}$  $C_{cross}$

- SA +

- SA +

- SA +

Does not reduce noise coupling, but couples same
noise into both bitlines so appears as common node
noise and is rejected.

BR 8/99

---

# Capacitor Design

- To achieve density for Mb, Gb DRAMs, capacitors had to go 3D
- CUB – Capacitor under Bitline
  - Trench capacitor, sidewall of trench used to form capacitor plate
  - Used up to 16Mb DRAM
- COB – Capacitor Over Bitline – stacked capacitor
  - Used in 64Mb+ DRAMs
  - Forms a 3D cylindrical capacitor in which both inside and outside surfaces can be used
  - Memory cell array at a different (higher) height than surrounding surfaces, can cause metallization problems.
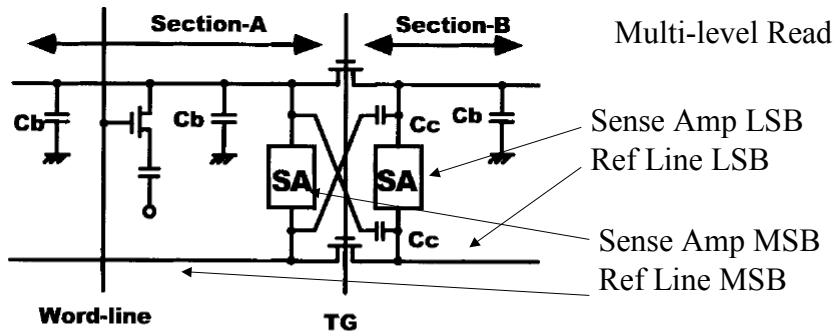
BR 8/99

# Multi-level DRAM Cells

- Multi-level DRAM cell are being investigated for 4 Gb DRAMs
- 4 states: 00, 01, 10, 11
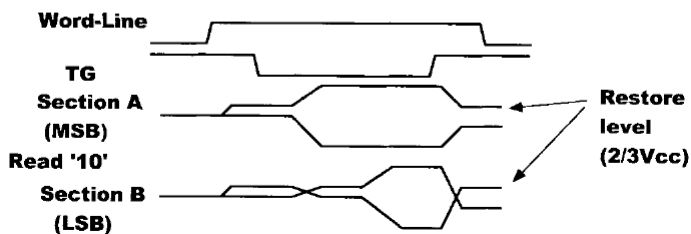- Multi-level SRAM cells have already been demonstrated for Flash RAM

TABLE I
FOUR-LEVEL STORAGE

| | Data | Storage | Reference | Signal |
|---|---|---|---|---|
| | 11 | Vcc | | 1/6 Vcc |
| 4-Level | 10 | 2/3 Vcc | 5/6 Vcc | |
| Storage | | | 3/6 Vcc | |
| | 01 | 1/3 Vcc | 1/6 Vcc | |
| | 00 | 0 (GND) | | |
| 2-Level | 1 | Vcc | | 1/2 Vcc |
| Storage | 0 | 0 (GND) | 1/2 Vcc | |

---

Multi-level Read



Sense Amp LSB
Ref Line LSB

Sense Amp MSB
Ref Line MSB

(a)

(b)

## Read of '10'

a. 4 line segments: Bl-A, Ref-A (msb), Bl-B, Ref-B (lsb). Bl-A coupled to Ref-B by Cc, Bl-B coupled to Ref-A by Cc. Define Vs as max voltage on bitlines. Coupling capacitor Cc designed to transfer 1/3 of swing to attached lines.

b. Transfer transistors can join segments Bl-A, Bl-B and segments Ref-A, Ref-B.

c. Ref-A used for MSB sensing, Ref-B used for LSB sensing.

d. Initial Conditions: Transfer transistors on, all segments precharged to 1/2Vs, cell contains 2/3 Vs

e. Word line asserted, Bl-A, Bl-B settle at about 2/3 Vs (a little under, Cs >> Cb). Via Cc, Ref-B raises by about 1/9Vs from 1/2Vs to 5.5/9Vs

f. Turn off Transfer transistors to isolate segments. MSB SA compares Bl-A (2/3Vs) to Ref-A(1/2Vs), senses a '1'. BL-A driven to Vcc, Ref-A driven to '0' by SA.

g. Via Cc, Bl-B drops from 2/3Vs to about 1/2Vs, while Ref-B raises another 1/9Vs to 6.5/9Vs.

h. LSB Sense amp compares Bl-B (1/2Vs) to Ref-B (6.5/9Vs) and senses a '0'. Bl-B driven to '0', Ref-B driven to '1'.

i. Turn on transfer transistors to join segments. Cap ratio of segments is 2 to 1, so Bl-A at Vs and Bl-B at 0v with 2/1 ratio (Bl-A twice cap of Bl-B) gives original cell vaue of 2/3 Vs for cell restore.

BR 8/99

---

## Read of '11'

a. 4 line segments: Bl-A, Ref-A (msb), Bl-B, Ref-B (lsb). Bl-A coupled to Ref-B by Cc, Bl-B coupled to Ref-A by Cc. Define Vs as max voltage on bitlines. Coupling capacitor Cc designed to transfer 1/3 of swing to attached lines.

b. Transfer transistors can join segments Bl-A, Bl-B and segments Ref-A, Ref-B.

c. Ref-A used for MSB sensing, Ref-B used for LSB sensing.

d. Initial Conditions: Transfer transistors on, all segments precharged to 1/2Vs, cell contains Vs

e. Word line asserted, Bl-A, Bl-B settle at about Vs (a little under, Cs >> Cb). Via Cc, Ref-B raises by about **1/6Vs** from 1/2Vs to **2/3Vs**

f. Turn off Transfer transistors to isolate segments. MSB SA compares Bl-A (**1.0Vs**) to Ref-A(1/2Vs), senses a '1'. BL-A driven to Vcc, Ref-A driven to '0' by SA.

g. Via Cc, Bl-B drops from 1.0Vs to about **5/6Vs**, while Ref-B remains stable at **2/3Vs** since Bl-A bitline was already at 1.0Vs.

h. LSB Sense amp compares Bl-B (**5/6Vs**) to Ref-B (**2/3Vs**) and senses a '1'. Bl-B driven to '1', Ref-B driven to '0'.

i. Turn on transfer transistors to join segments. **Both Bl-A and Bl-B are at Vs, which is original cell value for restore operation.**

BR 8/99

# DRAM Interface

Multiplexed Address bus (Row, Column).  RAS# (Row Address Strobe), CAS#(Column Address Strobe) used to latch in address.

**READ CYCLE**



**EDO-PAGE-MODE READ CYCLE**



Block transfer.  Access different bits on same row, change column address.

# Comments on Timings

- Typical times are Tras = 60 ns (RAS pulse width), Trc = 100 ns
  - Extra time on Read cycle (RAS high) is needed to recharge bitlines
- Block mode transfers (Page mode transfers) read bits from same row
  - Only change column address
  - Time to first bit on row = 50ns, time to successive bits = 25 ns (we have access to all bits on this row, just need to mux them out).
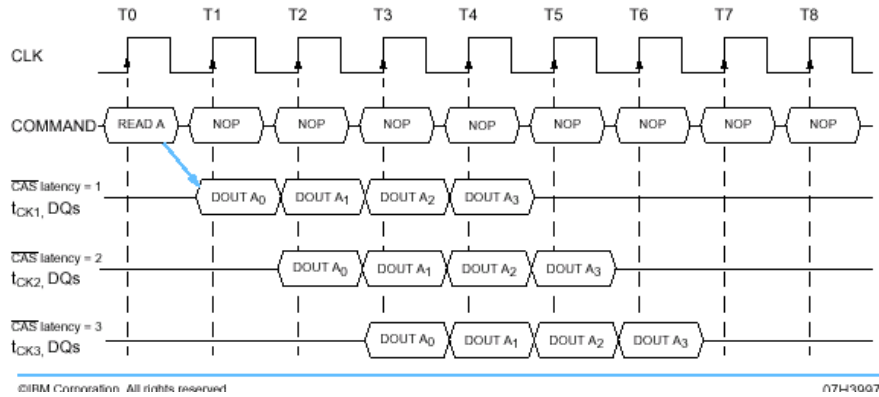
# Architectural Issues in DRAMs

- Need to support block transfers efficiently since DRAM used as main memory and reads/writes due to cache fills
- Add a clock to DRAM interface (SDRAM, DDR-SDRAM) to support burst mode operations for cache fills
  - Pentium burst mode is 2-1-1-1 (two clocks for first data, 1 clock for each sucessive data, address only provided for first data, internal counter on RAM used for address generation).
  - Pentium Pipelined burst mode is:
    2-1-1-1; 1*-1-1-1; 1*-1-1-1; ....
    Sucessive cycles pick up where the last cycle left off.

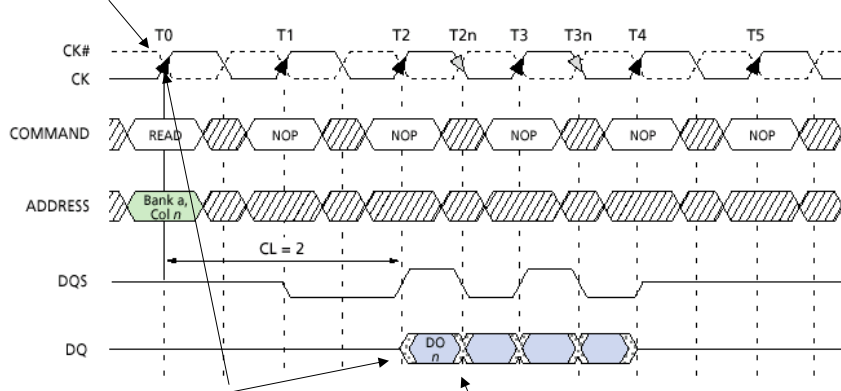# SDRAM - different CAS Latencies for burst operation

**Burst Read Operation** (Burst Length = 4, CAS latency = 1, 2, 3)

# Double Data Rate SDRAM (DDR-DRAM)

Differential Clocks



Two clock latency

Data transferred on each clock crossing

# Timing – SDRAM, DDR-SDRAM

- Clock Frequency – 133 Mhz, 100 Mhz
- Two clock latency to first data  (20 ns for 100 Mhz clock)
  - SDRAM -  10 ns per location afterwards.  For byte-wide, 100 MB/sec transfer rate.  400 MB/sec on 32-bit bus
  - DDR-SDRAM -  5 ns per location afterwards. For byte-wide, 200 MB/sec transfer rate.   On 32-bit bus, 800 MB/sec transfer rate.

# Rambus DRAM (RDRAM)

- DRAM with a high speed interface
- 400 Mhz differential clock, data transferred on each edge
- Reduced swing signaling about a reference voltage
  - Termination voltage is 1.5 V
  - Reference Voltage is 1.0 V
  - Signals swing +/- 200 mv about reference voltage
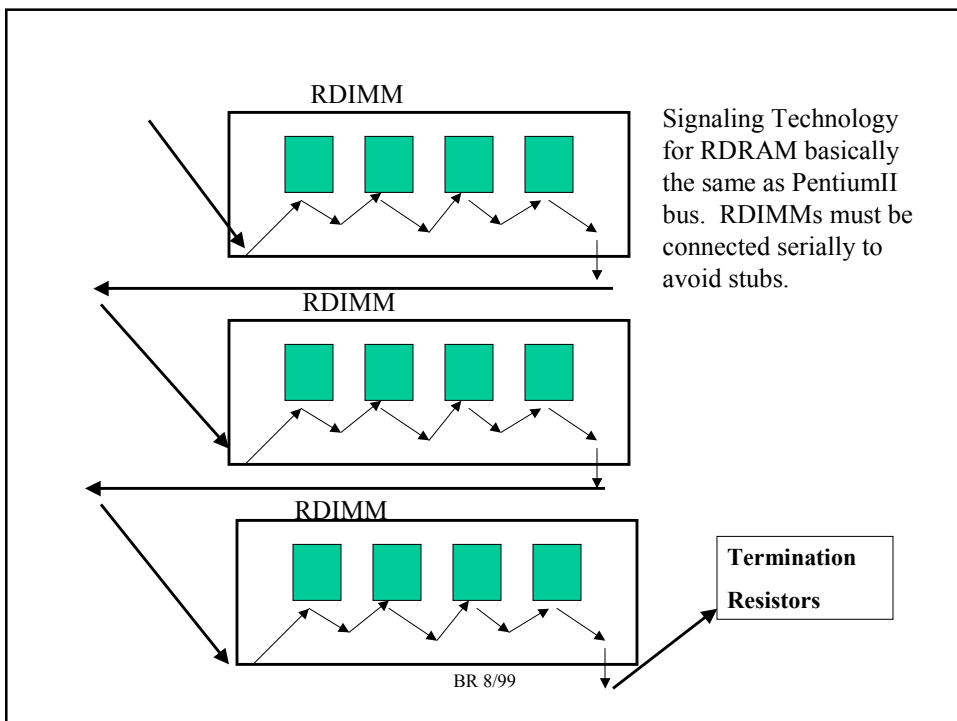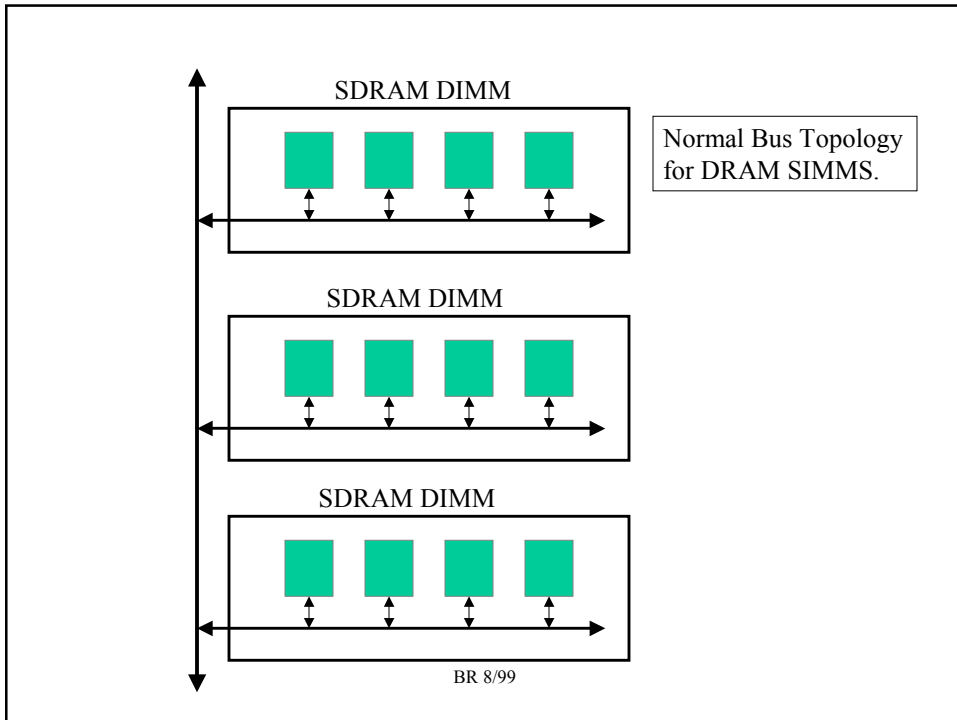  - All traces are transmission lines

# Bandwidth

- External bus is 18 bits wide (2 bytes + 2 parity bits)
- External clock cycle is 400 Mhz, but data is clocked on each edge
  - Actually, external clock is a differential pair and data is sampled at each crossing
- Total Bandwidth is 1.6 GBytes/s
  - 2 bytes * 400 Mhz * 2 edges => 1.6 Gbytes
  - Initial configurations are 4 M x 18 (72 Mbits)

# Maximum Bandwidth

- Note that maximum bandwidth with one RDRAM controller is 1.6GB/s.
  - Only one RDRAM chip can be active at a time on RDRAM bus.
  - More RDRAM chips increase capacity, not bandwidth.
    - With normal DRAM and SDRAM, can increase bandwidth by just adding more DRAM chips in parallel from same DRAM controller
  - To double the bandwidth, would need two separate RDRAM controllers

SDRAM DIMM

Normal Bus Topology
for DRAM SIMMS.

SDRAM DIMM

SDRAM DIMM

BR 8/99

RDIMM

Signaling Technology
for RDRAM basically
the same as PentiumII
bus. RDIMMs must be
connected serially to
avoid stubs.

RDIMM

RDIMM

**Termination**

**Resistors**

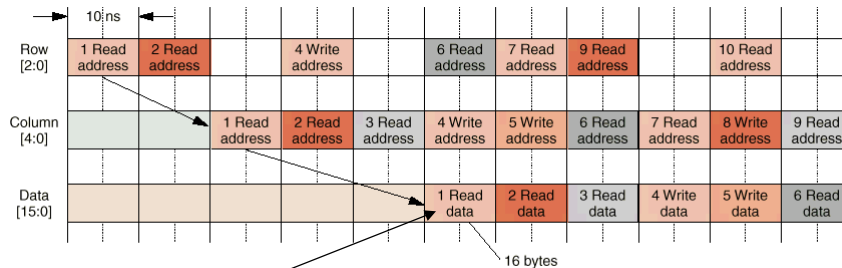BR 8/99

# Deep Pipelining => High Latency

Figure 7. Direct RDRAM interleaved memory transactions at full-memory bandwidth (16 bytes/10 ns).

16 bytes transferred because 4 clocks * 2 edges * 2 bytes/transfer

(external bus is 16 or 18 bits wide). 20 clock latency, 20 ns from column address)

BR 8/99

---

# Addressing

- 3-Bit Row bus used to give commands to RDRAM
- ROW Activate command used for read
  - 4 clocks transfers 8 groups of 3 bits over Row bus due to dual edge clocking (24 bits total)
  - 24 bits in Row Activate command split between device address (6 bits), bank select (4 bits), row select (9 bits), and reserved bits
- There are no chip select lines, internal register holds device address
  - All chips monitor bus - if bus device address matches internal id, then chip is selected.

BR 8/99

# Row Activate Command

10 ns

$T_0$     $T_1$     $T_2$     $T_3$

CTM/CFM

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| ROW2 | DR4T | DR2 | BR0 | BR3 | RsvR | R8 | R5 | R2 |
| ROW1 | DR4F | DR1 | BR1 | RsvB | RsvR | R7 | R4 | R1 |
| ROW0 | DR3 | DR0 | BR2 | RsvB | AV=1 | R6 | R3 | R0 |

ROWA Packet

R bits = row select

DR bits = device address

BR bits = bank select

BR 8/99

---

Controller    RDRAM 1    RDRAM 2    RDRAM *n*

IEEE Micro Nov/Dec 1997

INIT    INITo

$V_{TERM}$

Bus data [18:0]
RC[7:0]
RClk[2]
TClk[2]
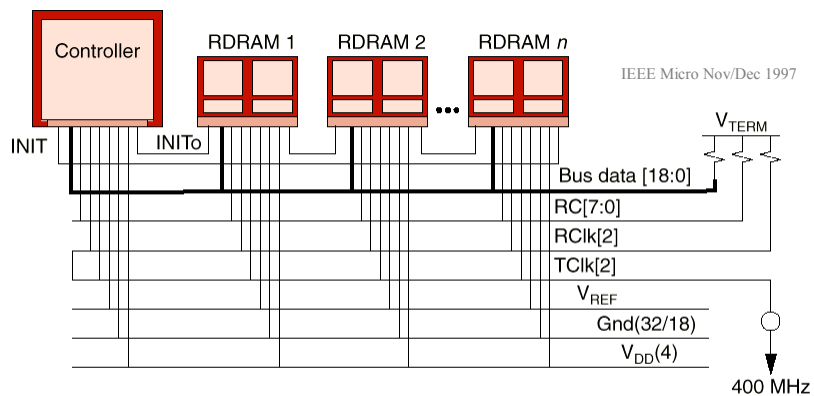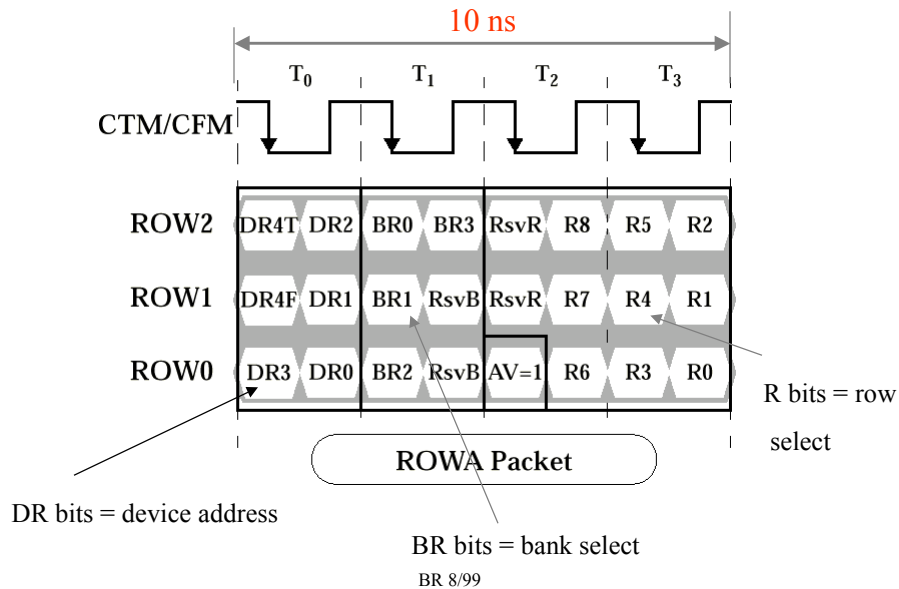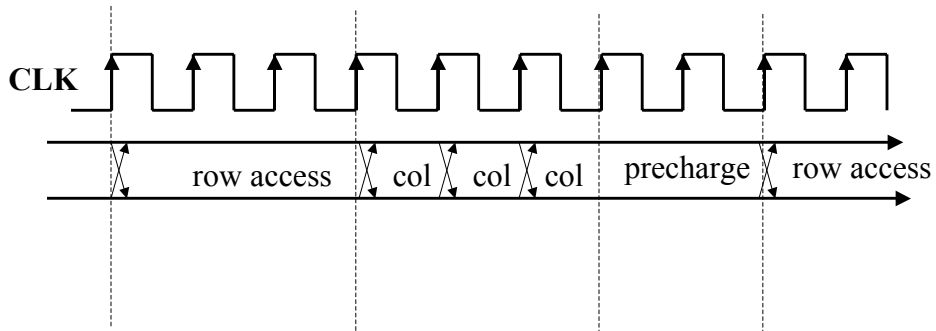$V_{REF}$
Gnd(32/18)
$V_{DD}(4)$

400 MHz

Figure 3. Direct RDRAM system.

18 bit wide external data bus which expands into 128 bit wide datapath internal to chip

BR 8/99

Portion of internal architecture ( 4M x 16  or 4M x 18)

16 banks of 512 rows of 64 dualocts (1 dualoct = 16 bytes = 128 bits)

$2^4$ (banks) * $2^9$ (rows) * $2^6$ (dualocts) * $2^7$ (one dualoct) = $2^{26}$ (64 Mbit)

A dualoct is the smallest addressable unit.

# Multiple DRAM Banks

- Multiple Banks are key to high throughput
- As one DRAM bank is recovering from read operation, next bank is being accessed
- Essentially on-chip memory interleaving
- Goal is to hide latency and bitline precharge time (recovery time)
  - Latency is access to first byte, critical path through row-decode and word line assertion
  - Bitline Precharge time (recovery time to next access) depends on number of bits in a column (number of rows)

# One Bank DRAM

**CLK**

row access | col | col | col | precharge | row access

---

# Multi-Bank DRAM

**CLK**

**Bank #1** row access | col | col | col | precharge | row access | col

**Bank #2** row access | col | col | col | precharge

Number of banks required to hire all row latency and precharge time depends on ratio of latency+precharge to column access time.
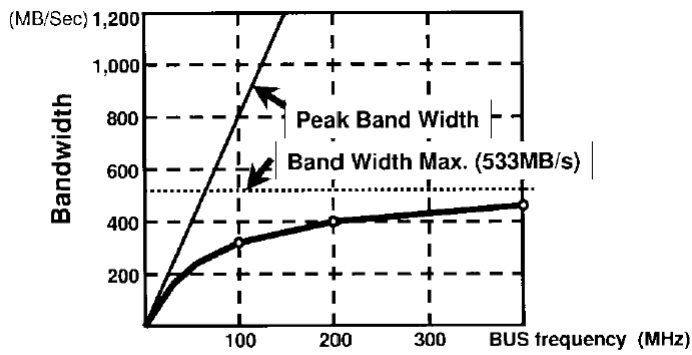
Fig. 10.   One-bank system bandwidth calculation.

Fig. 11.   Multibank system bandwidth calculations.

# Past Multi-bank



DRAM Banks

SRAM Cache

Mem Controller

---

# Architectural Issues in SRAMs

- Clock pin added to SRAMs to get Synchronous SRAMs for burst mode/pipelined burst capability
  - DDR-SRAMs available (clk = 300 Mhz)
  - Flowthru mode on SSRAMS to allow output in same clock as address to minimize latency
- Dual-port, Multi-port SRAMs are a big market
  - Multiprocessor systems
  - Telecommunications (networking hardware)
  - Any application with different speeds on 2 ports

# Dual Port SRAMS

- Separate left/right ports
- Independent read/write operation of each port for accesses to different locations or simultaneous read access to same location
- Asynchronous Dual Port
  - Need asynchronous arbitration circuit to determine 'winning' port in case of simultaneous write to same location – block losing port
- Synchronous Dual Port – simultaneous write to same location is undefined operation (results not guaranteed).