## 2.1 Buffer Memory Technology

- Memory Blocks On-Chip
  - On-chip SRAM area , power consumption, access rate
- Power Consumption for chip-to-chip communication
- Memory Chips (commercially available)
  - Chip periphery interface: communication standards to memory chips and their off-chip throughput
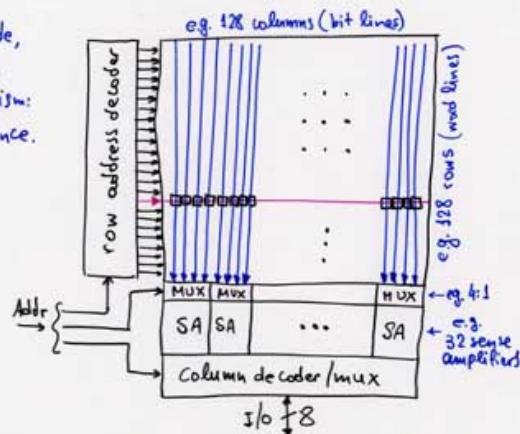  - DRAM chips, internal banks, Bank Interleaving

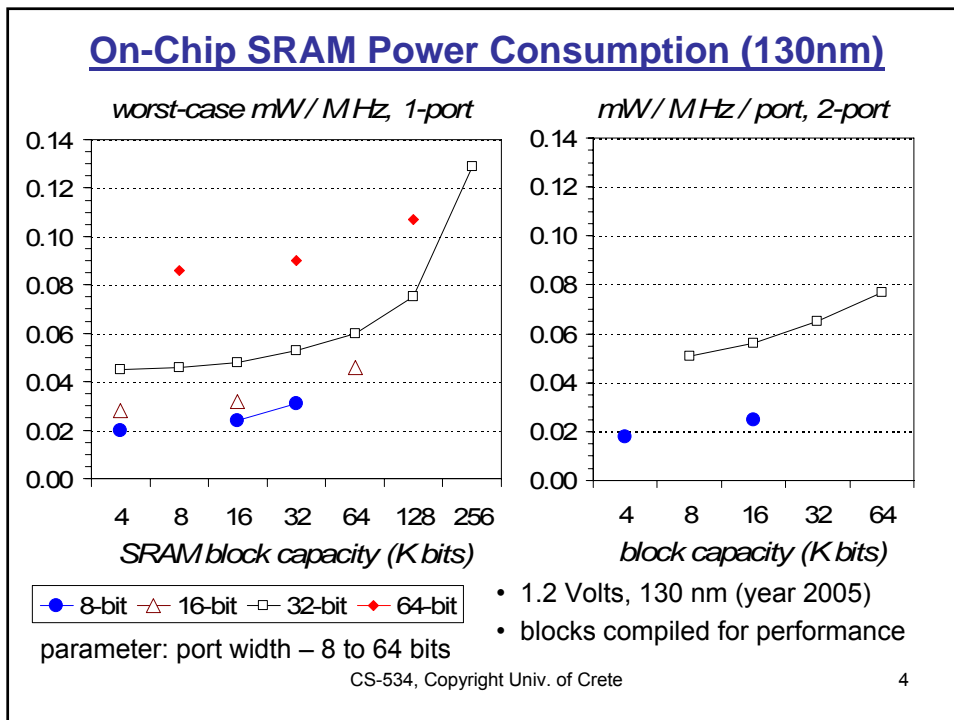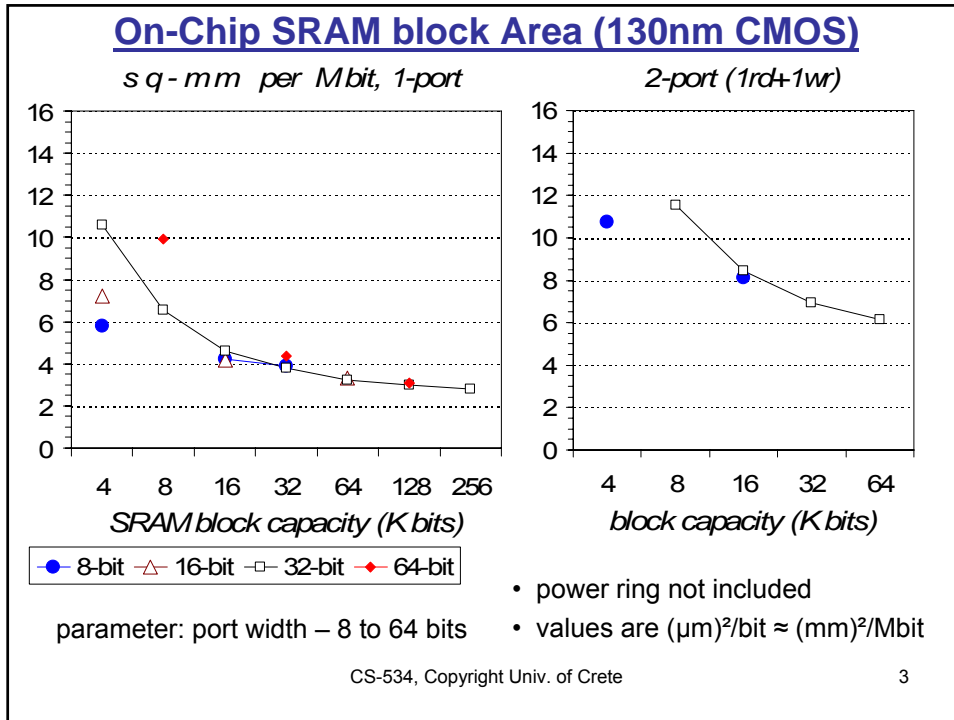CS-534, Copyright Univ. of Crete     1

---



## On-Chip SRAM

Memory blocks inherently provide, on-chip, very high throughputs, owing to their inherent parallelism: an entire row is accessed at once. This high throughput is available on-chip, due to the feasibility of very wide datapaths, running at high clock rates.
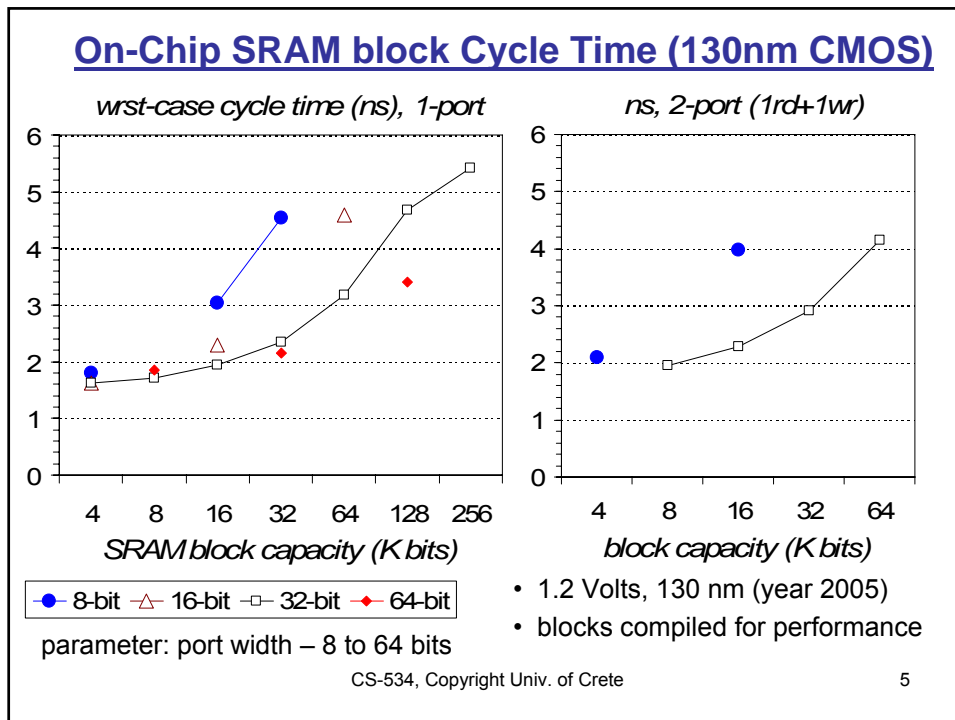
(Very wide or very large memories are made of several smaller memory blocks, to reduce capacitive loading on word lines and bit lines)

e.g. 128 columns (bit lines)

e.g. 128 rows (word lines)

MUX ← e.g. 4:1

SA ← e.g. 32 sense amplifiers

Column decoder/mux

I/o ⟂8

Example layout: 16 Kbit = 2k × 8

CS-534, Copyright Univ. of Crete     2

## On-Chip SRAM block Area (130nm CMOS)

*s q - m m  per M bit, 1-port*          *2-port (1rd+1wr)*

*SRAM block capacity (K bits)*     *block capacity (K bits)*

8-bit   16-bit   32-bit   64-bit

parameter: port width – 8 to 64 bits

• power ring not included
• values are (µm)²/bit ≈ (mm)²/Mbit

CS-534, Copyright Univ. of Crete          3

## On-Chip SRAM Power Consumption (130nm)

*worst-case mW / M Hz, 1-port*     *mW / M Hz / port, 2-port*

*SRAM block capacity (K bits)*     *block capacity (K bits)*

8-bit   16-bit   32-bit   64-bit

parameter: port width – 8 to 64 bits

• 1.2 Volts, 130 nm (year 2005)
• blocks compiled for performance

CS-534, Copyright Univ. of Crete          4

## On-Chip SRAM block Cycle Time (130nm CMOS)



*wrst-case cycle time (ns), 1-port*      *ns, 2-port (1rd+1wr)*

*SRAM block capacity (K bits)*      *block capacity (K bits)*

—●— 8-bit   —△— 16-bit   —□— 32-bit   —◆— 64-bit

parameter: port width – 8 to 64 bits

• 1.2 Volts, 130 nm (year 2005)
• blocks compiled for performance

CS-534, Copyright Univ. of Crete      5

---

## On-Chip SRAM block Cost, Performance

Area per Kbit:

• Area efficiency increases with block capacity: peripheral overhead (address decoders, column multiplexors, sense amplifiers) grows slower than core
• Port width costs significantly for small memories (more sense amp's, non-square aspect ratio)
• Two-port area ≈ 2 × one-port area
• Power ring: add 25 µm on each side of the block given in the above charts (width and heigth increase by 50 µm each)
• 1 sense amp / 8 col., usually
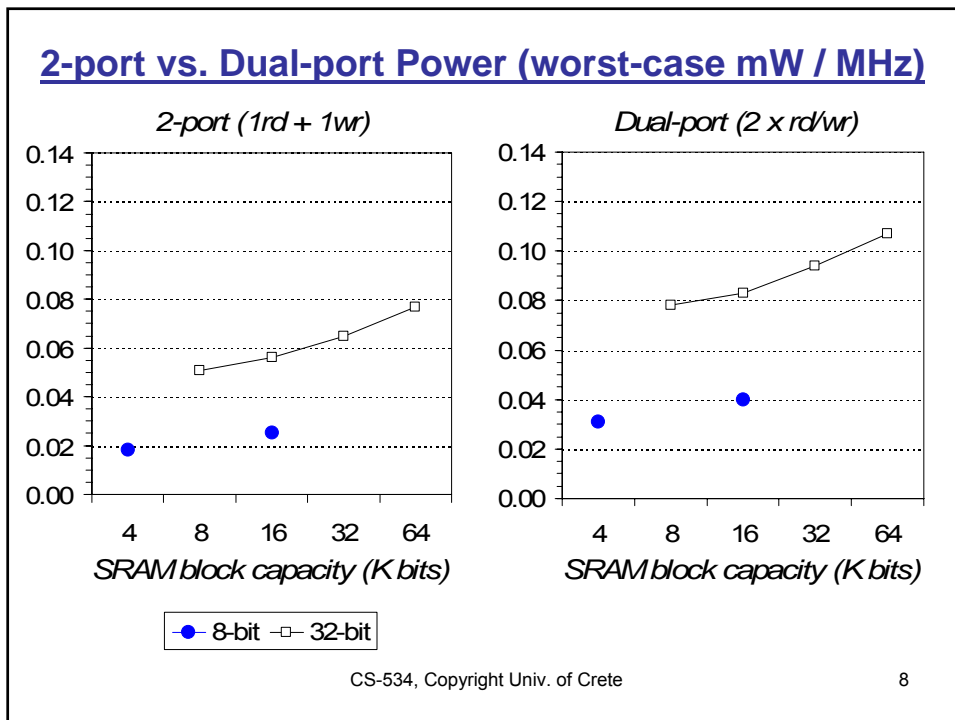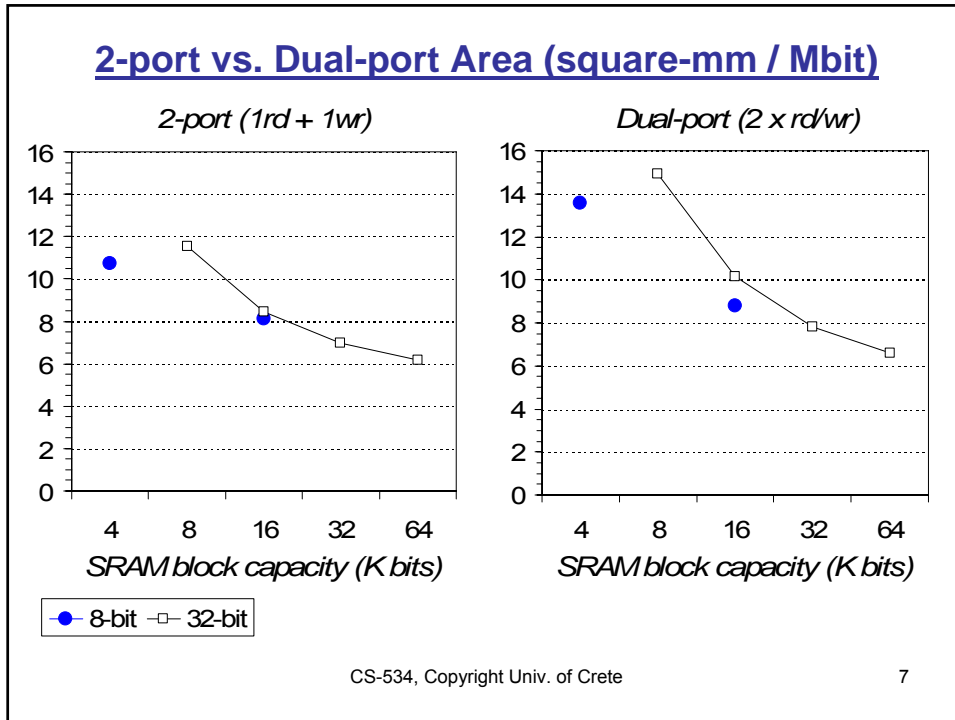• Quoted blocks have write-byte enable signals, except 8-bit ones

Power Consumption per MHz:

• Dominated by port-width for small mem's (sense amp. consumption)
• Dominated by block size for large mem's (word- & bit- line consum.)
• $P_{two-ports} \approx 2 \times P_{one-port}$

Access Rate (=1/cycle-time):

• Large blocks are quite slower than small ones, for sizes beyond the "knee" of the curve
• For large blocks, narrow ports reduce the speed, because of extra mux'es after sense amp's
• Two-port speed ≈ speed of 1-port block with twice the num. of bits

CS-534, Copyright Univ. of Crete      6

---

CS-534, 2.1: Memory Technology      3

## 2-port vs. Dual-port Area (square-mm / Mbit)

*2-port (1rd + 1wr)*    *Dual-port (2 x rd/wr)*

*SRAM block capacity (K bits)*    *SRAM block capacity (K bits)*

8-bit    32-bit

CS-534, Copyright Univ. of Crete    7

## 2-port vs. Dual-port Power (worst-case mW / MHz)

*2-port (1rd + 1wr)*    *Dual-port (2 x rd/wr)*

*SRAM block capacity (K bits)*    *SRAM block capacity (K bits)*

8-bit    32-bit

CS-534, Copyright Univ. of Crete    8

## 2-port vs. Dual-port Cycle Time (ns, worst-case)

*2-port (1rd + 1wr)*    *Dual-port (2 x rd/wr)*

SRAM block capacity (K bits)    SRAM block capacity (K bits)

● 8-bit  □ 32-bit

CS-534, Copyright Univ. of Crete                                    9

[intentionally left blank]

CS-534, Copyright Univ. of Crete                                    10

## On-Chip SRAM Buffer Example *(i):* 40-Byte wide

- Width = 1 min-size IP packet =
  = 40 Bytes = 320 bits = 5 blocks × 64 bits/block
- One-port, 2048 packets × 40 B = 80 KB = 640 Kb
- 130 nm CMOS, 1.2 Volts
- Area: 5 banks × 128 Kb/bank × 3 mm$^2$/Mb =
  = 0.64 Mb × 3 mm$^2$/Mb ≈ **2 mm$^2$**
- Throughput: 320 bits × 300 Macc/s ≈ **100 Gb/s**
- Power Consumption:
  5 banks × 0.11 mW/MHz × 300 MHz = **165 mW**

11

## On-Chip SRAM Buffer Example *(ii):* 256-Byte wide

- Width ≈ 1 average-size IP packet =
  = 256 Bytes = 2048 bits = 64 blocks × 32 bits/block
- Two-port (1rd+1wr), 2048 packets × 256 B = 512 KB = 4 Mb
- 130 nm CMOS, 1.2 Volts
- Area: 64 × 64 Kb × 6.1 mm$^2$/Mb = 4 M × 6.1 ≈ **25 mm$^2$**
- Throughput: 2 ports × 2048 b/port × 240 MHz ≈ **1 Tb/s**
  (500 Gb/s writes + 500 Gb/s reads)
- Power Consumption:
  64 banks × 2 ports × 0.08 mW/MHz × 240 MHz ≈ **2.4 W**
- Conclusion: "no problem" on-chip, except for small packets

12

## Power Consumption / Throughput: on-chip SRAM

- (1)  On-Chip Buffer Memories:

- 130 nm CMOS, "usual, medium" SRAM block sizes:
  - 1-port, ×16:  ≈ 0.03 mW/MHz = 0.03 mW / 16 Mbps ≈ 2.0 mW/Gbps
  - 1-port, ×32:  ≈ 0.05 mW/MHz = 0.05 mW / 32 Mbps ≈1.6 mW/Gbps
  - 1-port, ×64:  ≈ 0.10 mW/MHz = 0.10 mW / 64 Mbps ≈ 1.6 mW/Gbps
  - 2-port,  ×8:  ≈ 0.02 mW/MHz = 0.02 mW / 8 Mbps ≈ 2.5 mW/Gbps
  - 2-port, ×32:  ≈ 0.06 mW/MHz = 0.06 mW / 32 Mbps ≈ 2.0 mW/Gbps

- Conclusion:  1.5 to 2 mW / Gbps on-chip buffer memories

## Power Consumption / Throughput: Chip I/O

- (2)  Chip-to-Chip I/O Pin Power Consumption:

- both directions of a high-speed serial off-chip transceiver (without equalization –which consumes considerably)

- 130 nm CMOS:  10 to 25 mW / Gbps chip-to-chip comm

- copper cable power consumption is very small, by comparison

⇒ Chip-to-chip communication costs an order of magnitude more than on-chip buffering, in terms of power consumption

- Total chip power consumption (up to few tens of Watts) limits total chip throughput to about  1 Tbps/chip or less

Off-Chip Memory —or other networking/I/O chips:
How to Increase Chip-to-Chip Communication Throughput?

Old SRAM Read ("flow through"):

(1) Pipelined Reads (Synchronous, Registered Interface)

CS-534, Copyright Univ. of Crete    15



...Further Increasing the data pin throughput of chip-to-chip communication:

(2) DDR (Double Data Rate) Timing

Traditional Synchronous Intf.

DDR Interface!

Transmit and receive with a positive-edge-triggered register

Transmit with:
Receive with: two registers:
 • one positive-edge-tr. register
 • one negative-edge-tr. register

CS-534, Copyright Univ. of Crete    16

... further increasing the data pin throughput of chip-to-chip communication...

### (3) Source-Synchronous Data Clocking

when the clock frequency rises, the chip-to-chip (speed-of-light) delay becomes non-negligible wrt. pulse width

osc

ck1 — snd clock → CK2

Chip 1  snd data (addr)

Chip 2 (RAM or other)

ck1 domain

ck3 domain

return data

return clock

CK2

ck2

Synchronization — clock domain crossing

ck3 is a delayed version of ck1, i.e. has (exactly) the same frequency, but its delay (phase shift) may vary (slowly) with time...

CS-534, Copyright Univ. of Crete          17

---

## SRAM Data I/O Paths:

**Separate D (in) and Q (out) Paths:**        versus        **Shared "DQ" Data Bus:**

time-mux'ed write & read addresses

Addr

D → in write data → SRAM → out read data → Q

all addresses

Addr

SRAM

DQ

all data

⊖: data path underutilization when inbalanced ( ≠ 50% − 50% ) read/write transactions

⊖: bus turn-around overhead: data bus underutilization when frequently switching between read & write transactions

CS-534, Copyright Univ. of Crete          18

---

modern SRAM chip technology w. separate D(in) & Q(out) paths:

"QDR" (Quad Data Rate) SRAM

Other Version:
"burst-of-4":
- addr. path is plain (NOT DDR)
- each addr. refers to 4 data words.

CS-534, Copyright Univ. of Crete    19



Example QDR SRAM (2001) Micron's MT54V512H18

9 Mbits = 512 K × 18 bits

Clock freq. up to 167 MHz

T ≥ 6ns   pulse, bit width ≥ 3ns

peak write throughput = 167 MHz × 2 × 18 bits = 6 Gb/s /chip

peak read throughput = 167 MHz × 2 × 18 bits = 6 Gb/s /chip

Peak total throughput, when fully balanced 50-50 reads/writes } = 6 + 6 = 12 Gb/s /chip

2.5 Volt power supply; Power Consumption ≈ 1 Watt @ 167 MHz

⇒ power per throughput = $\frac{1 W}{12 Gbps}$ ≈ 0.08 $\frac{Watt}{Gbps}$

CS-534, Copyright Univ. of Crete    20

## Shared "DQ" Data Bus Timing:

### Naive Timing:



ck

R/W#  write / read   i d l e !!  write

Addr  A1  A'2  i d l e !!  A3

write / read / write

DQ  D1  D2  D3

Underutilization on every
read-to-write transition

### "ZBT" (Zero Bus Turn-around) Timing:

ck

R/W#  write / read  write  read  read

Addr  A1  A2  A3  A4  A5

write / read / write / read

DQ  D1  D2  D3

D1 has not yet been written at M[A1]
when reading from M[A2] starts...
...need to bypass mem. when A2==A1

CS-534, Copyright Univ. of Crete                          21

---

Example Shared Bus SRAM at the top current performance (2001):

Micron's MT57 V256 H36   **DDR SRAM**

9 Mbits = 256K × 36 bits

Clock frequ. up to 300 MHz (!) ⇒

$T \geq 3.3 \, ns$, bit pulse width $\geq 1.6 \, ns$

Burst-of-4 accesses only
(one address every 2 clock cycles)

Peak Throughput = 300 MHz × 2 × 36 b = 21.6 Gb/s

Throughput with alternating read/writes $= \frac{2}{3} \times peak = 14.4 \frac{Gbps}{chip}$

2.5 Volts Power Supply; Consumption = 1.6 W

⇒ ~ 0.1 $\frac{Watt}{Gbps}$

Although the ZBT concept is used,
due to the high clock frequency
and the unavoidable bus
turn around overhead (multiple
drivers on the same wire, each
using its own clock
(source-synchronous
timing)), 1 to 2 clock
cycles (= 2 to 4 word burst)
are lost on every
read-to-write
transition.

CS-534, Copyright Univ. of Crete                          22

---

## DRAM Basics: Row Address, Column Address, Precharge



CS-534, Copyright Univ. of Crete                    23

## Fast DRAM Example (2001)
Micron MT46 V2 M32
### DDR SDRAM
(Synchronous DRAM)
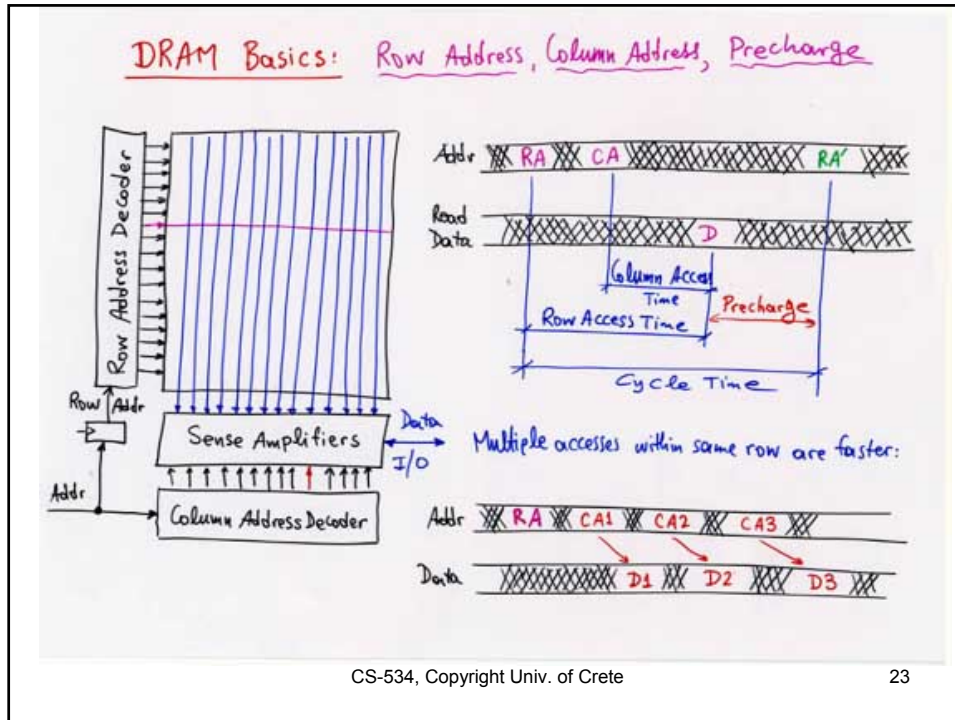
- **200 MHz** max. clock frequency
- 64 Mbits = $\underline{2M \times 32\ bits}$ = 
  = 512k $\times$ 32b $\times$ $\underline{4\ Banks}$

- 32-bit (shared DQ) databus, DDR timing ⇒ ⇒ 2 words × 32 bits each per clock cycle peak databus throughput

- ≈1 Watt at peak access rate, using one bank only, 2.5 Volt. (No number given for multibank op.)

- Row Address - to - Column Address: ----------- $t_{RCD} \geqslant 20ns$ (@200MHz: 4~)
- Column Address - to - Read Data (CAS latency): --- $CL \geqslant 15ns$ (@200MHz: 3~)
- Write Recovery Time (write data to precharge): ... $t_{WR} \geqslant$ ----------- 2~
- Precharge Time: - - - - - - - - - - - - - - $t_{RP} \geqslant 20ns$ (@200MHz: 4~)
- Cycle Time (same bank): - - - - - - - - - - $t_{RC} \geqslant 60ns$ (@200MHz: 12~)
- Bank - to - Bank Activation (other bank Row-to-Row): $t_{RRD}$ - - - - - - 2~
- Read - to - Write bus turn-around lost cycles: ----------- 3~
- Write - to - Read same bank lost cycles (write recovery time): ........ 2~
- Write - to - Read other bank lost cycles: - - - - - - - - - - - - - - - - - - $\emptyset$~

CS-534, Copyright Univ. of Crete                    24

Single-Bank Read Access

ACT = Activate
Ba = Bank #a
Ra = Row # Ra Address

RD = Read (the predefined burst size)
Ba = from the active Row within Bank #a
Ca = at Column Address #Ca

$D_a$ = $i^{th}$ word of burst from Ba, Ra, Ca

CS-534, Copyright Univ. of Crete                    25



Single-Bank Write Access

ACT = Activate
Ba = Bank #a
Ra = Row Address Ra

WR = Write (the predefined burst size)
Ba = into the active Row of Bank #a
Ca = at Column Address Ca

$D_a$ = $i^{th}$ word of burst destined to Ba, Ra, Ca

CS-534, Copyright Univ. of Crete                    26

Multiple Accesses to Different Columns in the same Row of a Bank

All transactions shown are to the same bank #a, and to the same activated row Ra in that bank.
The transactions shown are:
- Read from column Ca → a1, a2
- Read from column Cb → b1, b2
- Write c1, c2 at column Cc
- Write d1, d2 at column Cd
- Read from column Ce → e1, e2

CS-534, Copyright Univ. of Crete          27



**Multi-Bank Operation:  Memory Interleavin**

- burst length set to 8;  each successive READ command interrupts the preceding burst, resulting in net bursts of 6.

CS-534, Copyright Univ. of Crete          28