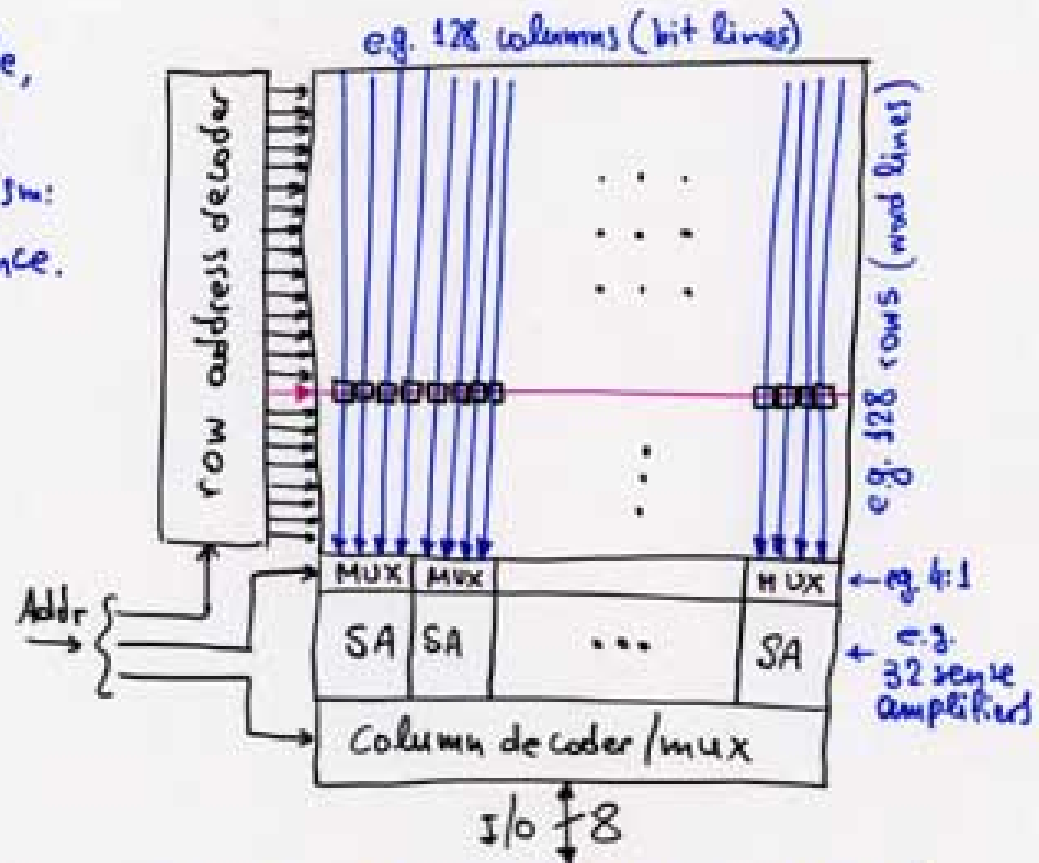# 2.1 Buffer Memory Technology

- ## 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

# 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.
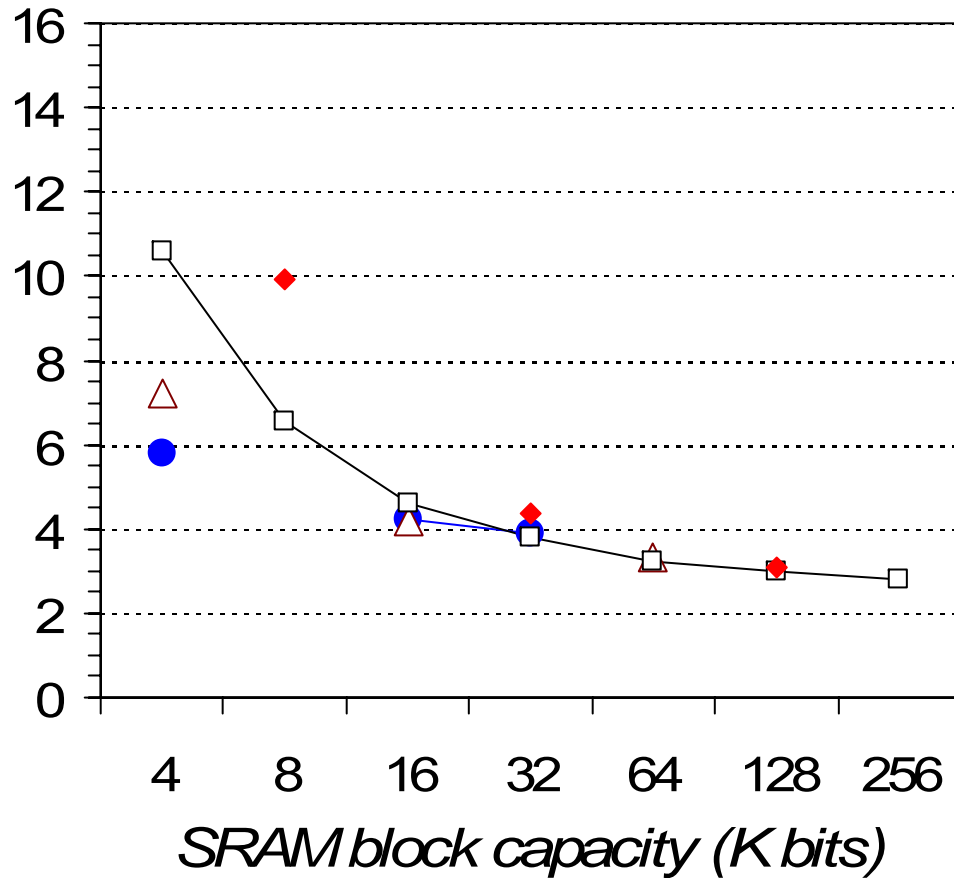
(Very wide × very large memories are made of several smaller memory blocks, to reduce capacitive loading on word lines and bit lines)
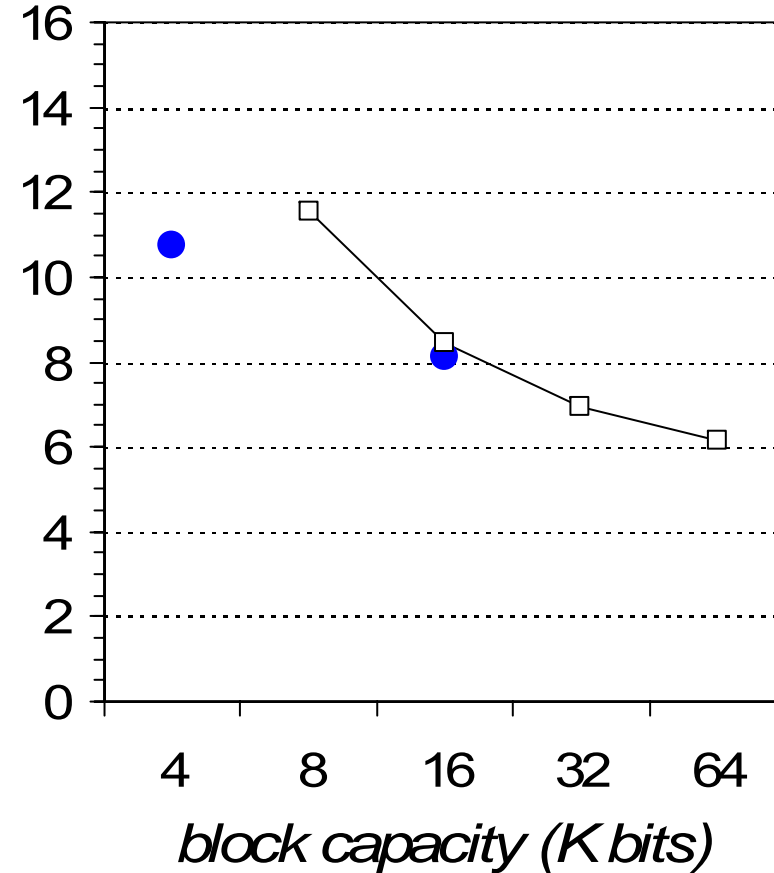
e.g. 128 columns (bit lines)

row address decoder

e.g. 128 rows (word lines)

Addr

MUX  MUX ... MUX  ← e.g. 4:1

SA  SA ... SA  ← e.g. 32 sense amplifiers

Column decoder/mux

I/O ↕ 8

Example layout: 16 Kbit = 2K × 8

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



*sq-mm per Mbit, 1-port*

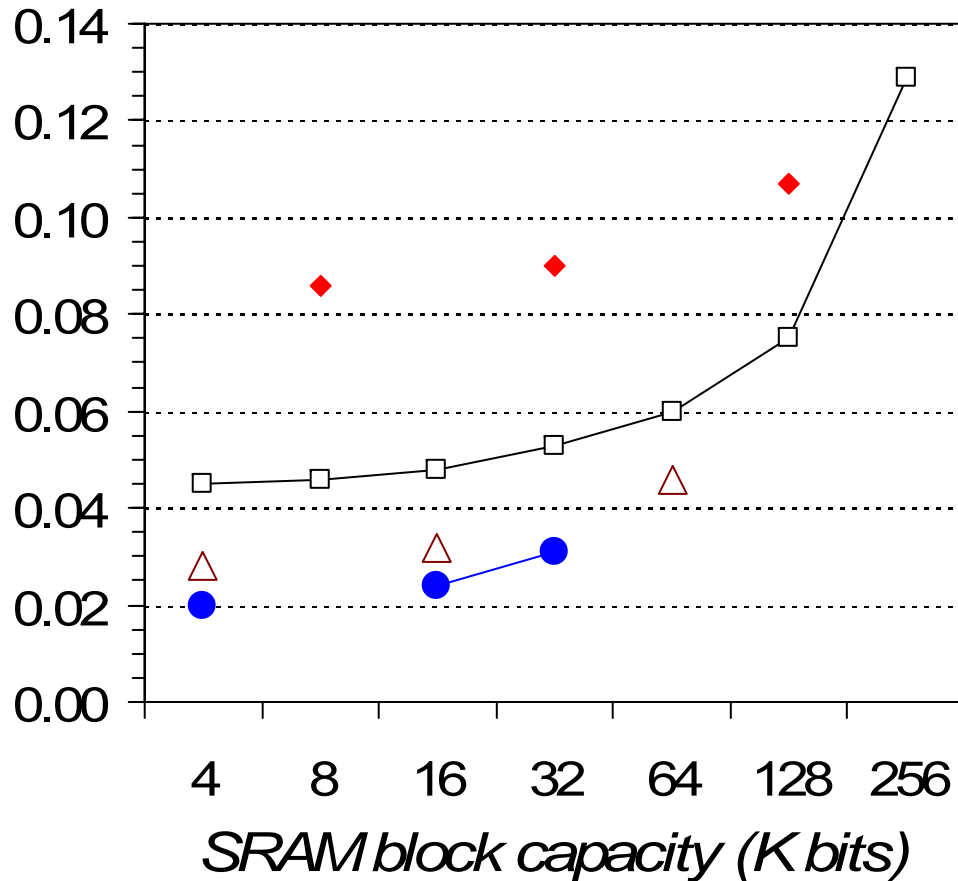*2-port (1rd+1wr)*

*SRAM block capacity (K bits)*

*block capacity (K bits)*

Legend: ● 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

# On-Chip SRAM Power Consumption (130nm)



*worst-case mW / MHz, 1-port*

*mW / MHz / port, 2-port*

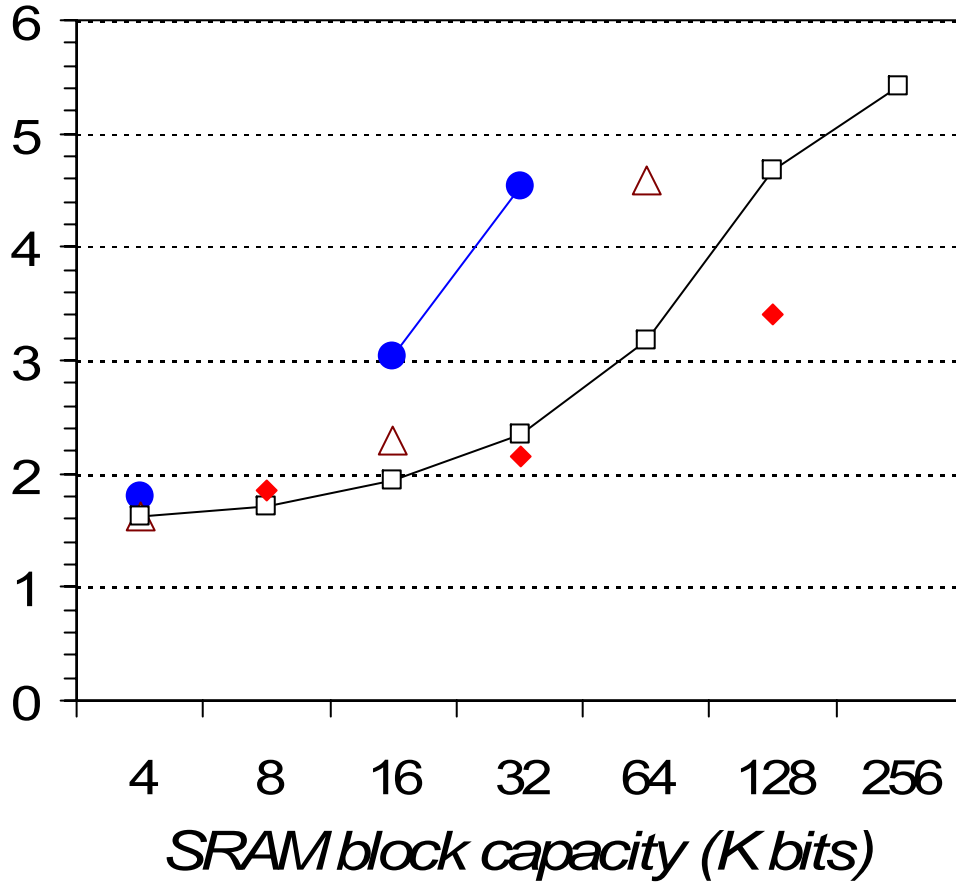*SRAM block capacity (K bits)*

*block capacity (K bits)*

● 8-bit  △ 16-bit  □ 32-bit  ◆ 64-bit

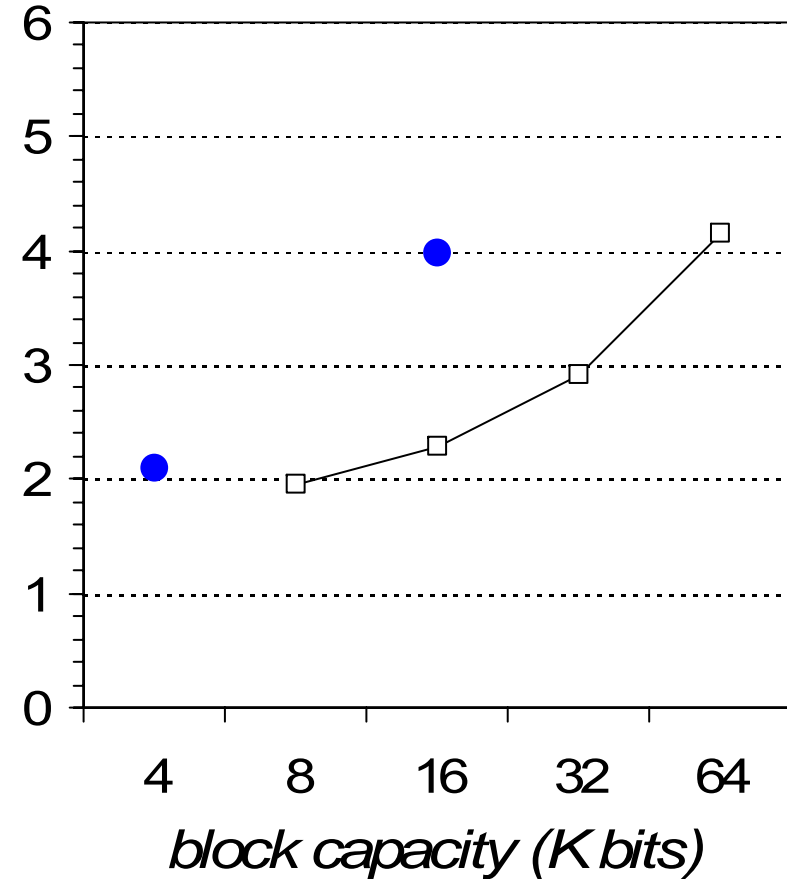parameter: port width – 8 to 64 bits

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

# 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)*

Legend: ●―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

# 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
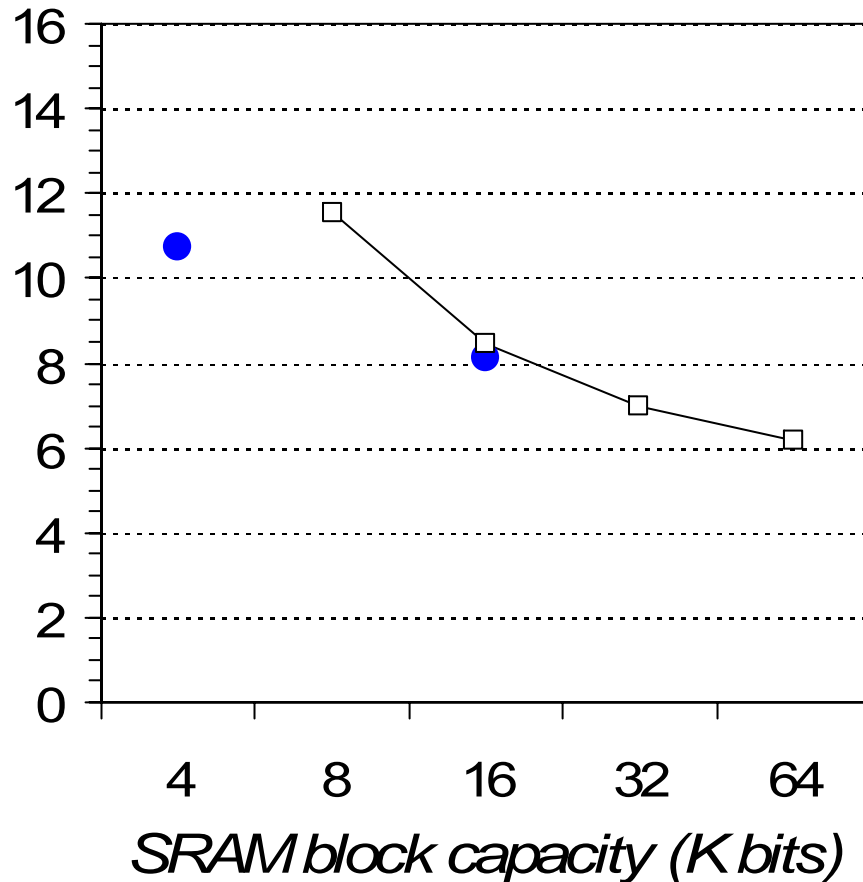
## Power Consumption per MHz:

- Dominated by port-width for small mem's (sense amp. consumption)

- Dominated by block size for large mem's (word- & bit- line consum.)

- $P_{\text{two-ports}} \approx 2 \times P_{\text{one-port}}$
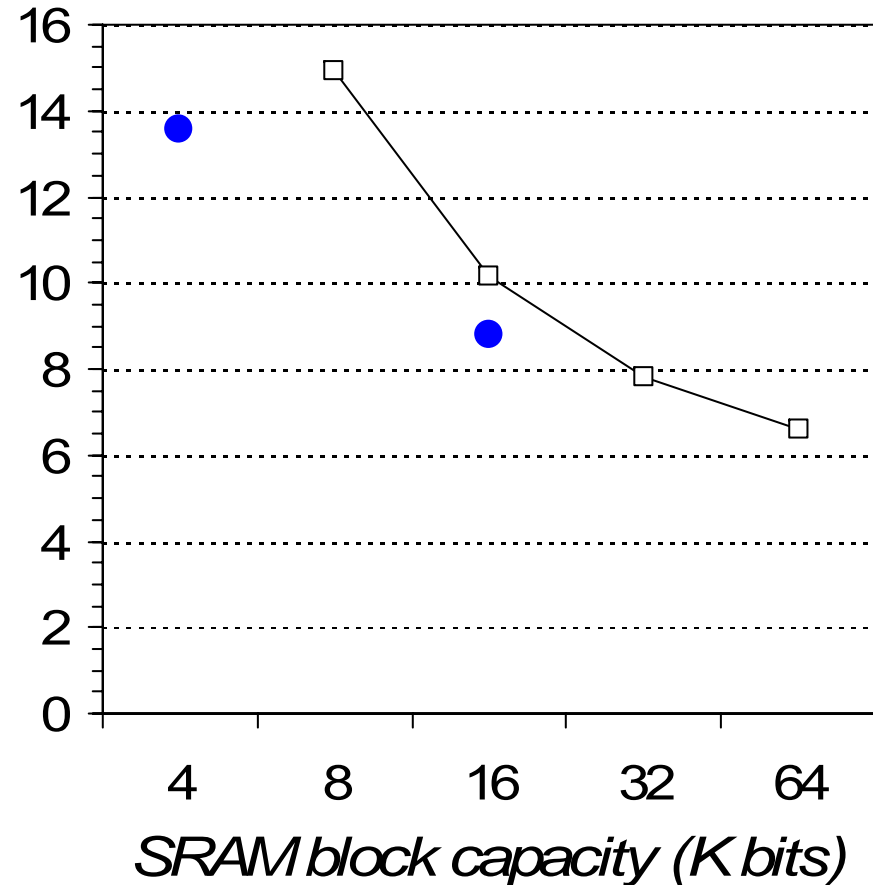
## Access Rate (=1/cycle-time):

- 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

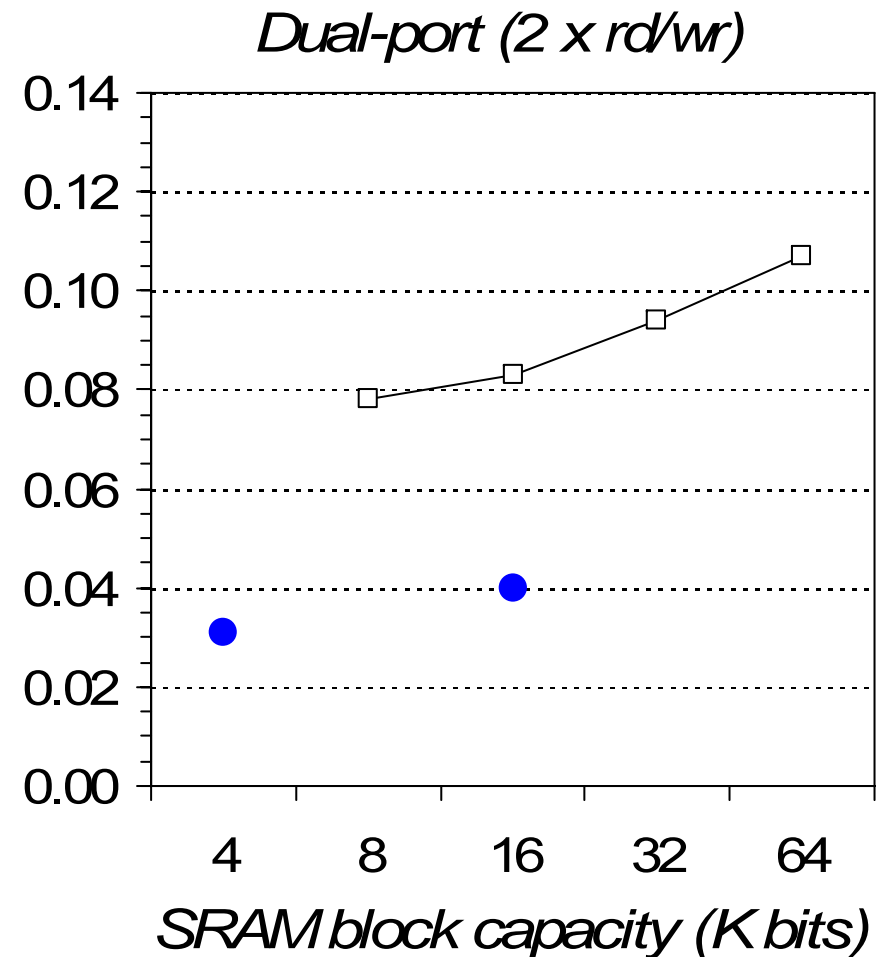# 2-port vs. Dual-port Area (square-mm / Mbit)
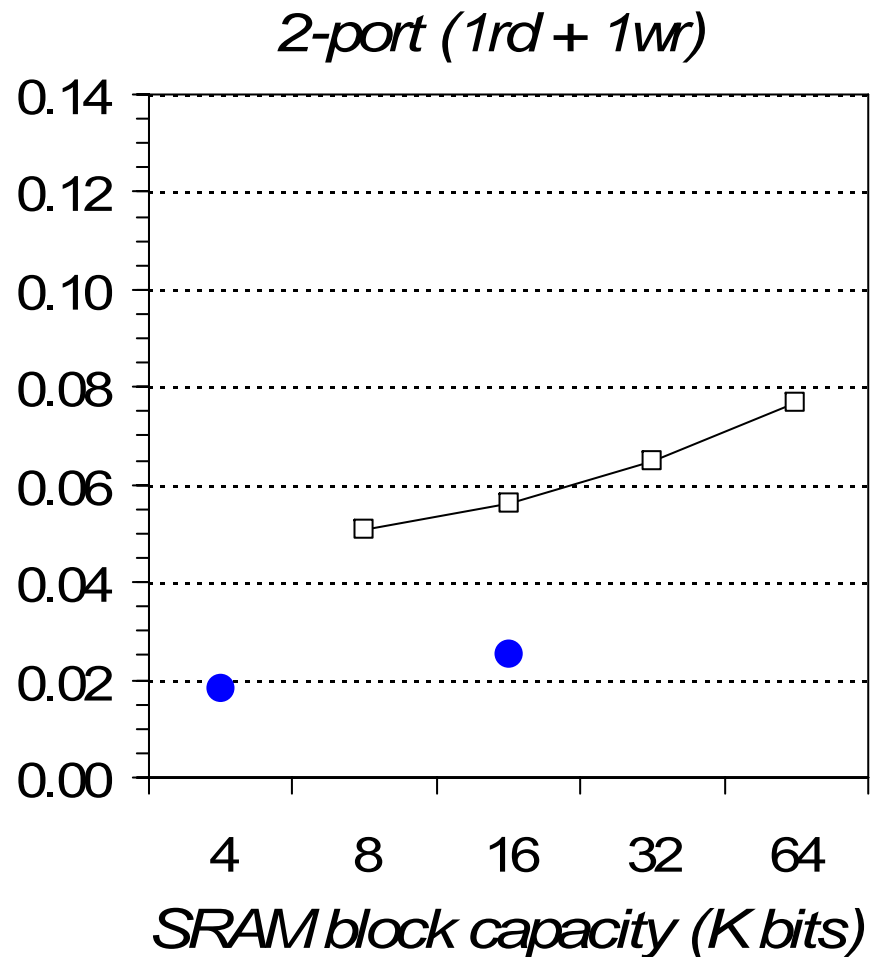
*2-port (1rd + 1wr)*



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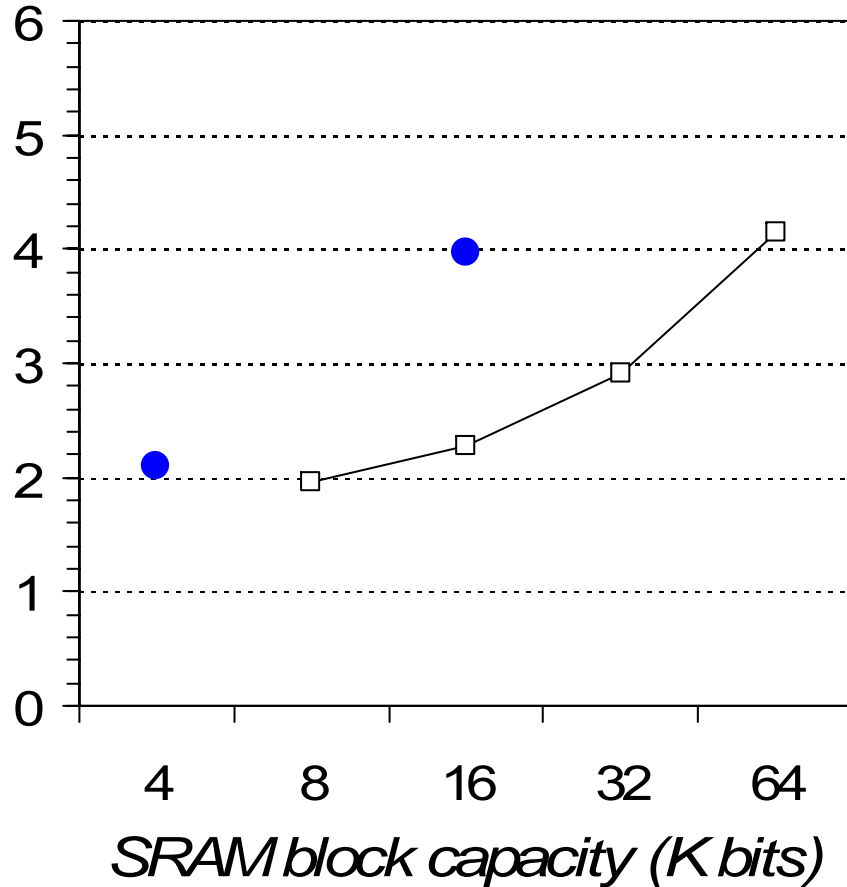


SRAM block capacity (K bits)

SRAM block capacity (K bits)

● 8-bit  –□– 32-bit

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

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

[intentionally left blank]

# 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**

# 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
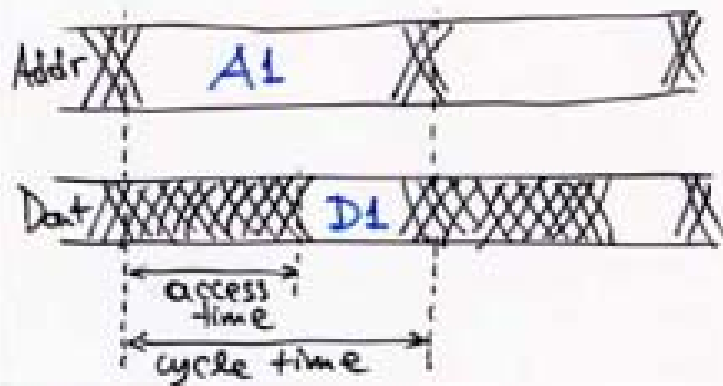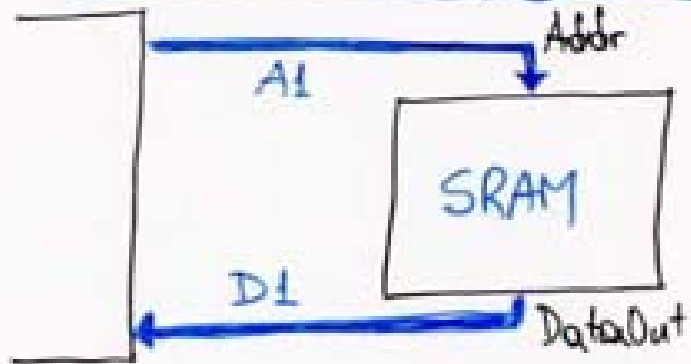
# Power Consumption / Throughput: on-chip SRAM

- (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

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

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

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

- copper cable power consumption is very small, by comparison

$\Rightarrow$ 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  <u>1 Tbps/chip or less</u>

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)

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

## (2) DDR (Double Data Rate) Timing

### Traditional Synchronous Intf!

$\approx T/2$    $\cong \frac{T}{2}$ ← minimum pulse width

ck

data

T

min. pulse width $\approx T$

Transmit and receive with a positive-edge-triggered register

### DDR Interface!

$T/2$

$T/2$

d1    ck

d2    data

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

... 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 w.r.t. pulse width

osc

Chip 1

ck1

send clock

send data (addr)

CK2

Chip2
(RAM
or
other)

ck1
domain

ck3
domain

return data

return clock

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...

# SRAM Data I/O Paths:

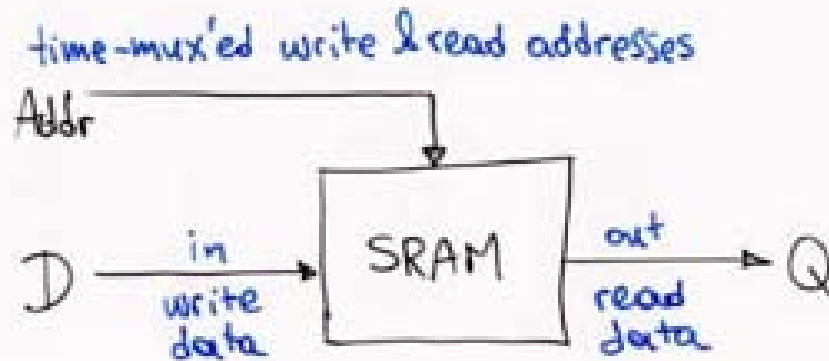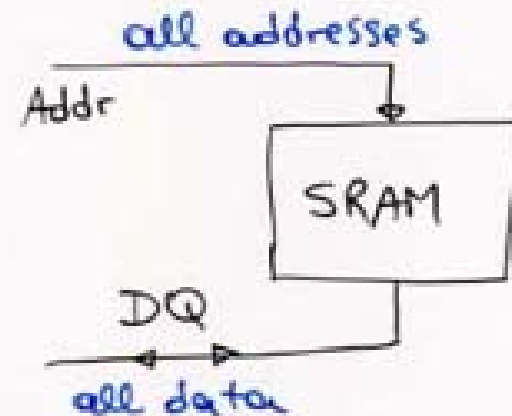## Separate D (in) and Q (out) Paths:                    Shared "DQ" Data Bus:

Versus

time-mux'ed write & read addresses

Addr ─────────────────┐                      all addresses
                       ▼                     Addr ─────────────┐
         ┌──────────────┐                                      ▼
         │              │  out                        ┌──────────────┐
 D ──in──▶    SRAM      ────read──▶ Q                  │    SRAM      │
   write │              │  data                        │              │
   data  └──────────────┘                              └──────────────┘
                                              DQ              │
                                             ◀────────▶───────┘
                                              all data

⊖: data path underutilization           ⊖: bus turn-around overhead:
   when inbalanced                          data bus underutilization
   ( ≠ 50% - 50% )                          when frequently switching
   read/write transactions                  between read & write
                                            transactions

modern SRAM chip technology w. separate D(in) & Q(out) paths:
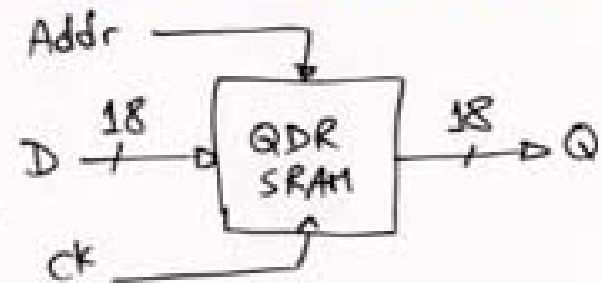
"QDR" (Quad Data Rate) SRAM

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

Example QDR SRAM (2001) Micron's MT54V512 H18

$9 \text{ Mbits} = \underline{512 \text{ k} \times 18 \text{ bits}}$

Clock freq. up to $\underline{167 \text{ MHz}}$

$T \geq 6ns$    pulse, bit width $\geq 3ns$

Addr ———

$D \xrightarrow{18}$ —□ | QDR SRAM | $\xrightarrow{18}$ —▷ Q

ck ———

peak write throughput = $167 \text{ MHz} \times \overset{\text{DDR}}{2} \times 18 \text{ bits} = 6 \text{ Gb/s /chip}$

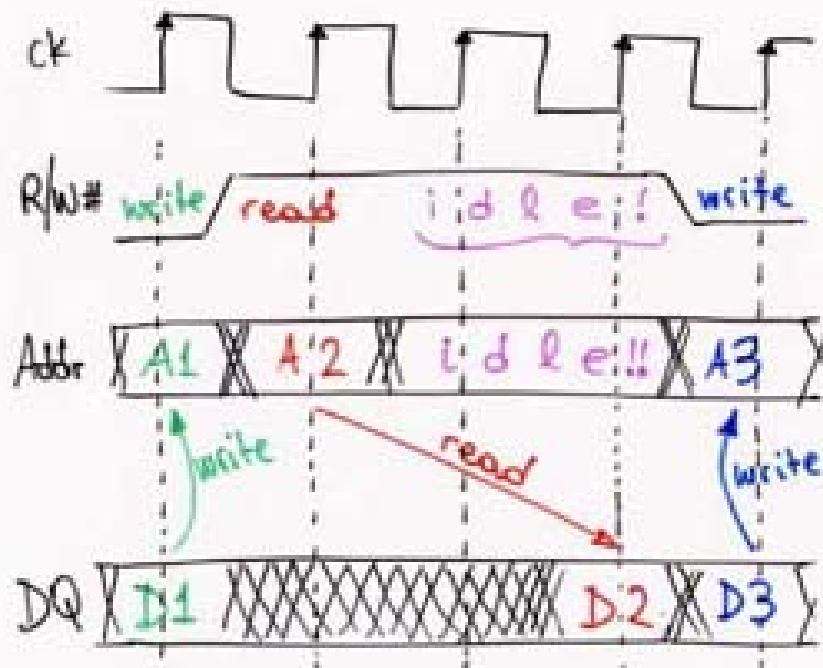peak read throughput = $167 \text{ MHz} \times 2 \times 18 \text{ bits} = 6 \text{ Gb/s /chip}$

Peak total throughput, when
fully balanced 50-50 reads/writes $\Big\} = 6 + 6 = \underline{12 \text{ Gb/s /chip}}$

$\underline{2.5 \text{ Volt}}$ power supply;   Power Consumption $\cong \underline{1 \text{ Watt}}$ @167 MHz

$\Rightarrow$ power per throughput = $\dfrac{1 \text{ W}}{12 \text{ Gbps}} \cong 0.08 \dfrac{\text{Watt}}{\text{Gbps}}$
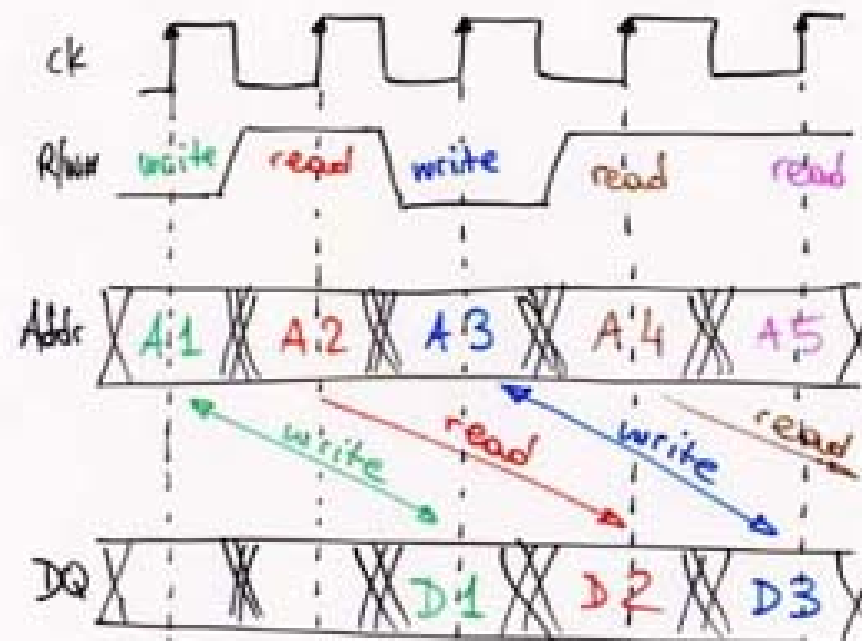
# Shared "DQ" Data Bus Timing:

## Naive Timing:



Underutilization on every
read-to-write transition

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



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

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 (!) → 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.

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

Burst-of-4 accesses only

(one address every 2 clock cycles)

$$\text{Peak Throughput} = 300 \text{ MHz} \times 2 \times 36 \text{ b} = \underline{21.6 \text{ Gb/s}}$$

DDR↑

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

2.5 Volts Power Supply; Consumption = $\underline{1.6 \text{ W}}$

$$\Rightarrow \sim 0.1 \frac{\text{Watt}}{\text{Gbps}}$$

DRAM Basics: Row Address, Column Address, Precharge

# Fast DRAM Example (2001)

Micron MT46 V2 M32

## DDR SDRAM
(Synchronous DRAM)

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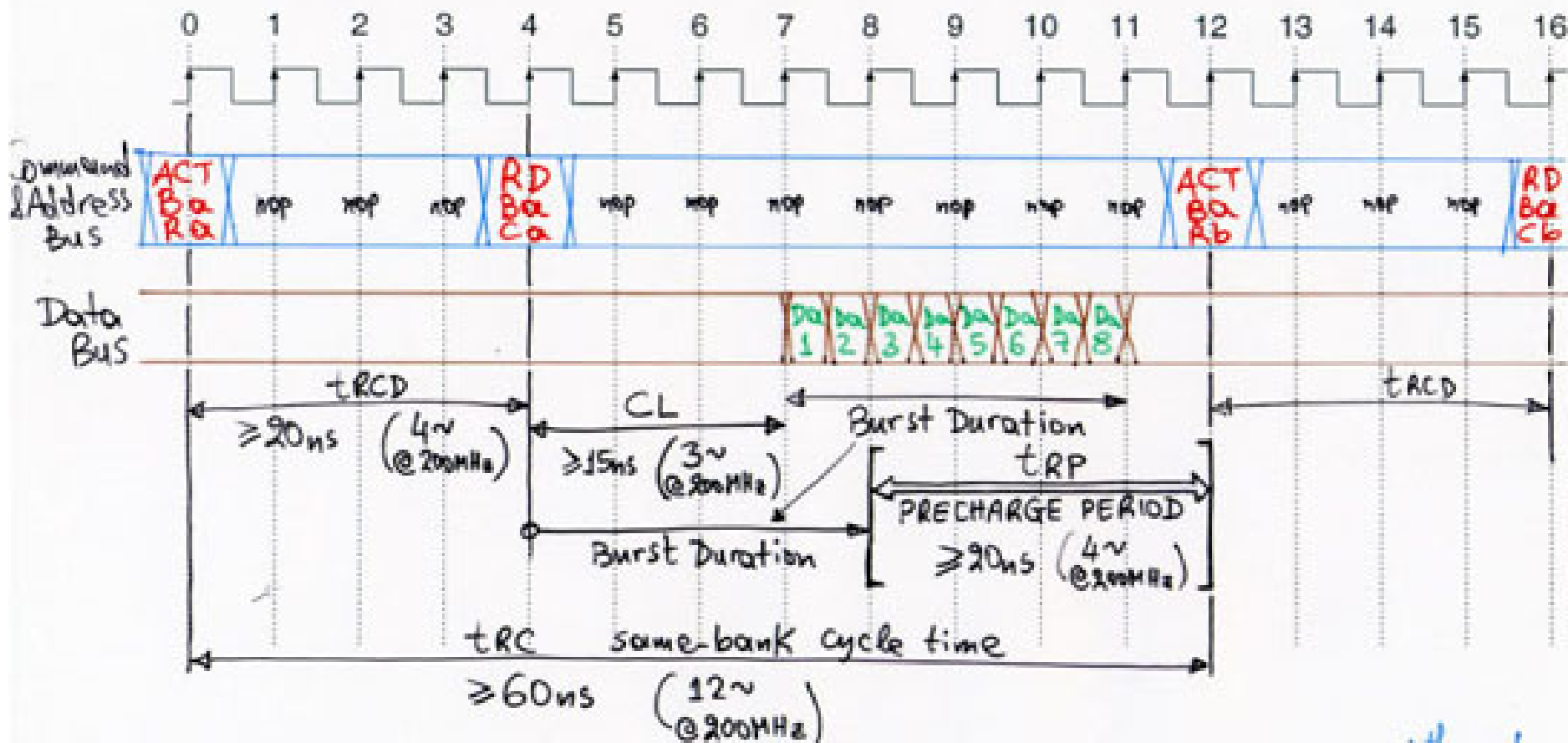  ⟹ 2 words × 32 bits each per clock cycle
  peak databus throughput

- 200 MHz max. clock frequency
- 64 Mbits = $2M \times 32$ bits =
  = $512k \times 32b \times 4$ Banks

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

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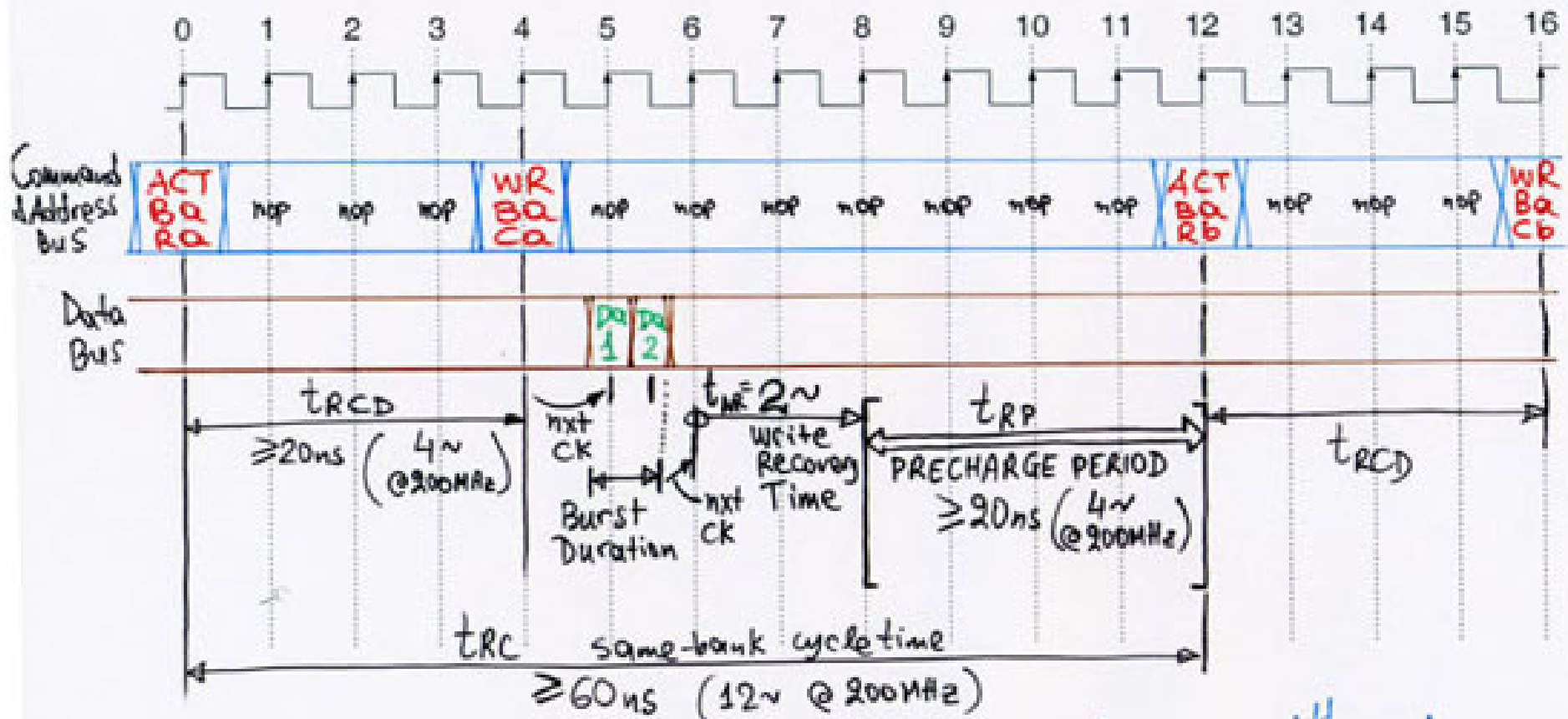
# Single-Bank Read Access



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

Da = i'th word of burst from Ba, Ra, Ca
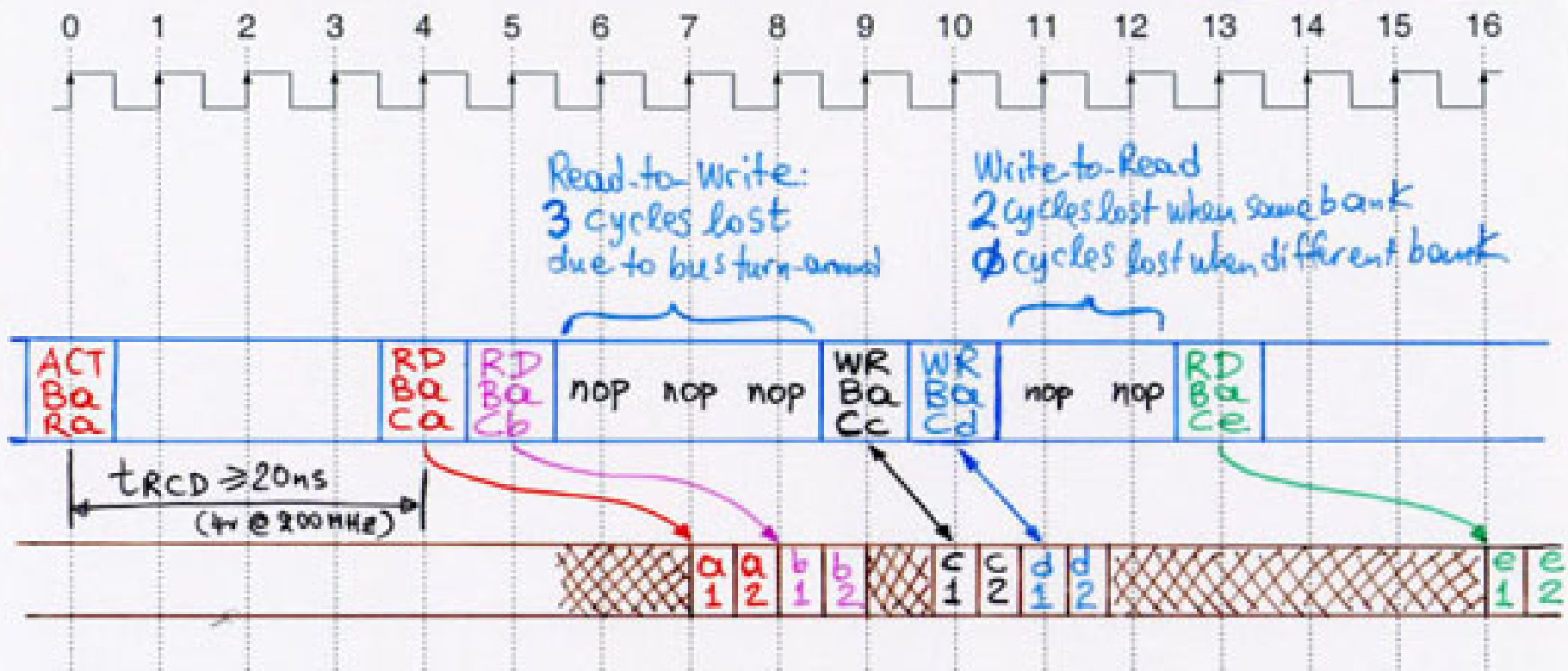
Single-Bank Write Access

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

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

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



Read-to-Write: 3 cycles lost due to bus turn-around

Write-to-Read: 2 cycles lost when same bank; 0 cycles lost when different bank

ACT Ba Ra | RD Ba Ca | RD Ba Cb | nop nop nop | WR Ba Cc | WR Ba Cd | nop nop | RD Ba Ce

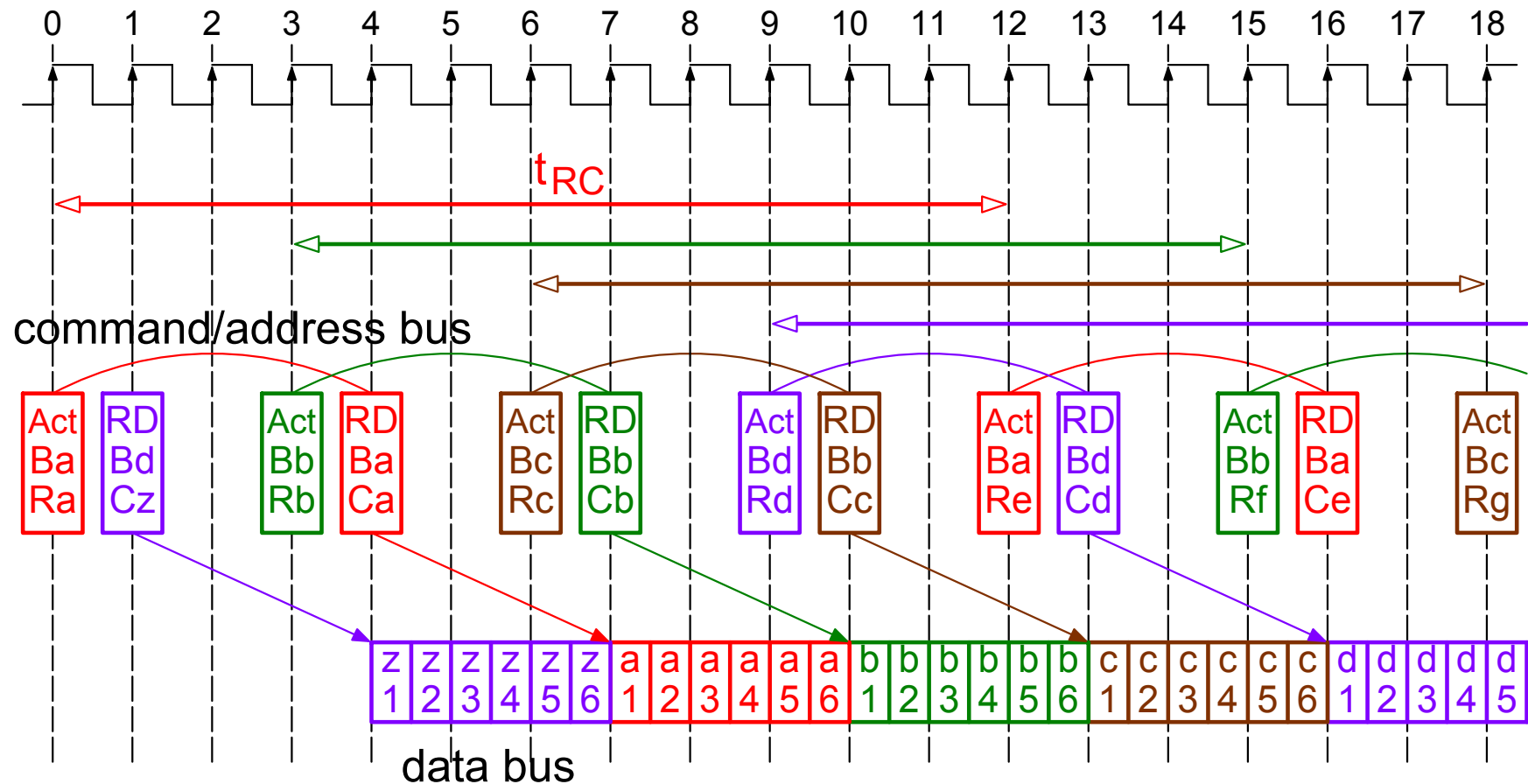$t_{RCD} \geq 20ns$ (4 @ 200 MHz)

a1 a2 b1 b2 | c1 c2 d1 d2 | e1 e2

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

# Multi-Bank Operation:  Memory Interleaving



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