# Memory System Design

## Chapter 16
## S. Dandamudi

# Outline

- Introduction
- A simple memory block
  - ∗ Memory design with D flip flops
  - ∗ Problems with the design
- Techniques to connect to a bus
  - ∗ Using multiplexers
  - ∗ Using open collector outputs
  - ∗ Using tri-state buffers
- Building a memory block

- Building larger memories
- Mapping memory
  - ∗ Full mapping
  - ∗ Partial mapping
- Alignment of data
- Interleaved memories
  - ∗ Synchronized access organization
  - ∗ Independent access organization
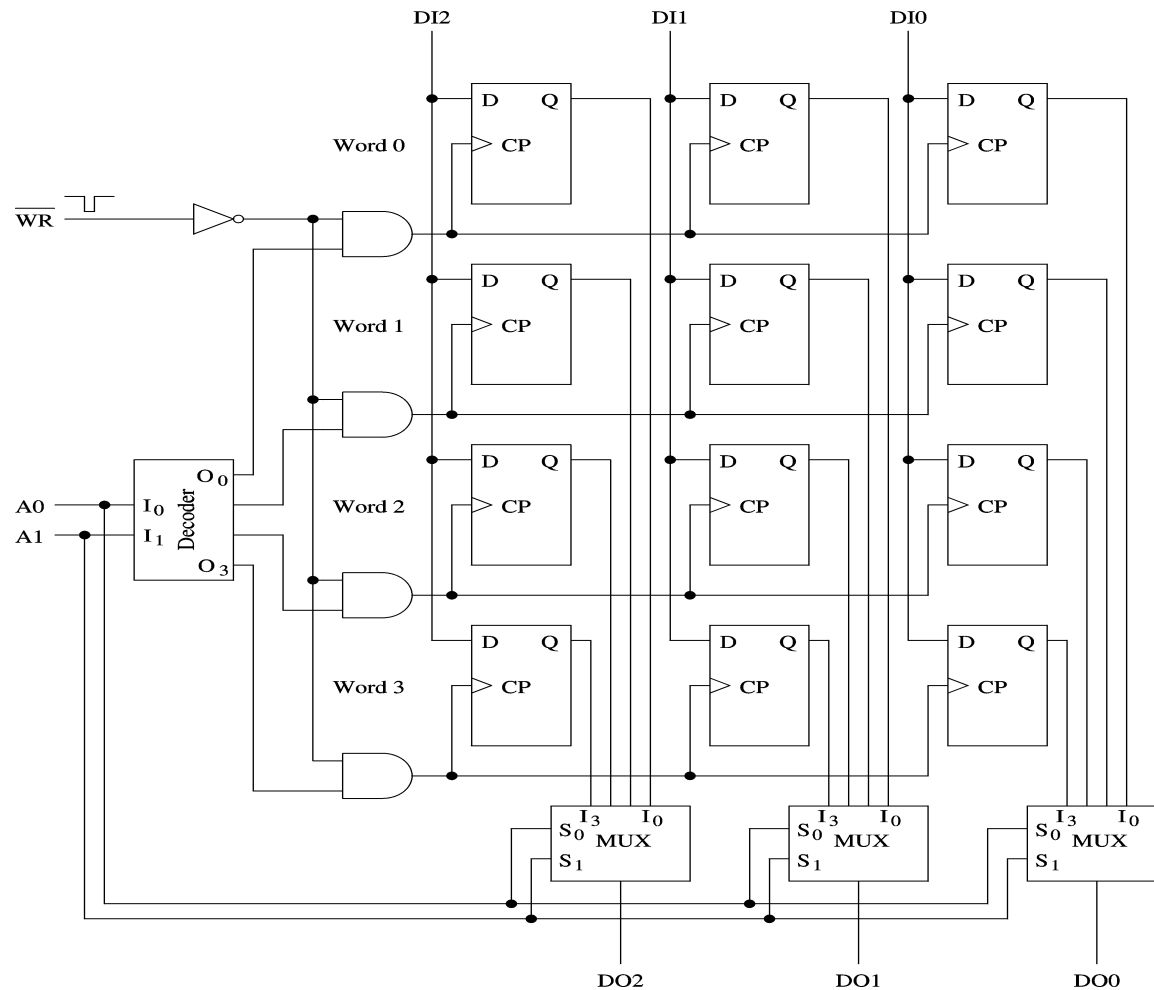  - ∗ Number of banks

# Introduction

- To store a single bit, we can use

    ∗ Flip flops or latches

- Larger memories can be built by

    ∗ Using a 2D array of these 1-bit devices

        » "Horizontal" expansion to increase word size

        » "Vertical" expansion to increase number of words

- Dynamic RAMs use a tiny capacitor to store a bit

- Design concepts are mostly independent of the actual technique used to store a bit of data

# Memory Design with D Flip Flops

- An example
  - ∗ 4X3 memory design
  - ∗ Uses 12 D flip flops in a 2D array
  - ∗ Uses a 2-to-4 decoder to select a row (i.e. a word)
  - ∗ Multiplexers are used to gate the appropriate output
  - ∗ A single WRITE (WR) is used to serve as a write and read signal
    - – zero is used to indicate write operation
    - – one is used for read operation
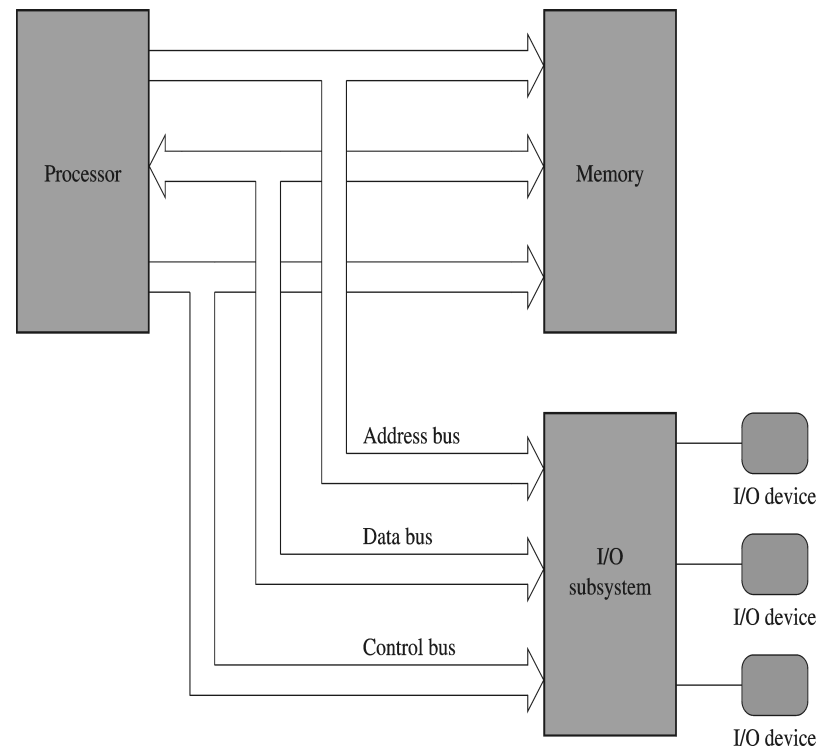  - ∗ Two address lines are needed to select one of four words of 3 bits each

# Memory Design with D Flip Flops (cont'd)

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.

# Memory Design with D Flip Flops (cont'd)

- Problems with the design
  - ∗ Requires separate data in and out lines
    - » Cannot use the bidirectional data bus
  - ∗ Cannot use this design as a building block to design larger memories
    - » To do this, we need a chip select input

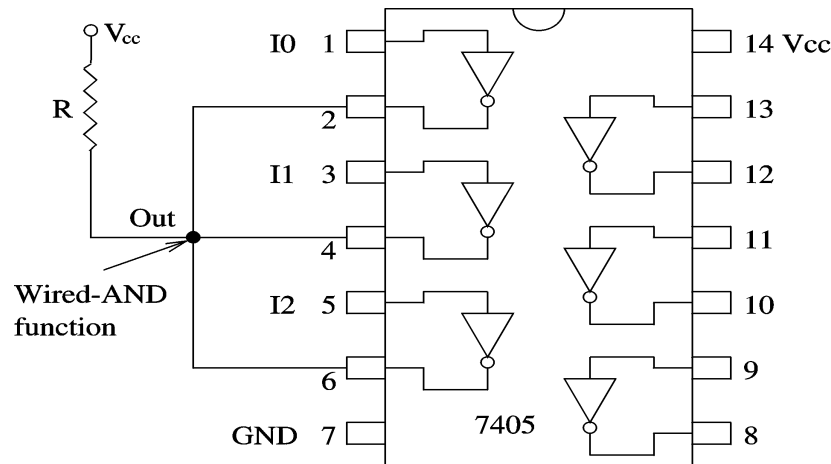- We need techniques to connect multiple devices to a bus



Processor

Memory

Address bus

Data bus

Control bus

I/O subsystem

I/O device

I/O device

I/O device

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.
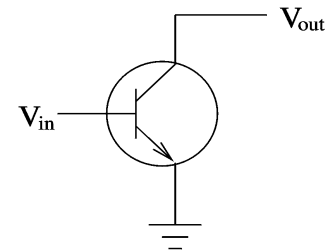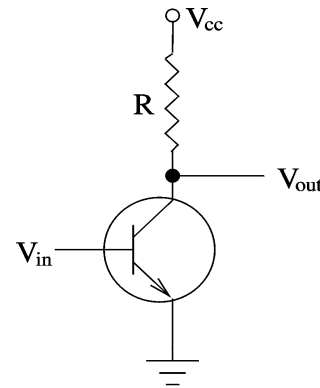
# Techniques to Connect to a Bus

- Three techniques
  - ∗ Use multiplexers
    - » Example
      - – We used multiplexers in the last memory design
    - » We cannot use MUXs to support bidirectional buses
  - ∗ Use open collector outputs
    - » Special devices that facilitate connection of several outputs *together*
  - ∗ Use tri-state buffers
    - » Most commonly used devices

# Techniques to Connect to a Bus (cont'd)

Open collector technique

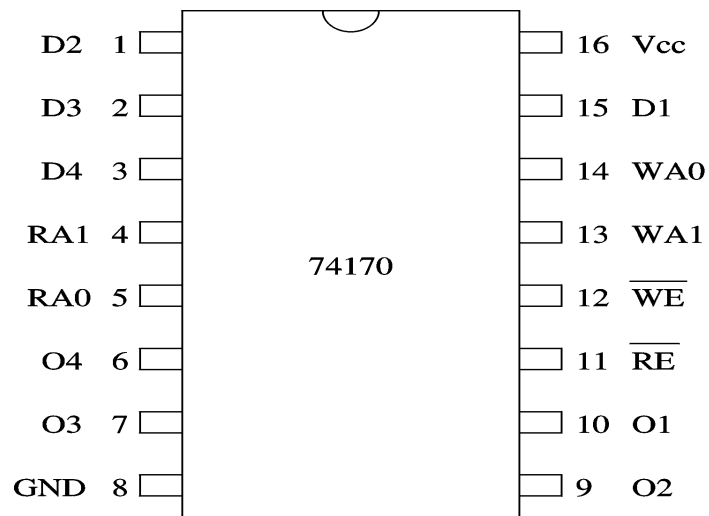| I2 | I1 | I0 | Out |
|----|----|----|-----|
| 0  | 0  | 0  | 1   |
| 0  | 0  | 1  | 0   |
| 0  | 1  | 0  | 0   |
| 0  | 1  | 1  | 0   |
| 1  | 0  | 0  | 0   |
| 1  | 0  | 1  | 0   |
| 1  | 1  | 0  | 0   |
| 1  | 1  | 1  | 0   |

(a) Connection diagram

(b) Truth table

# Techniques to Connect to a Bus (cont'd)

| D2 | 1 | | 16 | Vcc |
| D3 | 2 | | 15 | D1 |
| D4 | 3 | 74170 | 14 | WA0 |
| RA1 | 4 | | 13 | WA1 |
| RA0 | 5 | | 12 | $\overline{WE}$ |
| O4 | 6 | | 11 | $\overline{RE}$ |
| O3 | 7 | | 10 | O1 |
| GND | 8 | | 9 | O2 |

(a) Connection diagram

Open collector register chip

$\overline{WE}$  D1 D2 D3 D4
WA0
WA1
RA0
RA1
$\overline{RE}$  O1 O2 O3 O4

(b) Logic symbol

| $\overline{WE}$ | WA1 | WA0 | D inputs to |
|---|---|---|---|
| 0 | 0 | 0 | Word 0 |
| 0 | 0 | 1 | Word 1 |
| 0 | 1 | 0 | Word 2 |
| 0 | 1 | 1 | Word 3 |
| 1 | X | X | None |

(c) Write function table

| $\overline{RE}$ | RA1 | RA0 | Output from |
|---|---|---|---|
| 0 | 0 | 0 | Word 0 |
| 0 | 0 | 1 | Word 1 |
| 0 | 1 | 0 | Word 2 |
| 0 | 1 | 1 | Word 3 |
| 1 | X | X | None (Z) |

(d) Read function table

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.

# Techniques to Connect to a Bus (cont'd)

## Tri-State Buffers

DI ——•⟋• —— DO

(b) E = 0

| Inputs | | Output |
|---|---|---|
| E | DI | DO |
| 1 | 0 | 0 |
| 1 | 1 | 1 |
| 0 | X | Z |

E

DI ——▷—— DO

(a)

DI ——•—•• —— DO

(c) E = 1

(d)

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.

# Techniques to Connect to a Bus (cont'd)

Two example tri-state buffer chips



74367

74368

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.

# Techniques to Connect to a Bus (cont'd)

$\overline{OE}$  1 — 74373 — 20  Vcc

O0  2 — — 19  O7

D0  3 — — 18  D7

D1  4 — — 17  D6

O1  5 — — 16  O6

O2  6 — 74373 — 15  O5

D2  7 — — 14  D5

D3  8 — — 13  D4

O3  9 — — 12  O4

GND  10 — — 11  LE

8-bit tri-state register

D0  D1  D2  D3  D4  D5  D6  D7

$\overline{OE}$

LE

O0  O1  O2  O3  O4  O5  O6  O7

(a) Connection diagram

(b) Logic symbol

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.

# Building a Memory Block

A 4 X 3 memory design using D flip-flops

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.

# Building a Memory Block (cont'd)

## Block diagram representation of a 4x3 memory

Address lines

A0

A1

$\overline{WR}$   4 X 3

memory

$\overline{RD}$

$\overline{CS}$

Write

Read

Chip select

D0

D1

D2

Bidirectional
data lines

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.

# Building Larger Memories

2 X 16 memory module using 74373 chips

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.

# Designing Larger Memories

- Issues involved
  - ∗ Selection of a memory chip
    - » Example: To design a 64M X 32 memory, we could use
      - – Eight 64M X 4 in 1 X 8 array (i.e., single row)
      - – Eight 32M X 8 in 2 X 4 array
      - – Eight 16M X 16 in 4 X 2 array

- Designing M X N memory with D X W chips
  - ∗ Number of chips = $M \cdot N / D \cdot W$
  - ∗ Number of rows = $M/D$
  - ∗ Number of columns = $N/W$

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.

# Designing Larger Memories (cont'd)

64M X 32
memory using
16M X 16 chips

A2 — A25

A0 — A23
16M X 16
$\overline{\text{CS}}$
D0 — D15

A0 — A23
16M X 16
$\overline{\text{CS}}$
D0 — D15

D16 — D31    D0 — D15

D0 — D31

A2 — A25

Address bus

Decoder

A26    I0    O0
A27    I1    O1
              O2
              O3

A0 — A23
16M X 16
$\overline{\text{CS}}$
D0 — D15

A0 — A23
16M X 16
$\overline{\text{CS}}$
D0 — D15

D16 — D31    D0 — D15

A0 — A31

D0 — D31

A2 — A25

A0 — A23
16M X 16
$\overline{\text{CS}}$
D0 — D15

A0 — A23
16M X 16
$\overline{\text{CS}}$
D0 — D15

D16 — D31    D0 — D15

D0 — D31

D0 — D31

A2 — A25

A0 — A23
16M X 16
$\overline{\text{CS}}$
D0 — D15

A0 — A23
16M X 16
$\overline{\text{CS}}$
D0 — D15

D16 — D31    D0 — D15

D0 — D31

Data bus    D0 — D31

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.

# Designing Larger Memories (cont'd)

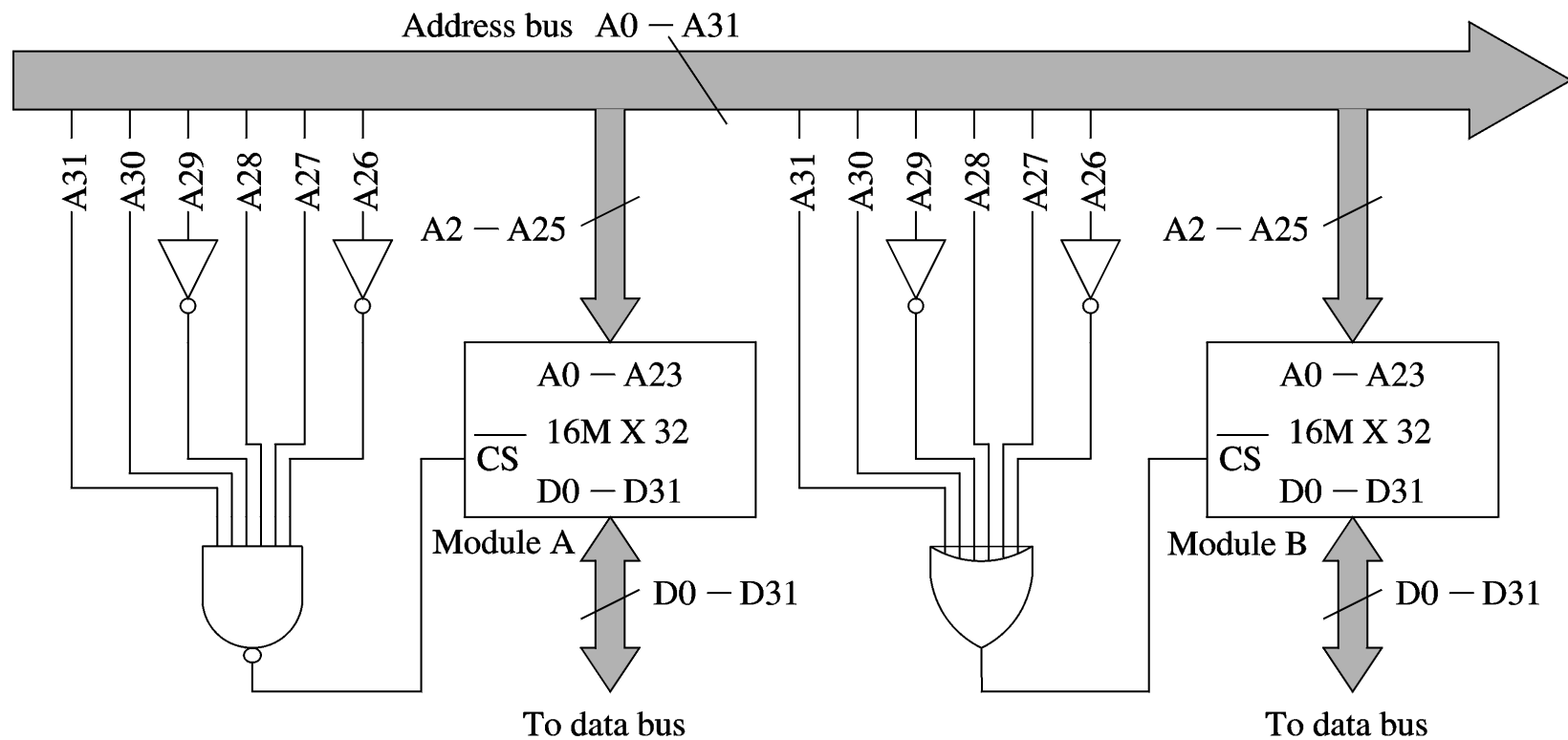- Design is simplified by partitioning the address lines (M X N memory with D X W memory chips)
  - ∗ Z bits are not connected ($Z = \log_2(N/8)$)
  - ∗ Y bits are connected to all chips ($Y = \log_2 D$)
  - ∗ X remaining bits are used to map the memory block
    - » Used to generate chip selects

MSB                           Address lines                        LSB

X                       Y                       Z

# Memory Mapping

## Full mapping

Address bus   A0 — A31

A31  A30  A29  A28  A27  A26

A2 — A25

A0 — A23

16M X 32

$\overline{CS}$

D0 — D31

Module A

D0 — D31

To data bus

A31  A30  A29  A28  A27  A26

A2 — A25

A0 — A23

16M X 32

$\overline{CS}$

D0 — D31

Module B

D0 — D31

To data bus

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.

# Memory Mapping (cont'd)

## Partial mapping

Address bus  A0 − A31

A31  A30  A29  A28  A27  A26

A2 − A25

A31  A30  A29  A28  A27  A26

A2 − A25

A0 − A23

16M X 32

$\overline{CS}$

D0 − D31

Module A

D0 − D31

A0 − A23

16M X 32

$\overline{CS}$

D0 − D31

Module B

D0 − D31

To data bus

To data bus

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.

# Alignment of Data



D24 — D31

3  7  11 15 19 23  byte address

D16 — D23

2  6  10 14 18 22  byte address

D8 — D15

1  5  9  13 17 21  byte address

D0 — D7

Data bus
D0 — D31

0  4  8  12 16 20  byte address

CPU

MEMORY

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.

# Alignment of Data (cont'd)

- Alignment
  - ∗ 2-byte data: Even address
    - » Rightmost address bit should be zero
  - ∗ 4-byte data: Address that is multiple of 4
    - » Rightmost 2 bits should be zero
  - ∗ 8-byte data: Address that is multiple of 8
    - » Rightmost 3 bits should be zero
  - ∗ Soft alignment
    - » Can handle aligned as well as unaligned data
  - ∗ Hard alignment
    - » Handles only aligned data (enforces alignment)

# Interleaved Memory

- In our memory designs
  - ∗ Block of contiguous memory addresses is mapped to a module
    - » One advantage
      - – Incremental expansion
    - » Disadvantage
      - – Successive accesses take more time
        - ⇥ Not possible to hide memory latency

- Interleaved memories
  - ∗ Improve access performance
    - » Allow overlapped memory access
    - » Use multiple banks and access all banks simultaneously
      - – Addresses are spread over banks
        - ⇥ Not mapped to a single memory module
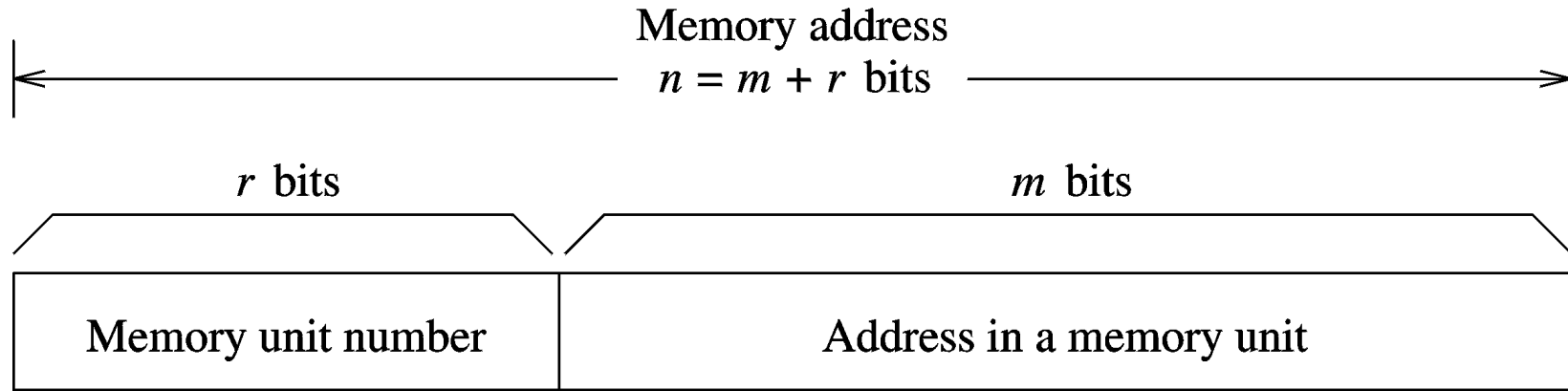
# Interleaved Memory (cont'd)

- The **n**-bit address is divided into two **r** and **m** bits:
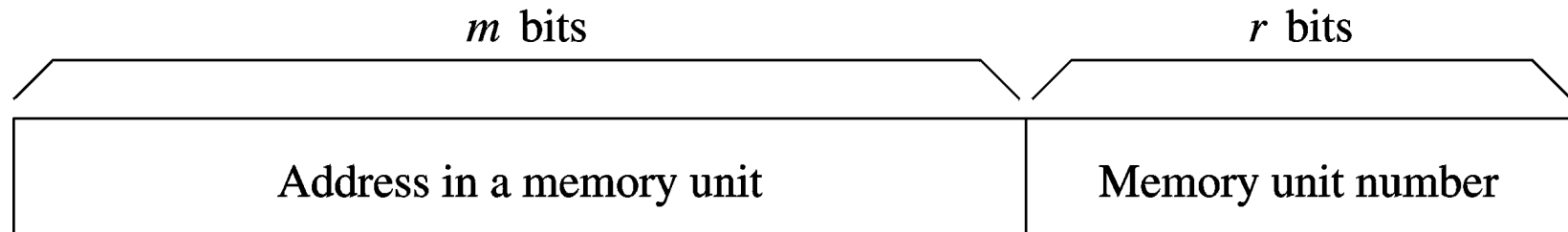
$$n = r + m$$

- Normal memory
  - ∗ Higher order **r** bits identify a module
  - ∗ Lower order **m** bits identify a location in the module
    - » Called high-order interleaving

- Interleaved memory
  - ∗ Lower order **r** bits identify a module
  - ∗ Higher order **m** bits identify a location in the module
    - » Called low-order interleaving
  - ∗ Memory modules are referred to as memory banks

# Interleaved Memory (cont'd)

Memory address
$n = m + r$ bits

| $r$ bits | $m$ bits |
|---|---|
| Memory unit number | Address in a memory unit |

(a) Normal memory address mapping

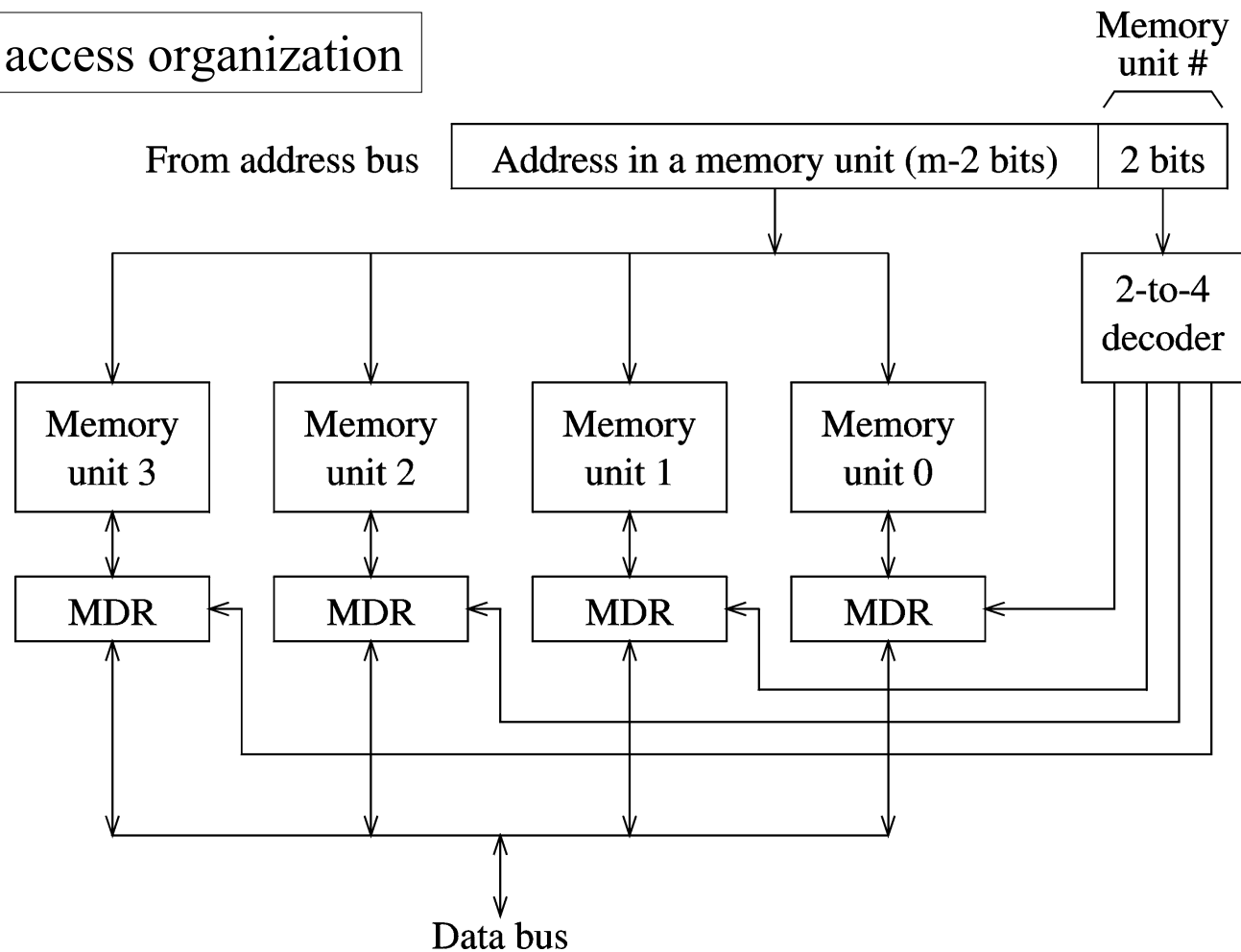| $m$ bits | $r$ bits |
|---|---|
| Address in a memory unit | Memory unit number |

(b) Interleaved memory address mapping
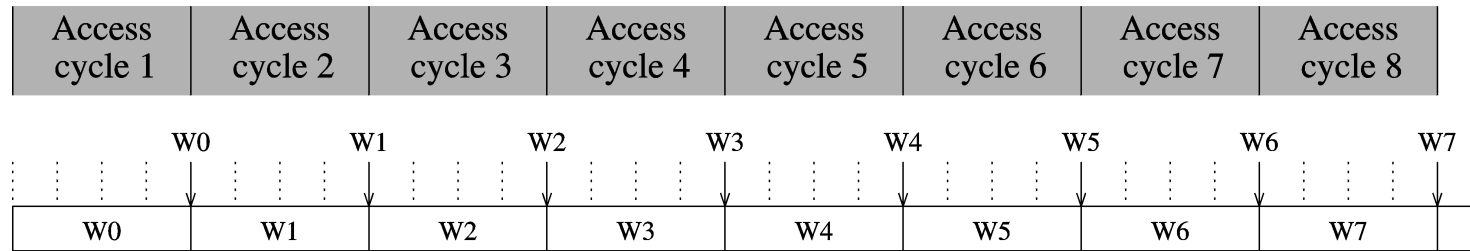
# Interleaved Memory (cont'd)

- Two possible implementations
  - ∗ Synchronized access organization
    - » Upper *m* bits are presented to all banks simultaneously
    - » Data are latched into output registers (MDR)
    - » During the data transfer, next *m* bits are presented to initiate the next cycle
  - ∗ Independent access organization
    - » Synchronized design does not efficiently support access to non-sequential access patterns
    - » Allows pipelined access even for arbitrary addresses
    - » Each memory bank has a memory address register (MAR)
      - – No need for MDR

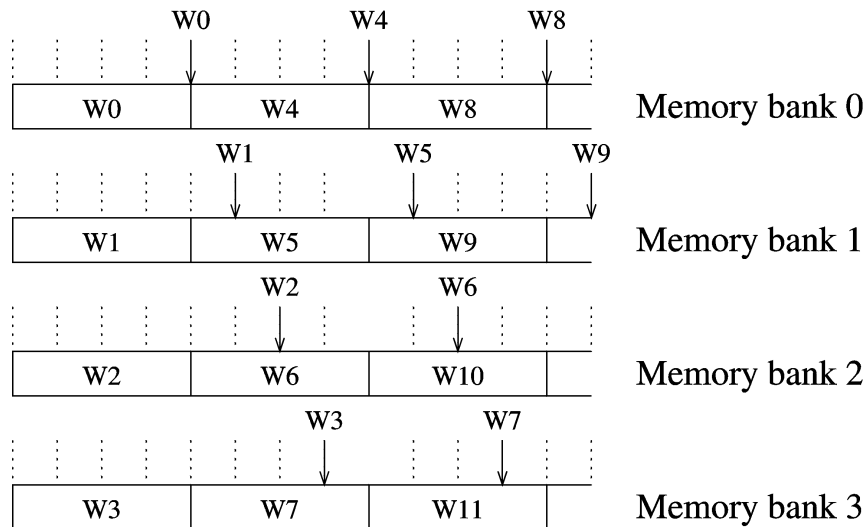# Interleaved Memory (cont'd)

Synchronized access organization

Memory unit #

From address bus | Address in a memory unit (m-2 bits) | 2 bits

2-to-4 decoder

Memory unit 3 | Memory unit 2 | Memory unit 1 | Memory unit 0

MDR | MDR | MDR | MDR

Data bus

# Interleaved Memory (cont'd)

| Access cycle 1 | Access cycle 2 | Access cycle 3 | Access cycle 4 | Access cycle 5 | Access cycle 6 | Access cycle 7 | Access cycle 8 |
|---|---|---|---|---|---|---|---|

W0    W1    W2    W3    W4    W5    W6    W7

| W0 | W1 | W2 | W3 | W4 | W5 | W6 | W7 |
|---|---|---|---|---|---|---|---|

(a) Noninterleaved memory access

W0    W4    W8

| W0 | W4 | W8 |
|---|---|---|
Memory bank 0

W1    W5    W9

| W1 | W5 | W9 |
|---|---|---|
Memory bank 1

W2    W6

| W2 | W6 | W10 |
|---|---|---|
Memory bank 2

W3    W7

| W3 | W7 | W11 |
|---|---|---|
Memory bank 3

Interleaved memory allows pipelined access to memory

(b) Interleaved memory access

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.

# Interleaved Memory (cont'd)

Independent access organization

Memory unit #

From address bus | Address in a memory unit (m-2 bits) | 2 bits

MAR | MAR | MAR | MAR | 2-to-4 decoder

Memory unit 3 | Memory unit 2 | Memory unit 1 | Memory unit 0

Data bus

To be used with S. Dandamudi, "Fundamentals of Computer Organization and Design," Springer, 2003.

# Interleaved Memory (cont'd)

- ## Number of banks
  - $*$ **$M$** = memory access time in cycles
  - $*$ To provide one word per cycle
    - » Number of banks $\geq$ **$M$**

- ## Drawbacks of interleaved memory
  - $*$ Involves complex design
    - » Example: Need MDR or MAR
  - $*$ Reduced fault-tolerance
    - » One bank failure leads to failure of the whole memory
  - $*$ Cannot be expanded incrementally