# CSEE 3827: Fundamentals of Computer Systems, Spring 2022

## Lecture 10

Prof. Dan Rubenstein (danr@cs.columbia.edu)

# Agenda (M&K 7.1-7.4)

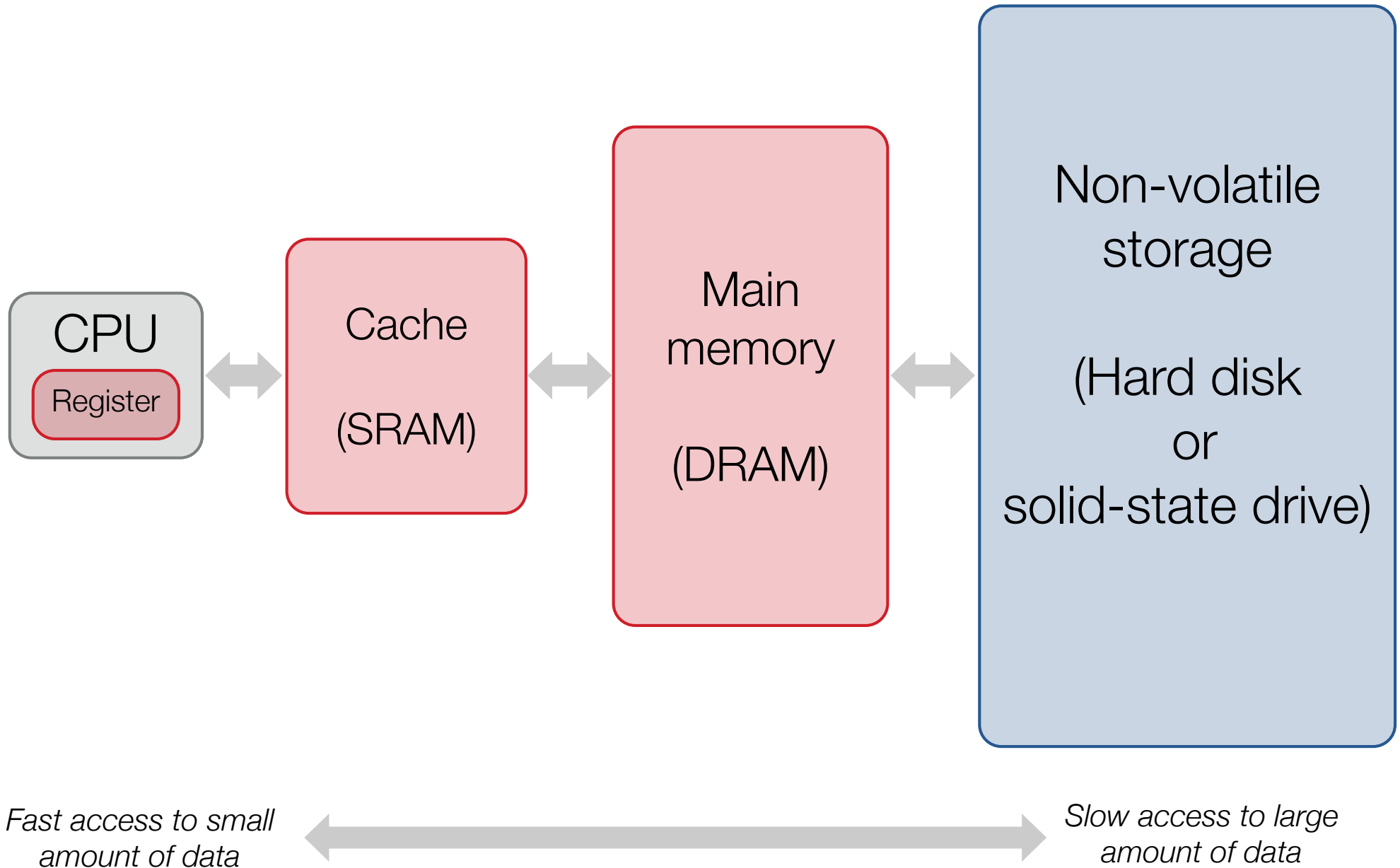- Memory

# Storage hierarchy



CPU

Register

Cache

(SRAM)

Main memory

(DRAM)

Non-volatile storage

(Hard disk or solid-state drive)

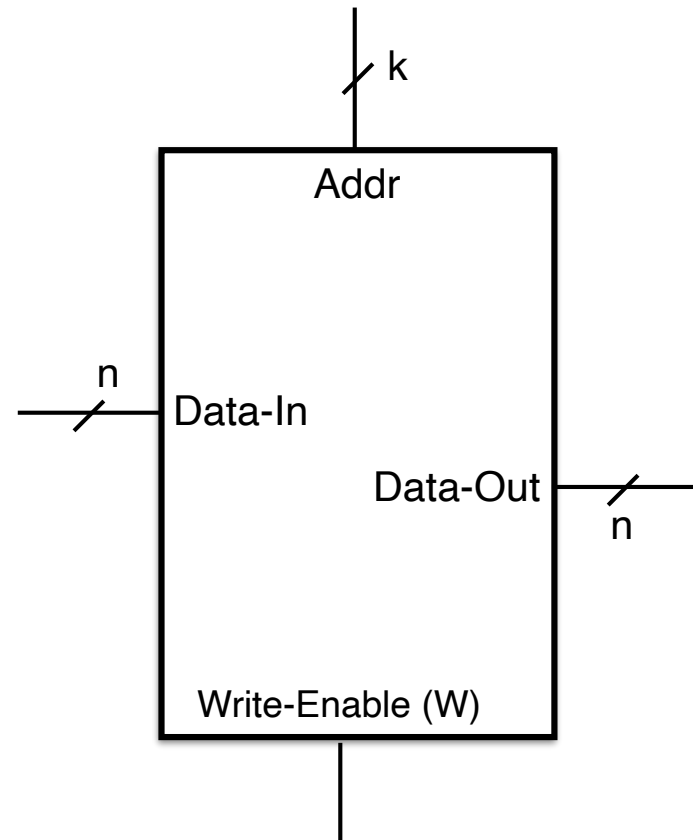*Fast access to small amount of data*

*Slow access to large amount of data*
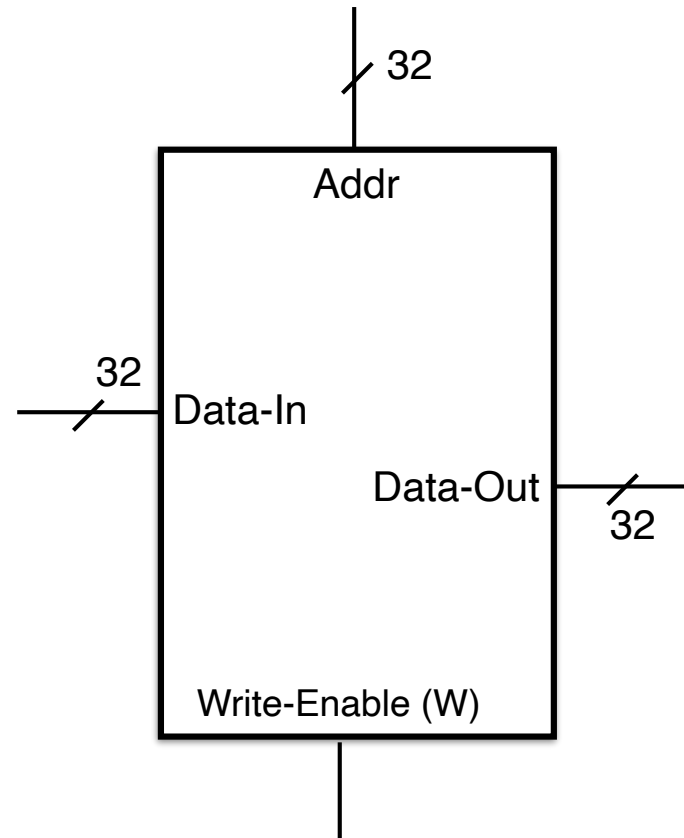
# Using Memory

# General Memory interface

- k-bit addresses ($2^k$ addresses), n-bit data

- Read from memory:

  - Specify a (k-bit) address, A

  - Data-out outputs the n-bit contents in memory at address, A

- Write to memory:

  - Specify a (k-bit) address, A

  - Through Data-In, supply n-bit data to write at address A

  - set Write-Enable to 1

- Careful, if wanting to read and Write-Enable set to 1, will also be writing.

```
            k
         ┌──┼──────────────┐
         │     Addr        │
         │                 │
   n     │                 │
 ───┼────┤ Data-In         │
         │                 │
         │        Data-Out ├───┼──── n
         │                 │
         │ Write-Enable (W)│
         └────────┼────────┘
                  │
```
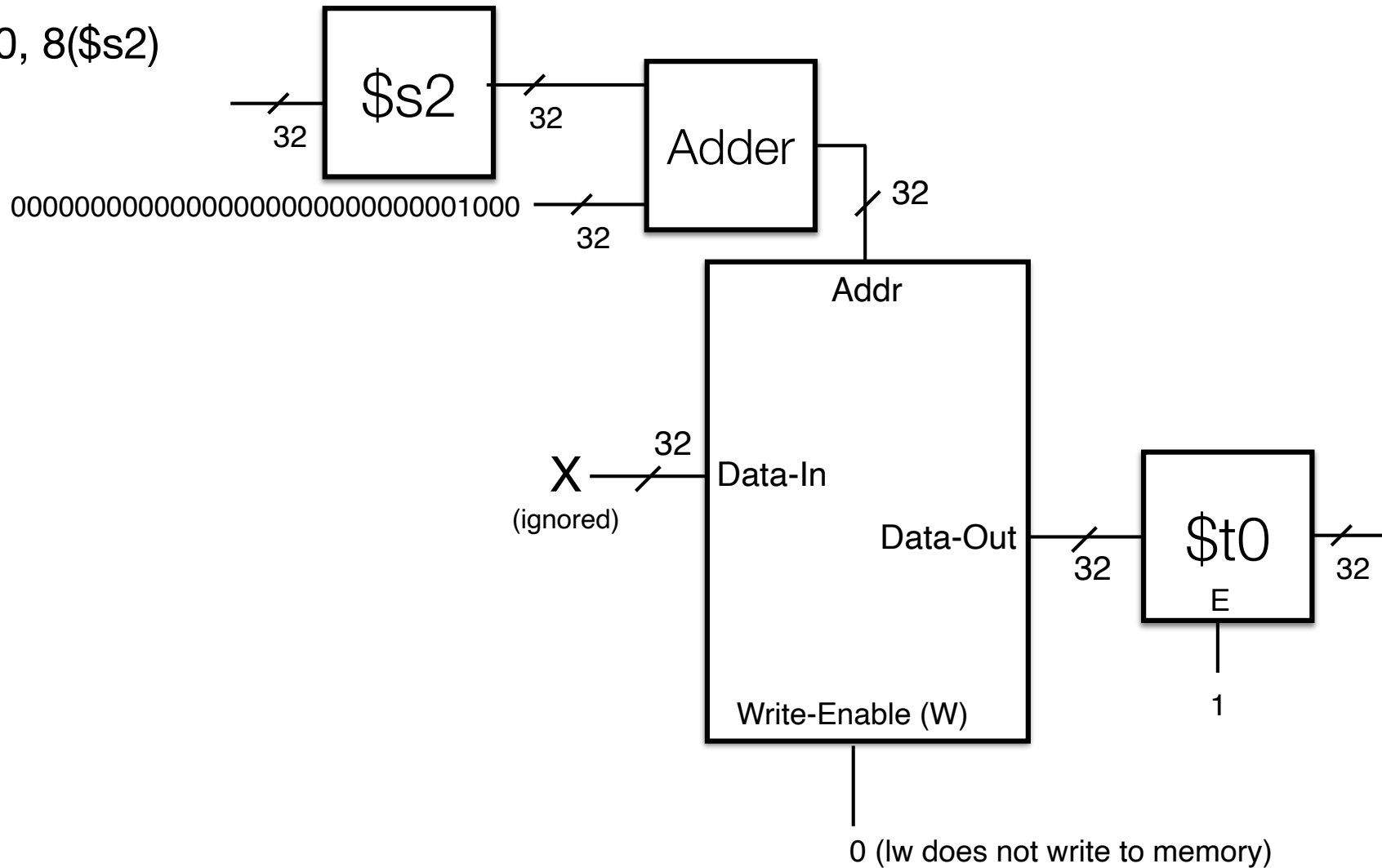
# MIPS-specific

- k=32: 32-bit addresses ($2^{32}$ byte-size addressable slots)

- n=32: we access memory 32-bits at a time

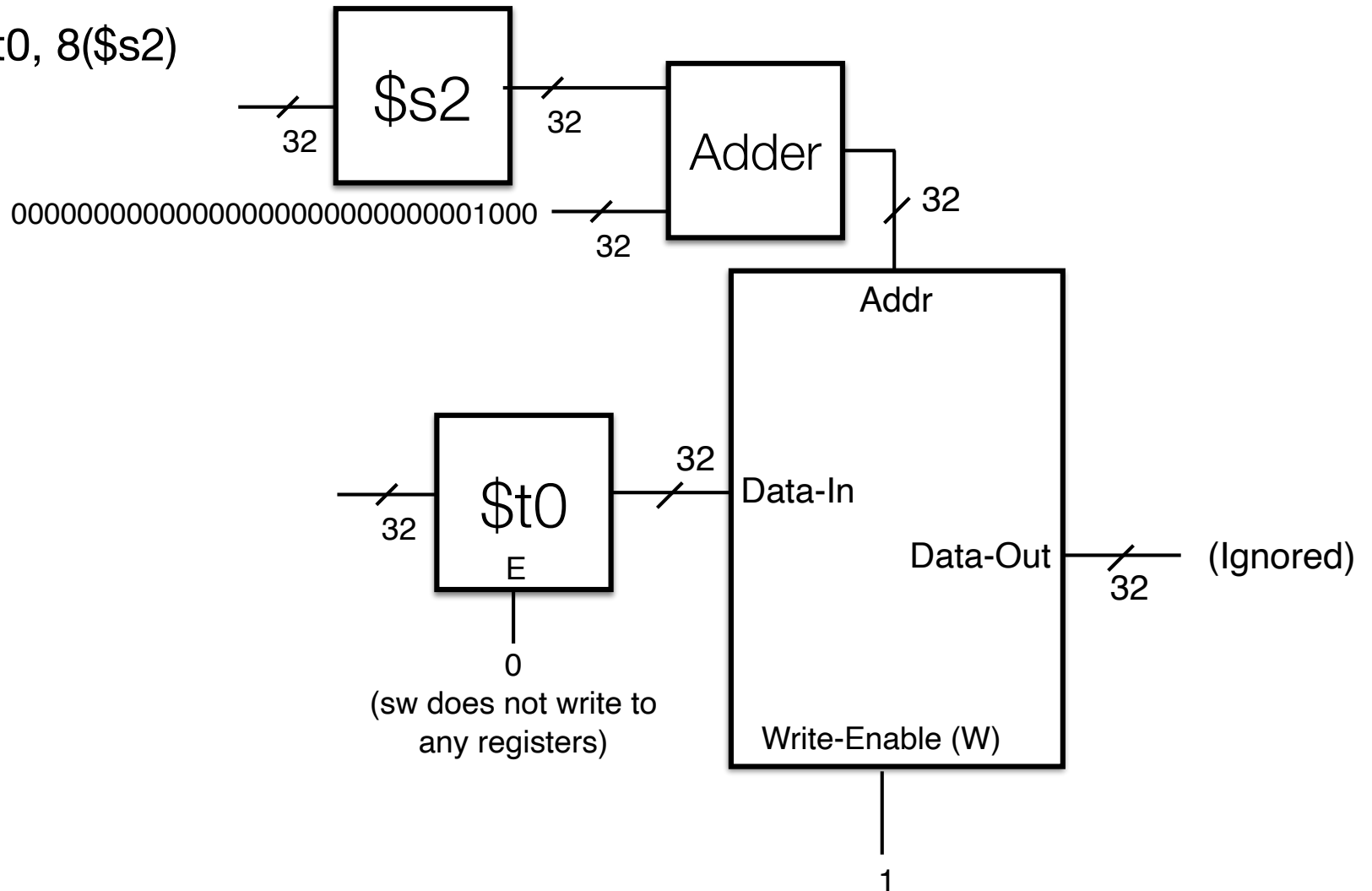- 32-bit address always multiple of 4 (2 least significant bits always "00")

32

Addr

32
Data-In

Data-Out
32

Write-Enable (W)

# MIPS-specific read example

- e.g., lw $t0, 8($s2)

$s2

32

32

Adder

00000000000000000000000000001000

32

32

Addr

32

X

(ignored)

Data-In

Data-Out

32

$t0

32

E

1

Write-Enable (W)

0 (lw does not write to memory)

# MIPS-specific read example

- e.g., sw $t0, 8($s2)

# Building Memory

# Recall: SR Latch

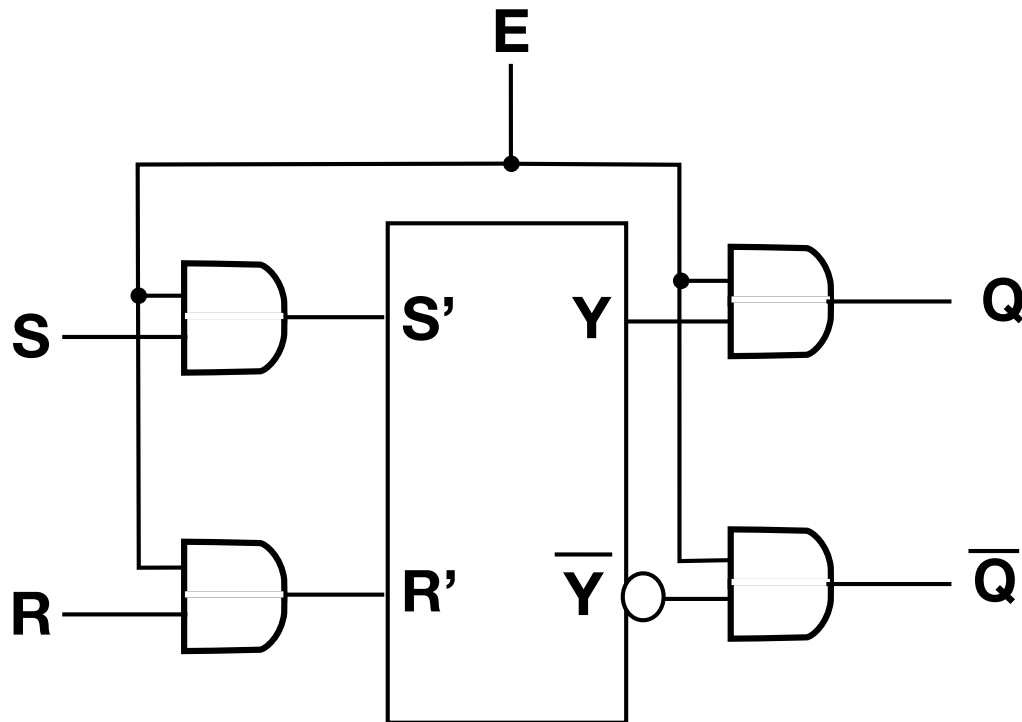| S | R | Q | $\overline{Q}$ |
|---|---|---|---|
| 0 | 0 | Q | $\overline{Q}$ |
| 0 | 1 | 0 | 1 |
| 1 | 0 | 1 | 0 |
| 1 | 1 | 0 | 0 |

Hold value

Reset

Set

# Single (1-bit) Memory Cell

# Memory "Cell"

- memory storage of a single bit

- built from SR-latch and some enablers (AND gates)



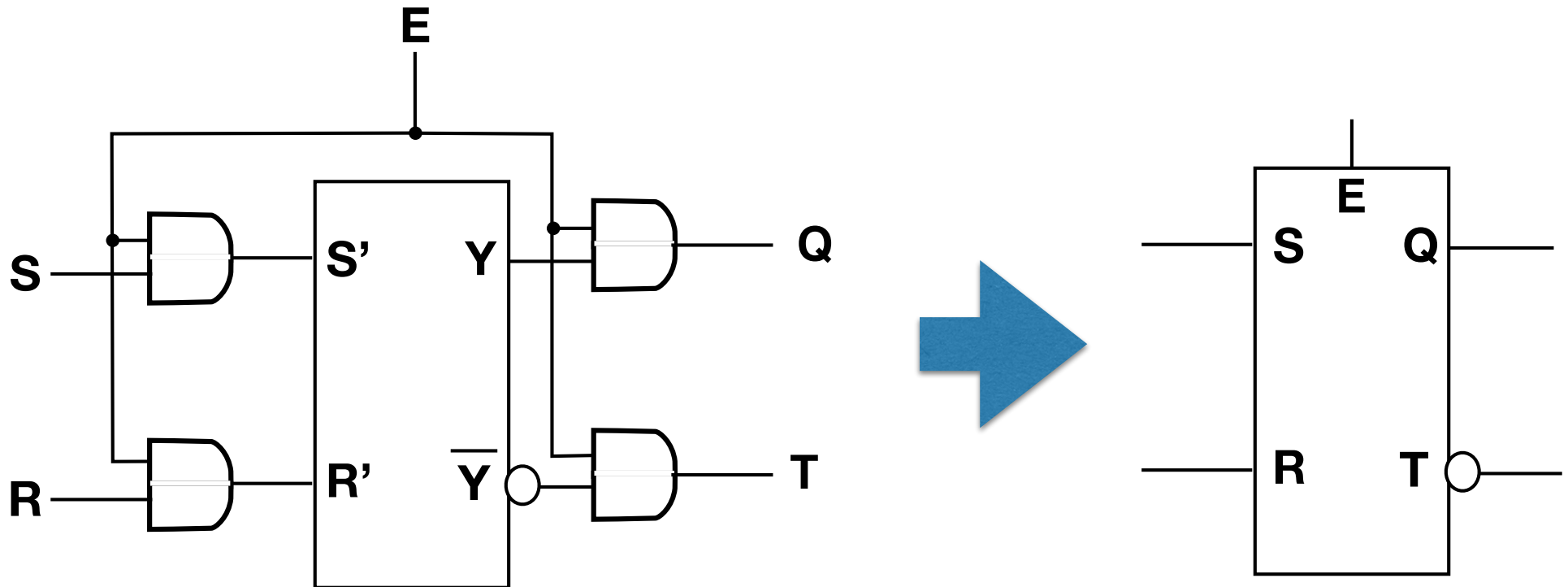| E | S | R | Y | $\overline{Y}$ | Q | $\overline{Q}$ | |
|---|---|---|---|---|---|---|---|
| 0 | X | X | Y | $\overline{Y}$ | 0 | 0 | Disable Inputs and Zero Outputs |
| 1 | 0 | 0 | Y | $\overline{Y}$ | Y | $\overline{Y}$ | Hold Value |
| 1 | 0 | 1 | 0 | 1 | 0 | 1 | Reset (to 0) |
| 1 | 1 | 0 | 1 | 0 | 1 | 0 | Set (to 1) |
| 1 | 1 | 1 | 0 | 0 | 0 | 0 | |

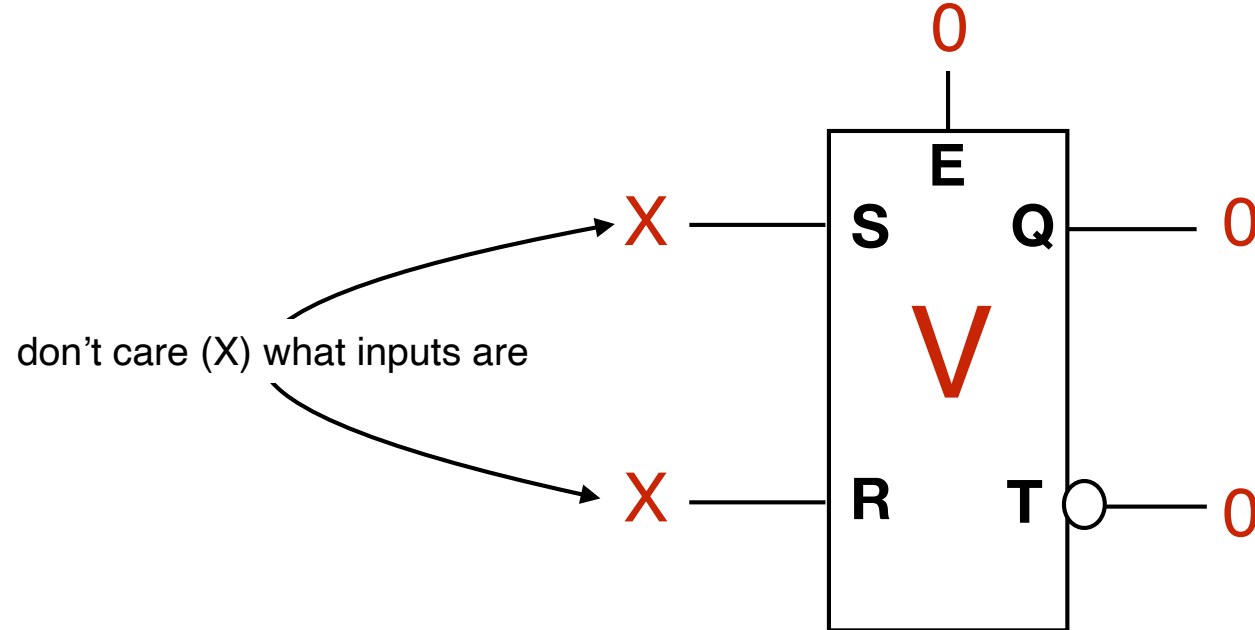Note: Setting E=0 does not erase value in latch, just "zeros" the value being output

# Memory "Cell"

- memory storage of a single bit

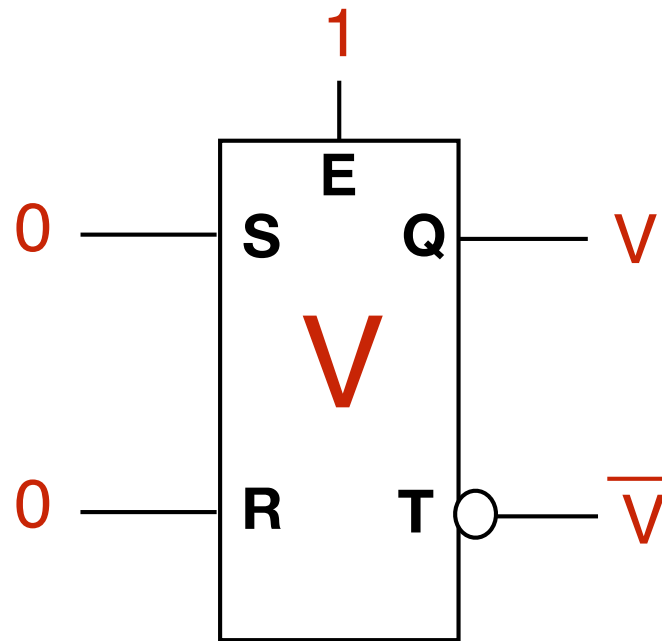- built from SR-latch and some enablers (AND gates)

# Disable a Cell

- Set E=0, S=X,R=X

- Output not 0'd, cell (1-bit) storage unchanged
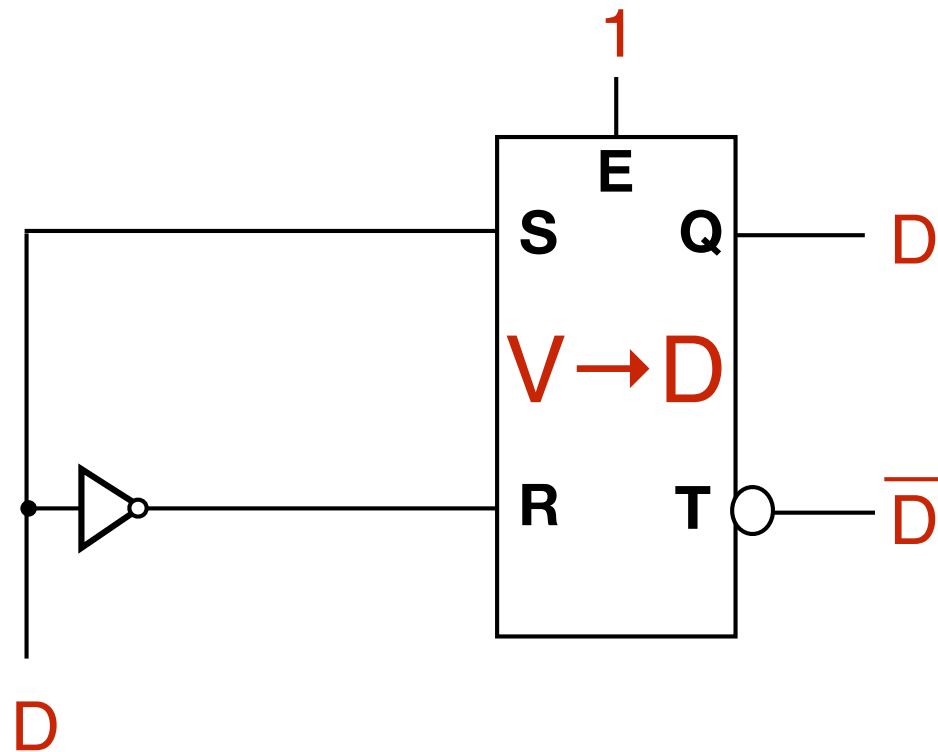


don't care (X) what inputs are

# Read from a Cell

- Set E=1, S=R=0

- Output not 0'd, cell (1-bit) storage unchanged

# Write to a Cell

- To write a data bit D into the cell

  - Enable the cell (E=1)

  - Set S=D, R=$\overline{D}$

# Reading or Writing

- E=0, ignore all inputs (outputs both 0)

- E=1

    - W=0: output bit stored in cell

    - W=1, write D into cell (also the output of the cell)

# Reading or Writing

- E=0, ignore all inputs (outputs both 0)

- E=1

  - W=0: output bit stored in cell

  - W=1, write D into cell (also the output of the cell)

# Reading or Writing

- E=0, ignore all inputs (outputs both 0)

- E=1

    - W=0: output bit stored in cell

    - W=1, write D into cell (also the output of the cell)

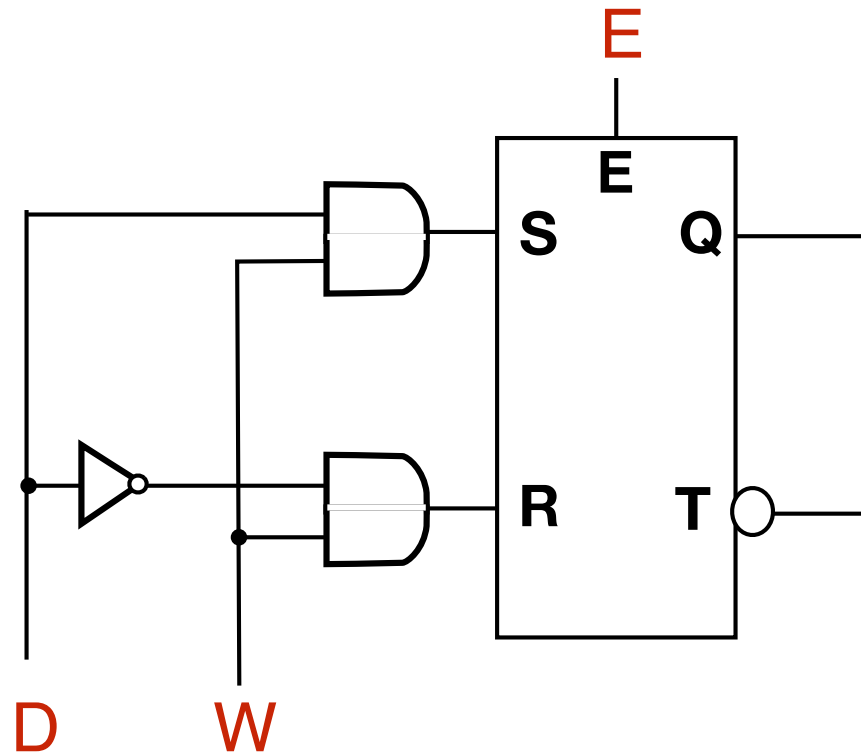# Reading or Writing
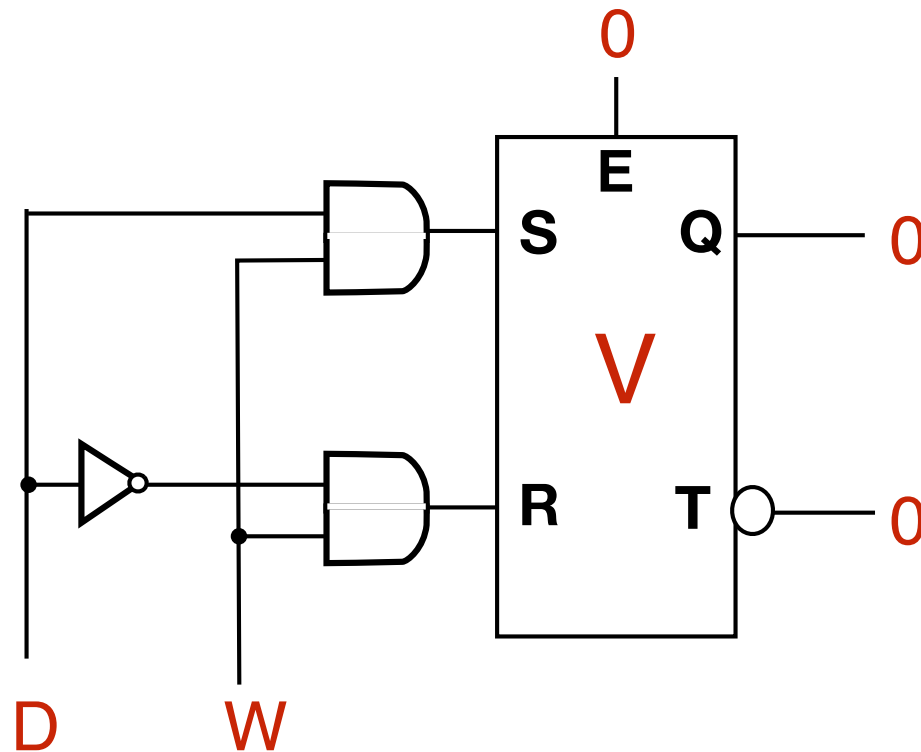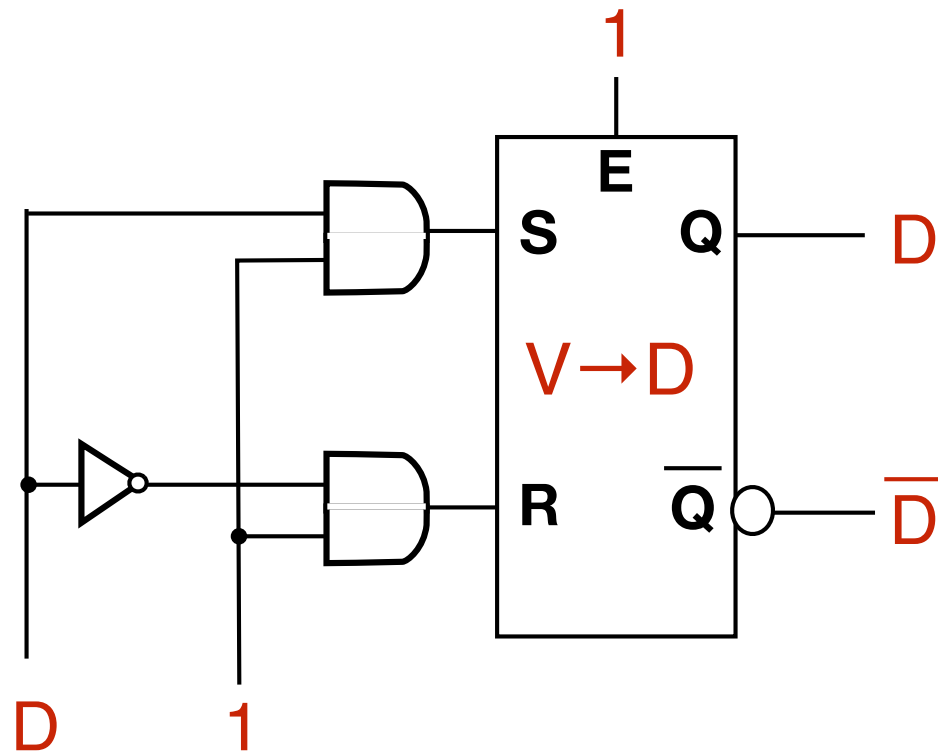
- E=0, ignore all inputs (outputs both 0)

- E=1

  - W=0: output bit stored in cell

  - W=1, write D into cell (also the output of the cell)

# Reading or Writing

- E=0, ignore all inputs (outputs both 0)

- E=1

  - W=0: output bit stored in cell

  - W=1, write D into cell (also the output of the cell)

# Memory Word
# (MIPS: 32 bits)

# Word of MIPS Memory

- memory storage of a single bit

- built from SR-latch and some enablers (AND gates)



## 32 latches w/ Enable

# Word of MIPS Memory



Attach Write Logic to each

# Word of MIPS Memory



Attach to same enable signal
(Enable all 32 or none)

# Word of MIPS Memory



Attach to same write signal
(write to all 32 (when enabled) or none)

# Word of MIPS Memory



Attach 32-bit Data-In signal, D, bit-by-bit
to corresponding latches

# Word of MIPS Memory



## Attach 32-bit Data-Out signal, R, bit-by-bit to corresponding latches

# Word of MIPS Memory

# Word of MIPS Memory

- E=0: don't write to this word of memory, output all 0's

0

X —/— **D**                                 **E**

32         **V** (32-bit) word of memory          **R** —/— 0

X — **W**                                     32

# Word of MIPS Memory

- E=1, W=0: read from this word of memory, but don't write

# Word of MIPS Memory

- E=1, W=1: write input at D into memory (read as well)

1

B /<sub>32</sub> **D**   **E**   **R** /<sub>32</sub> B
1 **W**   V→B (32-bit) word of memory

# $2^{32}$ bytes = $2^{30}$ words of Addressable Memory

**D**    32    **W**    Address 0    **E**    (32-bit) word of memory    **R**    32

**D**    32    **W**    Address 4    **E**    (32-bit) word of memory    **R**    32

**D**    32    **W**    Address 8    **E**    (32-bit) word of memory    **R**    32

■

■

■

**D**    32    **W**    Address $2^{32}$-4    **E**    (32-bit) word of memory    **R**    32

# $2^{32}$ bytes = $2^{30}$ words of Addressable Memory

[30 bits describing address] 00

- Each address is really described by 30 bits (since the 2 lowest-order bits of the 32-bit address are 00)

- Data Input to memory, D, fed into every word address

**D** | Address 0 | **E** (32-bit) word of memory | **R** 32
32 | | |
**W** | | |

**D** | Address 4 | **E** (32-bit) word of memory | **R** 32
32 | | |
**W** | | |

**D** | Address 8 | **E** (32-bit) word of memory | **R** 32
32 | | |
**W** | | |

**D** | Address $2^{32}$-4 | **E** (32-bit) word of memory | **R** 32
32 | | |
**W** | | |

D 32

W

# $2^{32}$ bytes = $2^{30}$ words of Addressable Memory

[30 bits describing address] 00

- Each address is really described by 30 bits (since the 2 lowest-order bits of the 32-bit address are 00)

- Data Input to memory, D, fed into word at each address

- Write Enable, W, fed into word at each addresss

# $2^{32}$ bytes = $2^{30}$ words of Addressable Memory

[30 bits describing address] 00

30

30-to-$2^{30}$ Decoder

| D | Address 0 | (32-bit) word of memory | E | R |
| W | | | | 32 |

32

| D | Address 4 | (32-bit) word of memory | E | R |
| W | | | | 32 |

32

| D | Address 8 | (32-bit) word of memory | E | R |
| W | | | | 32 |

32

■ ■ ■
■ ■ ■
■ ■ ■

| D | Address $2^{32}$-4 | (32-bit) word of memory | E | R |
| W | | | | 32 |

32

D 32

W

- Each address is really described by 30 bits (since the 2 lowest-order bits of the 32-bit address are 00)

- Data Input to memory, D, fed into word at each address

- Write Enable, W, fed into word at each addresss

- Decoder used to enable only 1 word

# $2^{32}$ bytes = $2^{30}$ words of Addressable Memory

[30 bits describing address] 00

30

30-to-$2^{30}$
Decoder

| D | | E | | R | |
|---|---|---|---|---|---|
| 32 | | | | | 32 |
| | Address 0 | (32-bit) word of memory | | | |
| W | | | | | |

| D | | E | | R | |
|---|---|---|---|---|---|
| 32 | | | | | 32 |
| | Address 4 | (32-bit) word of memory | | | |
| W | | | | | |

| D | | E | | R | |
|---|---|---|---|---|---|
| 32 | | | | | 32 |
| | Address 8 | (32-bit) word of memory | | | |
| W | | | | | |

■  ■        ■
■  ■        ■
■  ■        ■

| D | | E | | R | |
|---|---|---|---|---|---|
| 32 | | | | | 32 |
| | Address $2^{32}$-4 | (32-bit) word of memory | | | |
| W | | | | | |

D
32

1 ← Write data D into the word at enabled address
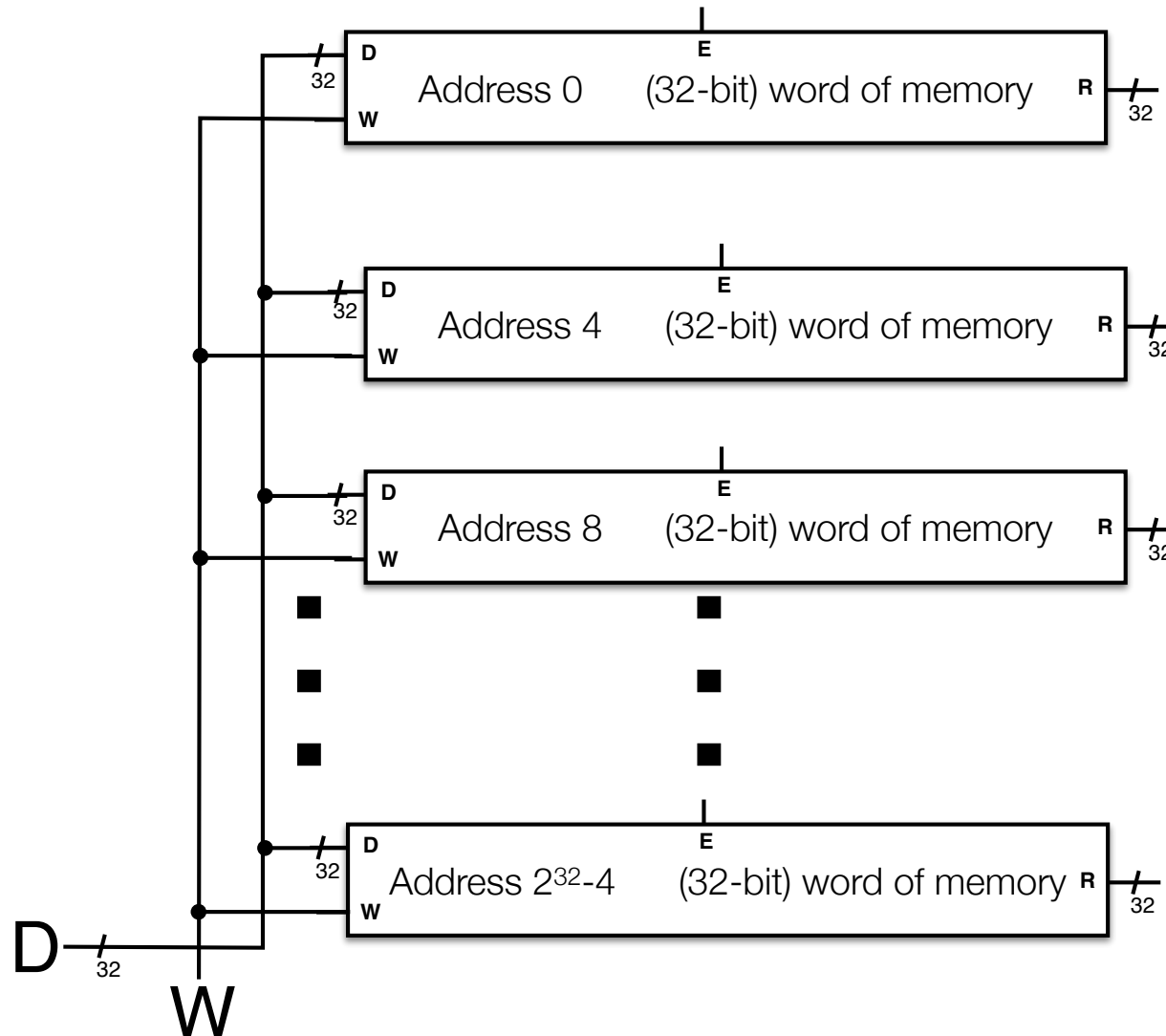
- Each address is really described by 30 bits (since the 2 lowest-order bits of the 32-bit address are 00)
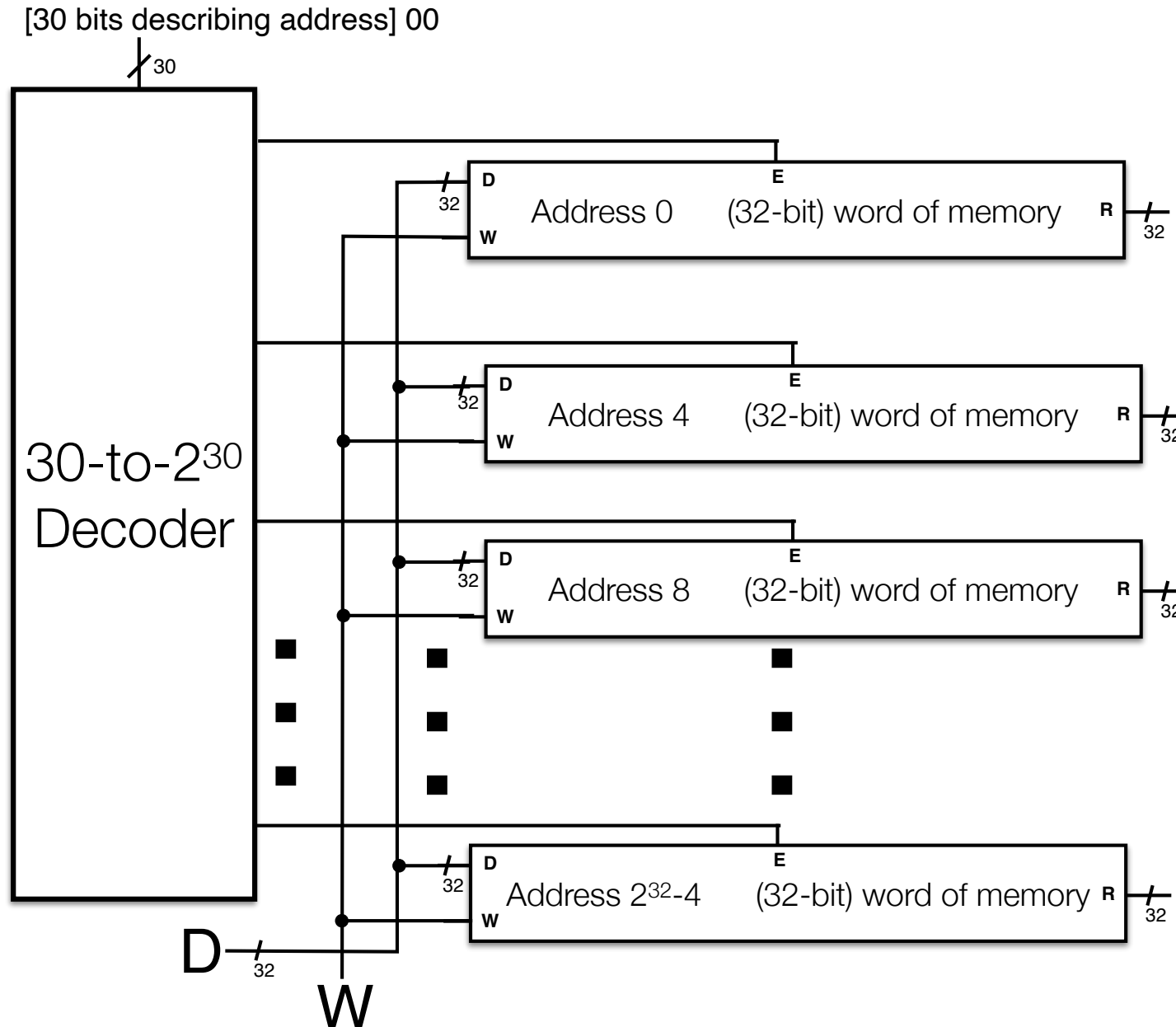
- Data Input to memory, D, fed into word at each address

- Write Enable, W, fed into word at each addresss

- Decoder used to enable only 1 word

# $2^{32}$ bytes = $2^{30}$ words of Addressable Memory

[30 bits describing address] 00

30-to-$2^{30}$ Decoder

| D | Address 0 | (32-bit) word of memory | E | R |
| W | | | | |

| D | Address 4 | (32-bit) word of memory | E | R |
| W | | | | |

| D | Address 8 | (32-bit) word of memory | E | R |
| W | | | | |

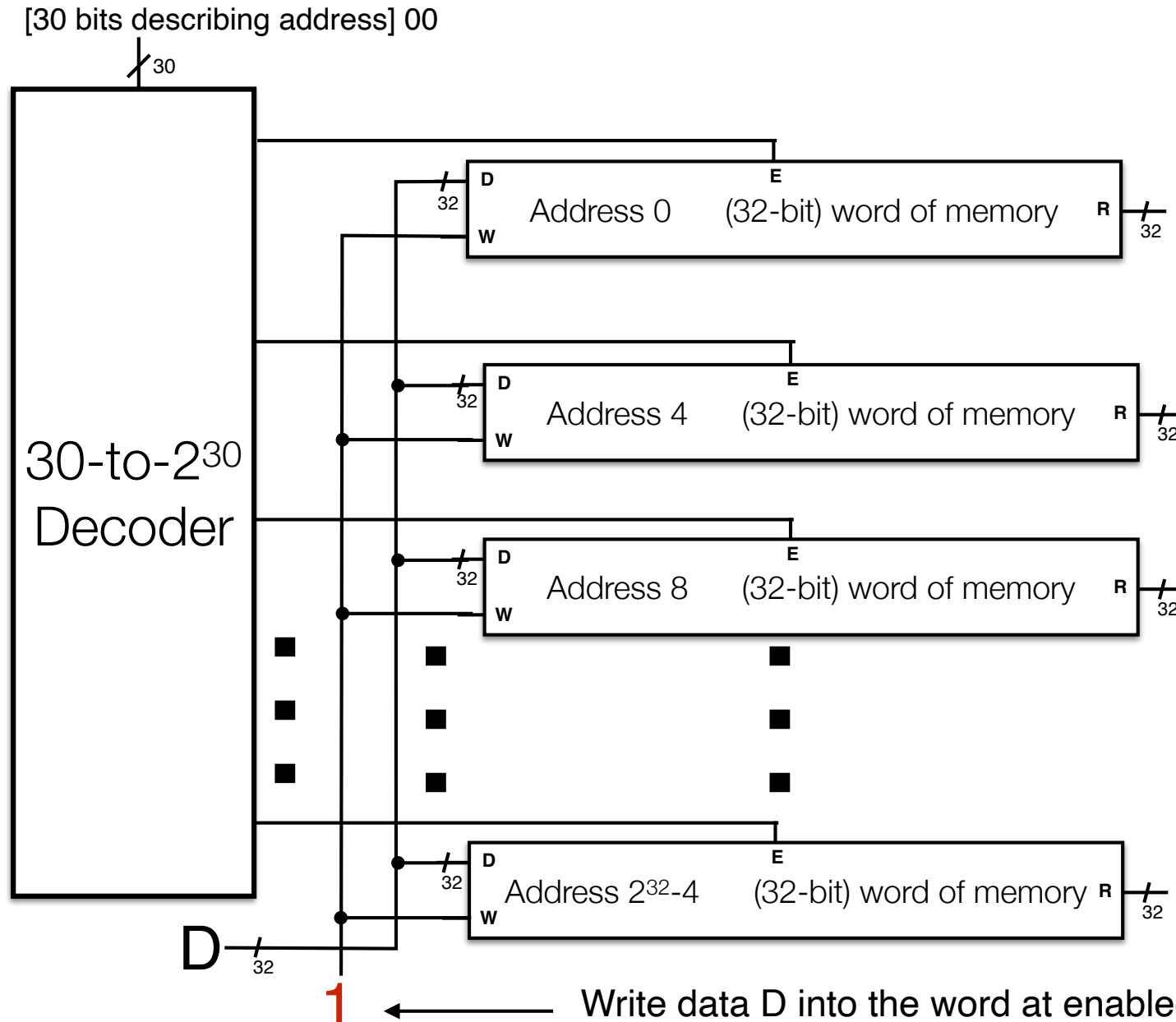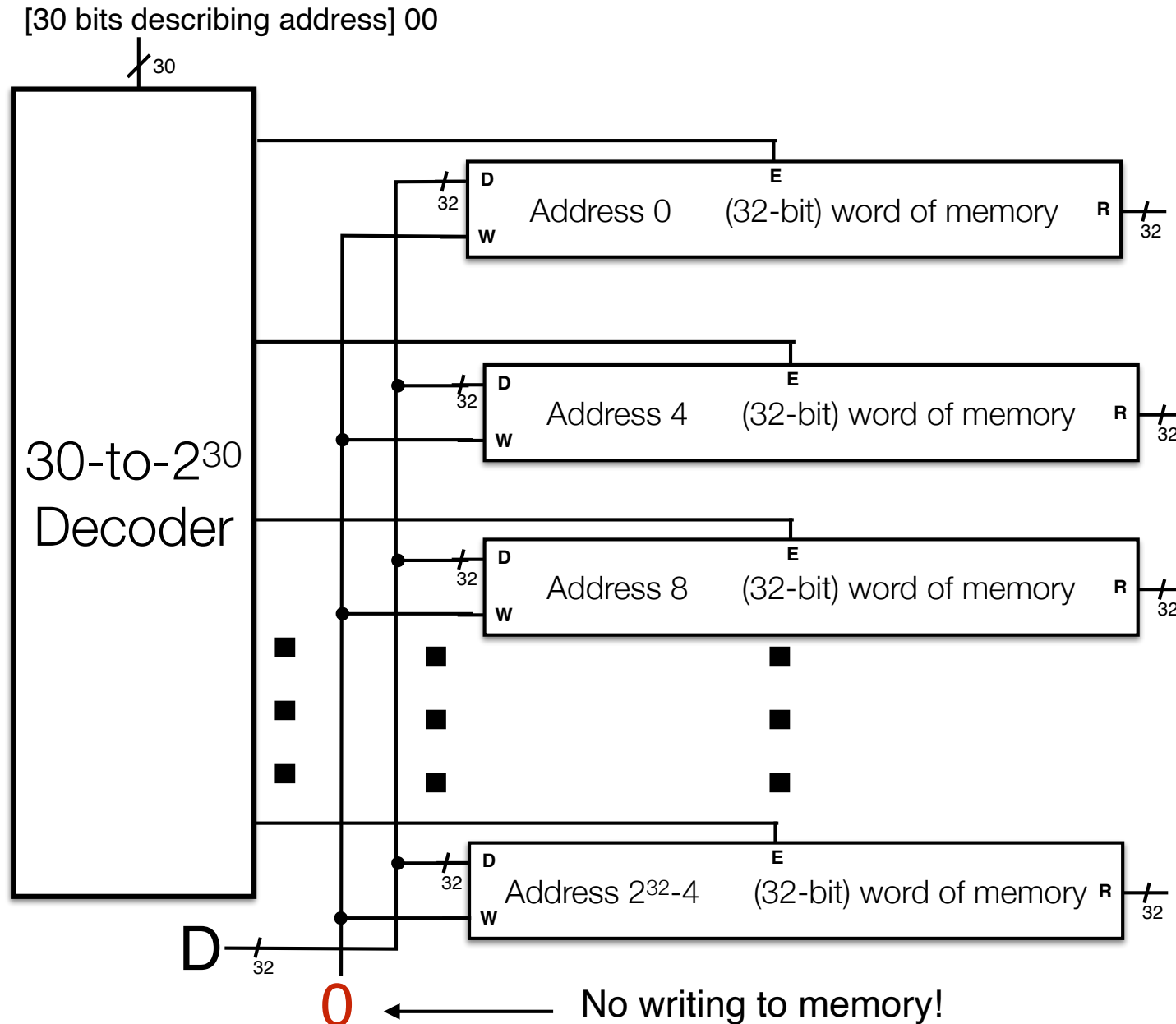| D | Address $2^{32}$-4 | (32-bit) word of memory | E | R |
| W | | | | |

D

0 ← No writing to memory!

- Each address is really described by 30 bits (since the 2 lowest-order bits of the 32-bit address are 00)

- Data Input to memory, D, fed into word at each address

- Write Enable, W, fed into word at each addresss

- Decoder used to enable only 1 word

# RECALL: Merge-with-known-0's

- Suppose have m input signals

- All of them = 0 except maybe one $X_j$.

- Circuit to retrieve the value of the one unknown

$$X_3$$
$$X_2$$
$$X_1$$
$$X_0$$

- If all $X_i=0$ for $i \neq j$, then the OR over all $X_i$ (including $X_j$) = $X_j$

# $2^{32}$ bytes $= 2^{30}$ words of Addressable Memory

[30 bits describing address] 00

30

30-to-$2^{30}$ Decoder

| | D | E | | |
|---|---|---|---|---|
| 32 | | Address 0 | (32-bit) word of memory | R |
| | W | | | 32 |

| | D | E | | |
|---|---|---|---|---|
| 32 | | Address 4 | (32-bit) word of memory | R |
| | W | | | 32 |

| | D | E | | |
|---|---|---|---|---|
| 32 | | Address 8 | (32-bit) word of memory | R |
| | W | | | 32 |

| | D | E | | |
|---|---|---|---|---|
| 32 | | Address $2^{32}$-4 | (32-bit) word of memory | R |
| | W | | | 32 |

D 32

W

32

- OR all word outputs together

- Only enabled word is not "zero'd out"

# Coincident Selection

# Recall: Decoder
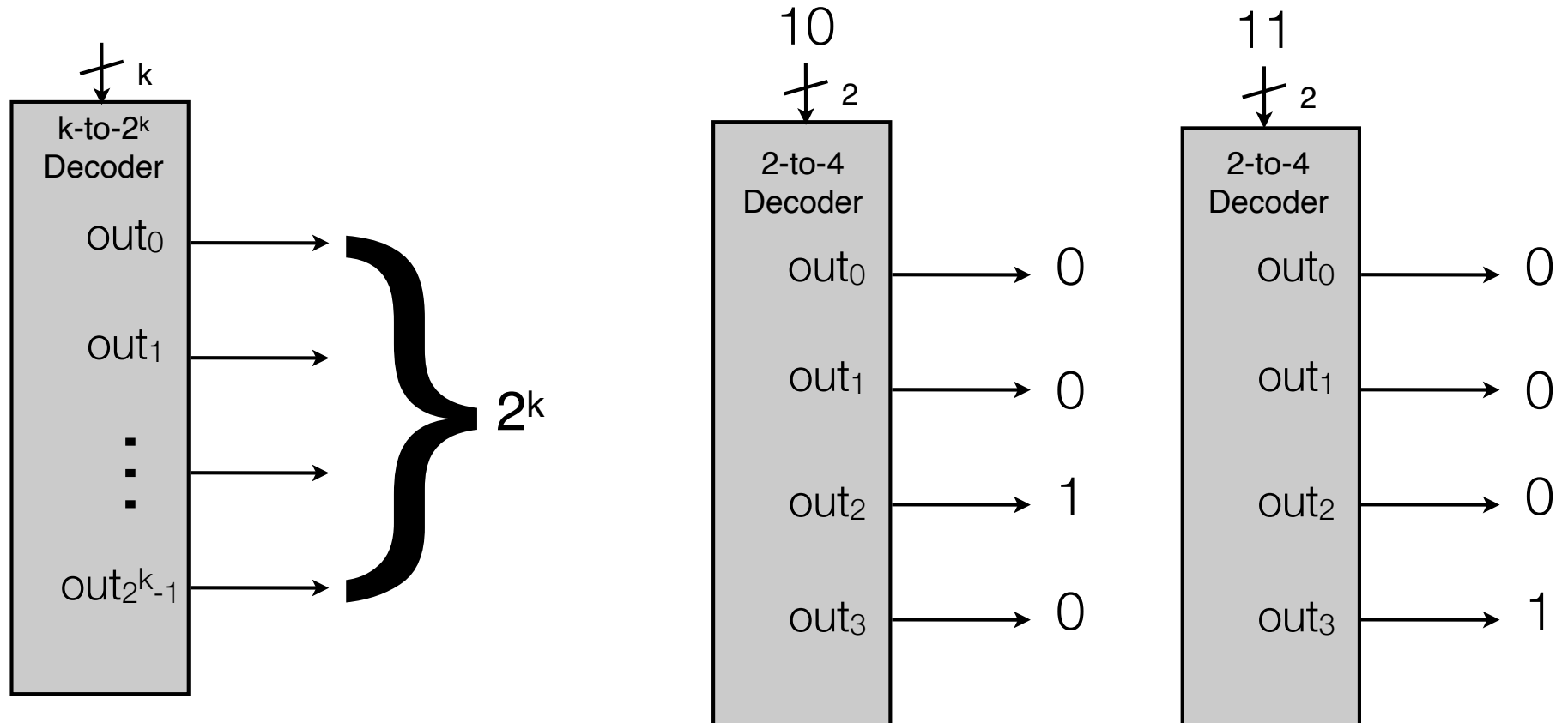
- No "DATA" inputs, just a k-bit "selector" input

- $2^k$ 1-bit outputs

- Selector input (i) chooses which output = 1, all other outputs = 0

k-to-$2^k$ Decoder
out$_0$
out$_1$
out$_{2^k-1}$
$2^k$

10
2-to-4 Decoder
out$_0$ → 0
out$_1$ → 0
out$_2$ → 1
out$_3$ → 0

11
2-to-4 Decoder
out$_0$ → 0
out$_1$ → 0
out$_2$ → 0
out$_3$ → 1

# Pick along 1 Dimension: Use a decoder

Suppose have $2^k$ "items" and want to choose 1
(here k=3)

# Pick along 1 Dimension: Use a decoder



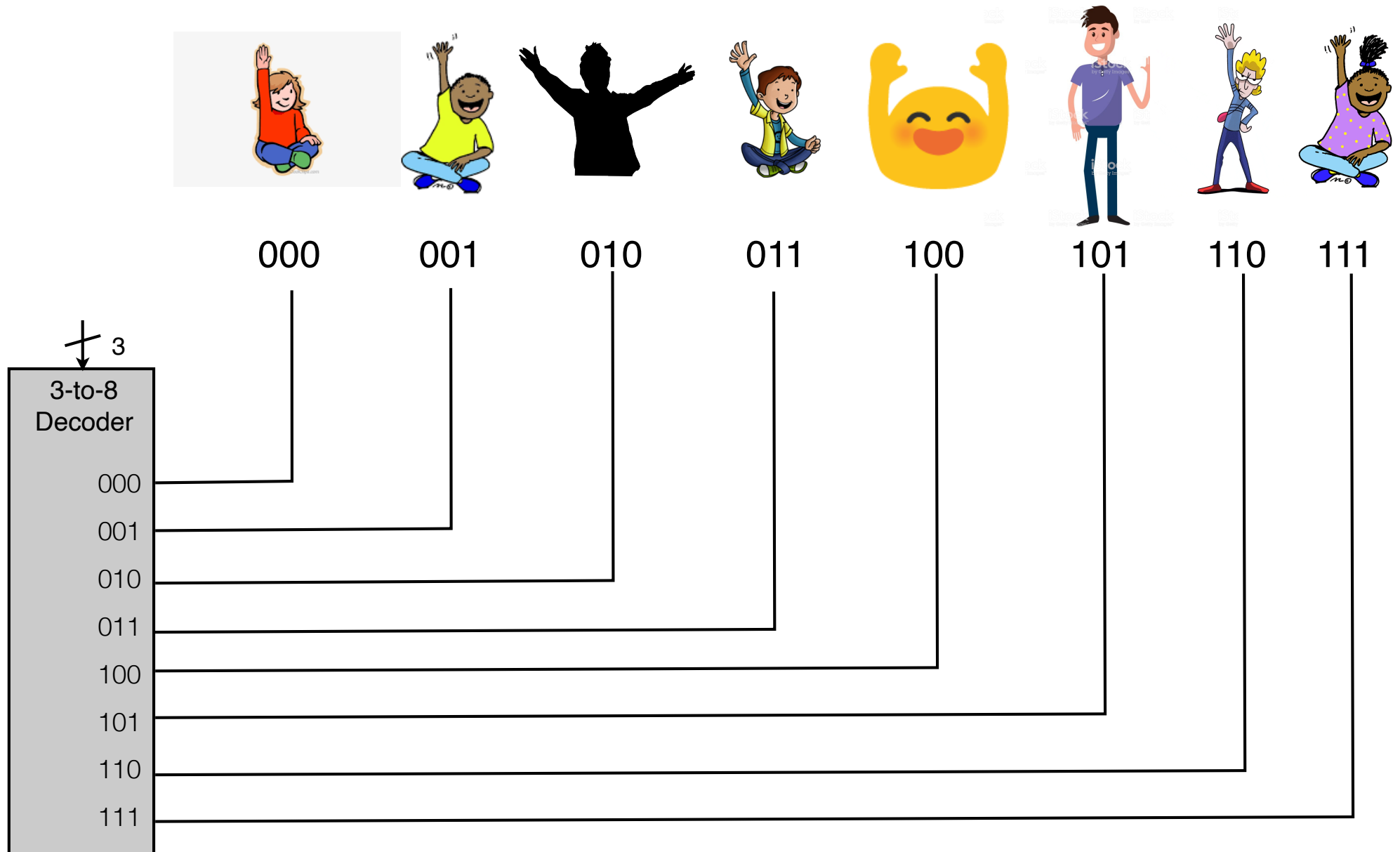| 000 | 001 | 010 | 011 | 100 | 101 | 110 | 111 |

Suppose have $2^k$ "items" and want to choose 1
(here k=3)

Assign a number to each "item"

# Pick along 1 Dimension: Use a decoder



000 001 010 011 100 101 110 111

3

3-to-8
Decoder

000
001
010
011
100
101
110
111

# 1-Dimensional Decoding Summary

- To enable 1 of $2^k$ addresses:

  - 1 Decoder

  - k selector "pins" on decoder

  - $2^k$ output "pins"

# Pick along 2 Dimensions: Use 2 decoders

Column 0   Column 1

Row 00

Row 01

Row 10

Row 11

# Pick along 2 Dimensions: Use 2 decoders
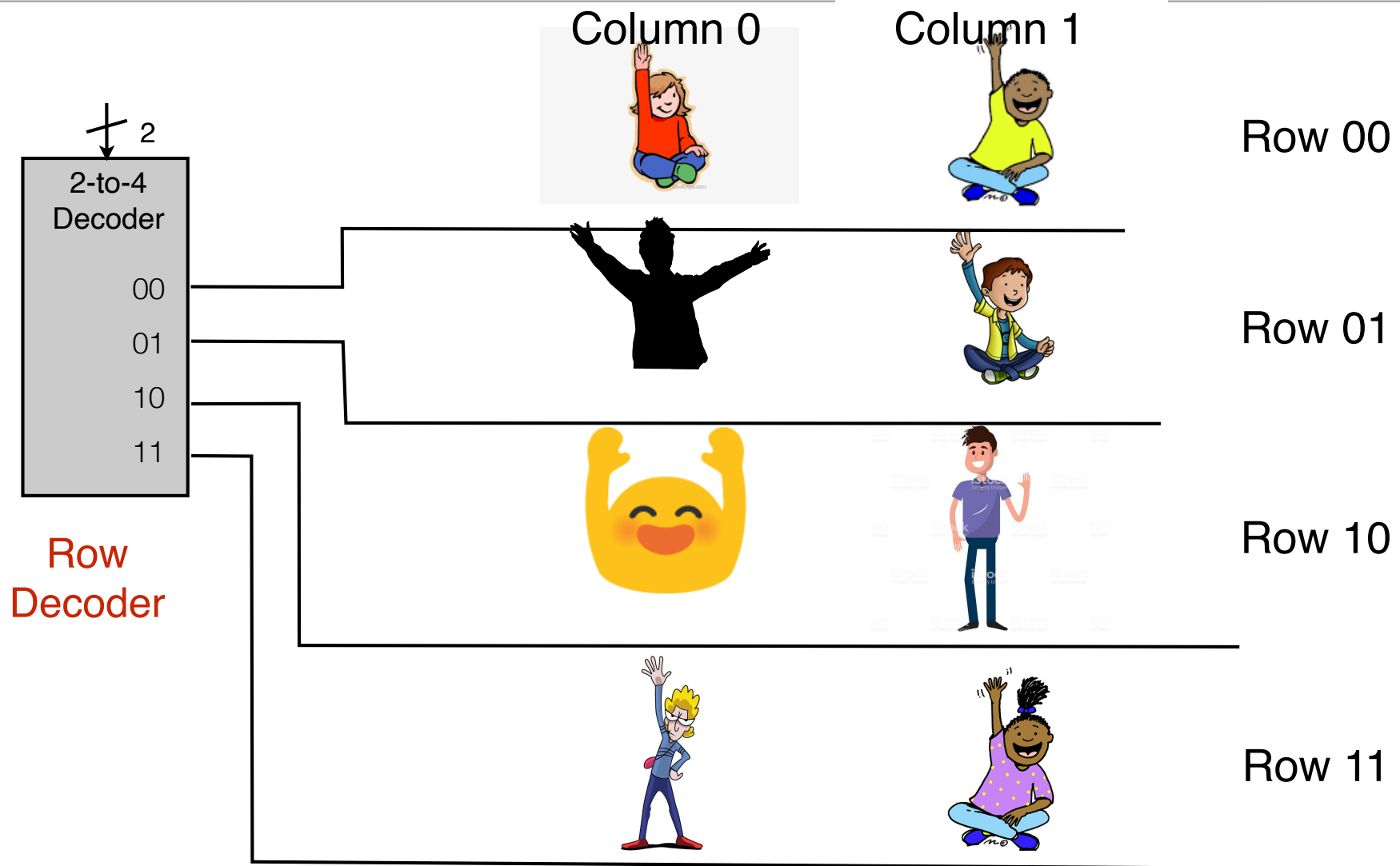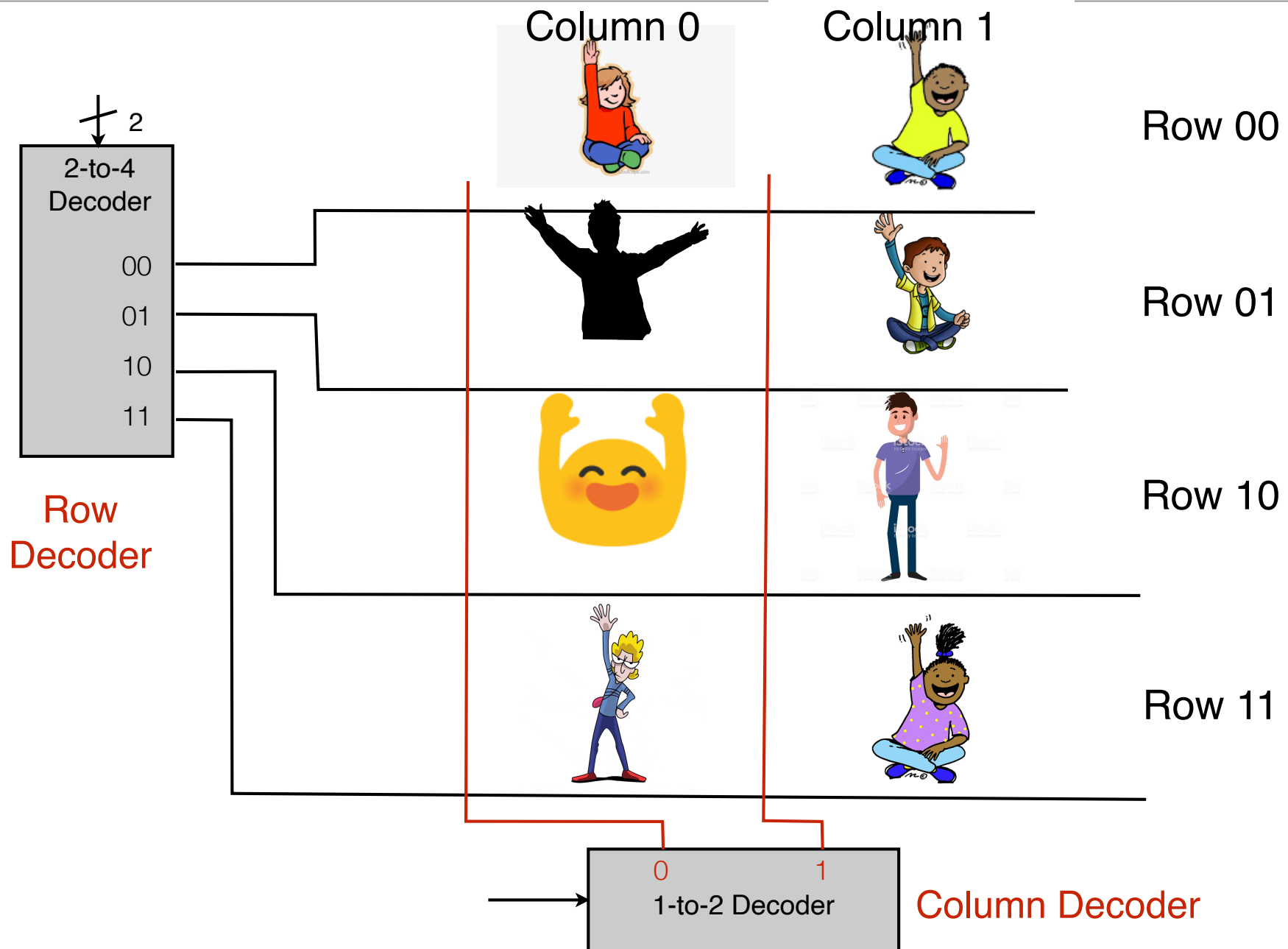
Column 0    Column 1

2

2-to-4
Decoder

00

01

10

11

Row
Decoder

Row 00

Row 01

Row 10

Row 11

# Pick along 2 Dimensions: Use 2 decoders

# Pick along 2 Dimensions: Use 2 decoders

# 1-D v. 2-D Decoding Comparison Summary

- 1D

- To enable 1 of $2^k$ addresses:

  - 1 Decoder

  - k selector "pins" on decoder

  - $2^k$ output "pins"

  - e.g., k=10, 1024 output pins

- 2D

- To enable 1 of $2^k$ addresses:

  - 2 Decoders

  - $k_1$ and $k_2$ selector "pins" on decoder with $k_1+k_2=k$

  - $2^{k_1} + 2^{k_2}$ output "pins"

  - e.g., $k_1=5, k_2=5$, 64 output pins

## 2D has Significantly Less Decoder Logic

# Modifying memory for coincident selection…

[30 bits describing address] 00

Recall - 1D structure



30-to-$2^{30}$ Decoder

D — E
W  Address 0    (32-bit) word of memory    R
32                                          32

D — E
W  Address 4    (32-bit) word of memory    R
32                                          32

D — E
W  Address 8    (32-bit) word of memory    R
32                                          32

D — E
W  Address $2^{32}$-4    (32-bit) word of memory    R
32                                          32

D
32

W

- OR all word outputs together

- Only enabled word is not "zero'd out"

52

# Layout words of memory into rows and columns

- Arrange memory words into 2-dimensional grid

# Layout words of memory into rows and columns

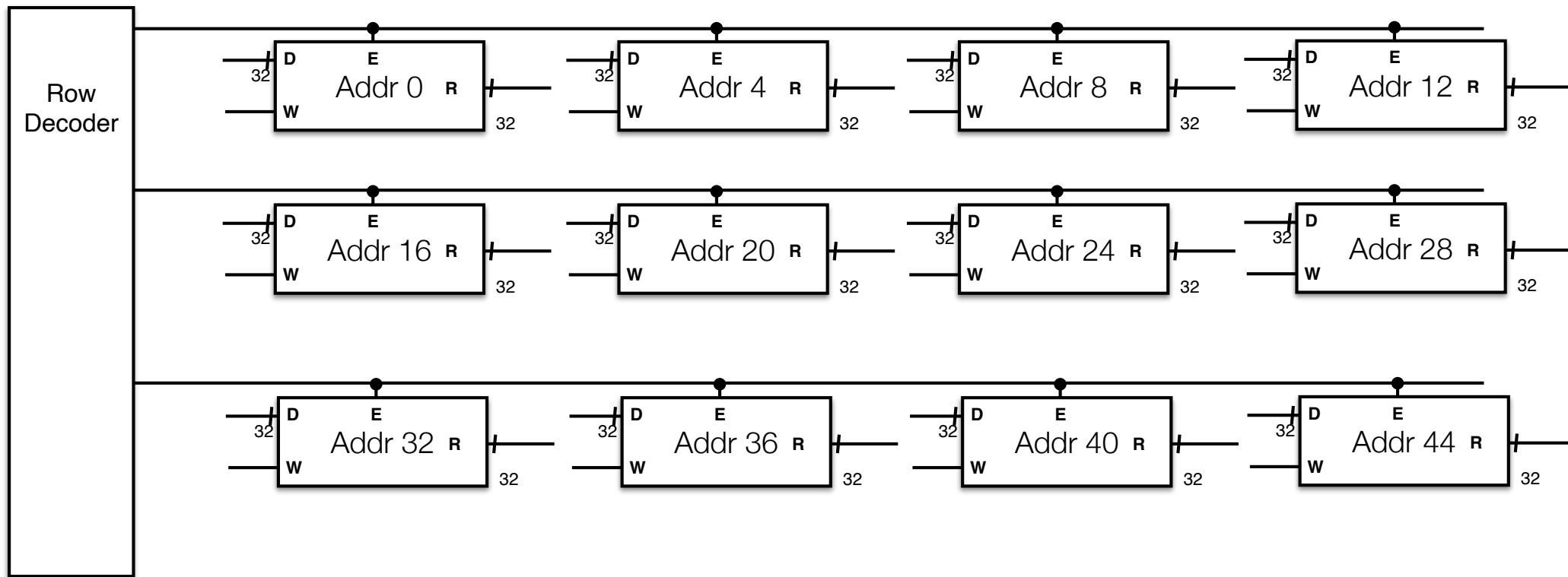- Row decoder enables 1 row of addresses for use

# Layout words of memory into rows and columns

- Column decoder will select a column, but how?

# Layout words of memory into rows and columns

- For write: AND decoder output with W input

# Layout words of memory into rows and columns

- For read, take output of each column (only 1 row enabled)…

# Layout words of memory into rows and columns

- Choose the column to enable also



Column
Decoder

# Using Multiple Memory Chips to Scale Memory

# Two ways to Scale Memory

- Given memory chips with $2^n$ addresses and k-bit word size

  - Suppose want more addresses (e.g., $2^{n+i}$)

  - Suppose want larger word size (e.g., mk for some integer m>1)

# Multi-chip memories: extend wordsize

- Extend wordsize by storing full word in parallel across multiple chips
- e.g., build 64-bit memory storing $2^n$ words from two 32-bit $2^n$-word chips

# Multi-chip memories: extend wordsize

- Extend wordsize by storing full word in parallel across multiple chips
- e.g., build 64-bit memory storing $2^n$ words from two 32-bit $2^n$-word chips

# Multi-chip memories: extend wordsize

- Extend wordsize by storing full word in parallel across multiple chips
- e.g., build 64-bit memory storing $2^n$ words from two 32-bit $2^n$-word chips

ADDR

$n$

ADDR

$2^n$ 32 bit words

32  Data In  W  Data Out  32

Store 32 high order bits of word

64-bit DATA-IN

$n$

ADDR

$2^n$ 32 bit words

32  Data In  W  Data Out  32

Store 32 low order bits of word

64-bit DATA-OUT

W

# Multi-chip memories: extend wordsize

- Extend wordsize by storing full word in parallel across multiple chips
- e.g., build 64-bit memory storing $2^n$ words from two 32-bit $2^n$-word chips

ADDR

$n$

ADDR

$2^n$ 32 bit words

32   Data      Data   32
     In    W   Out

Store 32 high order bits of word

64-bit DATA-IN

$n$

ADDR

$2^n$ 32 bit words

32   Data      Data   32
     In    W   Out

Store 32 low order bits of word
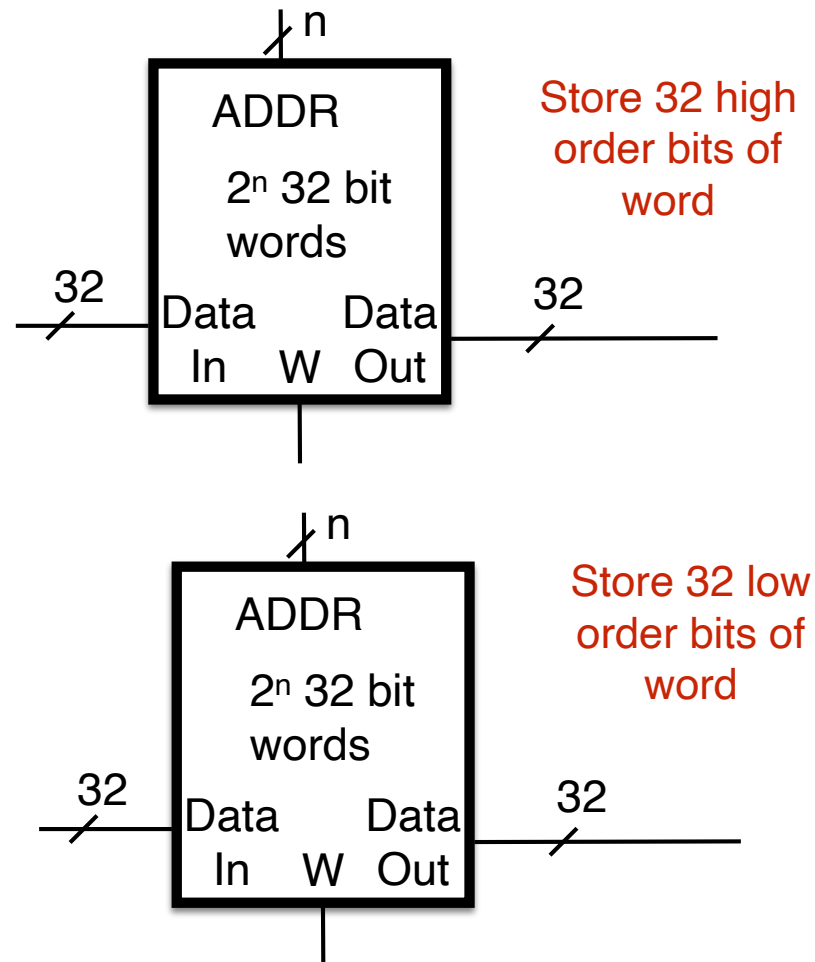
64-bit DATA-OUT

W

# Multi-chip memories: extend wordsize

- Extend wordsize by storing full word in parallel across multiple chips
- e.g., build 64-bit memory storing $2^n$ words from two 32-bit $2^n$-word chips

ADDR

$n$

ADDR

$2^n$ 32 bit words

32

Data In  W  Out  Data

32

Store 32 high order bits of word

64-bit DATA-IN

$n$

ADDR

$2^n$ 32 bit words

32

Data In  W  Out  Data

32

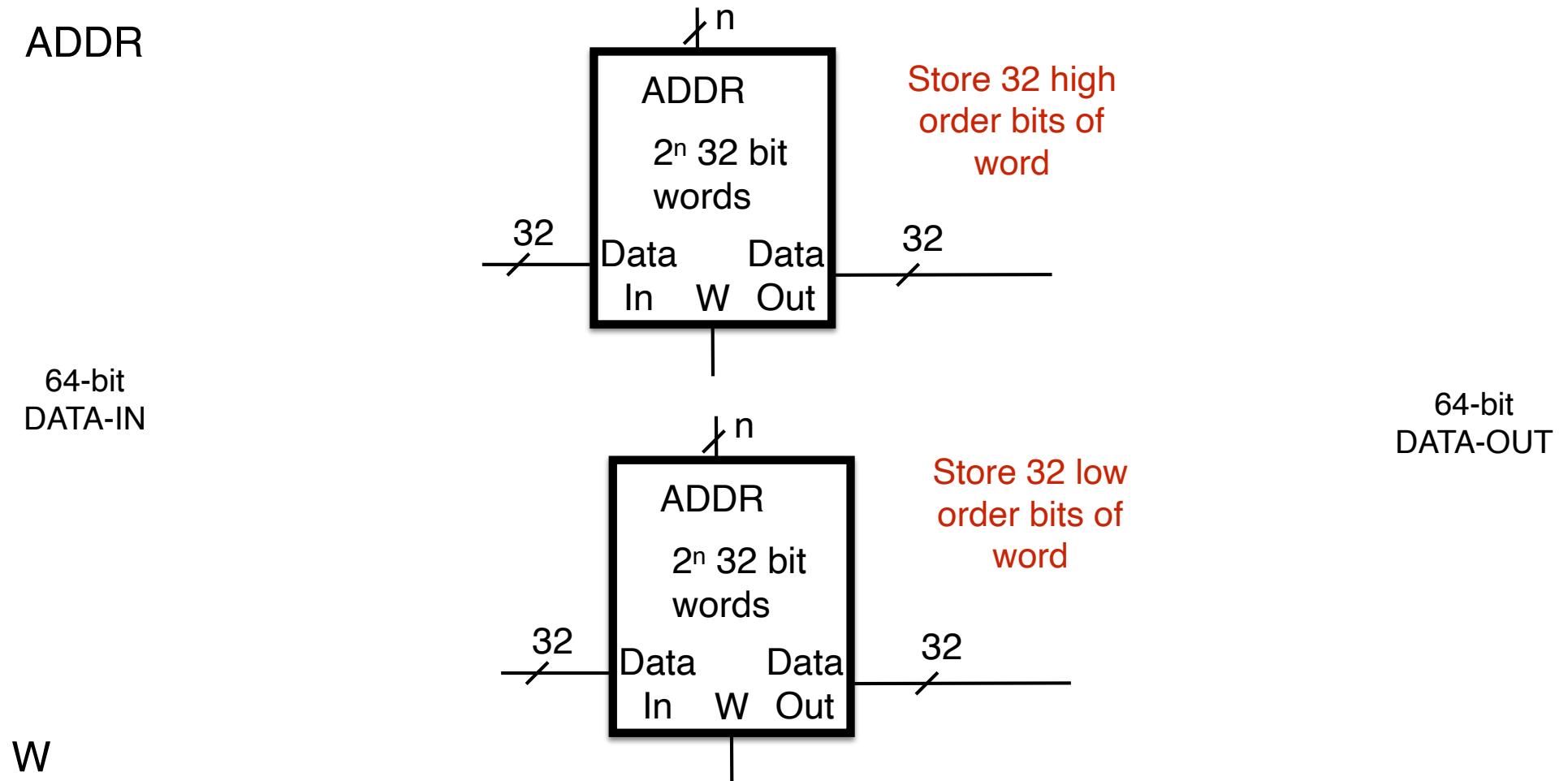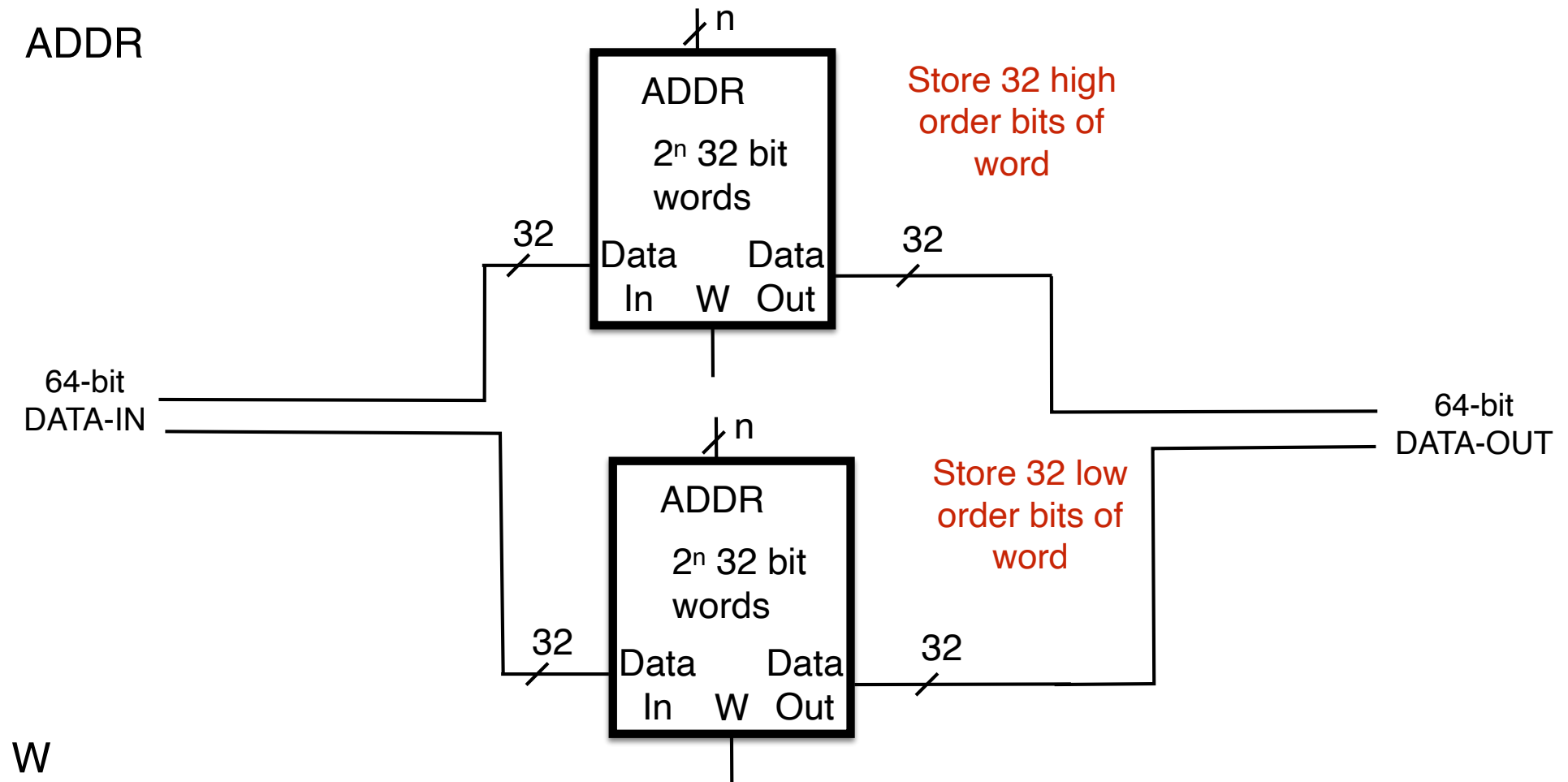Store 32 low order bits of word

64-bit DATA-OUT

W

# Multi-chip memories: extend wordsize
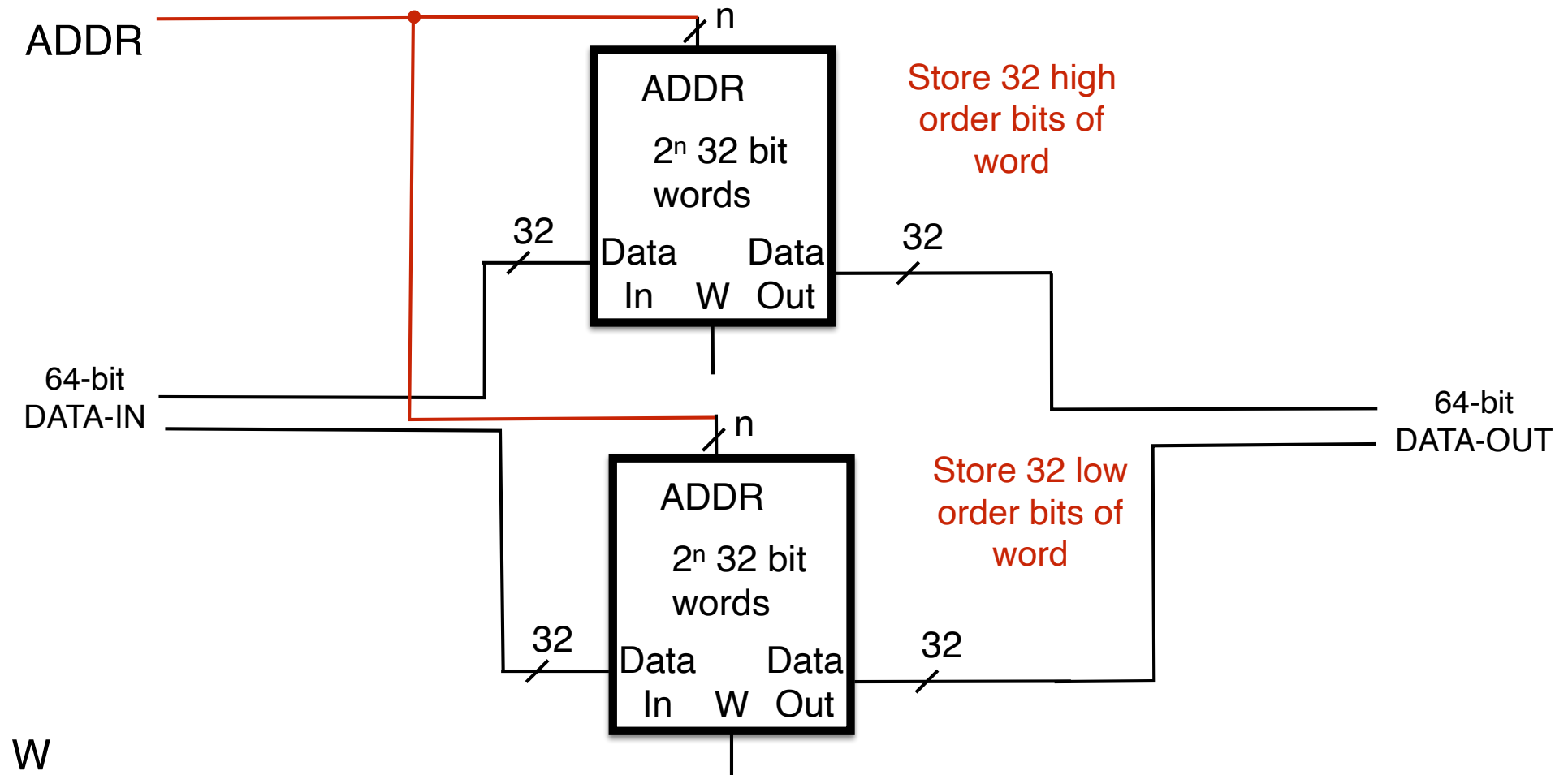
- Extend wordsize by storing full word in parallel across multiple chips
- e.g., build 64-bit memory storing $2^n$ words from two 32-bit $2^n$-word chips



ADDR

$n$

ADDR

$2^n$ 32 bit words

32

Data In   W   Out   Data

32

Store 32 high order bits of word

64-bit DATA-IN

$n$

ADDR

$2^n$ 32 bit words

32

Data In   W   Out   Data

32

Store 32 low order bits of word

64-bit DATA-OUT

W

# Multi-chip memories: extend address space

- Suppose each chip stores $2^n$ 32-bit words (has $2^n$ n-bit addresses)

$$n$$

```
          ADDR

          2ⁿ 32 bit
           words

32   Data      Data    32
      In    W  Out
```

# Multi-chip memories: extend address space

- Suppose each chip stores $2^n$ 32-bit words (has $2^n$ n-bit addresses)
- 2 of these chips can store $2^{n+1}$ 32-bit worse ($2^{n+1}$ (n+1)-bit addresses)

# Multi-chip memories: extend address space

- Suppose each chip stores $2^n$ 32-bit words (has $2^n$ n-bit addresses)
- 2 of these chips can store $2^{n+1}$ 32-bit worse ($2^{n+1}$ (n+1)-bit addresses)

$n$

ADDR

$2^n$ 32 bit words

Use for addresses where highest order bit of (n+1-bit) address is 0

n+1 bit address
$A_n$  $A_{n-1} \ldots A_1 A_0$

32

Data In    W    Data Out

32

$n$

ADDR

$2^n$ 32 bit words

32

Data In    W    Data Out

32

Use for addresses where highest order bit of (n+1-bit) address is 1

# Multi-chip memories: extend address space

- Suppose each chip stores $2^n$ 32-bit words (has $2^n$ n-bit addresses)
- 2 of these chips can store $2^{n+1}$ 32-bit worse ($2^{n+1}$ (n+1)-bit addresses)

$n$

**ADDR**

$2^n$ 32 bit words

Use for addresses where highest order bit of (n+1-bit) address is 0

n+1 bit address
$A_n$  $A_{n-1}...A_1 A_0$

32  Data       Data  32
In    W   Out

$n$

**ADDR**

$2^n$ 32 bit words

32  Data       Data  32
In    W   Out

Use for addresses where highest order bit of (n+1-bit) address is 1
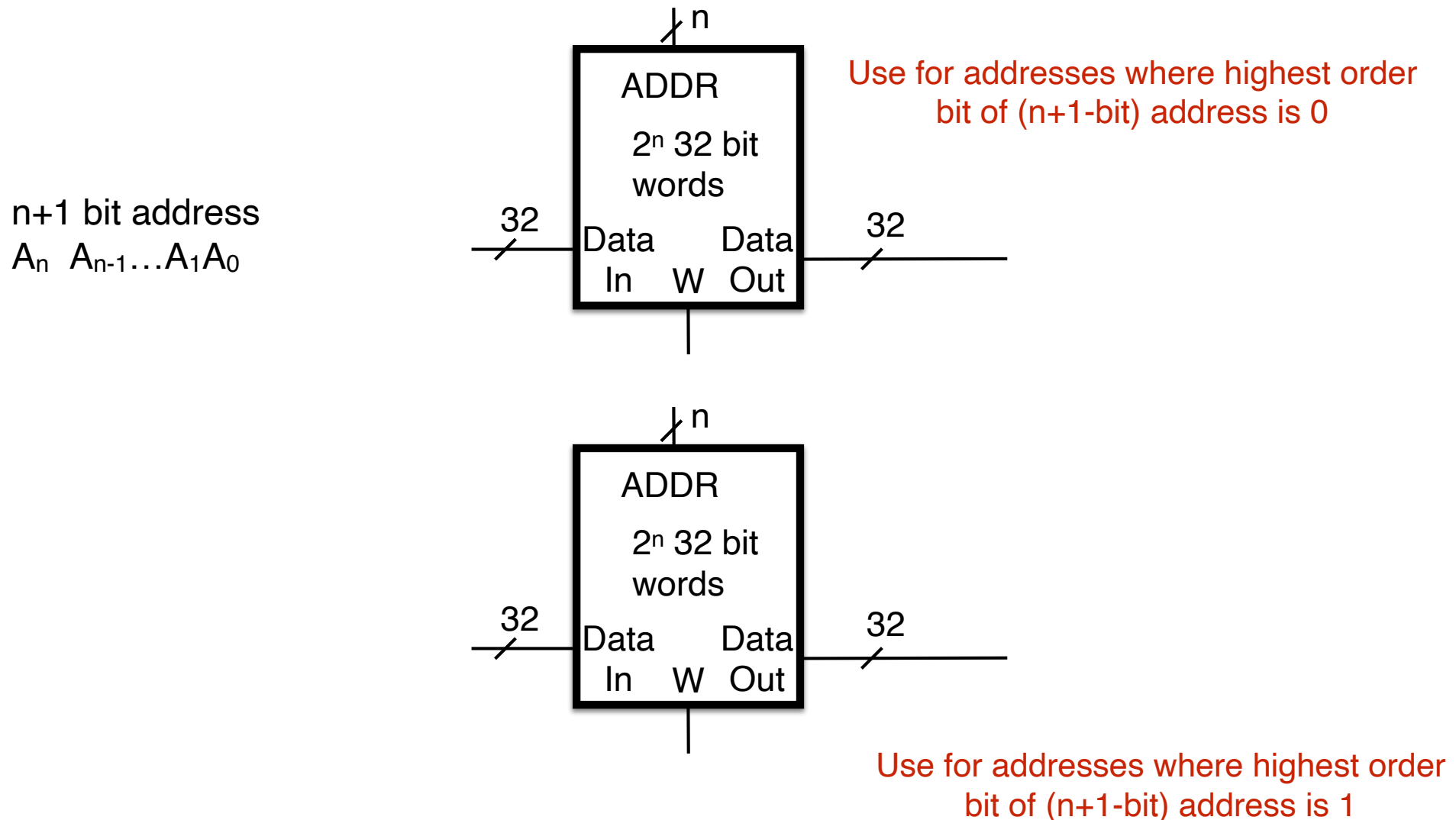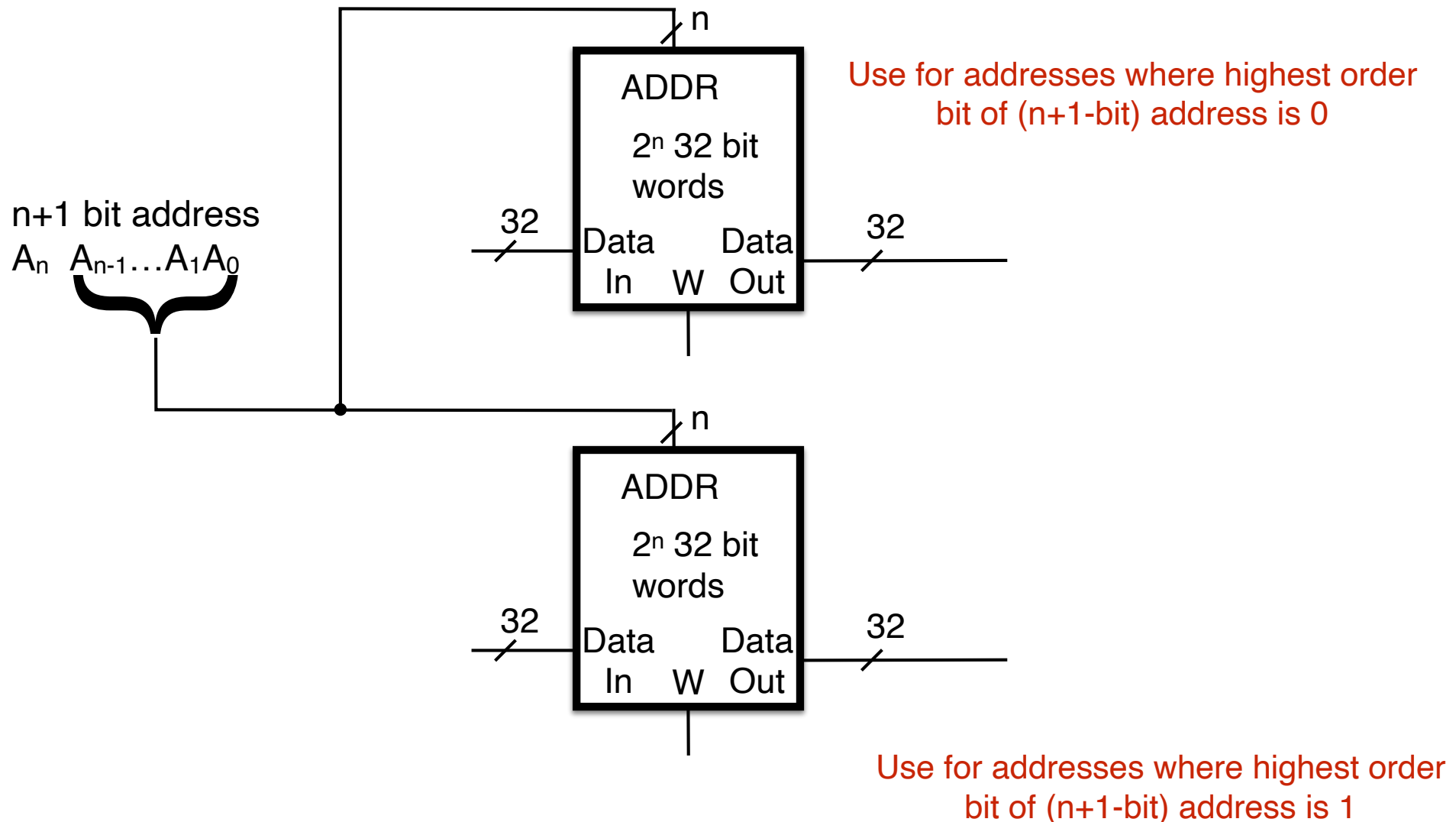
# Multi-chip memories: extend address space

- Suppose each chip stores $2^n$ 32-bit words (has $2^n$ n-bit addresses)
- 2 of these chips can store $2^{n+1}$ 32-bit worse ($2^{n+1}$ (n+1)-bit addresses)



Use for addresses where highest order bit of (n+1-bit) address is 0

n+1 bit address
$A_n$ $A_{n-1}\ldots A_1 A_0$

ADDR
$2^n$ 32 bit words
Data In   W   Data Out

2-to-1 MUX

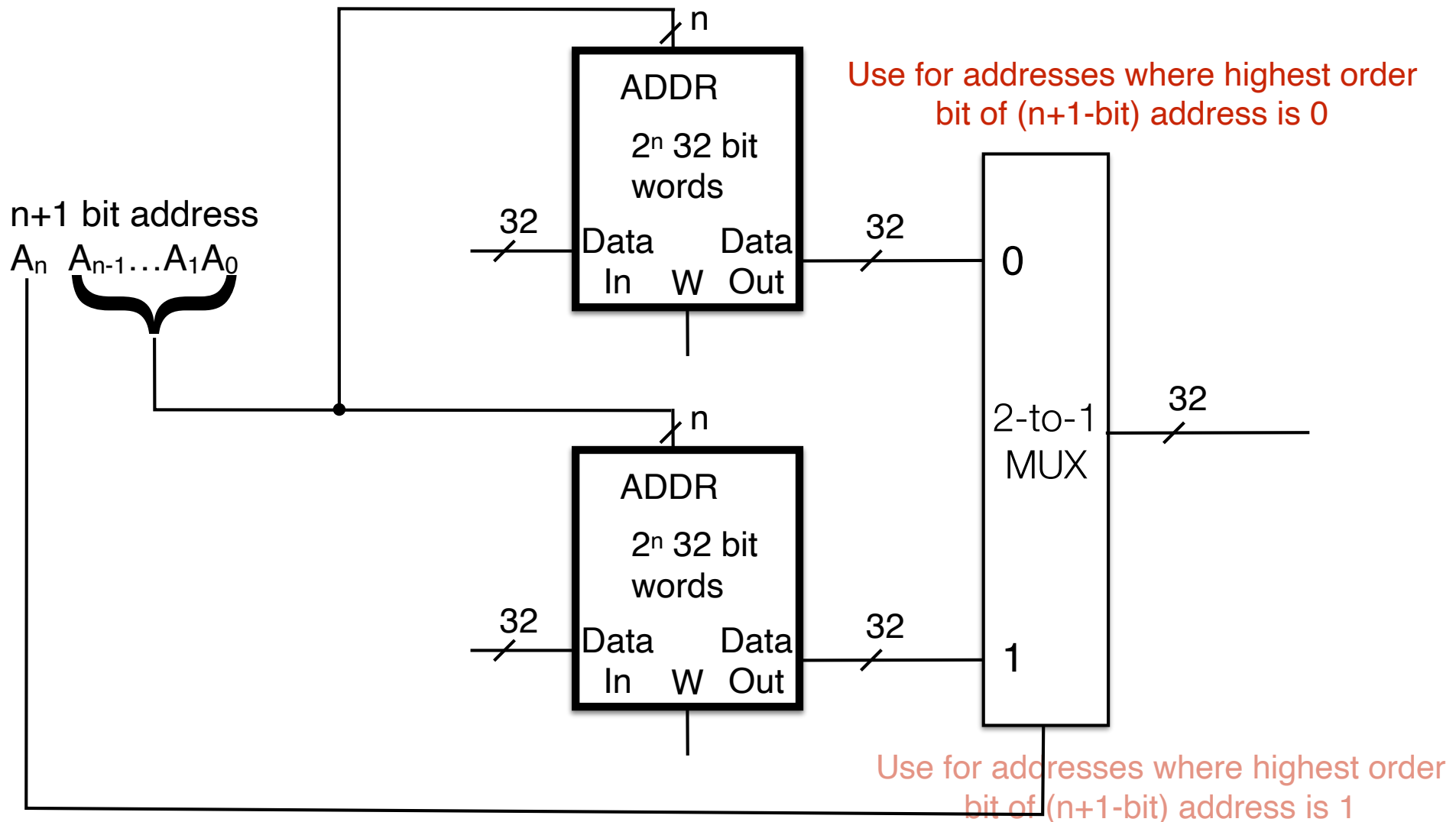Use for addresses where highest order bit of (n+1-bit) address is 1

# Multi-chip memories: extend address space

- Suppose each chip stores $2^n$ 32-bit words (has $2^n$ n-bit addresses)
- 2 of these chips can store $2^{n+1}$ 32-bit worse ($2^{n+1}$ (n+1)-bit addresses)



Data
(to write)

n+1 bit address
$A_n$ $A_{n-1}\ldots A_1 A_0$

W

**ADDR**

$2^n$ 32 bit words

Data In   W   Data Out

32        32

n

Use for addresses where highest order bit of (n+1-bit) address is 0

**ADDR**

$2^n$ 32 bit words

Data In   W   Data Out

32        32

n

0

1

2-to-1 MUX

32

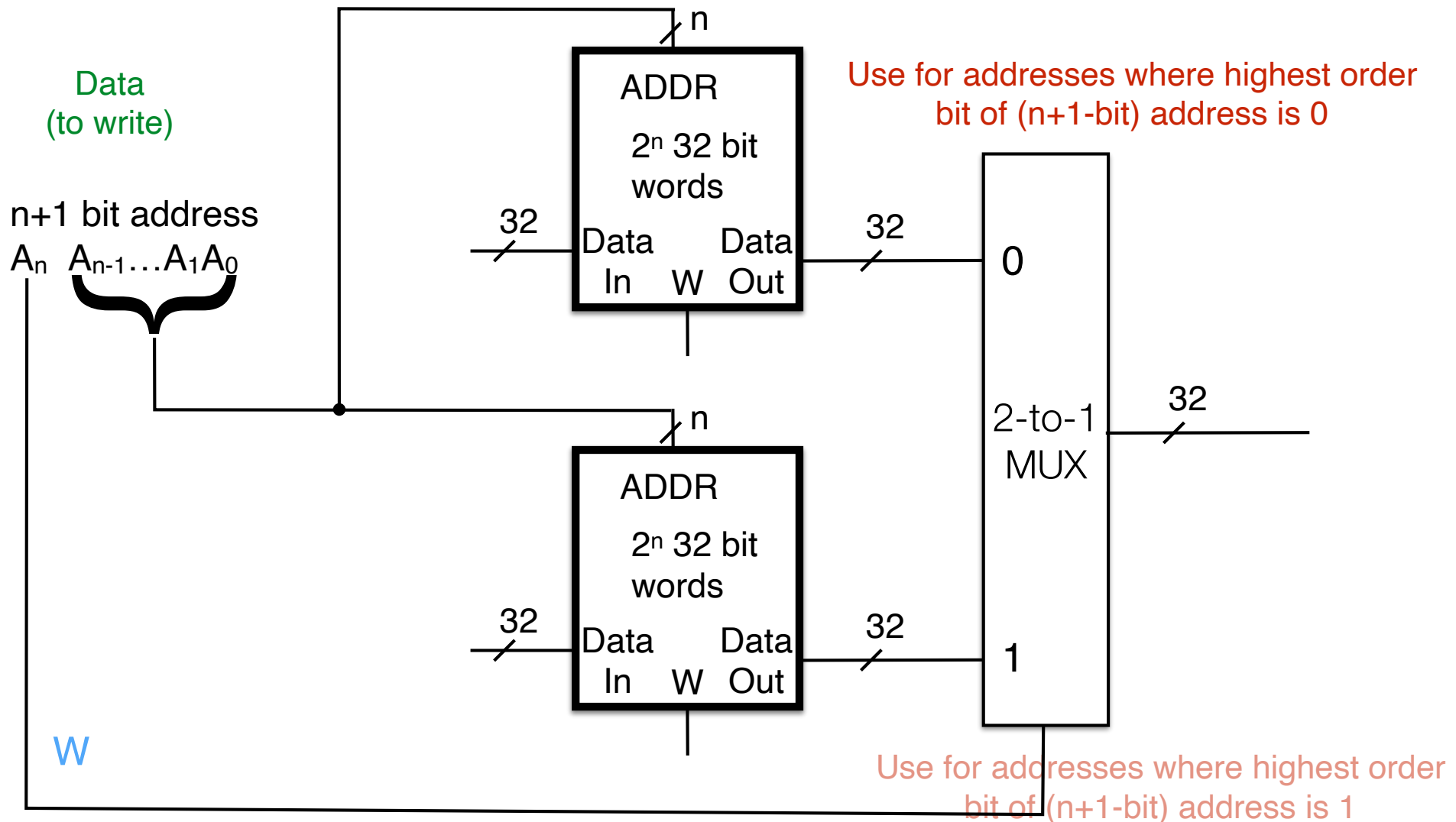Use for addresses where highest order bit of (n+1-bit) address is 1
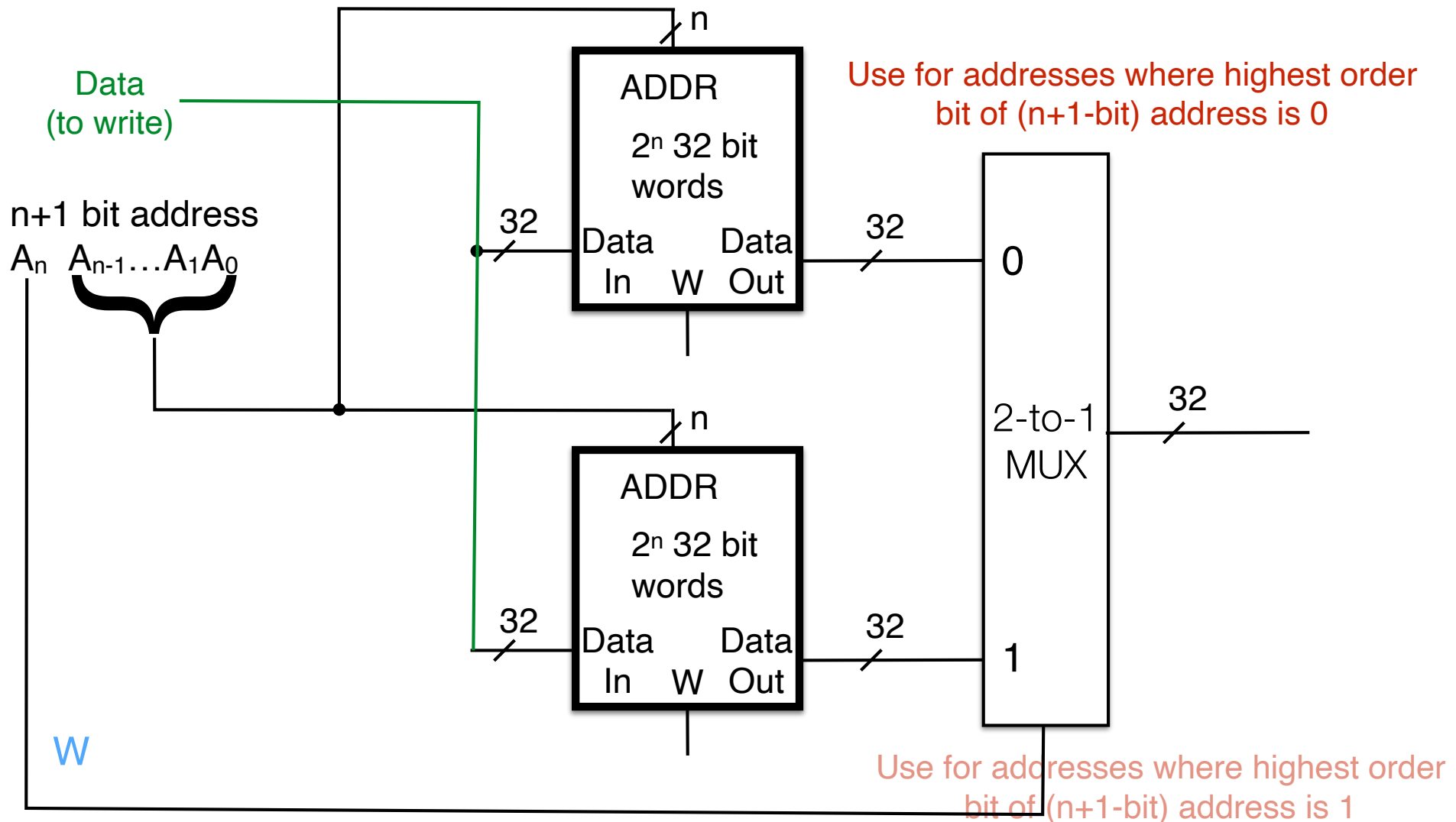
72

# Multi-chip memories: extend address space

- Suppose each chip stores $2^n$ 32-bit words (has $2^n$ n-bit addresses)
- 2 of these chips can store $2^{n+1}$ 32-bit worse ($2^{n+1}$ (n+1)-bit addresses)



Data (to write)

n+1 bit address
$A_n$ $A_{n-1}...A_1 A_0$

Use for addresses where highest order bit of (n+1-bit) address is 0

ADDR

$2^n$ 32 bit words

Data In    W    Data Out

2-to-1 MUX

W

Use for addresses where highest order bit of (n+1-bit) address is 1
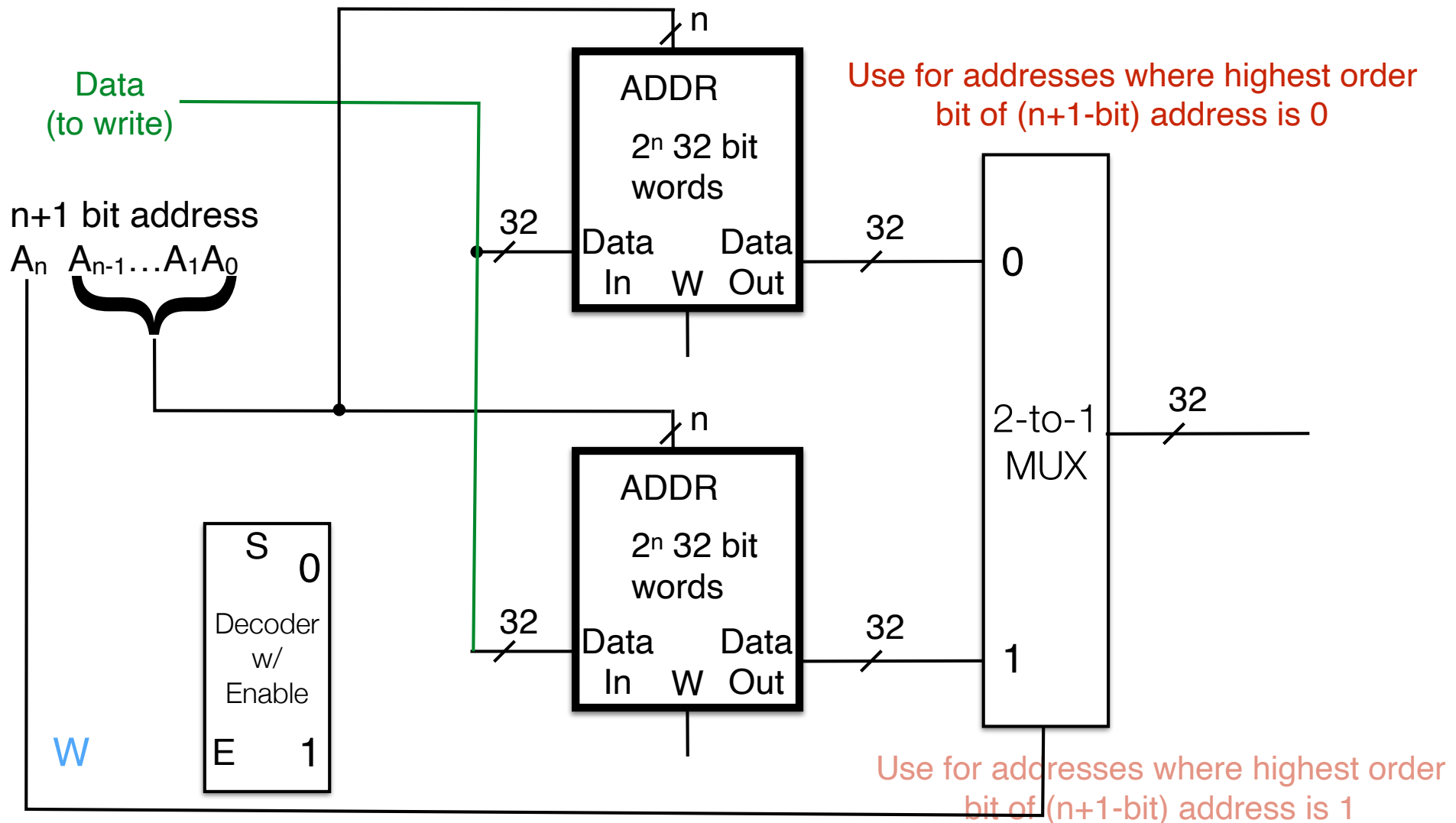
# Multi-chip memories: extend address space

- Suppose each chip stores $2^n$ 32-bit words (has $2^n$ n-bit addresses)
- 2 of these chips can store $2^{n+1}$ 32-bit worse ($2^{n+1}$ (n+1)-bit addresses)



Data
(to write)

n+1 bit address
$A_n$  $A_{n-1}\ldots A_1 A_0$

Use for addresses where highest order bit of (n+1-bit) address is 0

ADDR

$2^n$ 32 bit words

Data In    W    Data Out

32    0

2-to-1 MUX

32

S    0

Decoder w/ Enable

E    1

W

ADDR

$2^n$ 32 bit words

Data In    W    Data Out

32    1

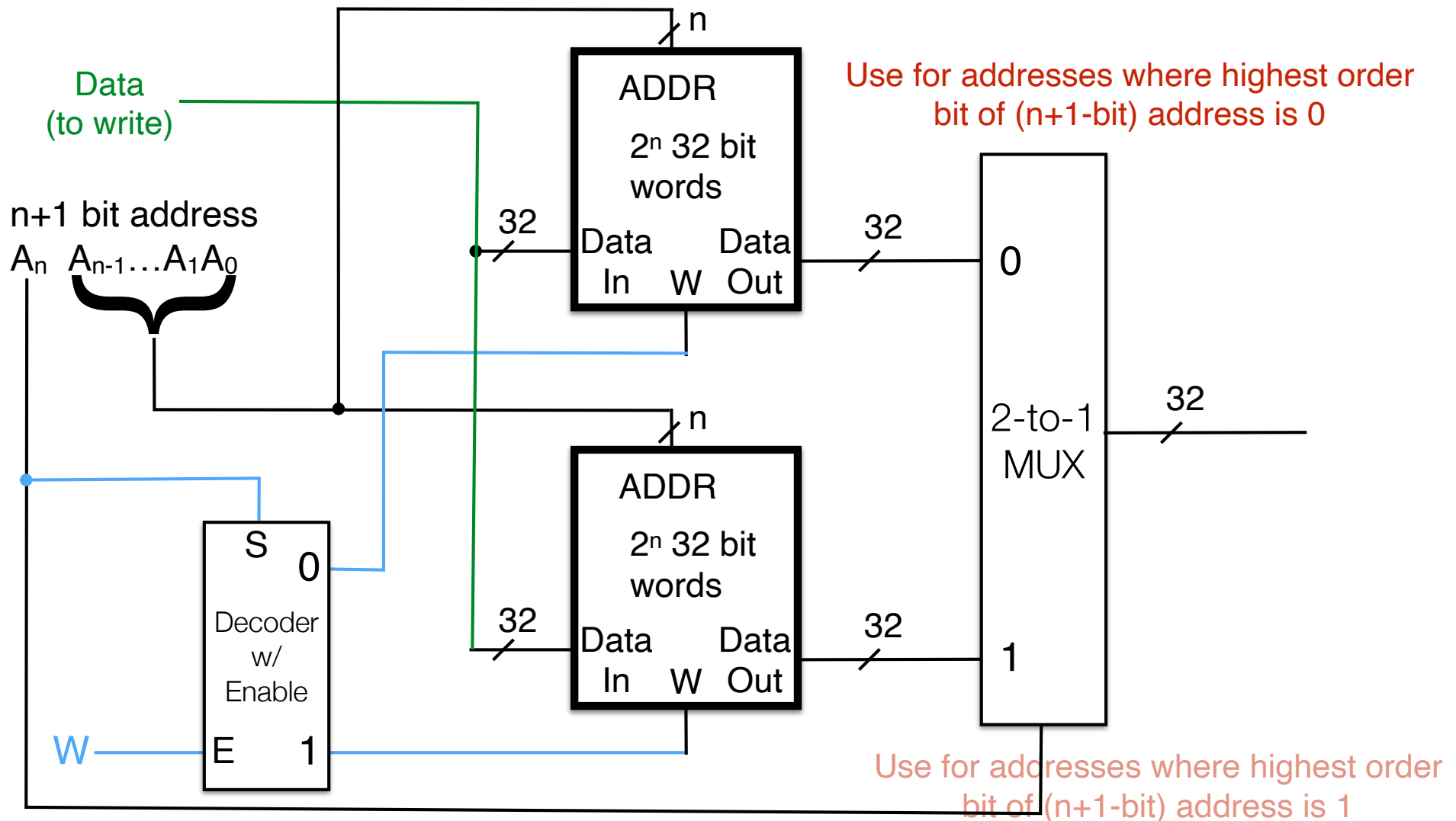Use for addresses where highest order bit of (n+1-bit) address is 1

# Multi-chip memories: extend address space

- Suppose each chip stores $2^n$ 32-bit words (has $2^n$ n-bit addresses)
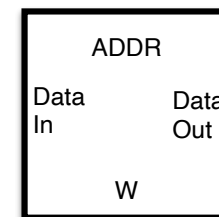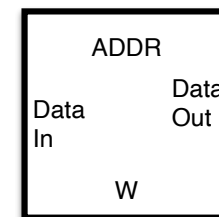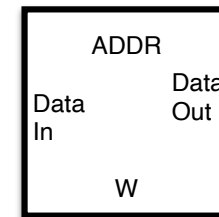- 2 of these chips can store $2^{n+1}$ 32-bit worse ($2^{n+1}$ (n+1)-bit addresses)

Data
(to write)

n+1 bit address
$A_n$  $A_{n-1}...A_1A_0$

n

ADDR

$2^n$ 32 bit words

Data In    W    Data Out

Use for addresses where highest order bit of (n+1-bit) address is 0

0

32

S    0

Decoder w/ Enable

E    1

W

ADDR

$2^n$ 32 bit words

Data In    W    Data Out

2-to-1 MUX

1

32

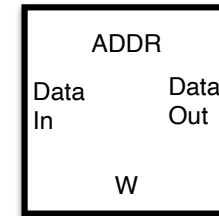Use for addresses where highest order bit of (n+1-bit) address is 1

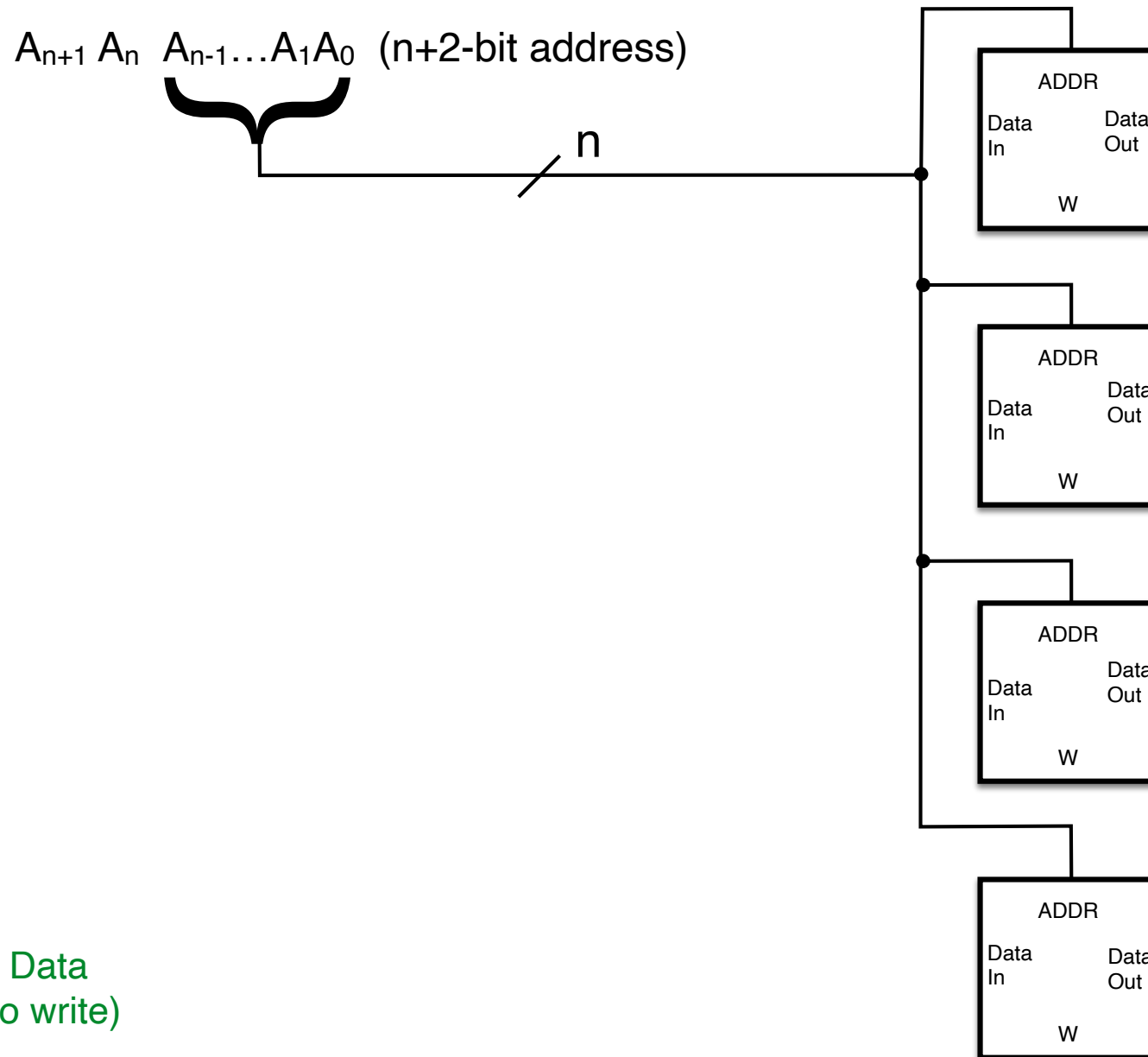One more example: quadruple address space

# Extending address beyond 2 chips

$A_{n+1} A_n A_{n-1} \ldots A_1 A_0$ (n+2-bit address)

W

Data
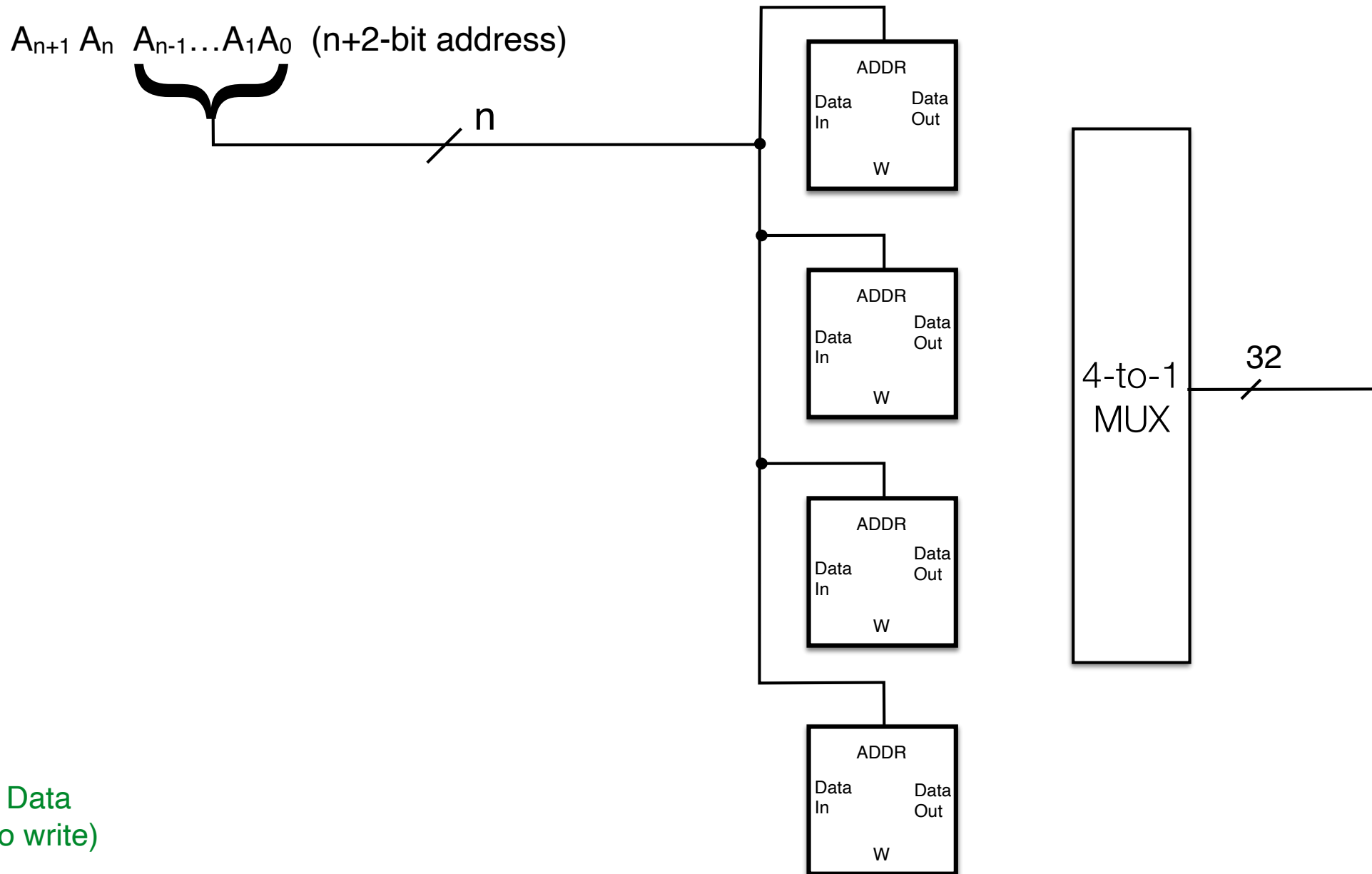(to write)

| ADDR | |
|------|------|
| Data In | Data Out |
| | W |

| ADDR | |
|------|------|
| Data In | Data Out |
| | W |

| ADDR | |
|------|------|
| Data In | Data Out |
| | W |

| ADDR | |
|------|------|
| Data In | Data Out |
| | W |

# Extending address beyond 2 chips

$A_{n+1}$ $A_n$ $A_{n-1} \ldots A_1 A_0$ (n+2-bit address)

n

W

Data
(to write)

ADDR

Data In    Data Out

W

ADDR

Data In    Data Out

W

ADDR

Data In    Data Out

W

ADDR

Data In    Data Out

W

# Extending address beyond 2 chips



$A_{n+1} A_n A_{n-1}...A_1 A_0$ (n+2-bit address)

n

W

Data
(to write)

ADDR
Data In
Data Out
W

ADDR
Data In
Data Out
W

ADDR
Data In
Data Out
W

ADDR
Data In
Data Out
W

4-to-1 MUX

32

# Extending address beyond 2 chips



$A_{n+1}$ $A_n$ $A_{n-1}\ldots A_1 A_0$ (n+2-bit address)

2

n

ADDR

Data In    Data Out    32

W

32

32

ADDR

Data In    Data Out    32

W

W

ADDR

Data In    Data Out    32

W

4-to-1 MUX    32

ADDR

Data In    Data Out    32

W

Data (to write)
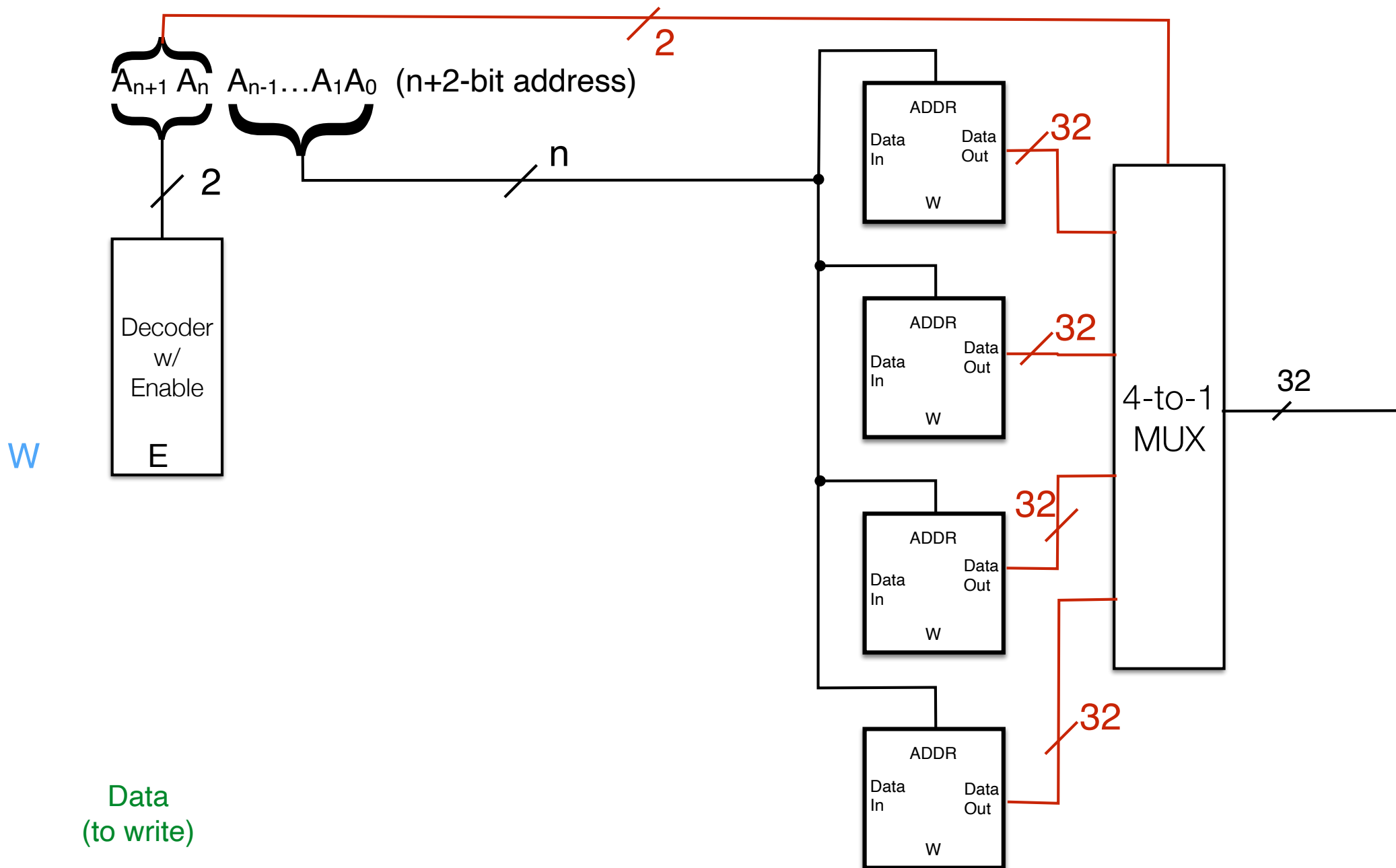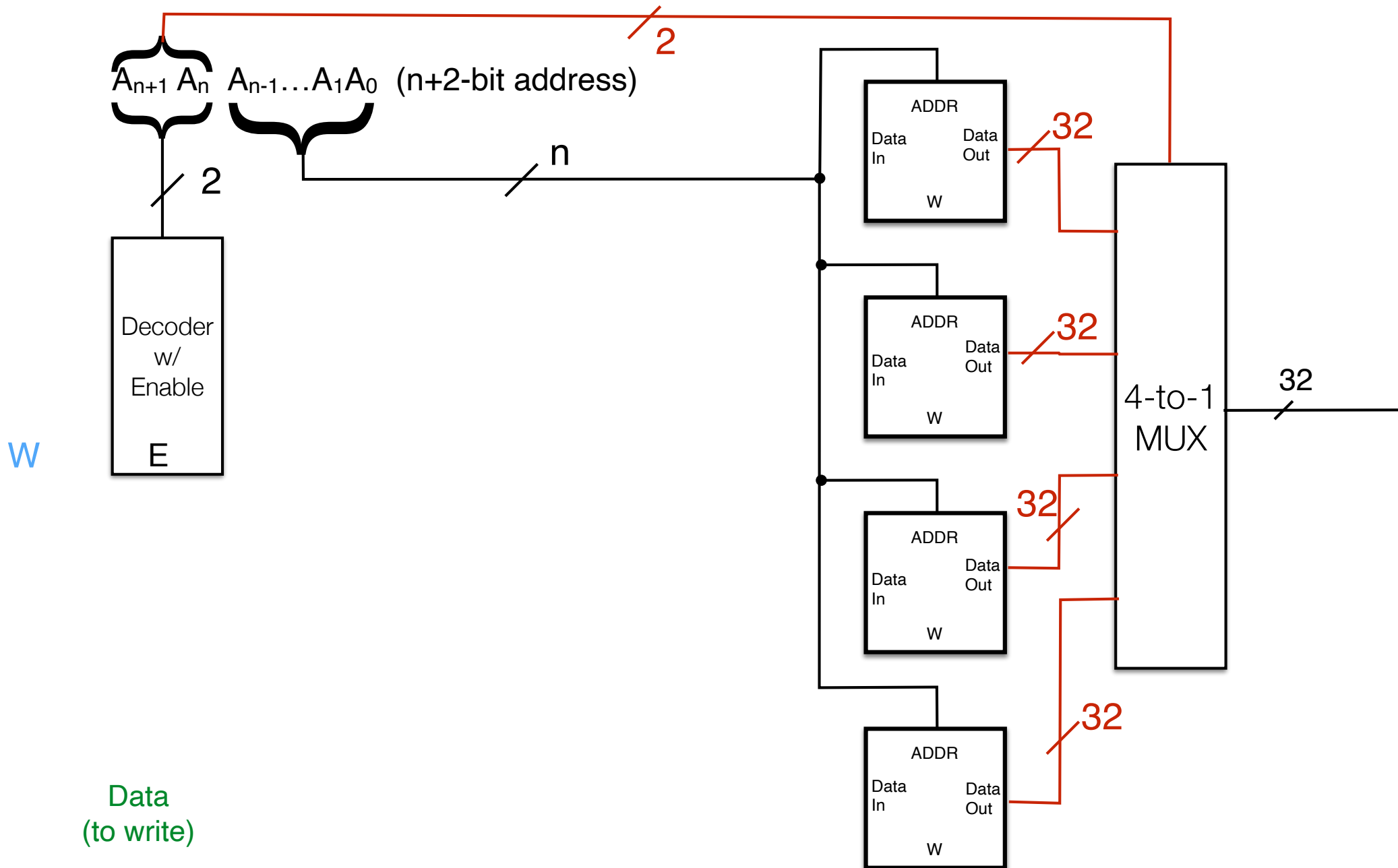
# Extending address beyond 2 chips

# Extending address beyond 2 chips



$A_{n+1} A_n \; A_{n-1} \ldots A_1 A_0$ (n+2-bit address)

2

n

2

Decoder w/ Enable

E

W

Data (to write)

ADDR
Data In
Data Out
W
32

ADDR
Data In
Data Out
W
32

ADDR
Data In
Data Out
W
32

ADDR
Data In
Data Out
W
32

4-to-1 MUX

32

# Extending address beyond 2 chips



$A_{n+1} A_n$  $A_{n-1}\ldots A_1 A_0$  (n+2-bit address)

2

n

Decoder w/ Enable

E

W

Data (to write)

ADDR
Data In    Data Out    32
W

ADDR
Data In    Data Out    32
W

ADDR
Data In    Data Out    32
W

ADDR
Data In    Data Out    32
W

4-to-1 MUX    32

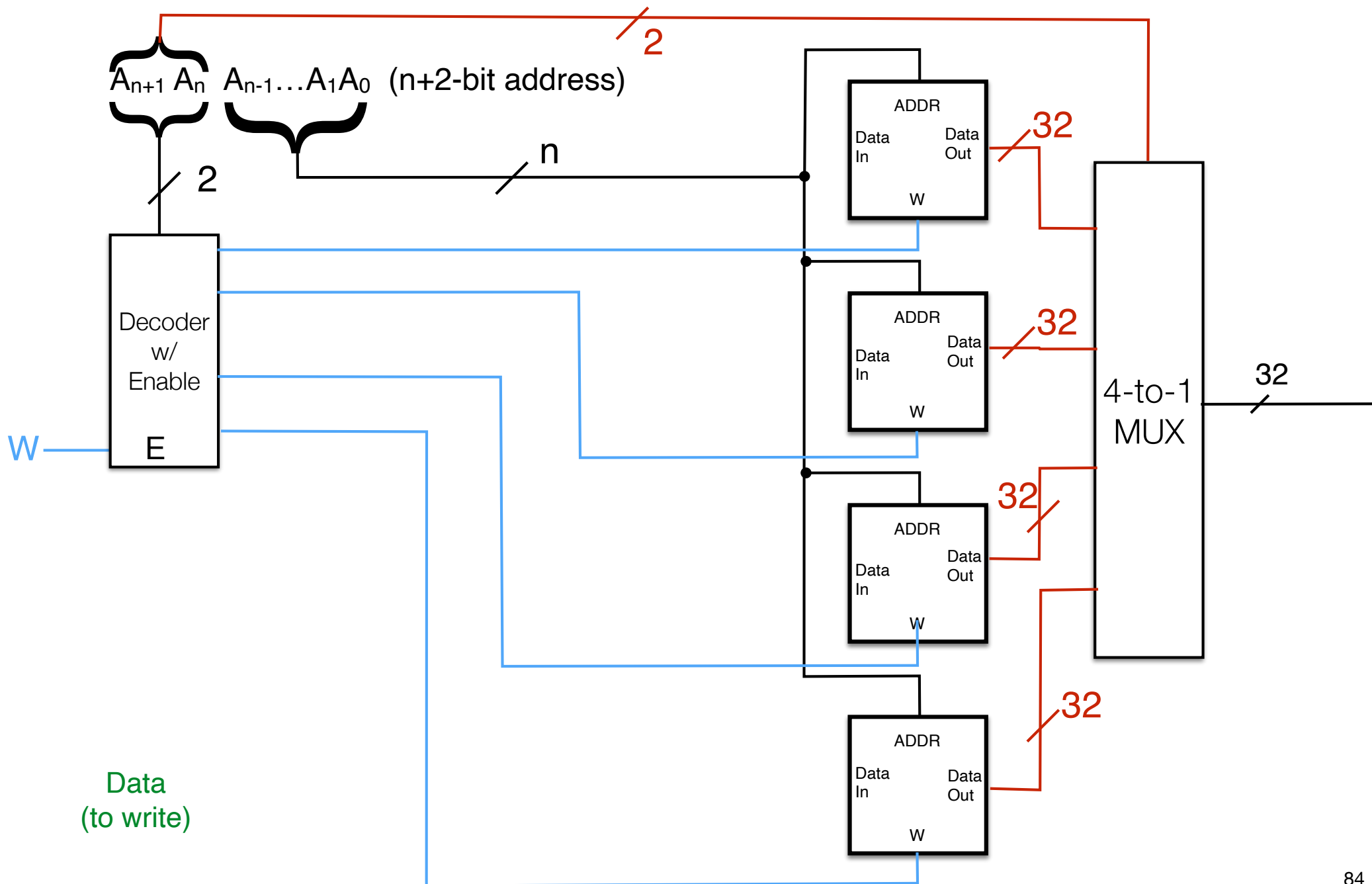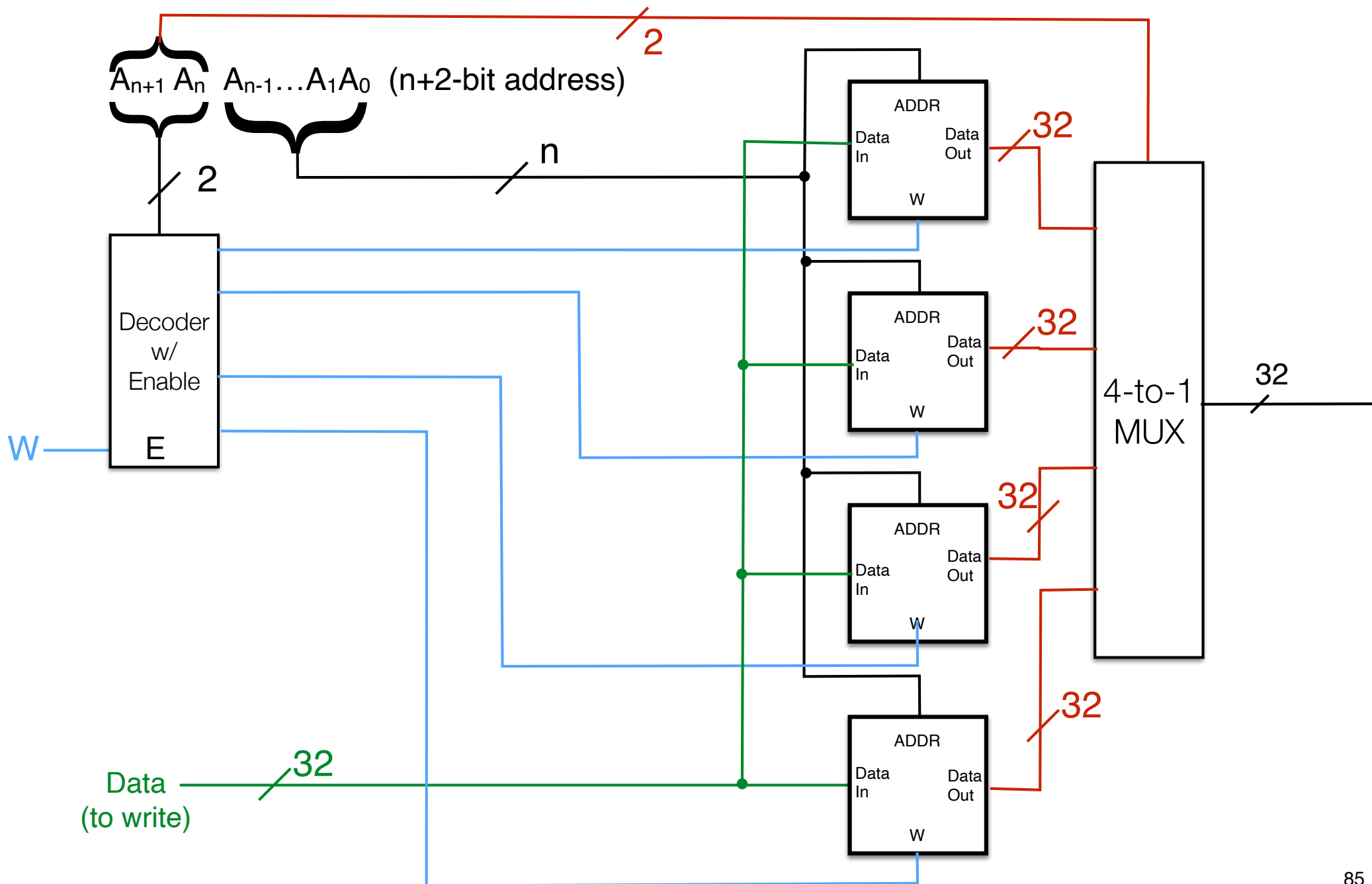# Extending address beyond 2 chips

# Extending address beyond 2 chips

# Some Timing Issues with Memory

# Memory Timing: Write example

- Even though memory not "on the clock", timing still an issue:

  - inputs to memory are "on the clock"

  - Real memory takes several clock cycles to access (will deal with this later via cache)

  - must first be properly enabled for reading or writing before data is transferred

    - Write:

      - Appropriate address chosen & Data to write fed in

      - Wait for address & Data to stabilize

      - Activate write

    - Read:

      - Appropriate address chosen

      - Wait for address & corresponding data out to stabilize
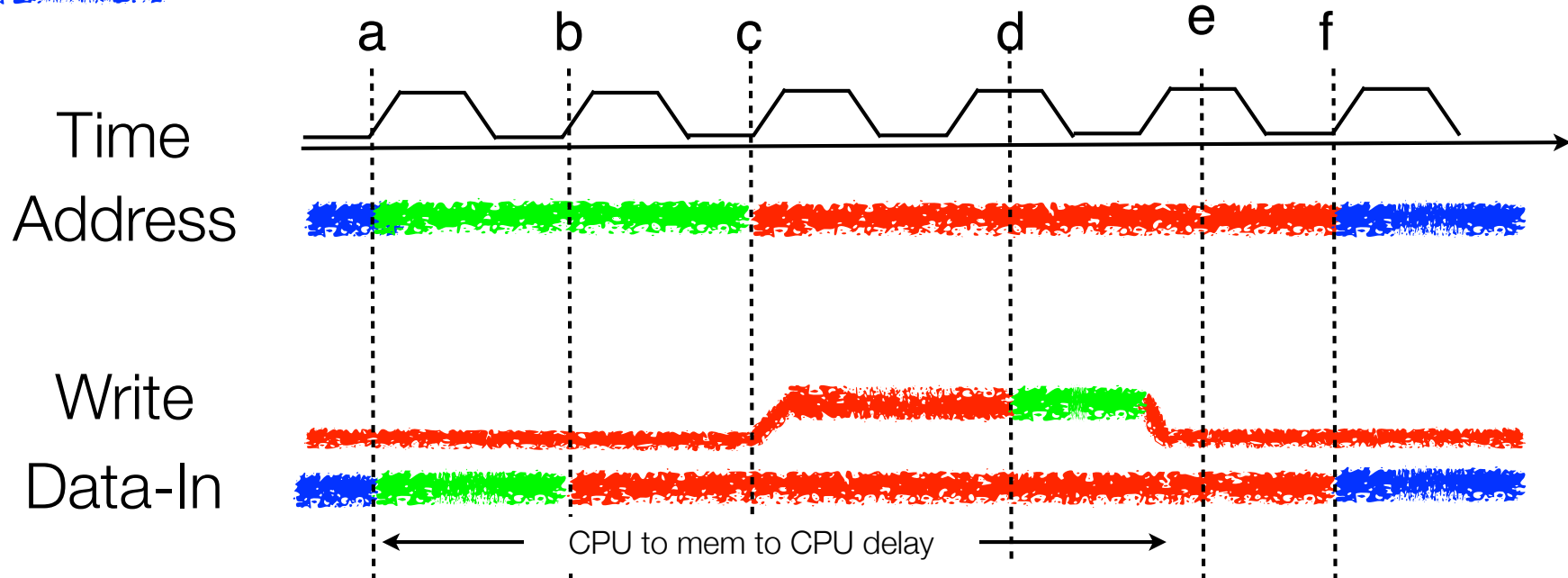
# Memory Timing: Write example

- Suppose:

  - 2 clock cycles to set address

  - 1 clock cycle to prepare data for writing

  - Write needs 1.5 clock cycles to ensure writing complete

a. Write request (address + data to write) flows to chip

b. Data to write stabilized

c. Address stabilized

d. Memory with new data stabilized

e. Memory "frozen" (no writing)

f. New instruction, new data

Correct Value arrives

Correct value "stabilized" in circuit

Not correct value for current write



CPU to mem to CPU delay

# Memory Timing: Read example

- Suppose:

  - 2 clock cycles to set address

  - 1 clock for data to flow from address to Data Out

a. Read request (address) flows to chip

b. Address stable, data from address flowing out

c. Data Out stable

d. New instruction, new data (read further at own risk…)