

Computer Architecture

Daniel Page

Department of Computer Science,
University Of Bristol,
Merchant Venturers Building,
Woodland Road,
Bristol, BS8 1UB. UK.
(csdsp@bristol.ac.uk)

October 14, 2024

Keep in mind there are *two* PDFs available (of which this is the latter):

1. a PDF of examinable material used as lecture slides, and
2. a PDF of non-examinable, extra material:
 - ▶ the associated notes page may be pre-populated with extra, written explanation of material covered in lecture(s), plus
 - ▶ anything with a “grey’ed out” header/footer represents extra material which is useful and/or interesting but out of scope (and hence not covered).

Notes:

Notes:

- **Problem:** an n -bit register based on latches (resp. flip-flops) is limited in that
 1. each latch (resp. flip-flop) in the register needs a relatively large number of transistors, which limits the viable capacity (i.e., n), and
 2. the register is not addressable, i.e.,
 - an **address** (or **index**) allows dynamic rather than static reference to some stored datum, so
 - by analogy, in a C program

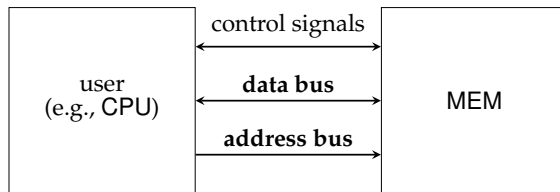
Listing	
1	<code>int A0, A1, A2, A3;</code>
2	
3	<code>A0 = 0;</code>
4	<code>A1 = 0;</code>
5	<code>A2 = 0;</code>
6	<code>A3 = 0;</code>

Listing	
1	<code>int A[4];</code>
2	
3	<code>A[0] = 0;</code>
4	<code>A[1] = 0;</code>
5	<code>A[2] = 0;</code>
6	<code>A[3] = 0;</code>

we *have* the left-hand side, but we *want* the right-hand side.

Notes:

- **Solution:** a **memory** component, i.e.,



such that

- MEM has a capacity of $n = 2^{n'}$ addressable words, and
- each such word is w bits (where $n \gg w$).

Notes:

► Agenda:

1. memory *cells*,
2. memory *devices*,

noting there are various ways to classify memories, e.g.,

1. volatility:
 - **volatile**, meaning the content is lost when the component is powered-off, or
 - **non-volatile**, meaning the content is retained even after the component is powered-off.
2. interface type:
 - **synchronous**, where a clock or pre-determined timing information synchronises steps, or
 - **asynchronous**, where a protocol synchronises steps.
3. access type:
 - random versus constrained (e.g., sequential) access to content,
 - **Random Access Memory (RAM)** which we can read from *and* write to, and
 - **Read Only Memory (ROM)** which, as suggested by the name, supports reads only.
4. ...

but we'll focus exclusively on a volatile, synchronous RAM.

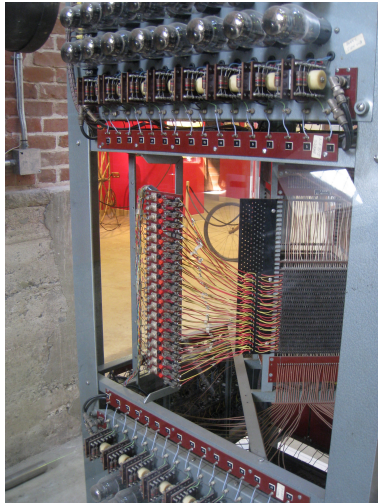
Notes:

An Aside: some history



- The EDSAC used **delay line** memory, where the rough idea is:
 - Each “line” is a tube of mercury (or something else in which sound waves propagate fairly slowly).
 - Put a speaker at one end to store sound waves into the line, and a microphone at the other to read them out.
 - Values are stored in the sense the corresponding waves take time to propagate; when they get to one end they are either replaced or fed back into the other.
- This is **sequential access** (cf. **random access**): you need to *wait* for the data you want to appear!

Notes:



- ▶ The Whirlwind used **magnetic-core** memory, where the rough idea is:
 - ▶ The memory is a matrix of small magnetic rings, or “cores”, which can be magnetically polarised to store values.
 - ▶ Wires are threaded through the cores to control them, i.e., to store or read values.
 - ▶ The magnetic polarisation is retained, so core memory is non-volatile!
- ▶ You might still hear main memory termed **core memory** (cf. **core dump**) which is a throw-back to this technology.

https://en.wikipedia.org/wiki/File:Project_Whirlwind_-_core_memory,_circa_1951_-_detail_1.JPG

Part 1: memory cells (1)

Comparison	Comparison
Static RAM (SRAM) is <ul style="list-style-type: none">▶ manufacturable in lower densities (i.e., smaller capacity),▶ more expensive to manufacture,▶ fast(er) access time (resp. lower access latency),▶ easy(er) to interface with,▶ ideal for latency-optimised contexts, e.g., as cache memory.	Dynamic RAM (DRAM) is <ul style="list-style-type: none">▶ manufacturable in higher densities (i.e., larger capacity),▶ less expensive to manufacture,▶ slow(er) access time (resp. higher access latency),▶ hard(er) to interface with,▶ ideal for capacity-optimised contexts, e.g., as main memory.

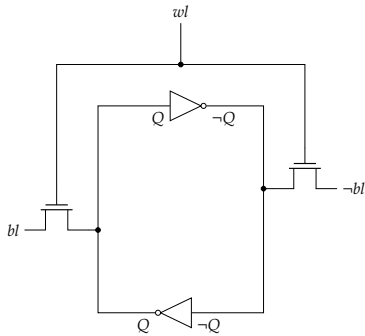
Notes:

Notes:

Part 1: memory cells (2)

An SRAM cell

Circuit



► Idea:

- internally, the cell is essentially two NOT gates,
- bl and $\neg bl$ are the **bit lines** (via which the state is accessed),
- wl is the **word line** (which controls access to the state),
- a “6T SRAM cell” requires 6 transistors (cf. ~ 20 or so for a D-type latch).

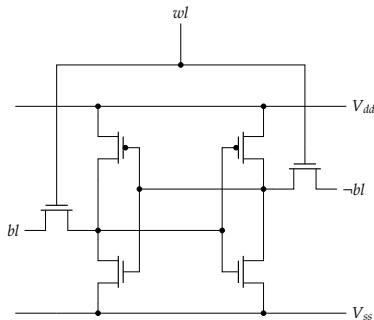
Notes:

- The initial NOT-based circuit might look odd, but clearly has two stable states: either $Q = 1$ and $\neg Q = 0$, or $Q = 0$ and $\neg Q = 1$. The transistors re-enforce each other; the state is maintained as long as the cell is powered-on.
- To read the cell we pre-charge $bl = 1$ and $\neg bl = 1$ then set $wl = 1$, after which $\neg bl$ (resp. bl) is discharged if state is 1 (resp. 0). To write x into the cell we pre-charge $bl = x$ and $\neg bl = \neg x$ then set $wl = 1$, after which the state matches x .
- The pre-charging steps are managed by extra **bit line conditioning**, the detail of which we ignore.

Part 1: memory cells (2)

An SRAM cell

Circuit



► Idea:

- internally, the cell is essentially two NOT gates,
- bl and $\neg bl$ are the **bit lines** (via which the state is accessed),
- wl is the **word line** (which controls access to the state),
- a “6T SRAM cell” requires 6 transistors (cf. ~ 20 or so for a D-type latch).

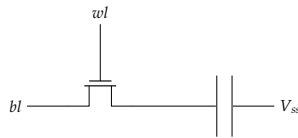
Notes:

- The initial NOT-based circuit might look odd, but clearly has two stable states: either $Q = 1$ and $\neg Q = 0$, or $Q = 0$ and $\neg Q = 1$. The transistors re-enforce each other; the state is maintained as long as the cell is powered-on.
- To read the cell we pre-charge $bl = 1$ and $\neg bl = 1$ then set $wl = 1$, after which $\neg bl$ (resp. bl) is discharged if state is 1 (resp. 0). To write x into the cell we pre-charge $bl = x$ and $\neg bl = \neg x$ then set $wl = 1$, after which the state matches x .
- The pre-charging steps are managed by extra **bit line conditioning**, the detail of which we ignore.

Part 1: memory cells (3)

A DRAM cell

Circuit



► Idea:

- internally, the cell is essentially one transistor and one capacitor,
- *bl* is the **bit line** (via which the state is accessed),
- *wl* is the **word line** (which controls access to the state),
- the capacitor
 1. discharges and charges (relatively) slowly,
 2. discharges when the cell is read, *and* also over time even if it's not read; this implies a need to **refresh** it.

Notes:

- To read the cell we set $wl = 1$, after which current flows (resp. does not flow) on *bl* if the capacitor is charged (resp. not charged) meaning the state is 1 (resp. 0). To write x into the cell we set $bl = x$ then $wl = 1$, after which the capacitor charges (resp. discharges) and the state matches x .
- The capacitor holds a tiny charge: this must be amplified to use as a driver in whatever circuit uses the cell.

Part 2: memory cells \leadsto memory devices (1)

► Concept: a **memory device** is constructed from (roughly) three components

1. a **memory array** (or matrix) of replicated cells with
 - r rows, and
 - c columnsmeaning a $(r \cdot c)$ -cell capacity,
 2. a **row decoder** which given an address (de)activates associated cells in that row, and
 3. a **column decoder** which given an address (de)selects associated cells in that column
- plus additional logic to allow use (depending on cell type), e.g.,
1. **bit line conditioning** to, e.g., ensure the bit lines are strong enough to be effective, and
 2. **sense amplifiers** to, e.g., ensure output from the array is usable.

Notes:

► (Typical, or exemplar) design: an SRAM device.

1. interface:
- auxiliary pin(s) for power and so on,

► D , a single 1-bit **data pin**,

► A , a collection of n' **address pins** where A_i is the i -th such pin,

► a **Chip Select (CS)** pin, which enables the device,

► a **Output Enable (OE)** pin, which signals the device is being read from,

► a **Write Enable (WE)** pin, which signals the device is being written to.

Notes:

- In some cases, D might be split into two separate D_{in} and D_{out} pins dedicated to input and output respectively.

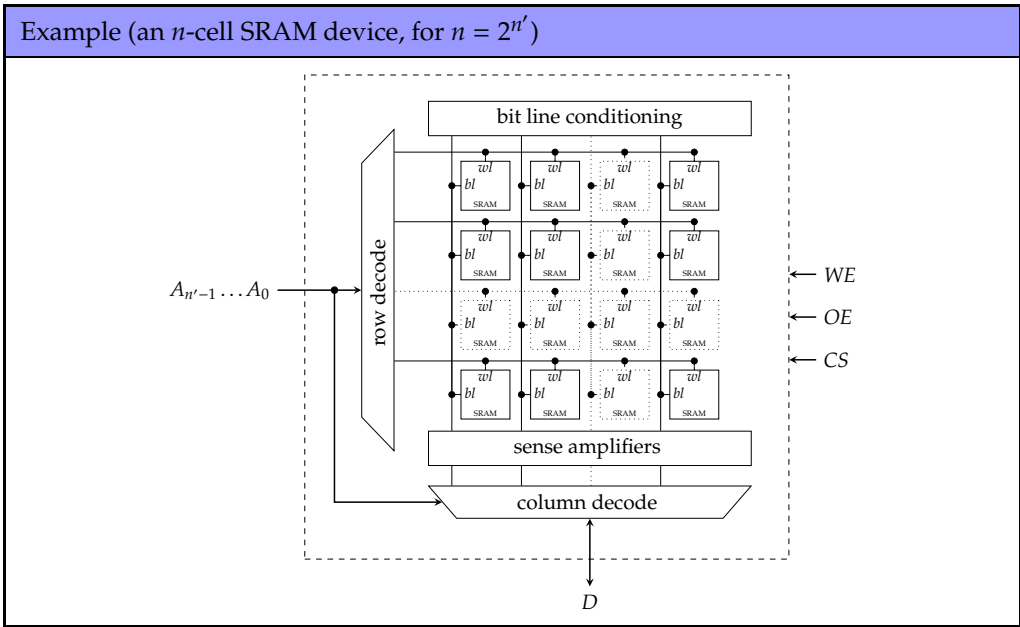
► (Typical, or exemplar) design: an SRAM device.

2. usage:

Algorithm (SRAM-READ)	Algorithm (SRAM-WRITE)
<div>Having performed the following steps</div> <div><div>► drive the address onto A,</div><div>► set $WE = \text{false}$, $OE = \text{true}$ and $CS = \text{true}$,</div></div> <div>1-bit of data is read and made available on D, then we set $CS = \text{false}$.</div>	<div>Having performed the following steps</div> <div><div>► drive the data onto D,</div><div>► drive the address onto A,</div><div>► set $WE = \text{true}$, $OE = \text{false}$ and $CS = \text{true}$,</div></div> <div>1-bit of data is written, then we set $CS = \text{false}$.</div>

Notes:

- In some cases, D might be split into two separate D_{in} and D_{out} pins dedicated to input and output respectively.



Notes:

ASI
Austin Semiconductor, Inc. **SRAM MT5C1001 Limited Availability**

1M x 1 SRAM
SRAM MEMORY ARRAY
AVAILABLE AS MILITARY SPECIFICATIONS
• MS0166-01/02/03
• MIL-STD-883C

FEATURES

- High Speed: 20, 25, 35, and 45
- Low power standby
- Single +5V (±0.7V) Power Supply
- Easy memory expansion with CE1 and CE0 options.
- All inputs and outputs are TTL compatible
- Three-state output

OPTIONS

Options	MARKING
• Timing	
20ns access	-20
25ns access	-25
35ns access	-35
45ns access	-45
55ns access	-55
70ns access	-70

Package(s)

Package	C	No.
Ceramic DIP (90mm)	C	No. 100
Ceramic LCC	LC	No. 207
Ceramic Flipchip	FC	No. 303
Ceramic Ball	CB	No. 303

Operating Temperature Ranges

Range	Marking
Industrial (-40°C to +85°C)	IT
Military (-55°C to +125°C)	M2

• 2V data retention/standby power L

GENERAL DESCRIPTION

The MT5C1001 employs low power, high performance silicon gate CMOS technology. Static design eliminates the need for external clocks or timing modes while CMOS circuitry reduces power consumption and provides for greater activity.

For flexibility in high-speed memory applications, ASI offers three output CE1 and output enable (OE) signals. These enhancements can place the outputs in High-Z for additional flexibility in system design. When these devices are reassembled when write enable (WE) and CE0 inputs are both LOW, reading is accomplished when WE is enabled (HIGH) while CE1 and CE0 go LOW. The devices offer a reduced power standby mode when disabled. The active system design to achieve low standby power requirements.

The "T" series provides an approximately 50 percent reduction in CMOS standby current (I_{CC1}) over the standard version.

All devices operate from a single +5V power supply and all inputs and outputs are fully TTL compatible.

For more products and information please visit our web site at www.austinsemiconductor.com

Rev. 2.0 1/00

ASI
Austin Semiconductor, Inc. **SRAM MT5C1001 Limited Availability**

FUNCTIONAL BLOCK DIAGRAM

TRUTH TABLE

MODE	CE1	WE	OUTPUT	POWER
READ	H	L	DATA	ACTIVE
WRITE	L	H	DATA	ACTIVE

PIN ASSIGNMENTS

PIN	ASSIGNMENT
A0-A10	Address Inputs
WE1	Write Enable
CE0	Chip Enable
D	Data Input
Q	Data Output
PD	Power Down
VCC	+5V Power Supply
VDD	Ground

Rev. 2.0 1/00

Notes:

► (Typical, or exemplar) design: a DRAM device.

1. interface:
- auxiliary pin(s) for power and so on,

► D , a single 1-bit **data pin**,

► A , a collection of $n'/2$ **address pins** where A_i is the i -th such pin,

► a **Chip Select (CS)** pin, which enables the device,

► a **Output Enable (OE)** pin, which signals the device is being read from,

► a **Write Enable (WE)** pin, which signals the device is being written to,

► a **Row Address Strobe (RAS)**, which controls the row buffer, and

► a **Column Address Strobe (CAS)**, which controls the column buffer.

Notes:

- In some cases, D might be split into two separate D_{in} and D_{out} pins dedicated to input and output respectively.
- An address is basically just an unsigned integer: if you have an n' -bit address, you can therefore use it to identify any one of $n = 2^{n'}$ addressable elements. An SRAM device will typically have n' address pins, so one can provide the address in a single step. An DRAM device will typically have less than n' address pins, so one must provide the address in multiple steps: here we see the use of $\frac{n'}{2}$ address pins, for example, which are used to provide an n' -bit address in 2 steps. The rationale for the latter is to cope with density. A DRAM device will typically has many addressable elements, but there is a physical limit on how many address pins are reasonable: by multiplexing, or sharing them, this issue is mitigated to some extent. Note that saying “typical” here is an important caveat, because it is not easy to give one definitive answer: there *are* other approaches and exceptions, for example which might be valid, even if less useful.

► (Typical, or exemplar) design: a DRAM device.

2. usage:

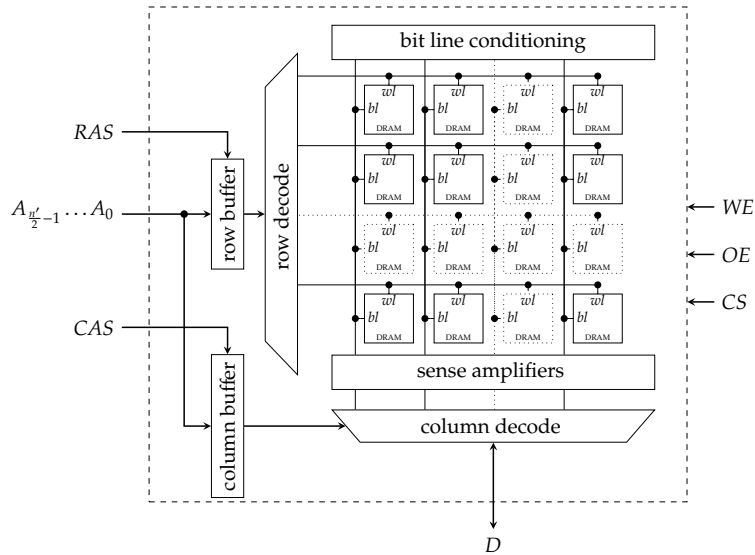
Algorithm (DRAM-READ)	Algorithm (DRAM-WRITE)
<div>Having performed the following steps</div> <div><div>► drive the row address onto A,</div><div>► set $RAS = \text{true}$ to latch row address,</div><div>► drive the column address onto A,</div><div>► set $CAS = \text{true}$ to latch column address,</div><div>► set $WE = \text{false}$, $OE = \text{true}$ and $CS = \text{true}$,</div></div> <div>1-bit of data is read and made available on D, then we set $CS = RAS = CAS = \text{false}$.</div>	<div>Having performed the following steps</div> <div><div>► drive the data onto D,</div><div>► drive the row address onto A,</div><div>► set $RAS = \text{true}$ to latch row address,</div><div>► drive the column address onto A,</div><div>► set $CAS = \text{true}$ to latch column address,</div><div>► set $WE = \text{false}$, $OE = \text{true}$ and $CS = \text{true}$,</div></div> <div>1-bit of data is written, then we set $CS = RAS = CAS = \text{false}$.</div>

Notes:

- In some cases, D might be split into two separate D_{in} and D_{out} pins dedicated to input and output respectively.
- An address is basically just an unsigned integer: if you have an n' -bit address, you can therefore use it to identify any one of $n = 2^{n'}$ addressable elements. An SRAM device will typically have n' address pins, so one can provide the address in a single step. An DRAM device will typically have less than n' address pins, so one must provide the address in multiple steps: here we see the use of $\frac{n'}{2}$ address pins, for example, which are used to provide an n' -bit address in 2 steps. The rationale for the latter is to cope with density. A DRAM device will typically has many addressable elements, but there is a physical limit on how many address pins are reasonable: by multiplexing, or sharing them, this issue is mitigated to some extent. Note that saying “typical” here is an important caveat, because it is not easy to give one definitive answer: there *are* other approaches and exceptions, for example which might be valid, even if less useful.

A DRAM device: implementation

Example (a n -cell DRAM memory device, for $n = 2^{n'}$)

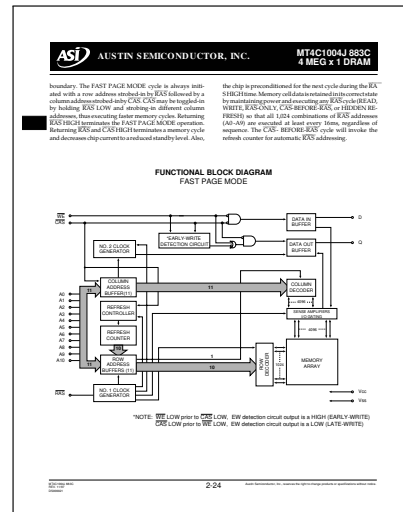
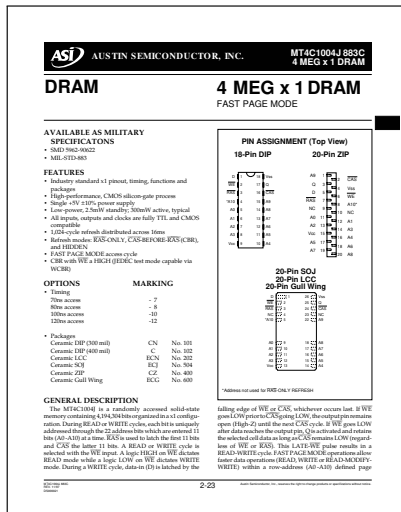


Notes:

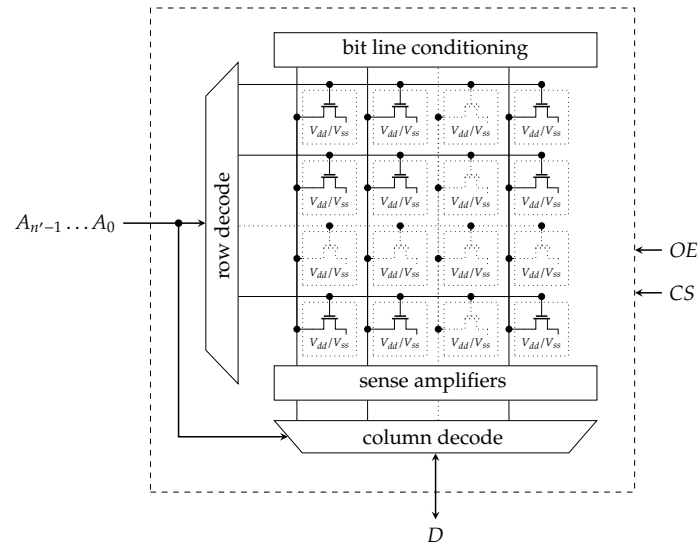
- We also need some logic, which is not shown here, to refresh the cells At a high level, the idea is that to cope with decay of content, we periodically read

Part 2: memory cells \leadsto memory devices (7)

A DRAM device: implementation



Notes:

Example (an n -cell ROM device, for $n = 2^{n'}$)

Notes:

Part 2: memory cells \leadsto memory devices (9)

► Concept:

- *externally*, the configuration of a device is described as something like

$$\delta \times \omega \times \beta$$

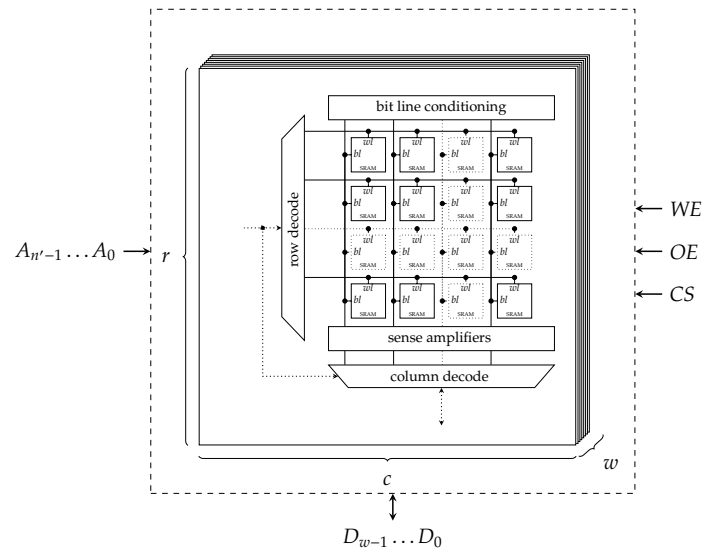
(plus maybe some timing information) where

- δ relates to capacity, usually measured in (large multiples of) bits,
 - ω describes the width of words, measured in bits, and
 - β is the number of internal **logical banks**.
- *internally*, this implies some organisational choices: for example,
1. for $\omega > 1$, we replicate the memory device internally to give ω arrays (each copy relates to one bit of a ω -bit word),
 2. for the arrays, r and c can be selected to match physical requirements (e.g., to get “square” or “thin” arrays), and
 3. for $\beta > 1$, each array is split into logical **banks**.

Notes:

Part 2: memory cells \leadsto memory devices (10)

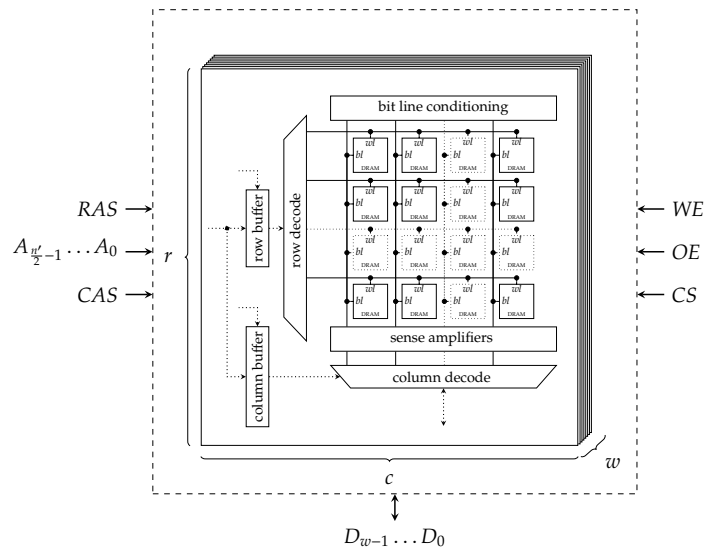
Example (from a 1-bit to w -bit SRAM device via replication)



Notes:

Part 2: memory cells \leadsto memory devices (11)

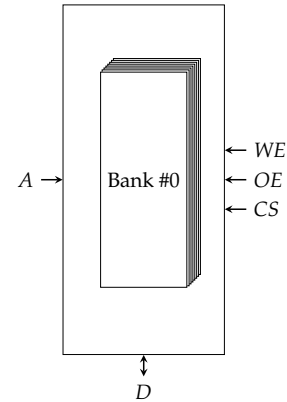
Example (from a 1-bit to w -bit DRAM device via replication)



Notes:

Example (a 512Mbit, 8-bit Device)

- ▶ A 64Mbit \times 8 \times 1 internal configuration:
 - ▶ The total capacity is 64Mbit \cdot 8 \cdot 1 = 512Mbit.
 - ▶ There is $\beta = 1$ logical bank: the bank consists of $\omega = 8$ arrays, each of which has $\delta = r \cdot c = 64 \cdot 1024 \cdot 1024$ cells.
 - ▶ Address x refers to cell x of each array; for example, $x = 42_{(10)}$ refers to cell 42 of each array, and therefore to an 8-bit word overall.



Notes:

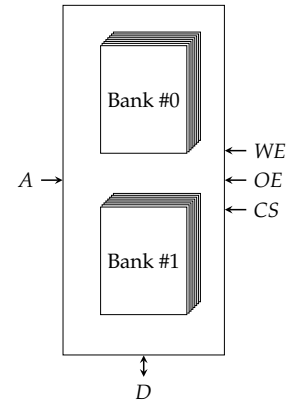
Example (a 512Mbit, 8-bit Device)

- ▶ A 32Mbit \times 8 \times 2 internal configuration:
 - ▶ The total capacity is 32Mbit \cdot 8 \cdot 2 = 512Mbit.
 - ▶ There are $\beta = 2$ logical banks: each bank consists of $\omega = 8$ arrays, each of which has $\delta = r \cdot c = 32 \cdot 1024 \cdot 1024$ cells.
 - ▶ Address x refers to

$$x \mapsto \begin{cases} \text{bank} & : x \bmod \beta \\ \text{cell} & : x \div \beta \end{cases}$$

of each array; for example,

- ▶ $x = 42_{(10)}$ refers to cell 21 of each array in bank 0, while
- ▶ $x = 43_{(10)}$ refers to cell 21 of each array in bank 1, each case therefore referring to 8-bit word overall.



Notes:

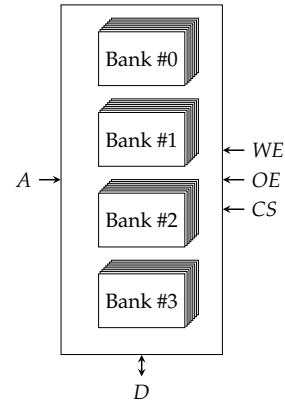
Example (a 512Mbit, 8-bit Device)

- ▶ A 16Mbit \times 8 \times 4 internal configuration:
 - ▶ The total capacity is 16Mbit \cdot 8 \cdot 4 = 512Mbit.
 - ▶ There are $\beta = 4$ logical banks: each bank consists of $\omega = 8$ arrays, each of which has $\delta = r \cdot c = 16 \cdot 1024 \cdot 1024$ cells.
 - ▶ Address x refers to

$$x \mapsto \begin{cases} \text{bank} & : x \bmod \beta \\ \text{cell} & : x \div \beta \end{cases}$$

of each array; for example,

- ▶ $x = 42_{(10)}$ refers to cell 10 of each array in bank 2, while
 - ▶ $x = 43_{(10)}$ refers to cell 21 of each array in bank 3,
- each case therefore referring to 8-bit word overall.



Notes:

Part 3: memory devices \leadsto memory modules (1)

- ▶ **Concept:** a **memory module** is essentially the combination of
 1. one or more memory devices, plus
 2. an interface which controls access.
 with two (physical) package types dominating:
 1. **Single Inline Memory Module (SIMM)**, which is (roughly) 1-sided, has less pins and a narrower word size, and
 2. **Dual Inline Memory Module (DIMM)**, which is (roughly) 2-sided, has more pins and a wider word size
 and are capable of housing different device types (e.g., EDO or SDRAM devices).

Notes:

► **Concept:**

- *externally*, the configuration of a module is described as something like

$$\Delta \times \Omega$$

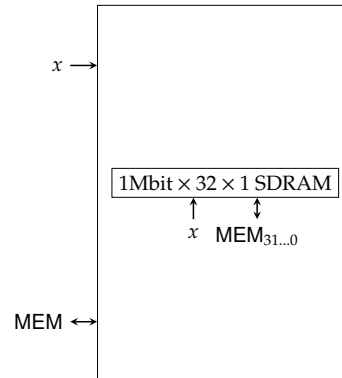
(plus maybe some timing information) where

- Δ relates to capacity, usually measured in (large multiplies of) bytes, and
 - Ω describes the width of words, measured in bits.
- *internally*, this implies some organisational choices: for example,
1. usually $\Omega > \omega$ so the module is “filled” using multiple devices to form one **physical bank**, and
 2. depending on the module type, the devices are organised into one or more **ranks**.

Notes:

Example (a 4MB, 32-bit SIMM)

- A $4\text{MB} \times 32$, 1-device configuration:
 - The device has a capacity of $1\text{Mbit} \cdot 32 \cdot 1 = 32\text{Mbit} = 4\text{MB}$, arranged into 32-bit words.
 - The bank, which has a 32-bit data-path, is filled by the single device: address x refers to the word MEM .



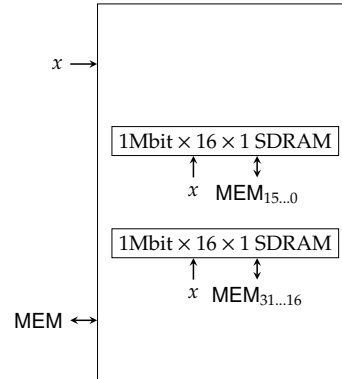
Notes:

Example (a 4MB, 32-bit SIMM)

- ▶ A $4\text{MB} \times 32$, 2-device configuration:
 - ▶ Each device has a capacity of $1\text{Mbit} \cdot 16 \cdot 1 = 16\text{Mbit} = 2\text{MB}$, arranged into 16-bit words.
 - ▶ Since there are 2 devices, the total capacity is $2 \cdot 2\text{MB} = 4\text{MB}$.
 - ▶ The bank, which has a 32-bit data-path, is filled by merging the devices: address x refers to

$$x \mapsto \begin{cases} \text{bits } 15 \dots 0 & \text{of MEM from device 0} \\ \text{bits } 31 \dots 16 & \text{of MEM from device 1} \end{cases}$$

merging each 16-bit part into the word MEM.



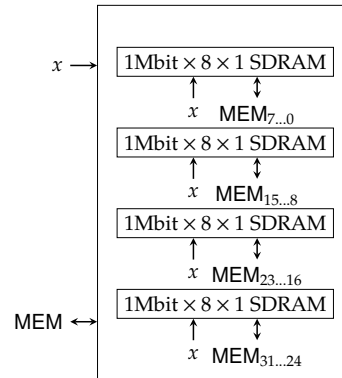
Notes:

Example (a 4MB, 32-bit SIMM)

- ▶ A $4\text{MB} \times 32$, 4-device configuration:
 - ▶ Each device has a capacity of $1\text{Mbit} \cdot 8 \cdot 1 = 8\text{Mbit} = 1\text{MB}$, arranged into 8-bit words.
 - ▶ Since there are 4 devices, the total capacity is $4 \cdot 1\text{MB} = 4\text{MB}$.
 - ▶ The bank, which has a 32-bit data-path, is filled by merging the devices: address x refers to

$$x \mapsto \begin{cases} \text{bits } 7 \dots 0 & \text{of MEM from device 0} \\ \text{bits } 15 \dots 8 & \text{of MEM from device 1} \\ \text{bits } 23 \dots 16 & \text{of MEM from device 2} \\ \text{bits } 31 \dots 24 & \text{of MEM from device 3} \end{cases}$$

merging each 8-bit part into the word MEM.



Notes:

Conclusions

► Take away points:

1. The initial goal was an n -element memory of w -bit words; the final solution is motivated by divide-and-conquer, i.e.,
 - 1.1 one or more channels, each backed by
 - 1.2 one or more physical banks, each composed from
 - 1.3 one or more devices, each composed from
 - 1.4 one or more logical banks, of
 - 1.5 one or more arrays, of
 - 1.6 many cells
2. The major complication is a large range of increasingly detailed options:
 - lots of parameters mean lots of potential trade-offs (e.g., between size, speed and power consumption),
 - need to take care of detail: there are so many cells, any minor change can have major consequences!
3. Even so, there is just one key concept: we have some cells, and however they are organised we just need to identify and use the right cells given some address.

Notes:

Additional Reading

- [Wikipedia: Computer Memory](https://en.wikipedia.org/wiki/Category:Computer_memory). URL: https://en.wikipedia.org/wiki/Category:Computer_memory.
- D. Page. “Chapter 8: Memory and storage”. In: *A Practical Introduction to Computer Architecture*. 1st ed. Springer, 2009.
- A.S. Tanenbaum and T. Austin. “Section 3.3.5: Memory chips”. In: *Structured Computer Organisation*. 6th ed. Prentice Hall, 2012.
- A.S. Tanenbaum and T. Austin. “Section 3.3.6: RAMs and ROMs”. In: *Structured Computer Organisation*. 6th ed. Prentice Hall, 2012.
- W. Stallings. “Chapter 5: Internal memory”. In: *Computer Organisation and Architecture*. 9th ed. Prentice Hall, 2013.

Notes:

References

- [1] [Wikipedia: Computer Memory](https://en.wikipedia.org/wiki/Category:Computer_memory). URL: https://en.wikipedia.org/wiki/Category:Computer_memory (see p. 67).
- [2] D. Page. “Chapter 8: Memory and storage”. In: *A Practical Introduction to Computer Architecture*. 1st ed. Springer, 2009 (see p. 67).
- [3] W. Stallings. “Chapter 5: Internal memory”. In: *Computer Organisation and Architecture*. 9th ed. Prentice Hall, 2013 (see p. 67).
- [4] A.S. Tanenbaum and T. Austin. “Section 3.3.5: Memory chips”. In: *Structured Computer Organisation*. 6th ed. Prentice Hall, 2012 (see p. 67).
- [5] A.S. Tanenbaum and T. Austin. “Section 3.3.6: RAMs and ROMs”. In: *Structured Computer Organisation*. 6th ed. Prentice Hall, 2012 (see p. 67).

Notes: