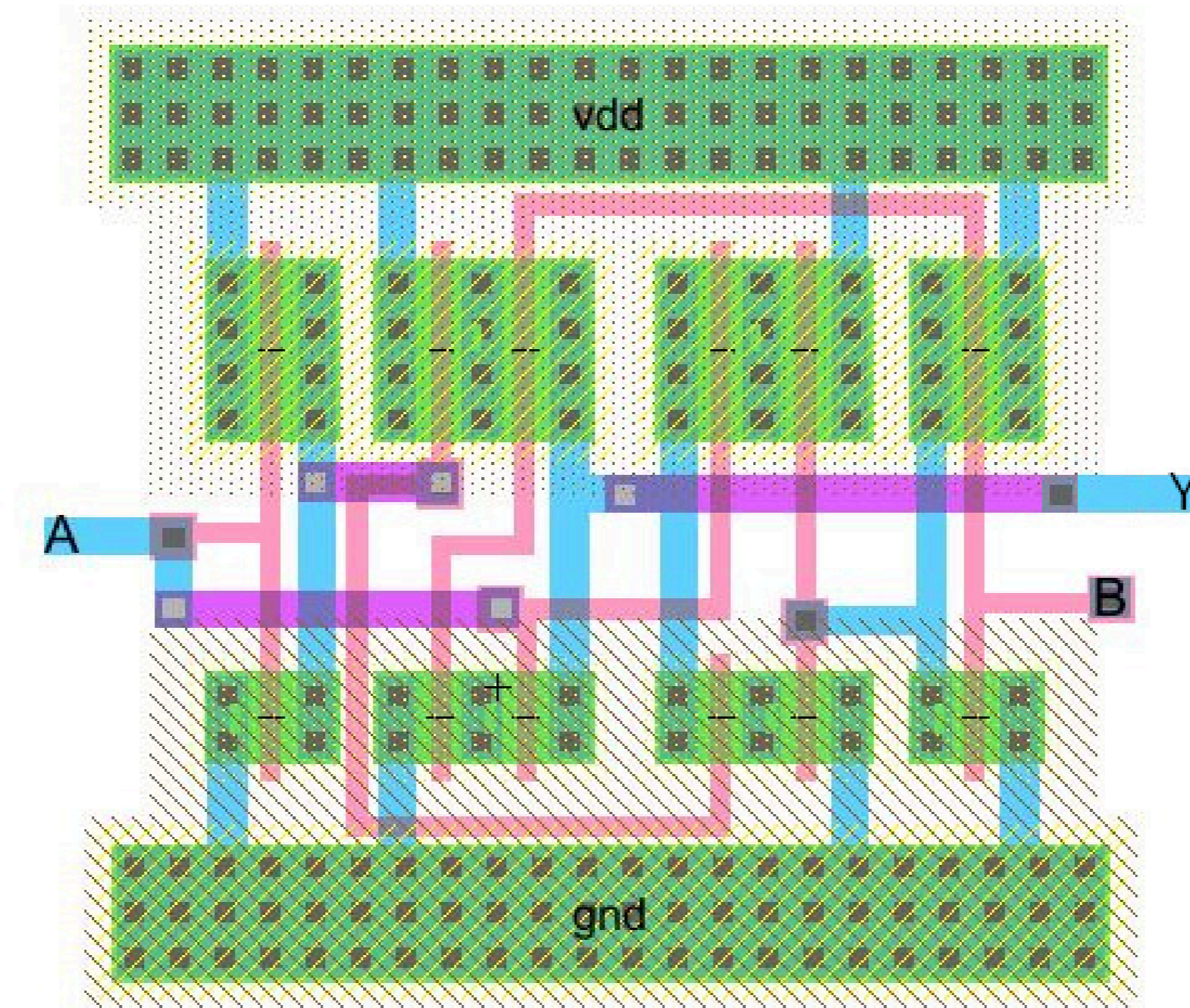# Hardware

# Levels of abstraction in IC design

# Levels of abstraction

0. integrated circuit (IC) layout

1. transistors

2. logic gates

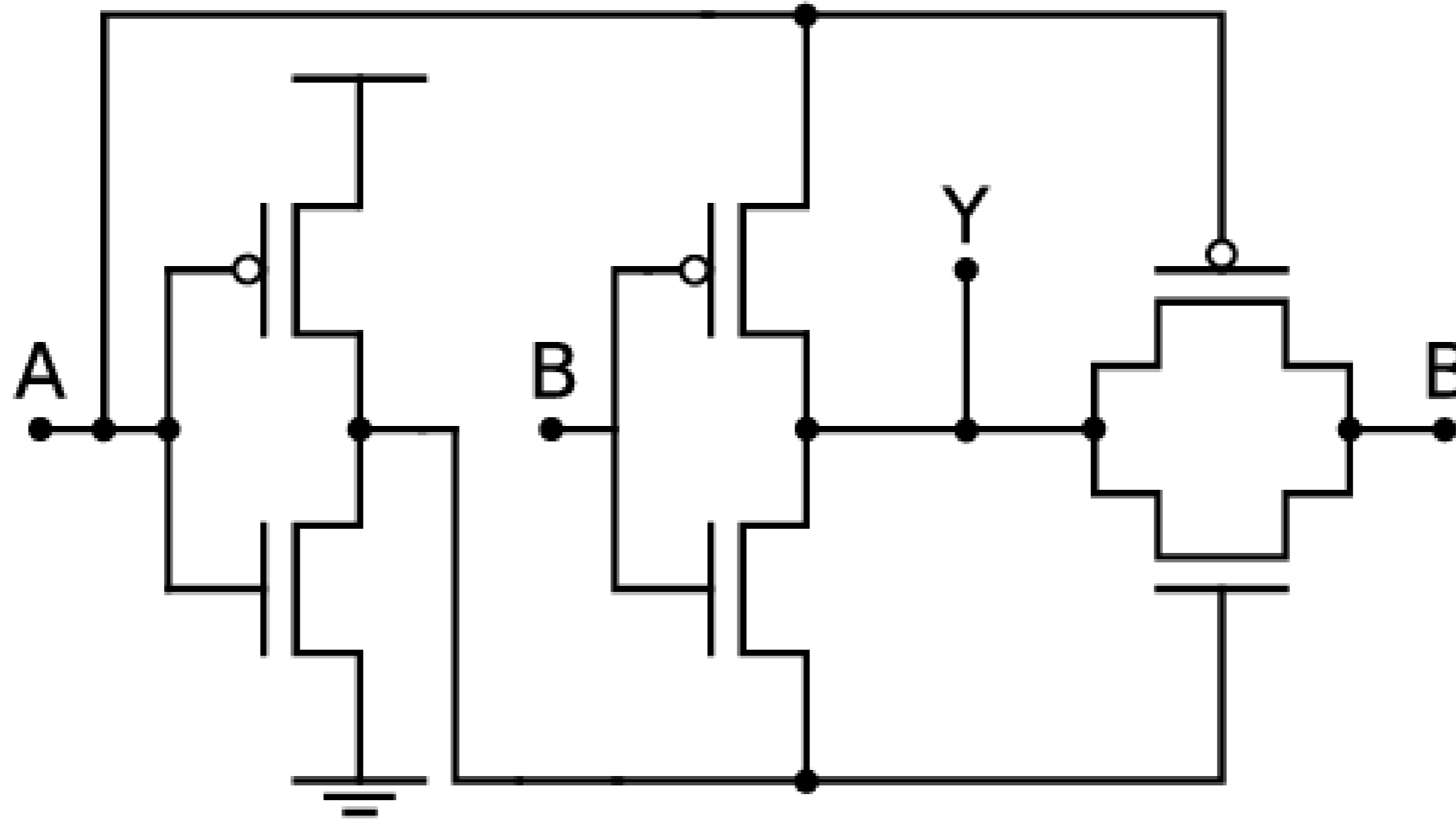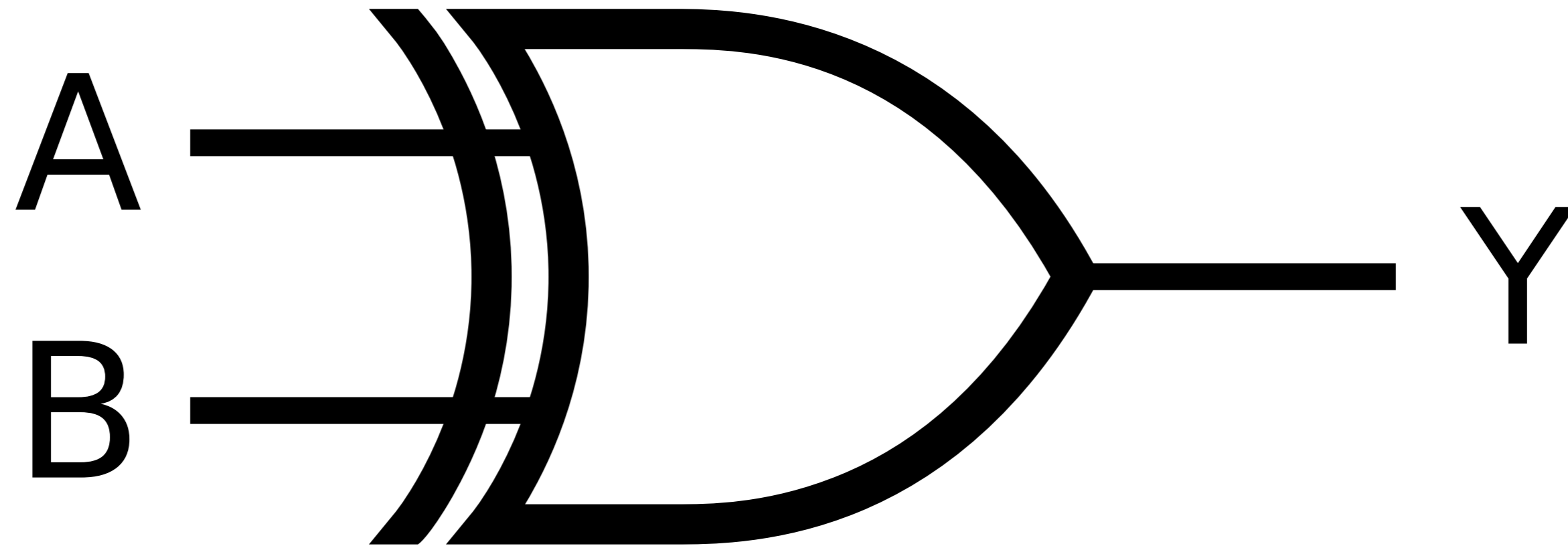3. "intellectual property" (IP) block

# 0. IC layout



Y := A xor B   (IC layout)

# 1. Transistors



Y := A xor B   (transistors)

# 2. Logic gates



Y := A xor B    (logic gate)

# 3. IP blocks

4-bit equality comparator:

$$\text{assuming} \quad \mathtt{a} \quad = \quad A_3 \times 2^3 + A_2 \times 2^2 + A_1 \times 2^1 + A_0 \times 2^0$$
$$\text{and} \quad \mathtt{b} \quad = \quad B_3 \times 2^3 + B_2 \times 2^2 + B_1 \times 2^1 + B_0 \times 2^0$$

```
if a == b then x := 1
if a != b then x := 0
```

# 3. IP blocks

4-bit equality comparator:

$$\text{assuming} \quad \mathtt{a} \;=\; A_3 \times 2^3 + A_2 \times 2^2 + A_1 \times 2^1 + A_0 \times 2^0$$
$$\text{and} \quad \mathtt{b} \;=\; B_3 \times 2^3 + B_2 \times 2^2 + B_1 \times 2^1 + B_0 \times 2^0$$

```
if a == b then x := 1
if a != b then x := 0
```

# Example: $n$-bit addition
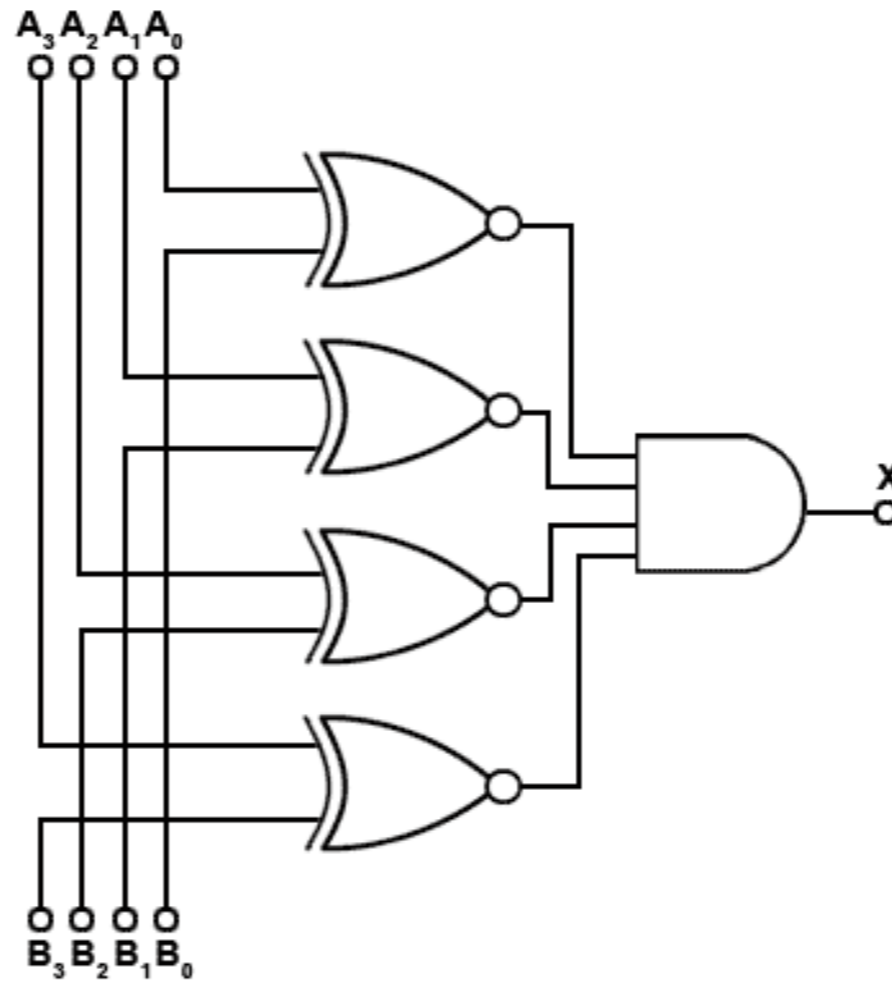
```
carry

a            …  0  1  1  0

b            …  0  1  1  1
```

a + b

# Example: $n$-bit addition

```
carry                    0

a              …  0  1  1  0

b              …  0  1  1  1
_____

a + b                       1
```

# Example: $n$-bit addition

```
carry              1 0
a          …  0  1  1  0
b          …  0  1  1  1
───────────────────────────
a + b               0  1
```

# Example: $n$-bit addition

```
carry           1  1  0
a           …   0  1  1  0
b           …   0  1  1  1
─────────────────────────
a + b              1  0  1
```

# Example: $n$-bit addition

```
carry      0  1  1  0
a        … 0  1  1  0
b        … 0  1  1  1
         ─────────────
a + b    … 1  1  0  1
```

# Example: $n$-bit addition

```
carry      0  1  1  0
a        … 0  1  1  0
b        … 0  1  1  1
─────────────────────
a + b    … 1  1  0  1
```

# The 1-bit "full adder"

- Input:
  - 1 bit of carry:   $C_{\text{in}}$
  - 1 bit of a:   $A$
  - 1 bit of b:   $B$

- Output:
  - 1 bit of carry:   $C_{\text{out}}$
  - 1 bit of the sum a + b:   $S$

- Operation:   Compute $C_{\text{out}}$ and $S$ such that

$$A + B + C_{\text{in}} = C_{\text{out}} \times 2 + S$$

```
carry      ...  1   1   0
a          ...  0   1   1   0
b          ...  0   1   1   1
─────────────────────────────
a + b      ...  1   1   0   1
```

# The 1-bit "full adder"

- Input:

  - 1 bit of carry: $C_{\text{in}}$
  - 1 bit of a: $A$
  - 1 bit of b: $B$

- Output:

  - 1 bit of carry: $C_{\text{out}}$
  - 1 bit of the sum a + b: $S$

- Operation: Compute $C_{\text{out}}$ and $S$ such that

$$A + B + C_{\text{in}} = C_{\text{out}} \times 2 + S$$

Truth table:

| $C_{\text{in}}$ | $A$ | $B$ | $C_{\text{out}}$ | $S$ |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 1 | 1 | 0 |
| 1 | 0 | 0 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 1 | 1 | 0 | 1 | 0 |
| 1 | 1 | 1 | 1 | 1 |

# The 1-bit "full adder"

Boolean expressions:

- $S := A \text{ xor } B \text{ xor } C_{\text{in}}$
- $C_{\text{out}} := (A \text{ and } B) \text{ or } ((A \text{ xor } B) \text{ and } C_{\text{in}})$

Logic diagram:



Truth table:

| $C_{\text{in}}$ | $A$ | $B$ | $C_{\text{out}}$ | $S$ |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 1 | 1 | 0 |
| 1 | 0 | 0 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 1 | 1 | 0 | 1 | 0 |
| 1 | 1 | 1 | 1 | 1 |

# The 1-bit "full adder"

Notes:

- Called "full" in contrast to the "1-bit half adder" which has no $C_{in}$.
- There can be multiple valid Boolean expressions (and logic diagrams)
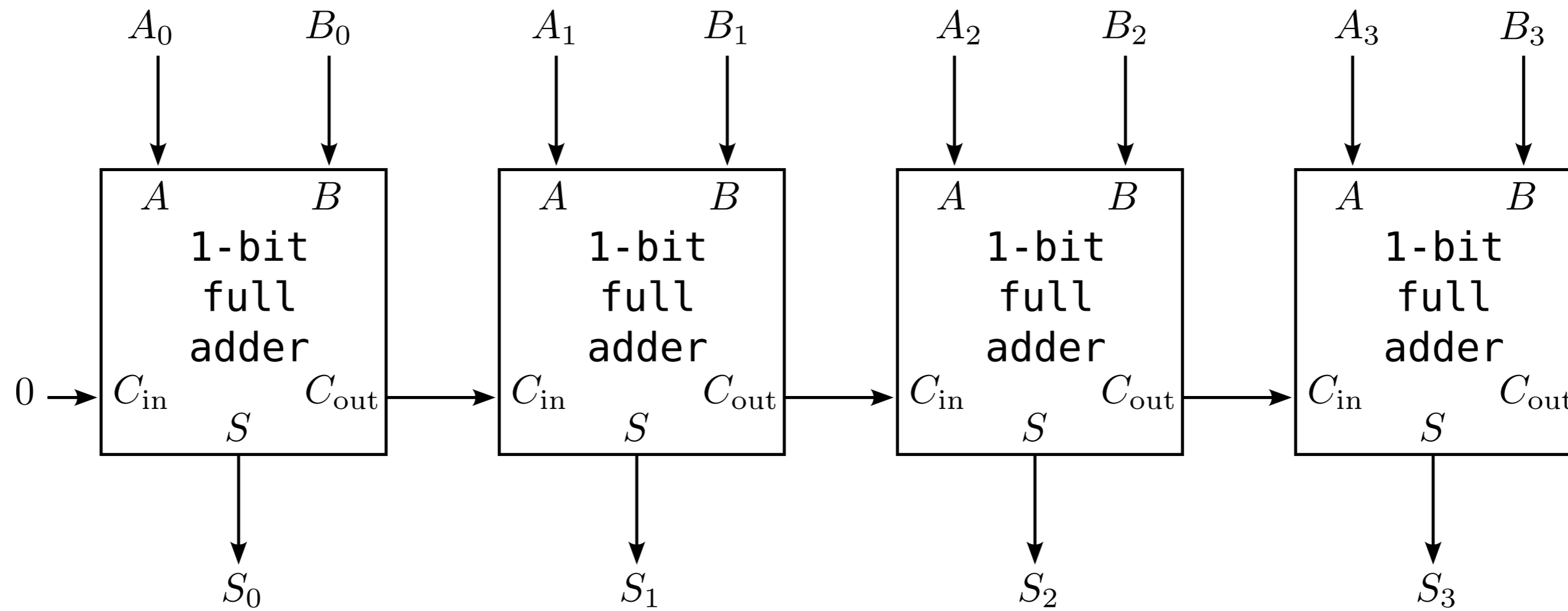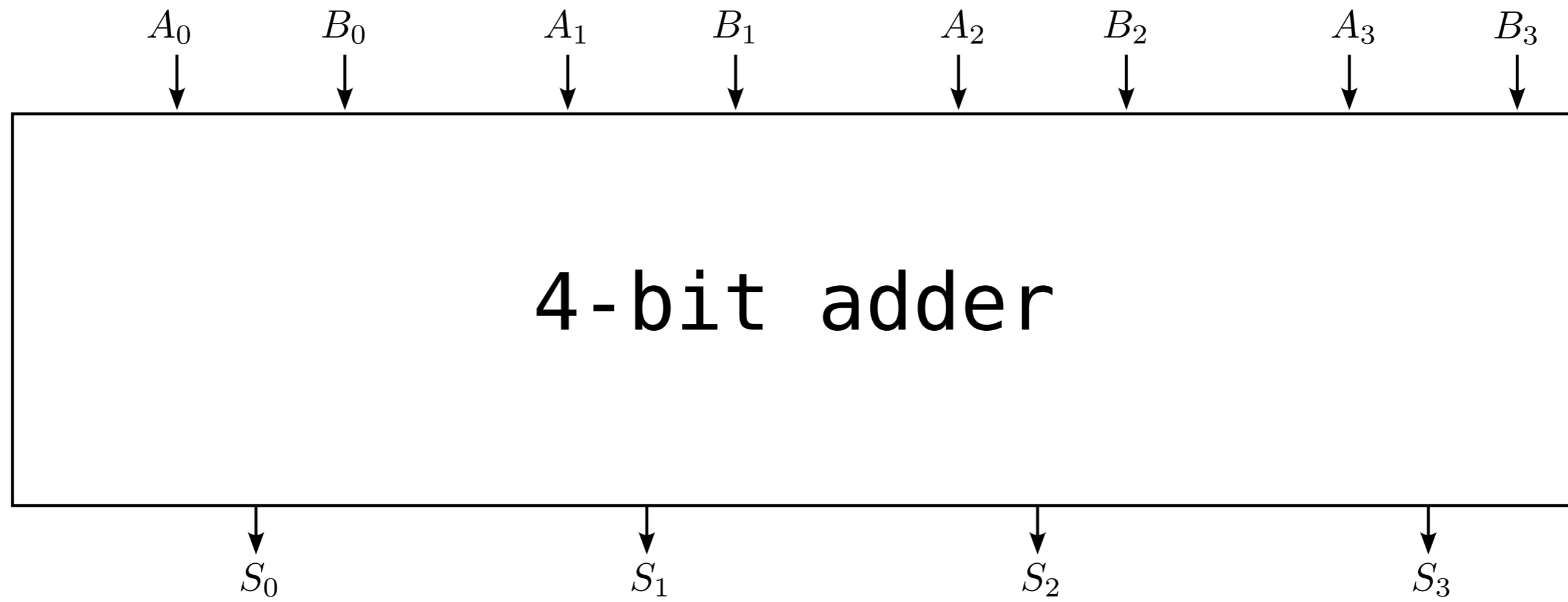
# The 1-bit "full adder"

Notes:

- Called "full" in contrast to the "1-bit half adder" which has no $C_{\text{in}}$.

- There can be multiple valid Boolean expressions (and logic diagrams)

# 4-bit adder

# 4-bit adder

# How IP blocks are designed

IP blocks are designed (and combined) in hardware description languages ("HDL"):
    Verilog, VHDL and derivatives

HDL is then translated into IC layouts by specialized tools.

```systemverilog
interface adder_if();
  logic       rstn;
  logic [7:0]  a;
  logic [7:0]  b;
  logic [7:0]  out;
  logic       carry;
endinterface

module adder(adder_if i);
  always_comb begin
    if (i.rstn) begin
      i.out <= 0;
      i.carry <= 0;
    end else begin
      {i.carry, i.out} <= i.a + i.b;
    end
  end
endmodule
```

SystemVerilog code for an 8-bit adder

# The IC design and manufacturing industry

# Foundries

- "**foundry**" or "fabrication plant" ("**fab**"):     plant in which ICs are manufactured, and by extension, companies who own such plants (e.g. Intel, TSMC, Samsung).

- "**process node**":     marketing name given by foundries to a particular version of their manufacturing process     –     usually a "**feature size**": the size of some parts of the transistors (e.g. 5nm, 7nm, etc.)     –     but not directly comparable across companies.

  Smaller tansistors means:
  - more transistors (per unit of area)
  - lower power consumption
  - faster ICs (propagation delay↓, power consumption↓, heat dissipation↓)

# Types of industry players

- **fabless**:    company that does not own fabrication plants (e.g. ARM, AMD, Apple, nVidia).

  Such companies either:

  - sell the designs of their IP blocks (ARM), or

  - subcontract foundries to manufacture their designs for them (AMD, Apple, nVidia).


- **"pure-play" foundries** (e.g. TSMC):    foundries who manufacture other companies' designs.


- **integrated device manufacturer (IDM)** (e.g. Intel):    designs and manufactures its own ICs.

# Recent history of the industry

- In the early 2000s, Intel (US), AMD (US) and IBM (US) have the best manufacturing technology. All three are IDMs – they design and manufacture in-house.

- Around 2008, AMD (US) spun off its foundry (as "GlobalFoundries") and became fabless.
  Note: first iPhone released in 2007 in the US

- In the 2010s, TSMC (Taiwan) emerges as a major foundry for Apple (US), nVidia (US) and AMD.

- In 2014, IBM sells its manufacturing business to GlobalFoundries.

- As of 2024, GlobalFoundries is still one of the largest pure-play foundries but it has fallen behind in terms of technology (by 5-10 years).

# State of the industry

- Since ~2018, TSMC (Taiwan) has had the best process node, ahead of Samsung (Korea) and Intel (US).

- TSMC's advantage in large part due to early bet on extreme ultra-violet ("EUV") technology.
  Now, Samsung and Intel also use EUV tech.
  ASML (Netherlands) is currently the only supplier of EUV-capable machines.

- Apple (US), AMD (US), nVidia (US), Qualcomm (US) all mostly subcontract TSMC to fabricate their top-of-the-line ICs.

# Microprocessors

- Logic gate circuits allow us to compute Boolean functions **very** fast

  limited by propagation delay in copper (nanoseconds per meter) and transistors (picoseconds)

- Boolean functions can model essentially anything we can compute today.

But

- we cannot design and manufacture a new IC for each algorithm or computing task

- we need **many** logic gates, even for simple things

  ~100k transistors for a 64-bit integer division

  for context, modern microprocessors have 1-100 billion transistors

$\rightarrow$ We break down complex algorithms into simple steps.

# Components in a microprocessor

- Logic gates

- A **clock**

- **Memory**

- Input and output devices

# A simple model

- Memory is $N$ bits $x \in \{0, 1\}^N$ (e.g. for 16 GB, $N \simeq 128 \times 10^9$)

- At every **clock cycle** (e.g. 1.2 GHz), we update the memory:

$$x_i' \leftarrow f_i(x) \qquad \forall i = 0, \dots, N$$

- To simplify the model

  - Some of the memory comes from input devices

  - Some of the memory is sent to output devices

# Issue with the simple model

In this model, we update the whole memory at every clock cycle:

- That would be $128 \times 10^9 \times 1.2 \times 10^9 = 153.6 \times 10^{18}$ b/s

$$\simeq 19,200,000,000 \text{ GB/s}$$

- As of 2024, memory maxes out at $\sim 800$ GB/s

Therefore, we cannot have too many different Boolean functions $f_i$

# A more realistic model

Instead, at each cycle, the computer executes one of a limited set of **instructions** in a **microprocessor**.    Ex.: "Central Processing Unit" (CPU), "Graphics Processing Unit" (GPU).

Instructions are read sequentially from memory and they can be:

- a memory read / write (a tiny amount, like 512 bits)

- 64-bit arithmetic (+, -, ×, /, …)

- a comparison (<, >, =, …)

- a branch (`if`, `while`, …) which alters the control flow of instructions

# Instruction Set Architectures (ISA)

An **ISA** specifies:

- How the machine is organized (memory, etc.)

- What instructions are available

- How instructions are encoded into bits

Two major ISAs in practice:

- x86_64 (aka. x64, x86_64, AMD64): Intel® and AMD® 64-bit CPUs

- AArch64 (aka. ARM64): ARM®-based 64-bits CPUs (most phones, Apple M1 – M4)

Many older or less-prominent ISAs:

x86, Itanium, ARMv7, RISC-V, PowerPC, …

```c
int f(int a, int b, int c)
{
    return (a * b) / c;
}
```

x86_64:

89 f8 89 d1 0f af c6 99 f7 f9 c3

```
f:
  mov eax, edi   # 89 f8
  mov ecx, edx   # 89 d1
  imul eax, esi  # 0f af c6
  cdq            # 99
  idiv ecx       # f7 f9
  ret            # c3
```

↑    assembly    ↑

AArch64:

1b 01 7c 00 1a c2 0c 00 d6 5f 03 c0

```
f:
  mul w0, w0, w1    # 1b 01 7c 00
  sdiv w0, w0, w2   # 1a c2 0c 00
  ret               # d6 5f 03 c0
```

# Assembly

- Assembly is the lowest-level programming language

- Usually in 1:1 correspondence with binary encoding of instructions

- Typically, one line per instruction

# Instructions (x86_64)

```
f:
  mov eax, edi    # 89 f8
  mov ecx, edx    # 89 d1
  imul eax, esi   # 0f af c6
  cdq             # 99
  idiv ecx        # f7 f9
  ret             # c3
```

| | | |
|---|---|---|
| $\text{mov}\ a, b$ | move | $a \leftarrow b$ |
| $\text{imul}\ a, b$ | signed integer multiply | $a \leftarrow a \times b$ |
| $\text{idiv}\ a$ | signed integer divide | $\text{eax} \leftarrow \text{eax}\ /b$ |
| $\text{cdq}$ | convert double-word (32 bits) to quad-word (64 bits) | sign-extend eax into edx:eax |
| $\text{ret}$ | return | return to calling function |

# Instructions (AArch64)

```
f:
  mul w0, w0, w1    # 1b 01 7c 00
  sdiv w0, w0, w2   # 1a c2 0c 00
  ret               # d6 5f 03 c0
```

$\texttt{mul}\ a, b, c$     multiply                $a \leftarrow b \times c$

$\texttt{sdiv}\ a, b, c$     signed integer divide    $a \leftarrow b/c$

$\texttt{ret}$             return                return to calling function

# Registers

x86_64:

```
f:
  mov eax, edi   # 89 f8
  mov ecx, edx   # 89 d1
  imul eax, esi  # 0f af c6
  cdq            # 99
  idiv ecx       # f7 f9
  ret            # c3
```

AArch64:

```
f:
  mul w0, w0, w1   # 1b 01 7c 00
  sdiv w0, w0, w2  # 1a c2 0c 00
  ret              # d6 5f 03 c0
```

- small, fixed set of variables that can be accessed instantly

- 16 (x86_64) or 31 (AArch64) general-purpose 64-bit registers

- plus special registers and flags (not accessible directly)

- plus larger registers for extended operations (e.g. non-integer numbers)

# General-purpose registers (x86_64)

- sixteen 64-bit registers:

  rax, rbx, rcx, rdx, rbp, rsp, rsi, rdi, r8, r9, r10, r11, r12, r13, r14, r15

- we can access the lower 32 bits separately:

  eax, ebx, ecx, edx, ebp, esp, esi, edi, r8d, r9d, r10d, r11d, r12d, r13d, r14d, r15d

- we can access the lower 16 bits separately:

  ax, bx, cx, dx, bp, sp, si, di, r8w, r9w, r10w, r11w, r12w, r13w, r14w, r15w

- we can access the lower 8 bits separately:

  al, bl, cl, dl, bpl, spl, sil, dil, r8b, r9b, r10b, r11b, r12b, r13b, r14b, r15b

- we can access bits 8-15 separately for some registers:

  ah, bh, ch, dh

Example:

| bits | 63…56 | 55…48 | 47…40 | 39…32 | 31…24 | 23…16 | 15…8 | 7…0 |
|------|-------|-------|-------|-------|-------|-------|------|-----|
| 64 | rax | | | | | | | |
| 32 | | | | | eax | | | |
| 16 | | | | | | | ax | |
| 8 | | | | | | | ah | al |

# General-purpose registers (AArch64)

- thirty-one 64-bit registers:

  x0, ..., x30

- we can access the lower 32 bits separately:

  w0, ..., w30

- register 31 (x31, w31) is read-only (zero in most cases)

Example:

| bits | 63…56 | 55…48 | 47…40 | 39…32 | 31…24 | 23…16 | 15…8 | 7…0 |
|------|-------|-------|-------|-------|-------|-------|------|-----|
| 64 | x0 | | | | | | | |
| 32 | | | | | w0 | | | |

## Note:

- In both cases, registers are treated as integer numbers

- We cannot (directly) access individual bits

- When it matters, the instruction specifies whether the register is signed or not:

x86_64:

```
 idiv ecx       # f7 f9   (signed)
 div ecx        # f7 f1   (unsigned)
```

AArch64:

```
 sdiv w0, w0, w2  # 1a c2 0c 00   (signed)
 udiv w0, w0, w2  # 1a c2 08 00   (unsigned)
```

# Memory

```c
int g(int *a, int *b)
{
    return *a + *b;
}
```

x86_64:

```
g:
  mov eax, DWORD PTR [rsi]
  add eax, DWORD PTR [rdi]
  ret
```

AArch64:

```
g:
  ldr w2, [x0]
  ldr w0, [x1]
  add w0, w2, w0
  ret
```

# Memory

- From a process' perspective, memory is seen as a single long array of <span style="color:red">bytes</span>

  (8 bits, treated as a single signed or unsigned integer)

- Like registers, memory can be accessed in larger chunks

  (e.g. 16, 32 or 64 bits integer)

- But the smallest addressable unit is the <span style="color:red">byte</span>

# Byte ordering

| address | 0 | 1 | 2 | 3 | ... | 239 | 240 | 241 | 242 | 243 | 244 | ... |
|---------|---|---|---|---|-----|-----|-----|-----|-----|-----|-----|-----|
| value (hex) | ef | cd | ab | 89 | ... | ff | a0 | a1 | a2 | a3 | 42 | ... |

- the byte at address 240 is (hex) a0 = (decimal) 160

- the byte at address 241 is (hex) a1 = (decimal) 161

- the byte at address 242 is (hex) a2 = (decimal) 162

- the byte at address 243 is (hex) a3 = (decimal) 163

Q: What is the value of the 32-bit integer at address 240?

A: It depends!

# Byte ordering / "Endianess"

| address | 0 | 1 | 2 | 3 | … | 239 | 240 | 241 | 242 | 243 | 244 | … |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| value (hex) | ef | cd | ab | 89 | … | ff | a0 | a1 | a2 | a3 | 42 | … |

- "**big-endian**" (BE): 32-bit int at 240 is (hex) a0 a1 a2 a3

$$= \text{(decimal)} \ 160 \times 2^{24} + 161 \times 2^{16} + 162 \times 2^8 + 163$$

$$= \text{(decimal)} \ 2{,}694{,}947{,}491$$

- "**little-endian**" (LE): 32-bit int at 240 is (hex) a3 a2 a1 a0

$$= \text{(decimal)} \ 163 \times 2^{24} + 162 \times 2^{16} + 161 \times 2^8 + 160$$

$$= \text{(decimal)} \ 2{,}745{,}344{,}416$$

- x86_64 is LE

- AArch64 is LE by default (LE-only on Windows, MacOS, Linux)

# Bit ordering

Because we cannot access individual bits on a CPU (smallest chunk is a byte),
bit ordering does not matter here.


However the same problem crops up in other contexts (USB, Ethernet, Wifi, …)

# Memory access notation

- In assembly, accessing memory is denoted using " [ " and " ] "

  - Moving the value 240 into a register:

    ```
    mov eax, 240 # eax = 240
    ldr w0, 240  # w0  = 240
    ```

  - Moving the 4 bytes of memory at address 240 into a register:

    ```
    mov eax, DWORD PTR [240] # eax = (hex) a3a2a1a0
    ldr w0, [240]            # w0  = (hex) a3a2a1a0
    ```

```c
int g(int *a, int *b)
{
    return *a + *b;
}
```

x86_64:

```
g:
  mov eax, DWORD PTR [rsi]
  add eax, DWORD PTR [rdi]
  ret
```

AArch64:

```
g:
  ldr w2, [x0]
  ldr w0, [x1]
  add w0, w2, w0
  ret
```