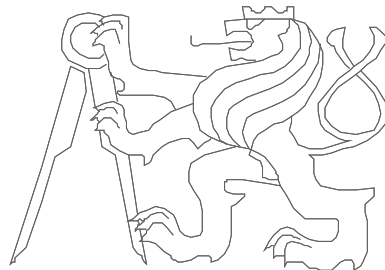


# Computer Architectures

Number Representation and Computer Arithmetics

Pavel Píša, Richard Šusta

Michal Štepanovský, Miroslav Šnorek



Czech Technical University in Prague, Faculty of Electrical Engineering

English version partially supported by:

European Social Fund Prague & EU: We invests in your future.





**Let health allows you and your dears to fulfill one's dreams,  
and nothing and nobody prevents you from sharing happiness  
and results of your work and creativity.**

## Important Introductory Note

- The goal is to understand the structure of the computer so you can make better use of its options to achieve its higher performance.
- It is also discussed interconnection of HW / SW
- **Webpages:**  
<https://cw.fel.cvut.cz/b192/courses/b35apo/>  
<https://dcenet.felk.cvut.cz/apo/> - *they will be opened*
- Some followup related subjects:
  - B4M35PAP - Advanced Computer Architectures
  - B3B38VSY - Embedded Systems
  - B4M38AVS - Embedded Systems Application
  - B4B35OSY - Operating Systems (OI)
  - B0B35LSP – Logic Systems and Processors (KyR + part of OI)
- Prerequisite: Šusta, R.: APOLOS , CTU-FEE 2016, 51 pg.

# Important Introductory Note

- The course is based on a world-renowned book of authors Paterson, D., Hennessey, V.: Computer Organization and Design, The HW/SW Interface. Elsevier, ISBN: 978-0-12-370606-5



**David Andrew Patterson**

University of California, Berkeley

Works: RISC processor Berkley RISC → SPARC, DLX, RAID, Clusters, RISC-V



**John Leroy Hennessy**

10th President of [Stanford University](#)

Works: RISC processors MIPS, DLX a MMIX

2017 Turing Award for pioneering a systematic, quantitative approach to the design and evaluation of computer architectures with enduring impact on the microprocessor industry. → A New Golden Age for Computer Architecture – RISC-V

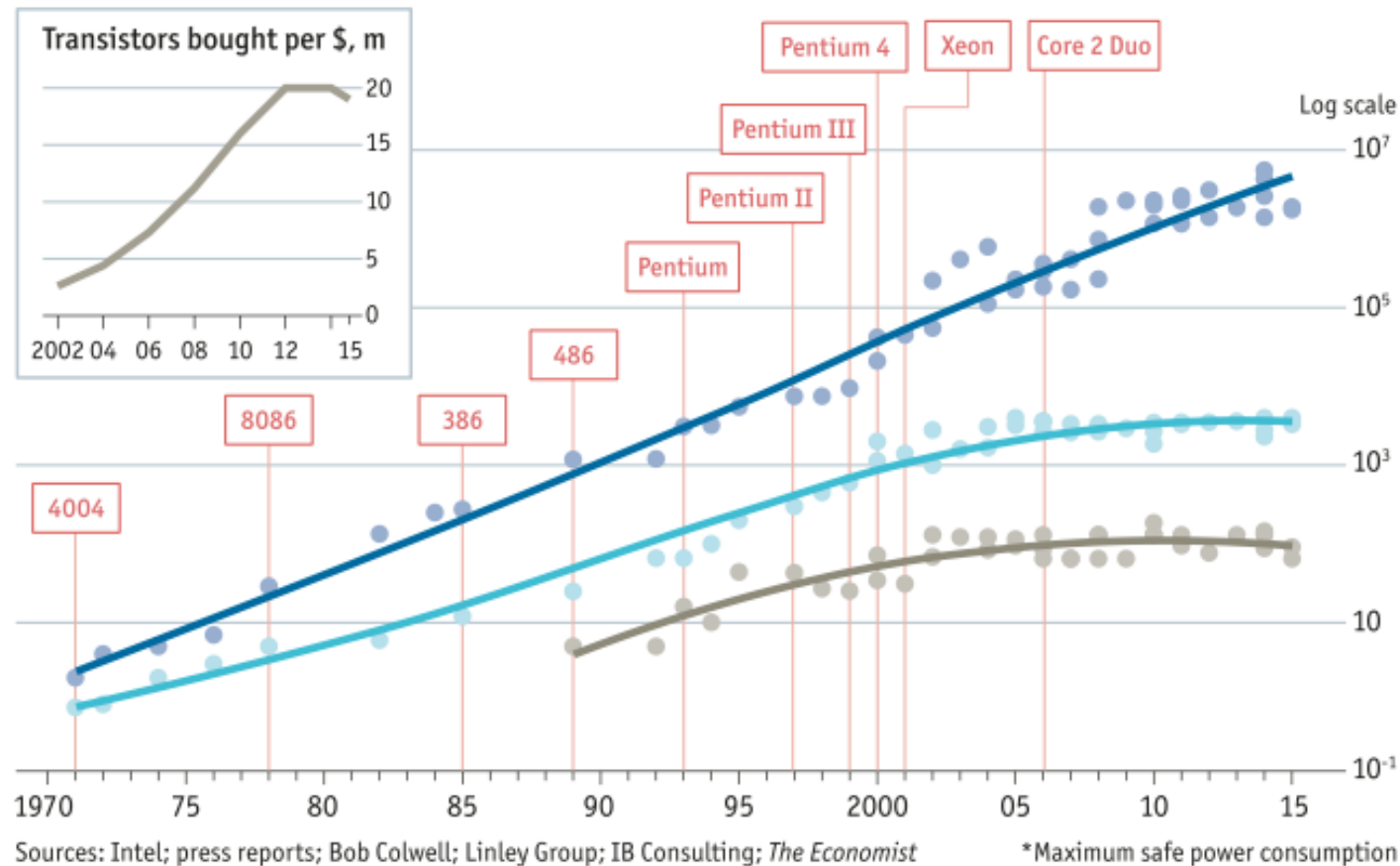


# Moore's Law

**Gordon Moore**, founder of Intel, in 1965: " *The number of transistors on integrated circuits doubles approximately every two years* "

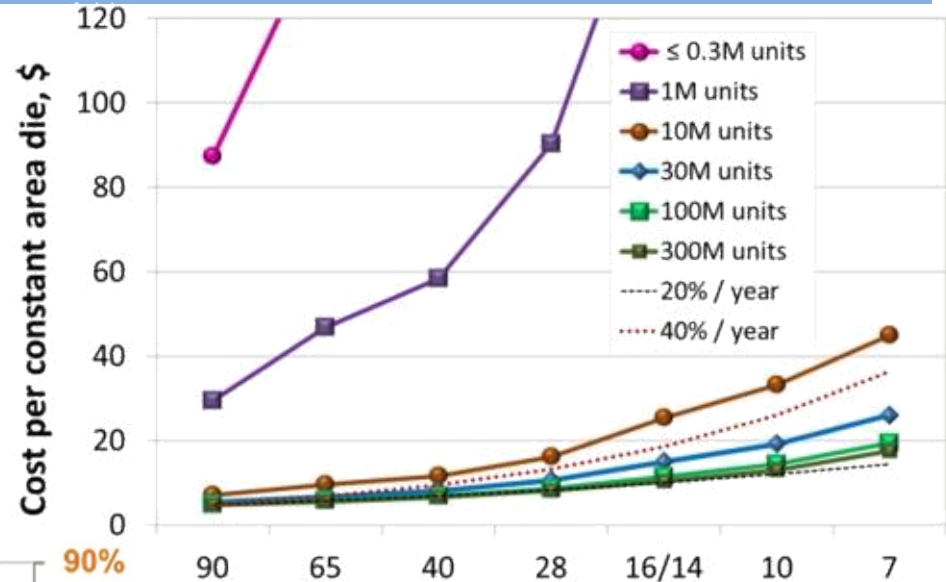
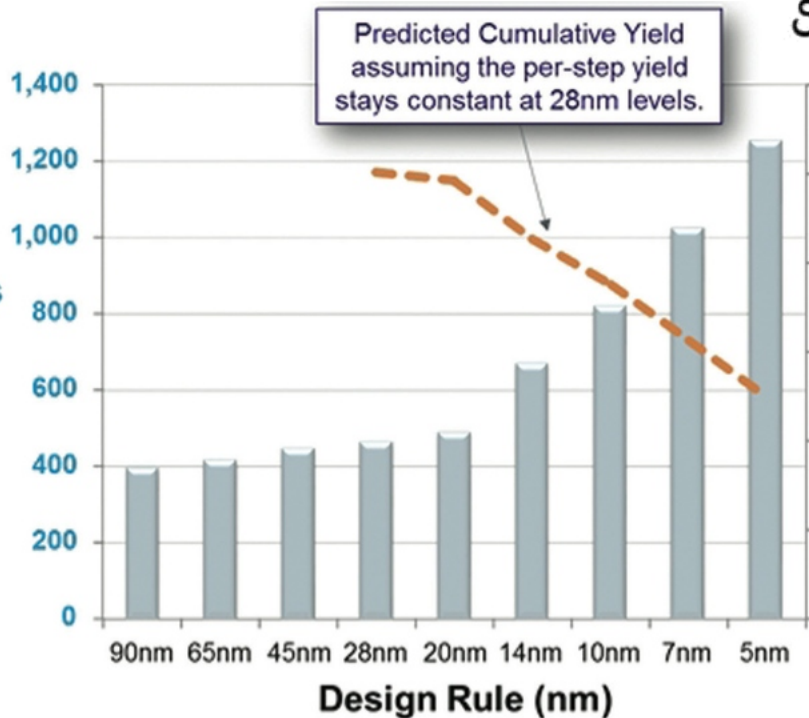
## Stuttering

● Transistors per chip, '000 ● Clock speed (max), MHz ● Thermal design power\*, w □ Chip introduction dates, selected



# The cost of production is growing with decreasing design rule

Moore's Law will be stopped by cost...



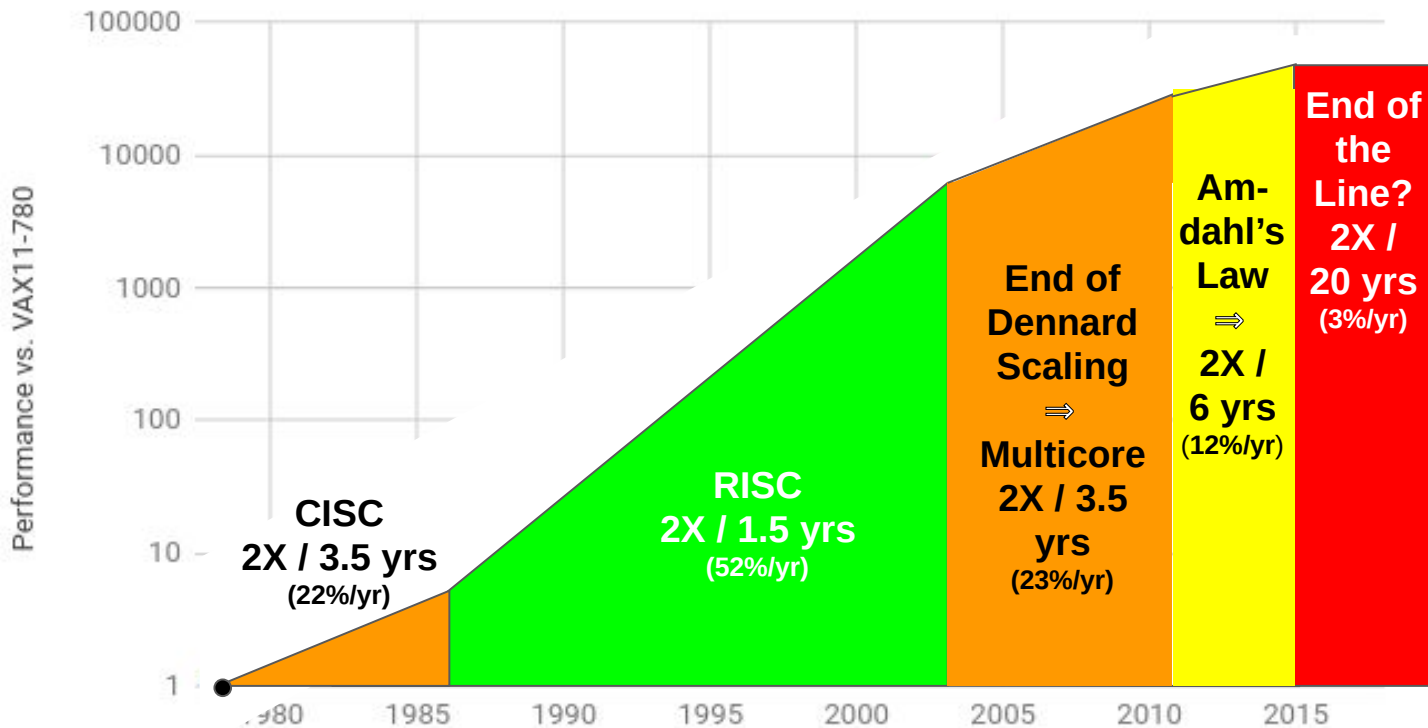
Source: <http://www.eetimes.com/>

Cumulative Yield

Source: <http://electroiq.com/>

# End of Growth of Single Program Speed?

## 40 years of Processor Performance



Based on SPECintCPU. Source: John Hennessy and David Patterson, Computer Architecture: A Quantitative Approach, 6/e. 2018

# Processors Architectures Development in a Glimpse

- 1960 – IBM incompatible families → IBM System/360 – one ISA to rule them all,

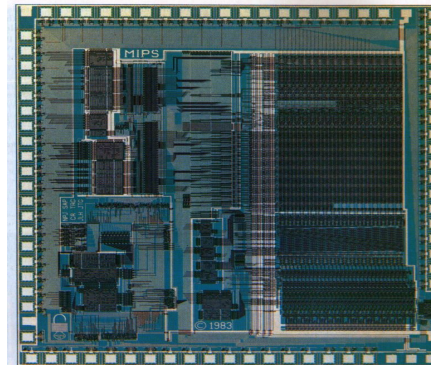
Model	M30	M40	M50	M65
Datapath width	8 bits	16 bits	32 bits	64 bits
Microcode size	4k x 50	4k x 52	2.75k x 85	2.75k x 87
Clock cycle time (ROM)	750 ns	625 ns	500 ns	200 ns
Main memory cycle time	1500 ns	2500 ns	2000 ns	750 ns
Price (1964 \$)	\$192,000	\$216,000	\$460,000	\$1,080,000
Price (2018 \$)	\$1,560,000	\$1,760,000	\$3,720,000	\$8,720,000

- 1976 – Writable Control Store, Verification of microprograms, David Patterson Ph.D., UCLA, 1976
- Intel iAPX 432: Most ambitious 1970s micro, started in 1975 – 32-bit capability-based object-oriented architecture, Severe performance, complexity (multiple chips), and usability problems; announced 1981
- Intel 8086 (1978, 8MHz, 29,000 transistors), “Stopgap” 16-bit processor, 52 weeks to new chip, architecture design 3 weeks (10 person weeks) assembly-compatible with 8 bit 8080, further i80286 16-bit introduced some iAPX 432 lapses, i386 paging



# CISC and RISC

- IBM PC 1981 picks Intel 8088 for 8-bit bus (and Motorola 68000 was out of main business)
- Use SRAM for instruction cache of user-visible instructions
- Use simple ISA – Instructions as simple as microinstructions, but not as wide, Compiled code only used a few CISC instructions anyways, Enable pipelined implementations
- Chaitin's register allocation scheme benefits load-store ISAs
- Berkeley (RISC I, II → SPARC) & Stanford RISC Chips (MIPS)



Stanford MIPS (1983) contains 25,000 transistors, was fabbed in 3  $\mu\text{m}$  & 4  $\mu\text{m}$  NMOS, ran at 4 MHz (3  $\mu\text{m}$ ), and size is 50 mm<sup>2</sup> (4  $\mu\text{m}$ ) (Microprocessor without Interlocked Pipeline Stages)

# CISC and RISC

- CISC executes fewer instructions per program ( $\approx 3/4X$  instructions), but many more clock cycles per instruction ( $\approx 6X$  CPI)  
⇒ RISC  $\approx 4X$  faster than CISC

## PC Era

- Hardware translates x86 instructions into internal RISC Instructions (Compiler vs Interpreter)
- Then use any RISC technique inside MPU
- $> 350M$  / year !
- x86 ISA eventually dominates servers as well as desktops

## PostPC Era: Client/Cloud

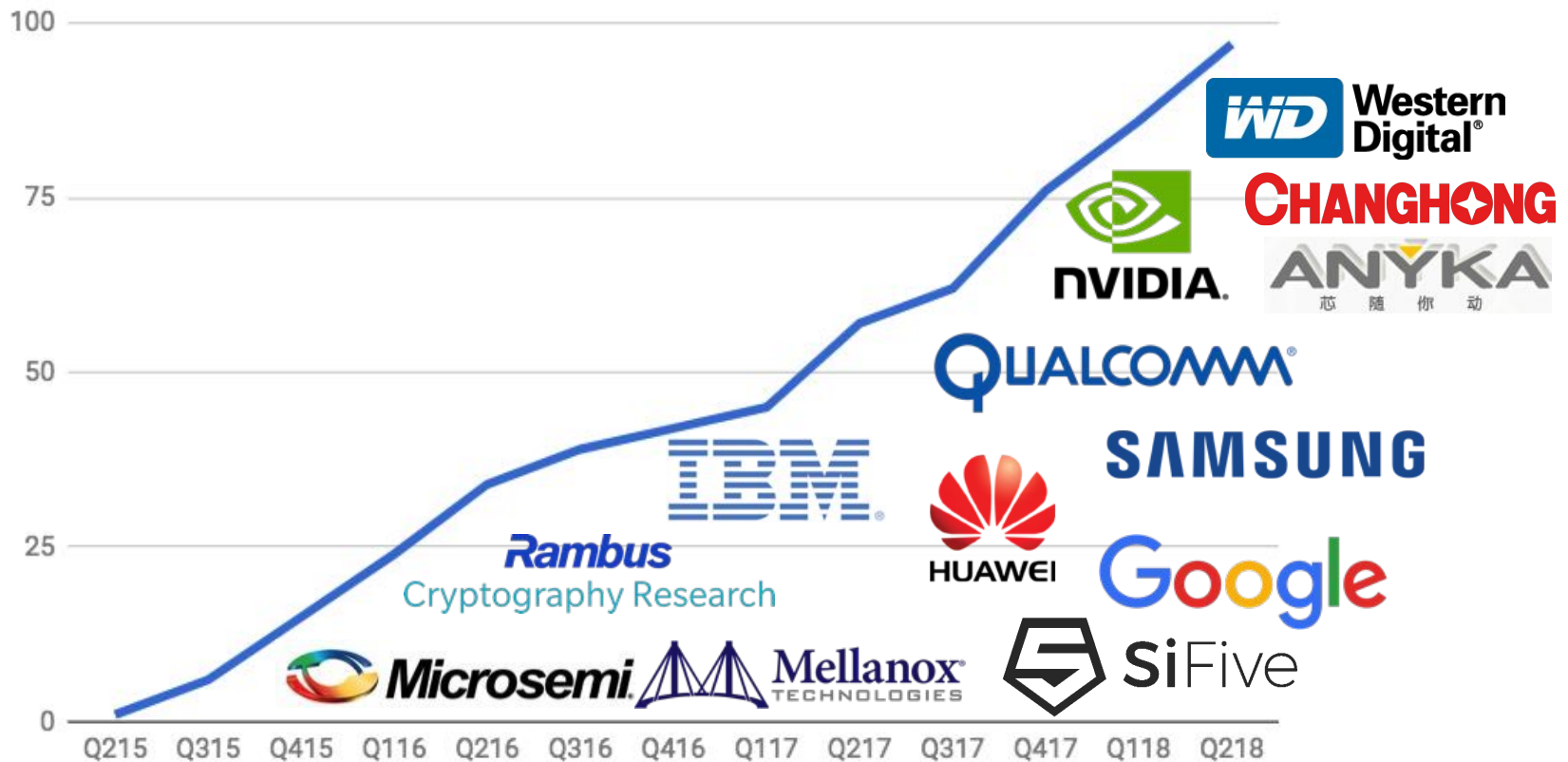
- IP in SoC vs. MPU
- Value die area, energy as much as performance
- $> 20B$  total / year in 2017
- 99% Processors today are RISC
- Marketplace settles debate

- Alternative, Intel Itanium VLIW, 2002 instead 1997
- *“The Itanium approach...was supposed to be so terrific –until it turned out that the wished-for compilers were basically impossible to write.”* - Donald Knuth, Stanford

# RISC-V

- ARM, MIPS, SPARC, PowerPC – Commercialization and extensions results in too complex CPUs again, with license and patents preventing even original investors to use real/actual implementations in silicon to be used for education and research
- Krste Asanovic and other prof. Patterson's students initiated development of new architecture (start of 2010), initial estimate to design architecture 3 months, but 3 years
- Simple, Clean-slate design (25 years later, so can learn from mistakes of predecessors, Avoids  $\mu$ architecture or technology-dependent features), Modular, Supports specialization, Community designed
- A few base integer ISAs (RV32E, RV32I, RV64I)
- Standard extensions (M: Integer multiply/divide, A: Atomic memory operations, F/D: Single/Double-precision FI-point, C: Compressed Instructions (<x86), V: Vector Extension for DLP (>SIMD\*\*))

# Foundation Members since 2015



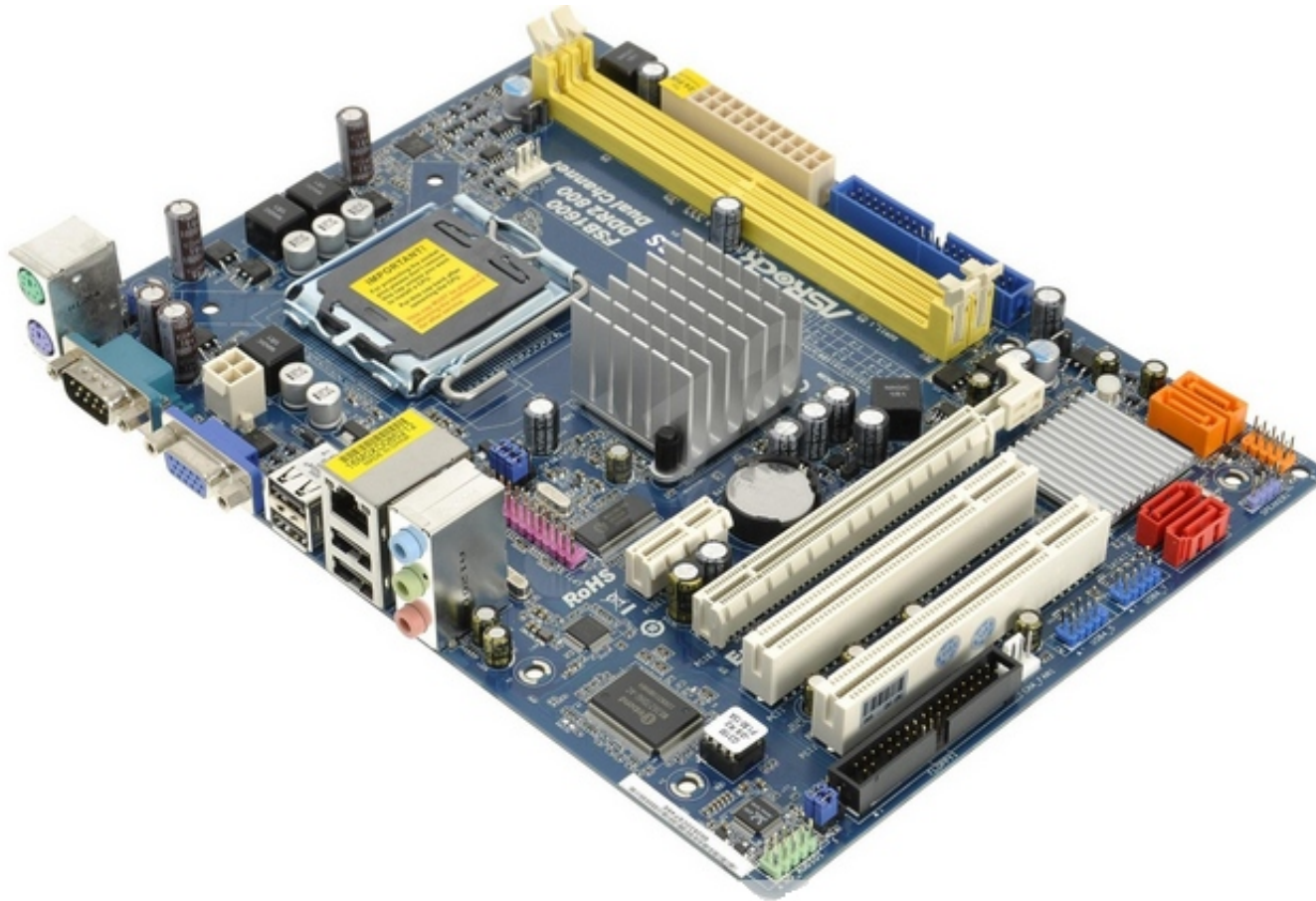
## Open Architecture Goal

Create industry-standard open ISAs for all computing devices

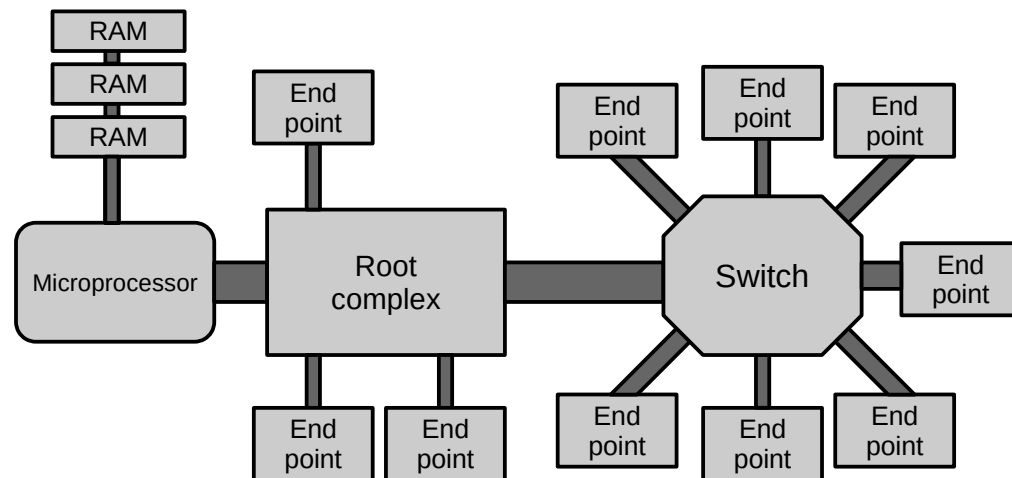
“Linux for processors”



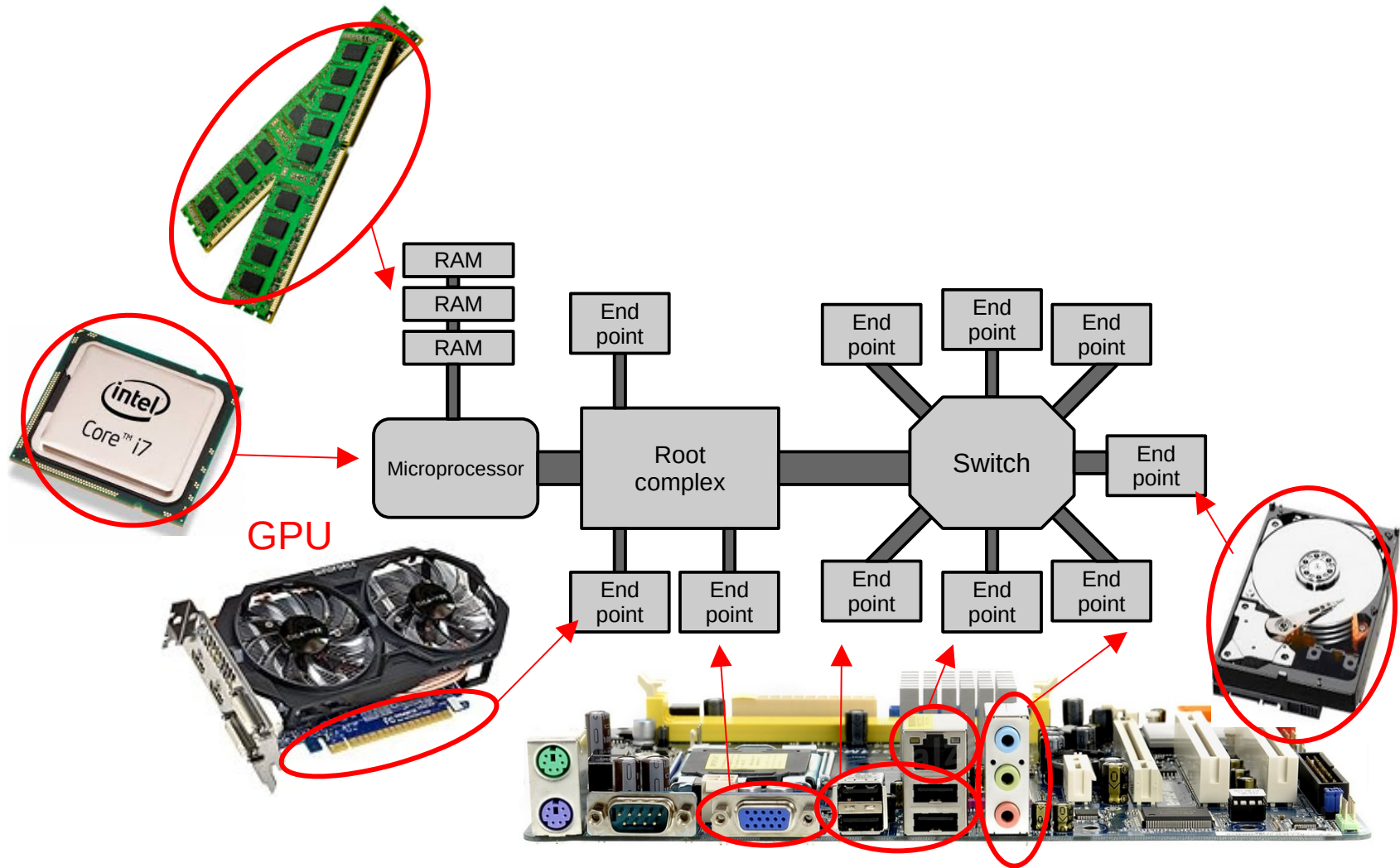
# Today PC Computer Base Platform – Motherboard



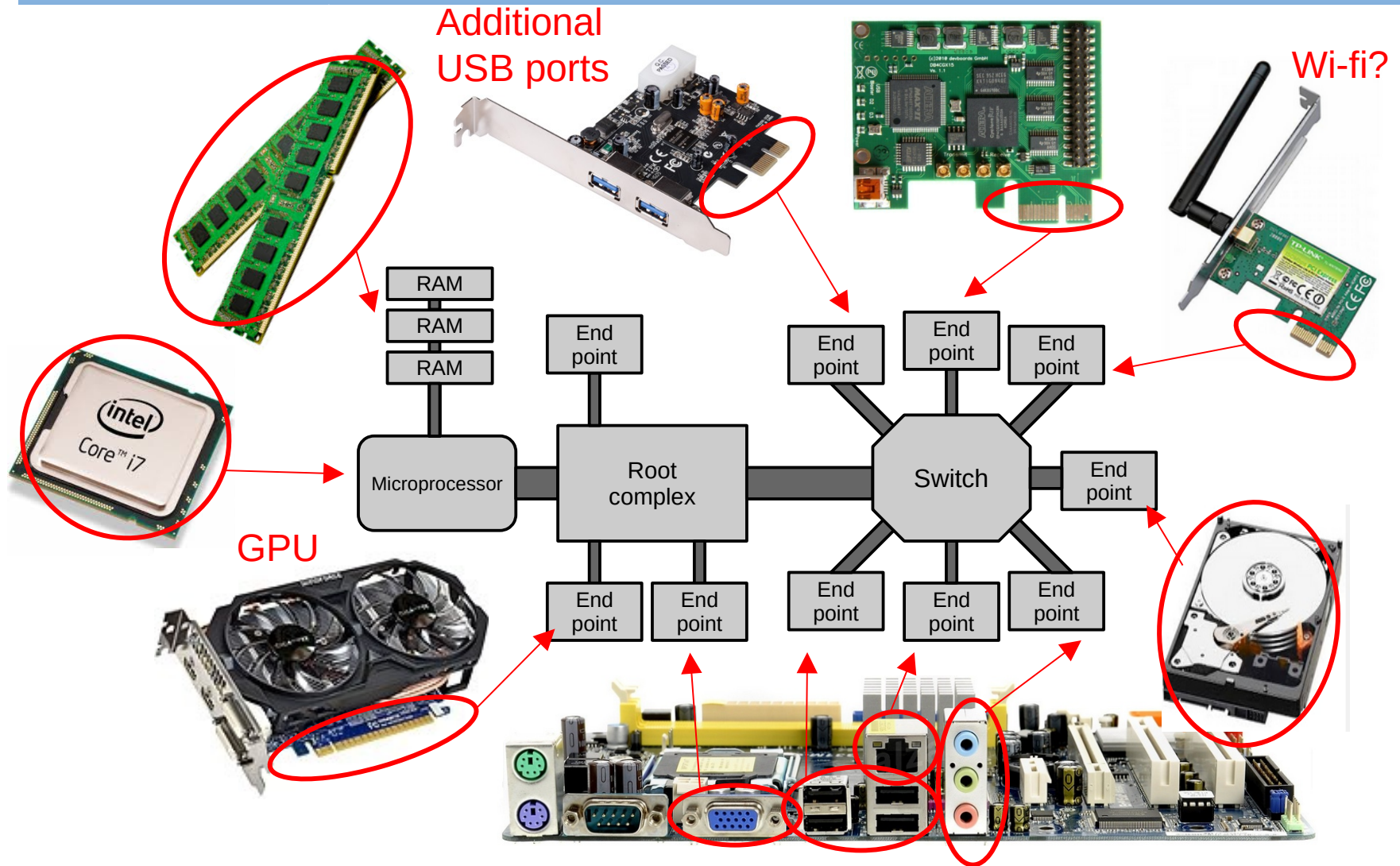
# Block Diagram of Components Interconnection



# Block Diagram of Components Interconnection

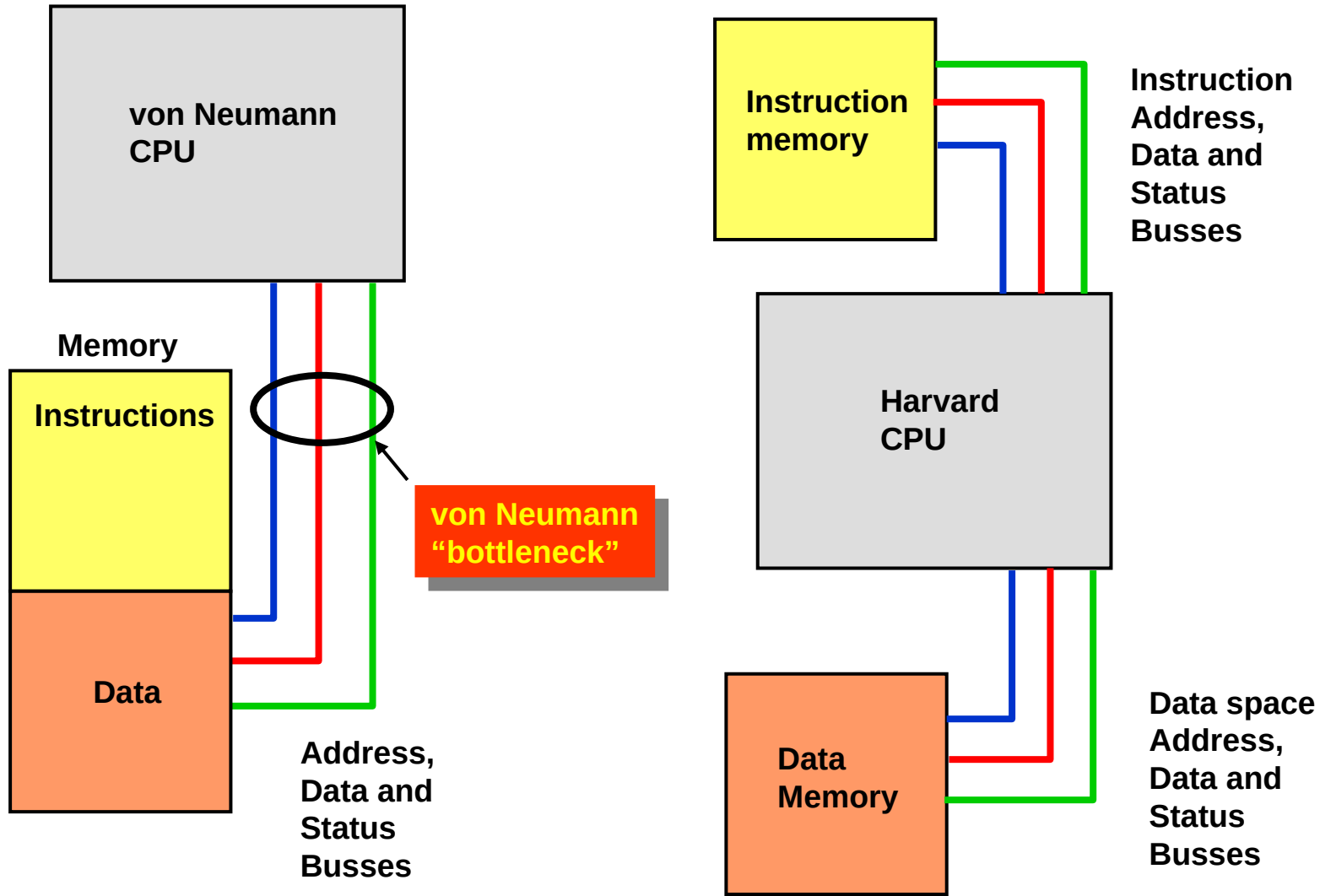


# Block Diagram of Components Interconnection





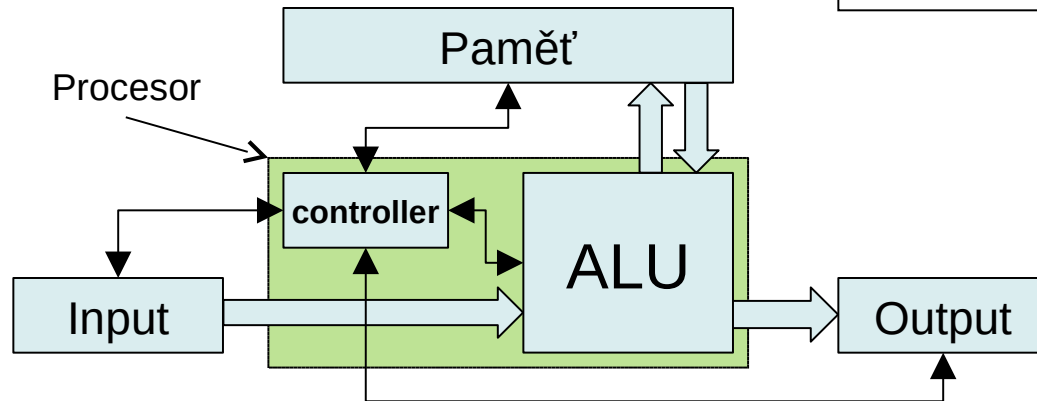
# Von Neumann and Harvard Architectures



[Arnold S. Berger: Hardware Computer Organization for the Software Professional]

# John von Neumann

Princeton Institute for Advanced Studies



28. 12. 1903 -  
8. 2. 1957



## 5 units:

- A processing unit that contains an arithmetic logic unit and processor registers;
- A control unit that contains an instruction register and program counter;
- Memory that **stores data and instructions**
- External mass storage
- Input and output mechanisms

# Samsung Galaxy S4 inside



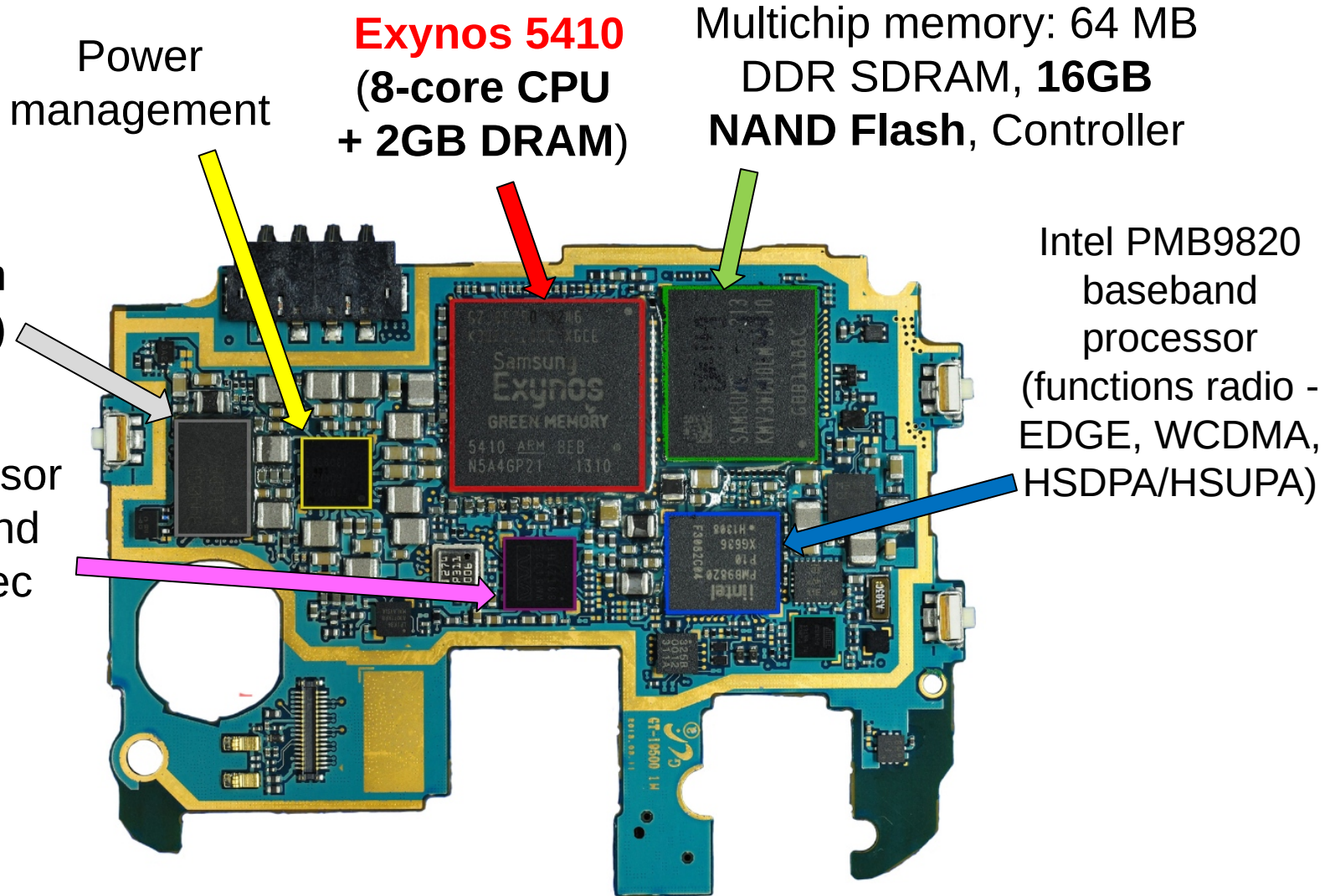
- **Android 5.0 (Lollipop)**
- **2 GB RAM**
- **16 GB user RAM user**
- **1920 x 1080 display**
- **8-core CPU (chip Exynos 5410):**
  - **4 cores 1.6 GHz ARM Cortex-A15**
  - **4 cores 1.2 GHz ARM Cortex-A7**

# Samsung Galaxy S4 inside



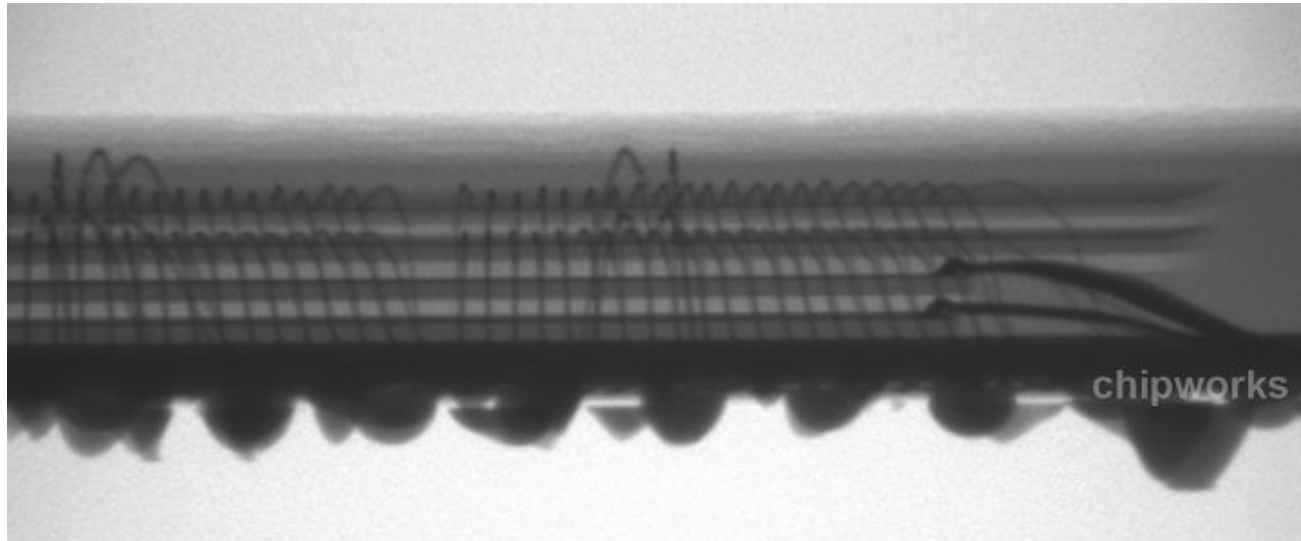


# Samsung Galaxy S4 inside



## Samsung Galaxy S4 inside

X-ray image of Exynos 5410 chip from the side :



We see that this is QDP (Quad die package)

*To increase capacity, chips have multiple stacks of dies.*

*A **die**, in the context of integrated circuits, is a small block of semiconducting material on which a given functional circuit is fabricated. [Wikipedia]*

# Samsung Galaxy S4 inside

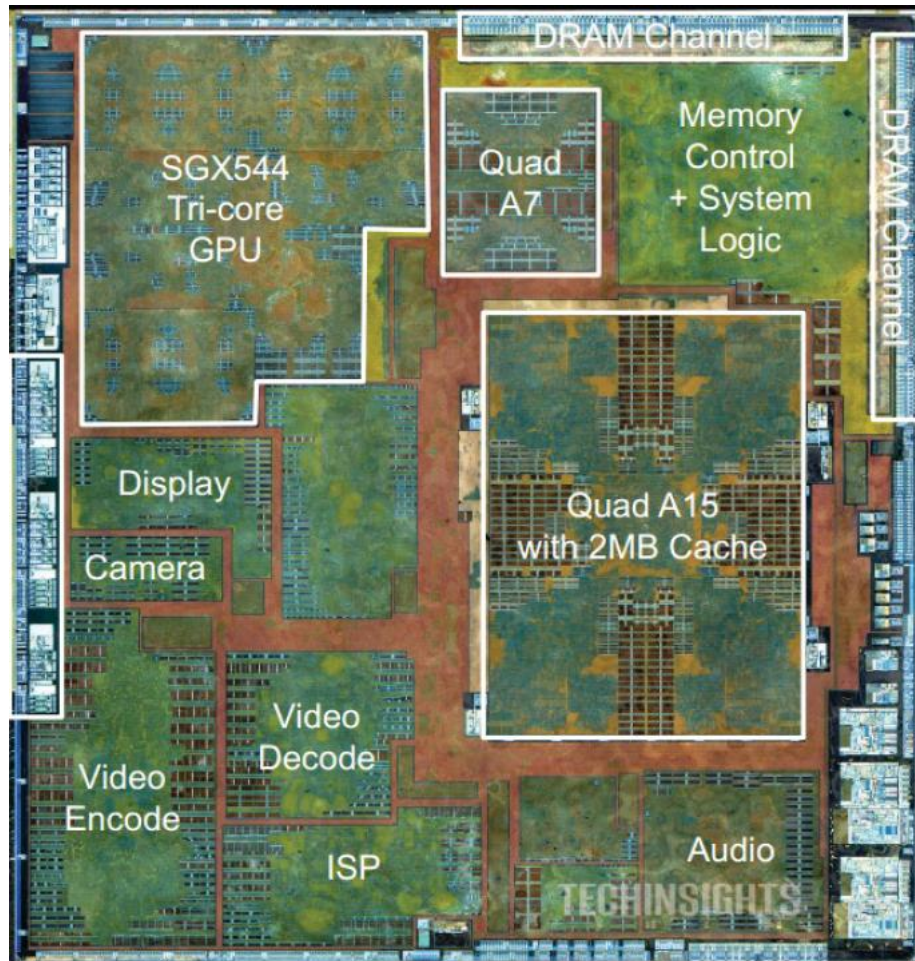
Chip Exynos 5410 – here, we see DRAM





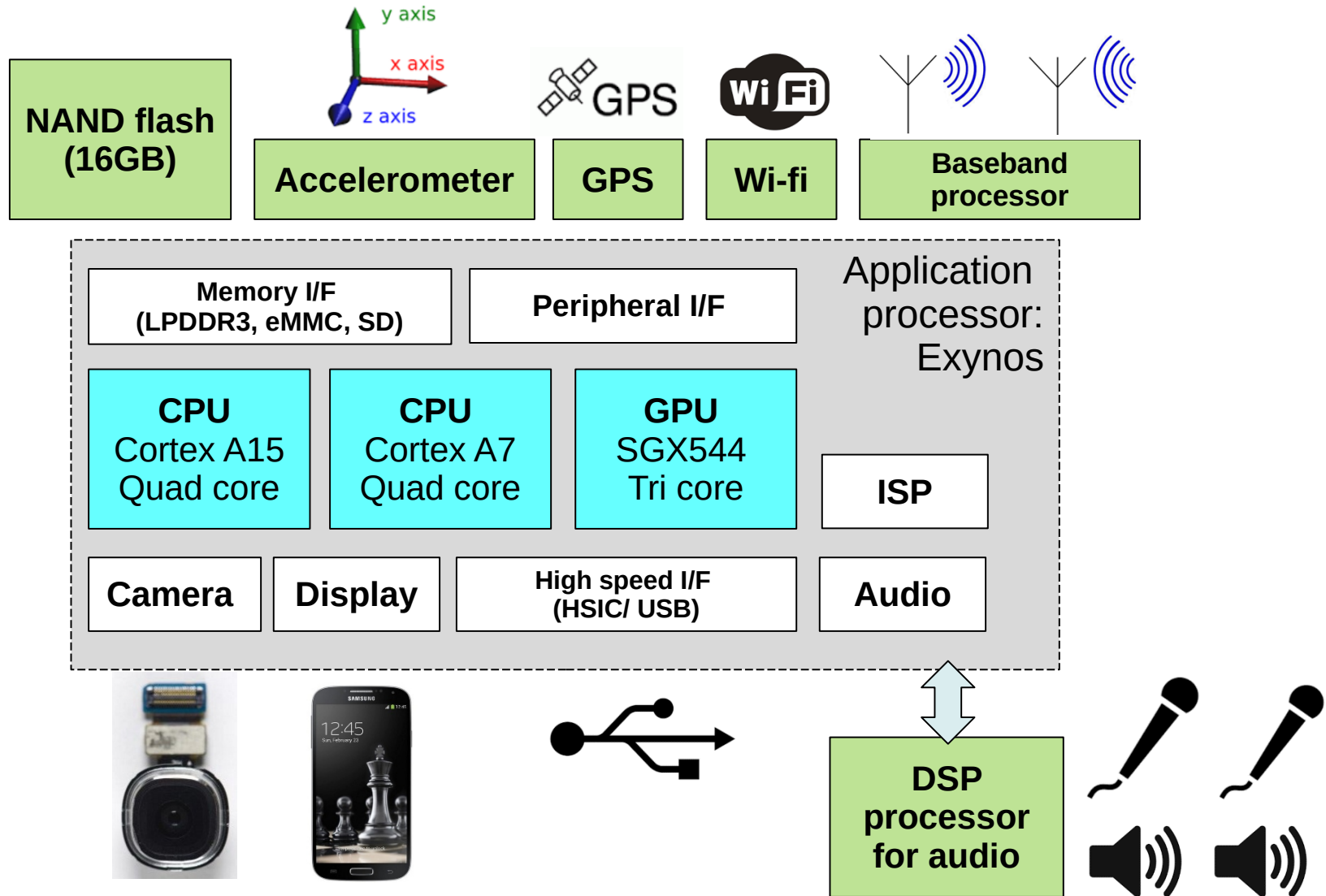
# Samsung Galaxy S4 inside

## Chip Exynos 5410

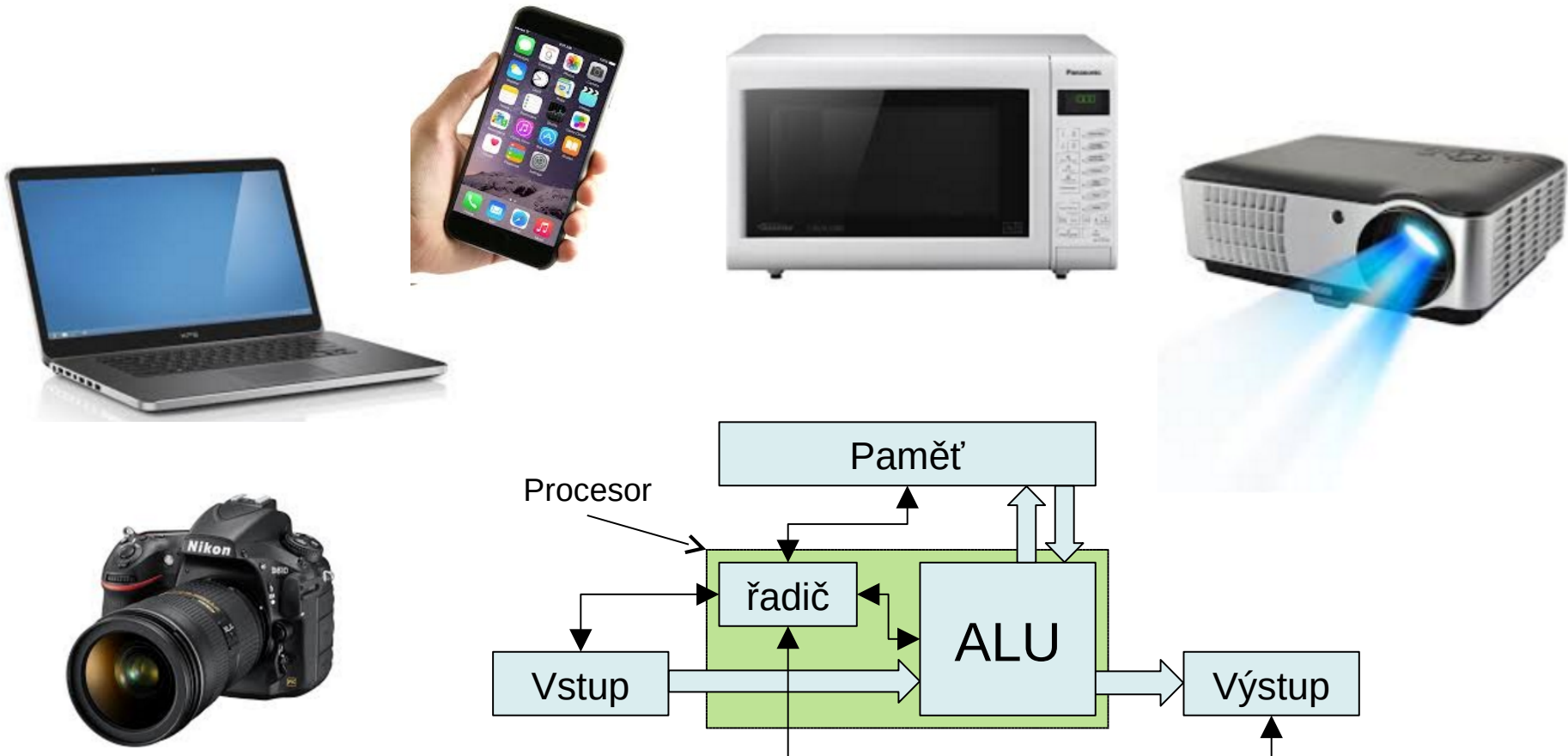


- Note the different sizes of 4 cores A7 and 4 cores A15
- On the chip, other components are integrated outside the processor: the GPU, Video coder and decoder, and more. This is SoC (System on Chip)

# Samsung Galaxy S4 inside



# Common concept

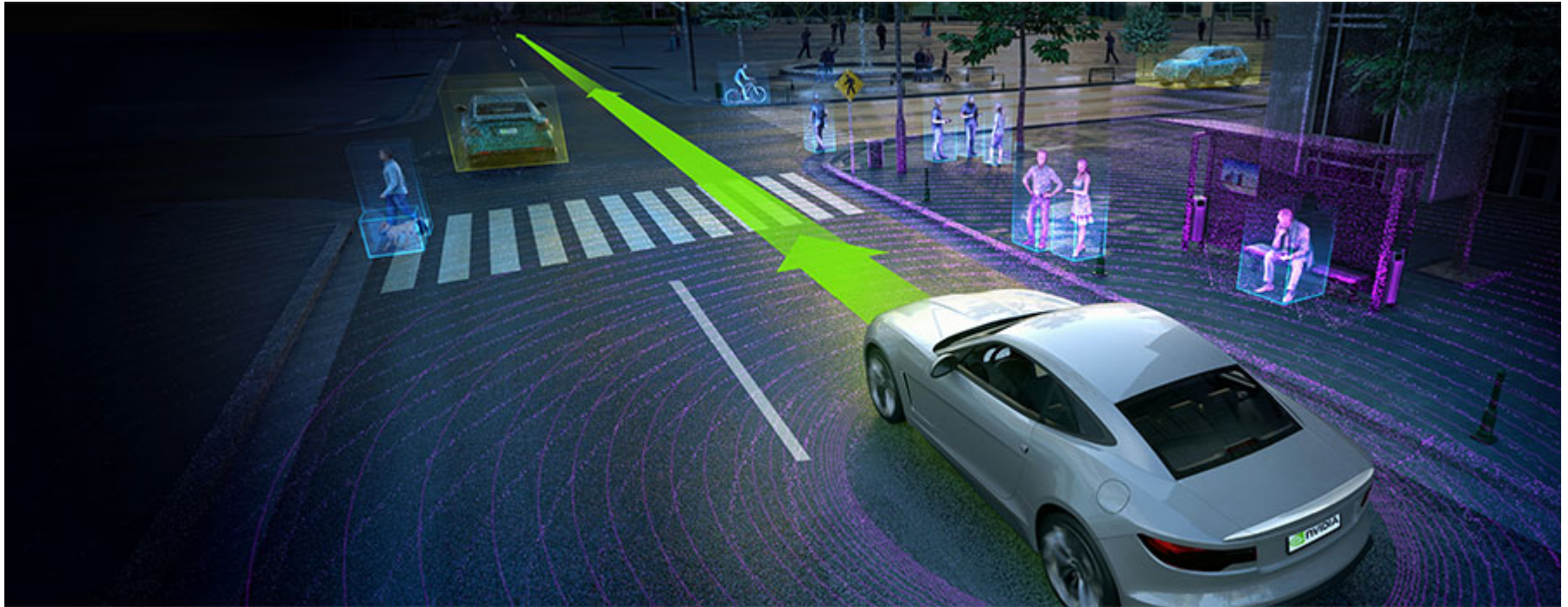


- The processor performs stored memory (ROM, RAM) instructions to operate peripherals, to respond to external events and to process data.



# Example of Optimization

## Autonomous cars



Source: <http://www.nvidia.com/object/autonomous-cars.html>

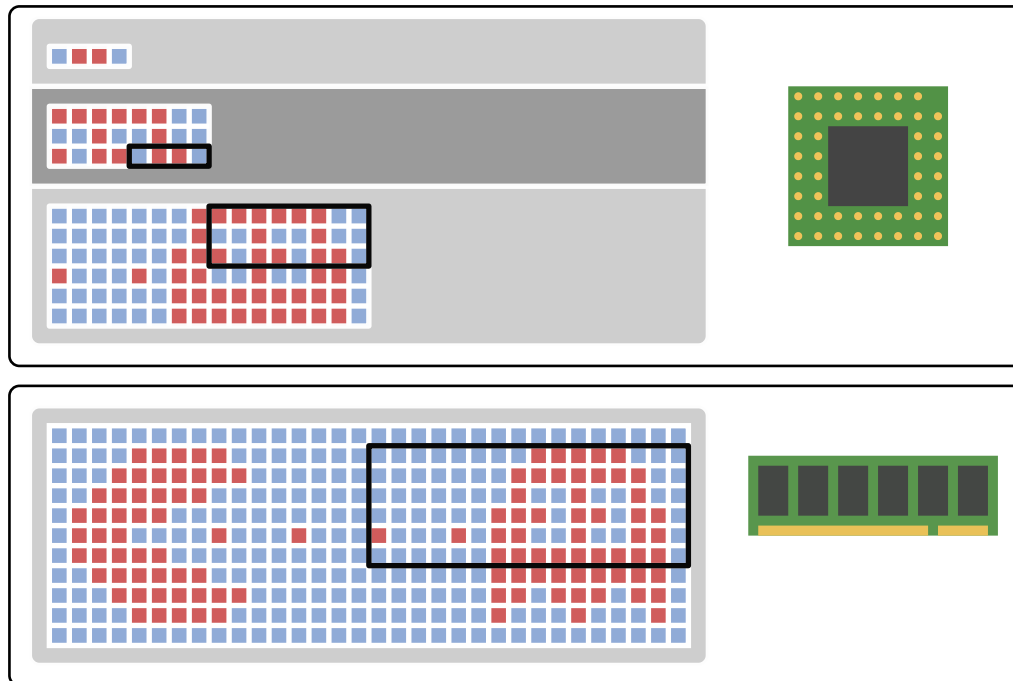
- Many artificial intelligence tasks are based on deep neural networks (deep neural networks)

# Neural network passage -> matrix multiplication

- How to increase calculation?
- The results of one of many experiments
  - Naive algorithm (3 × for) – 3.6 s = 0.28 FPS
  - Optimizing memory access – 195 ms = 5.13 FPS  
(necessary knowledge of HW)
  - 4 cores– 114 ms = 8.77 FPS  
(selection of a proper synchronization)
  - GPU (256 processors) — 25 ms = 40 FPS  
(knowledge of data transfer between CPU and coprocessors)
- Source: Naive algorithm, library Eigen (1 core), 4 cores (2 physical on i7-2520M, compiler flags -O3), GPU results Joela Matějka, Department of Control Engineering, FEE, CTU  
<https://dce.fel.cvut.cz/>
- How to speedup?

# Optimize Memory Accesses

- Algorithm modification with respect to memory hierarchy
- Data from the (buffer) memory near the processor can be obtained faster (but fast memory is small in size)



# Prediction of jumps / accesses to memory

- In order to increase average performance, the execution of instructions is divided into several phases => the need to read several instructions / data in advance
- Every condition (if, loop) means a possible jump - poor prediction is expensive
- It is good to have an idea of how the predictions work and what alternatives there are on the CPU / HW. (Eg vector / multimedia inst.)

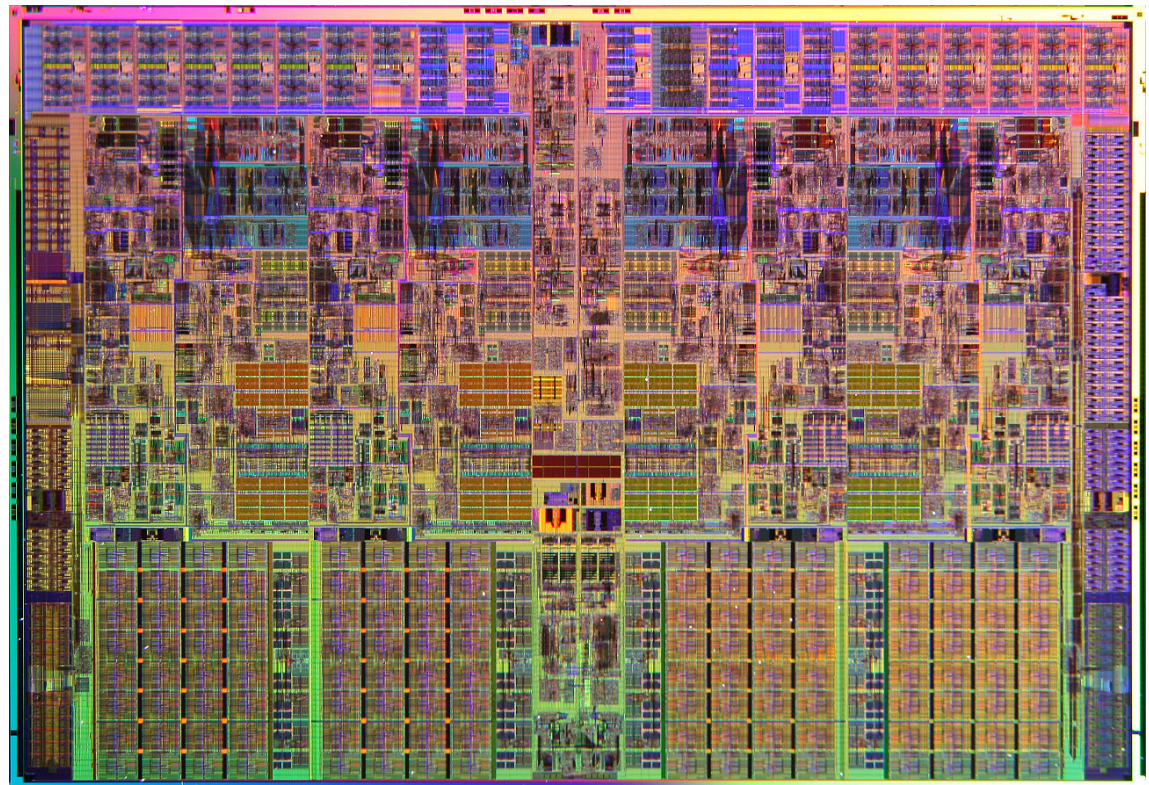


Source: [https://commons.wikimedia.org/wiki/File:Plektita\\_trakforke\\_14.jpeg](https://commons.wikimedia.org/wiki/File:Plektita_trakforke_14.jpeg)



# Parallelization - Multicore Processor

- Synchronization requirements
- Interconnection and communication possibilities between processors
- Transfers between memory levels are very expensive
- Improper sharing/access from more cores results in slower code than on a single CPU



Intel Nehalem Processor, Original Core i7

Source: [http://download.intel.com/pressroom/kits/corei7/images/Nehalem\\_Die\\_Shot\\_3.jpg](http://download.intel.com/pressroom/kits/corei7/images/Nehalem_Die_Shot_3.jpg)

# Computing Coprocessors - GPU

- Multi-core processor (hundreds)
- Some units and bclocks shared
- For effective use it is necessary to know the basic hardware features

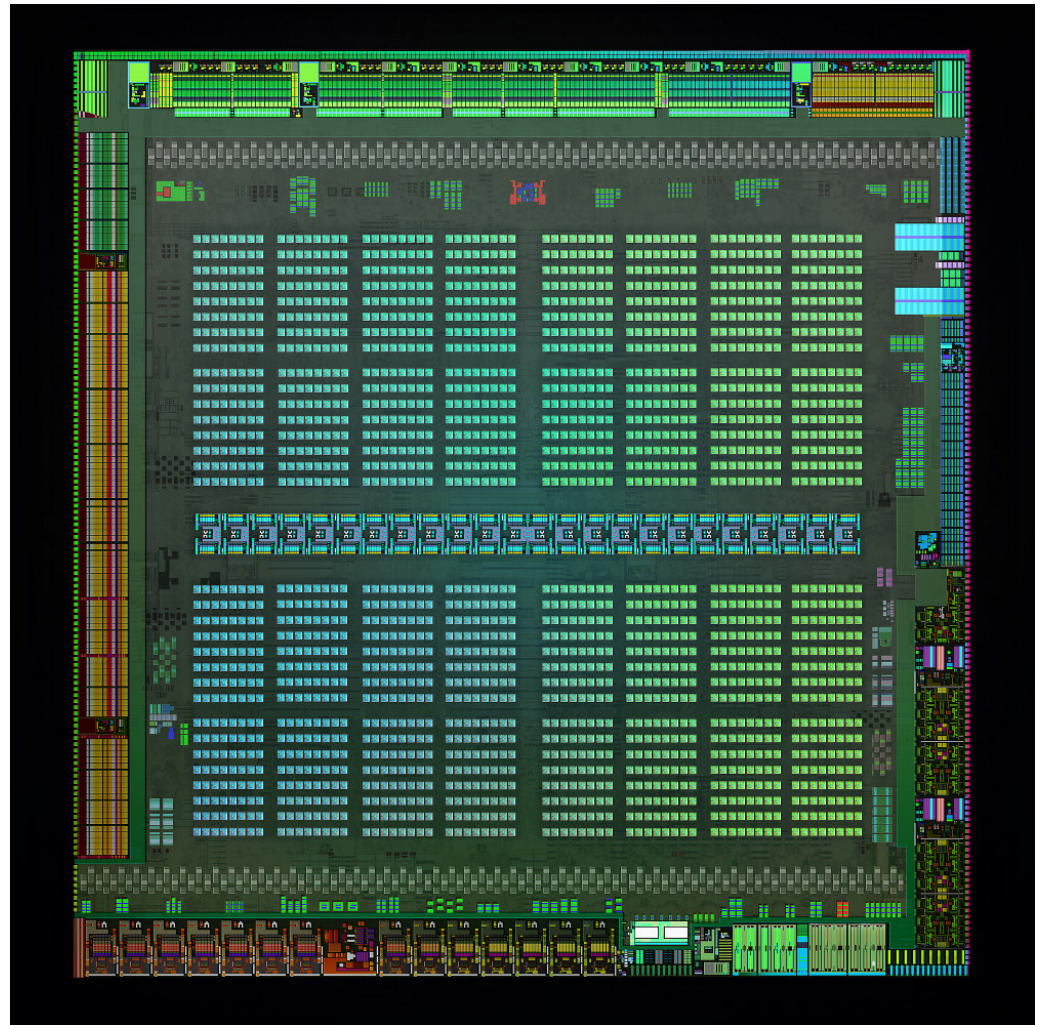


Source: <https://devblogs.nvidia.com/paralleforall/inside-pascal/>



# GPU – Maxwell

- GM204
- 5200 milins trasistors
- 398 mm<sup>2</sup>
- PCIe 3.0 x16
- 2048 computation units
- 4096 MB
- 1126 MHz
- 7010 MT/s
- 72.1 GP/s
- 144 GT/s
- 224 GB/s



Source: <http://www.anandtech.com/show/8526/nvidia-geforce-gtx-980-review/3>

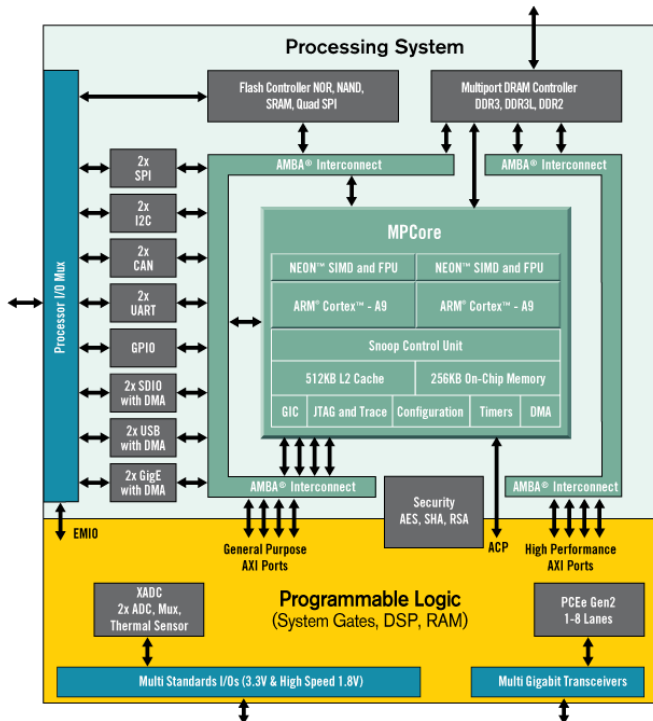
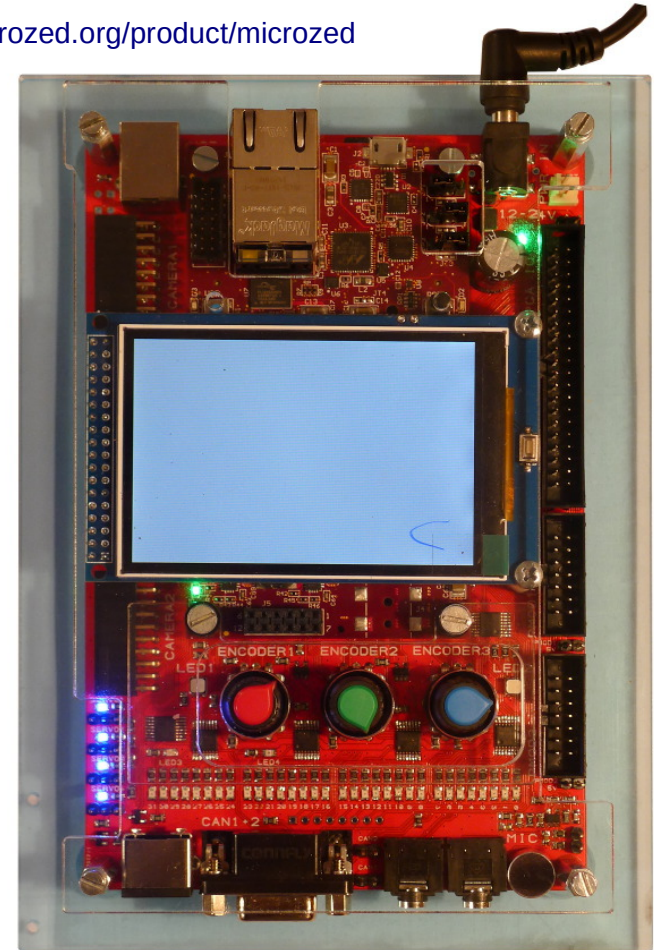
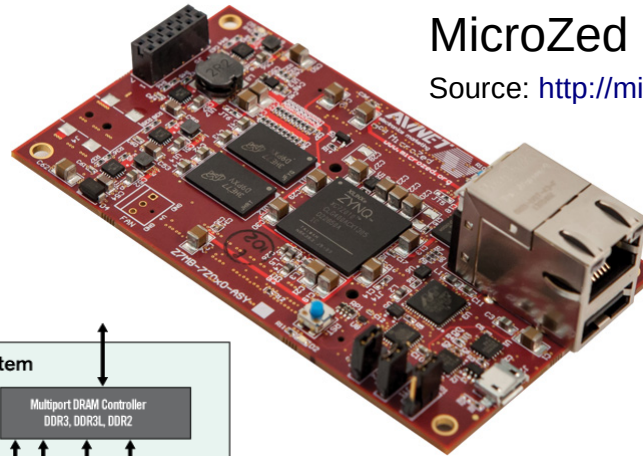
## FPGA – design/prototyping of own hardware

- Programmable logic arrays
- Well suited for effective implementation of some digital signal manipulation (filters – images, video or audio, FFT analysis, custom CPU architecture...)
- Programmer interconnects blocks available on the chip
- Zynq 7000 FPGA – two ARM cores equipped by FPGA – fast and simple access to FPGA/peripherals from own program
- (the platform is used for your seminars but you will use only design prepared by us, the FPGA programming/logic design is topic for more advanced courses)

# Xilinx Zynq 7000 a MicroZed APO

## MicroZed

Source: <http://microzed.org/product/microzed>



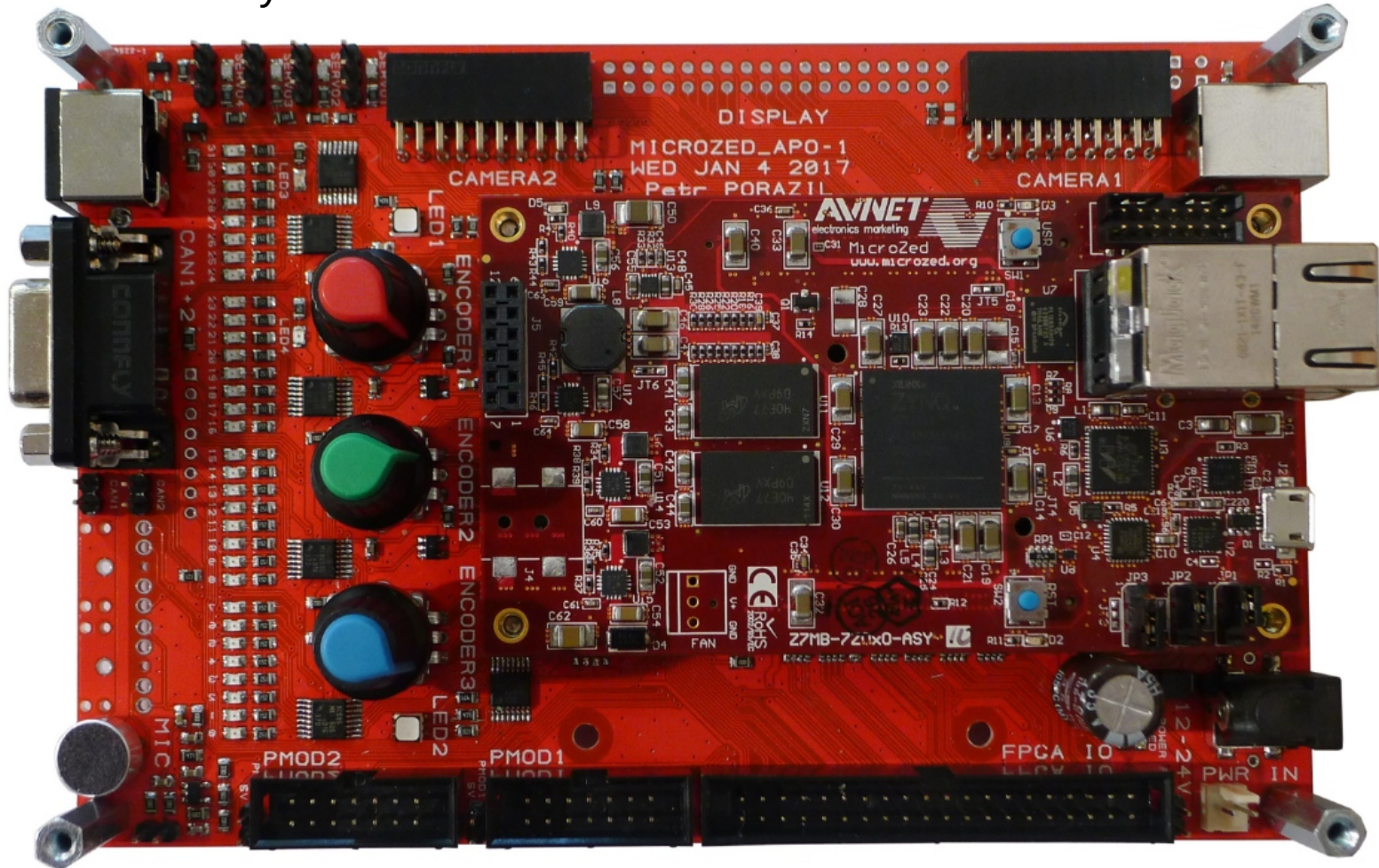
Source: <https://www.xilinx.com/products/silicon-devices/soc/zynq-7000.html>

Source: <https://cw.fel.cvut.cz/wiki/courses/b35apo/start>



# MZ\_APO board

you will later work with this board



## MZ\_APO – features

- The core chip: Zynq-7000 All Programmable SoC
- Typ: Z-7010, device XC7Z010
- CPU: Dual ARM® Cortex™-A9 MPCore™ @ 866 MHz (NEON™ & Single / Double Precision Floating Point)  
2x L1 32+32 kB, L2 512 KB
- FPGA: 28K Logic Cells (~430K ASIC logic gates, 35 kbit)
- Computational capability of FPGA DSP blocks: 100 GMACs
- Memory for FPGA design: 240 KB
- Memory on MicroZed board: 1GB
- Operating system: GNU/Linux
  - GNU LIBC (libc6) 2.19-18+deb8u7
  - Kernel: Linux 4.9.9-rt6-00002-ge6c7d1c
  - Distribution: Debian Jessie

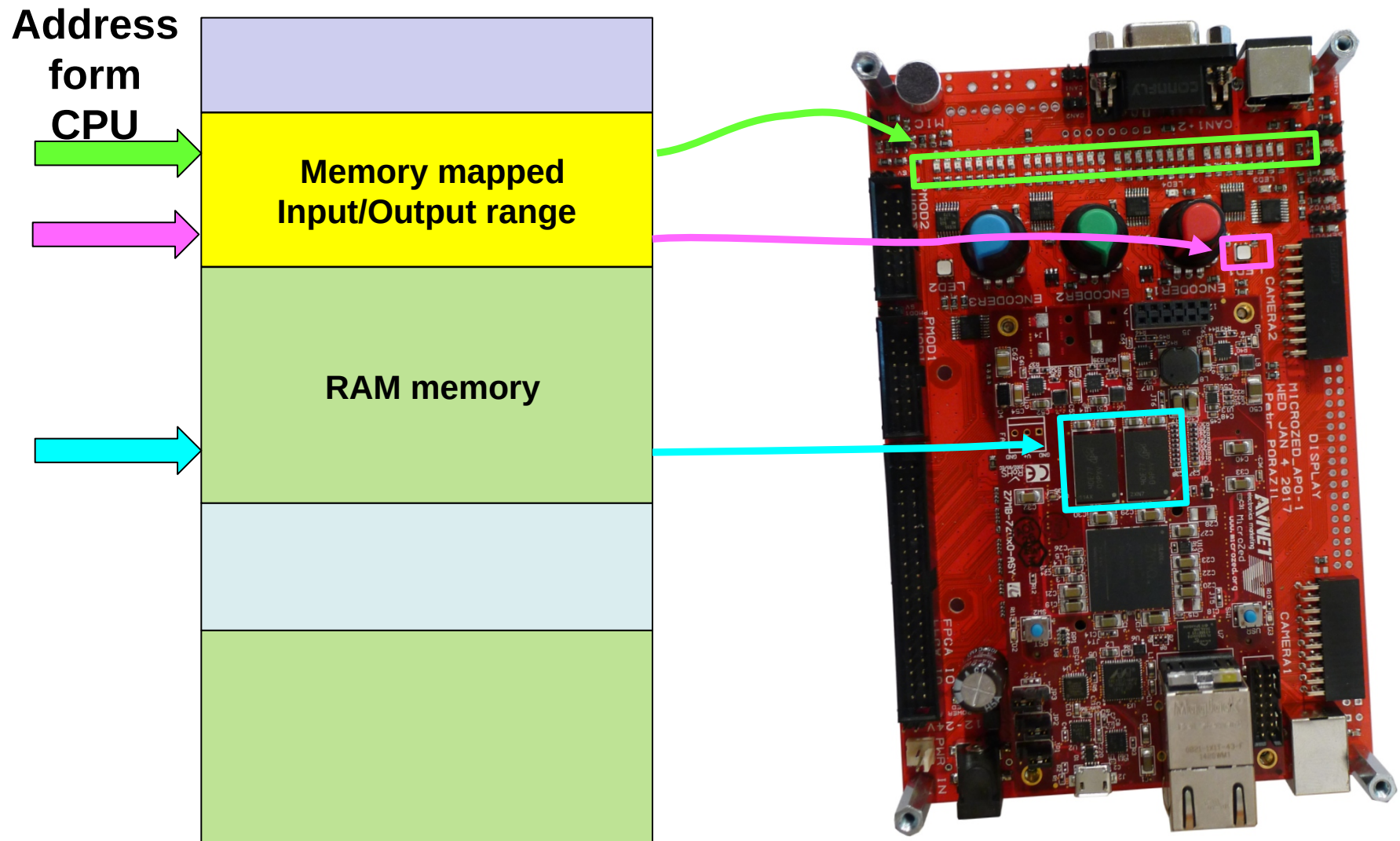
# MZ\_APO – Logic design done in Xilinx Vivado

The screenshot displays the Xilinx Vivado IDE interface for a logic design project named MZ\_APO. The main window shows a Block Design with a Hierarchy view on the left and a Diagram view on the right. The Hierarchy view shows a top-level wrapper structure containing various sub-blocks, including an AXI Interconnect, a display bus, and a processing system. The Diagram view shows the physical connections between these blocks and the ZYNQ7 Processing System. The Tcl Console at the bottom shows the following commands and output:

```
Tcl Console
Adding cell -- xilinx.com:ip:axi_protocol_converter:2.1 - auto_pc
Adding cell -- xilinx.com:ip:axi_protocol_converter:2.1 - auto_pc
Adding cell -- xilinx.com:ip:axi_protocol_converter:2.1 - auto_pc
Adding cell -- xilinx.com:ip:axi_protocol_converter:2.1 - auto_pc
Successfully read diagram <top> from BD file </home/pi/fpga/zynq/canbench-sw/system/src/top.bd>
open_bd_design: Time (s): cpu = 00:00:24 ; elapsed = 00:00:19 . Memory (MB): peak = 6008.051 ; gain = 153.621 ; free physical = 80 ; free virtual = 7868
set_property location {-22 483} [get_bd_ports CAN2_RXD]
set_property location {-26 1138} [get_bd_ports ENCDATA]
write_bd_layout -format pdf -orientation portrait /home/pi/mz_apo-v10-top.pdf
/home/pi/mz_apo-v10-top.pdf
```



# The first seminar – physical address space on MZ\_APO





# Linux – from tiny to supercomputers

- TOP500 <https://www.top500.org/> (<https://en.wikipedia.org/wiki/TOP500> )
  - Actual top one: Summit supercomputer – IBM AC922
  - June 2018, US Oak Ridge National Laboratory (ORNL),
  - 200 PetaFLOPS, 4600 “nodes”, 2× IBM Power9 CPU +
  - 6× Nvidia Volta GV100
  - 96 lanes of PCIe 4.0, 400Gb/s
  - NVLink 2.0, 100GB/s CPU-to-GPU,
  - GPU-to-GPU
  - 2TB DDR4-2666 per node
  - 1.6 TB NV RAM per node
  - 250 PB storage
  - POWER9-SO, Global Foundries 14nm FinFET,  
8×109 tran., 17-layer, 24 cores, 96 threads (SMT4)
  - 120MB L3 eDRAM (2 CPU 10MB), 256GB/s
- Other example: SGI SSI (single system image) Linux, 2048 Itanium CPU a 4TiB RAM



Source: <http://www.tomshardware.com/>

# Linux kernel and open-source

- Linux kernel project
  - 13,500 developers from 2005 year
  - 10,000 lines of code inserted daily
  - 8,000 removed and 1,500 till 1,800 modified
  - GIT source control system
- Many successful open-source projects exists
- Open for joining by everybody
- Google Summer of Code for university students
  - <https://developers.google.com/open-source/gsoc/>

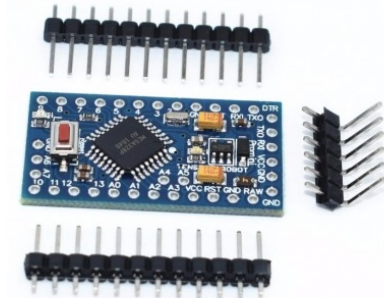
Zdroj: [https://www.theregister.co.uk/2017/02/15/think\\_different\\_shut\\_up\\_and\\_work\\_harder\\_says\\_linus\\_torvalds/](https://www.theregister.co.uk/2017/02/15/think_different_shut_up_and_work_harder_says_linus_torvalds/)



# Back to the Motivational Example of Autonomous Driving

The result of a good knowledge of hardware

- Acceleration (in our case  $18 \times$  using the same number of cores)
  - Reduce the power required
  - Energy saving
  - Possibility to reduce current solutions
  - Using GPUs, we process 40 fps.
- But in an **embedded** device, it is sometimes necessary to reduce its consumption and cost. There are used very simple processors or microcontrollers, sometimes without real number operations, and programmed with low-level C language.

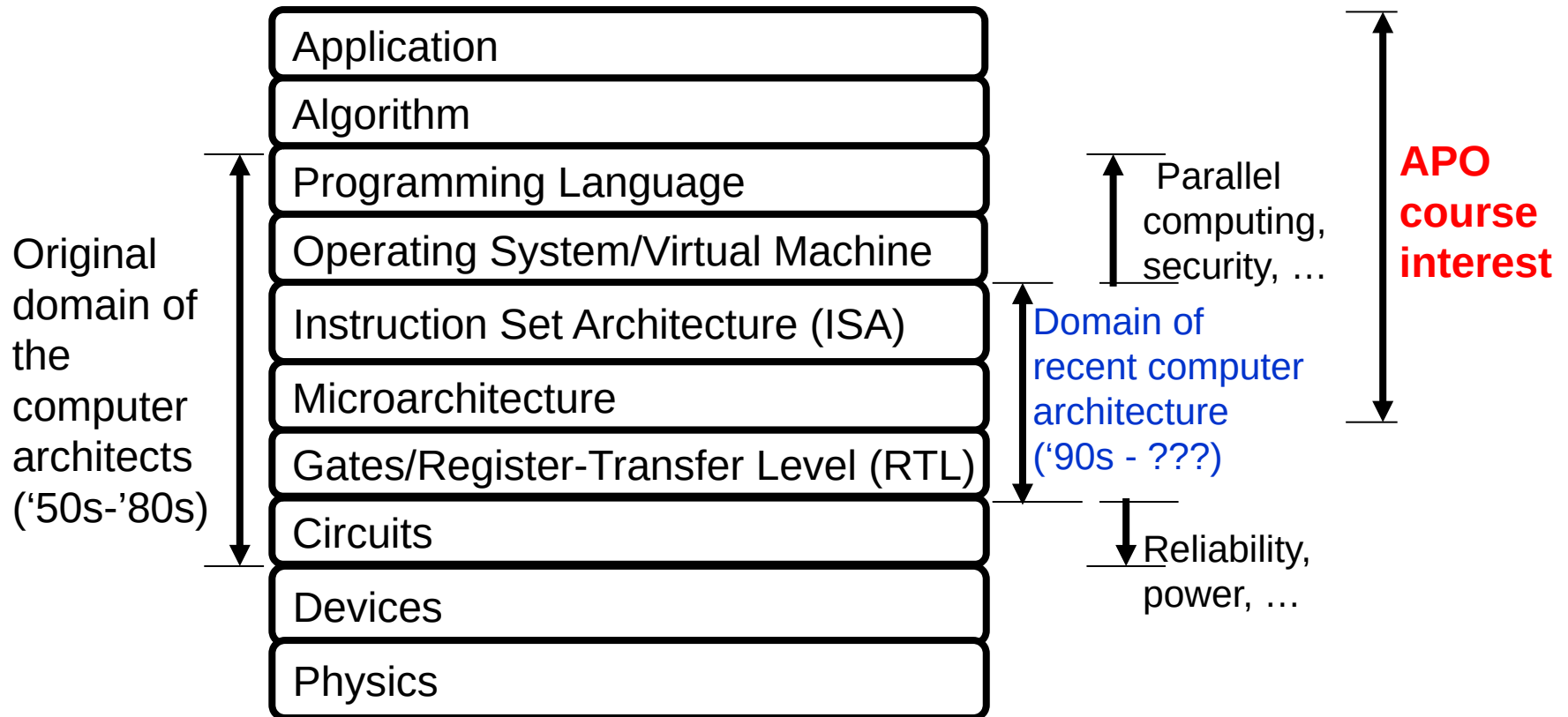


## Applicability of Knowledge and Techniques from the Course

- Applications not only in autonomous control
- In any embedded device - reduce size, consumption, reliability
- In data sciences - considerably reduce runtime and energy savings in calculations
- In the user interface - improving application response
- Practically everywhere...



# Computer



Reference: John Kubiatowicz: EECS 252 Graduate Computer Architecture, Lecture 1. University of California, Berkeley

## Reasons to study computer architectures

- To invent/design new computer architectures
- To be able to integrate selected architecture into silicon
- To gain knowledge required to design computer hardware/ systems (big ones or embedded)
- To understand generic questions about computers, architectures and performance of various architectures
- **To understand how to use computer hardware efficiently** (i.e. how to write good software)
  - It is not possible to efficiently use resources provided by any (especially by modern) hardware without insight into their constraints, resource limits and behavior
  - It is possible to write some well paid applications without real understanding but this requires abundant resources on the hardware level. But no interesting and demanding tasks can be solved without this understanding.

## More motivation and examples

- The knowledge is necessary for every programmer who wants to work with medium size data sets or solve little more demanding computational tasks
- No multimedia algorithm can be implemented well without this knowledge
- The 1/3 of the course is focussed even on peripheral access
- Examples
  - Facebook – HipHop for PHP → C++/GCC → machine code
  - BlackBerry (RIM) – our consultations for time source
  - RedHat – JAVA JIT for ARM for future servers generation
  - Multimedia and CUDA computations
  - Photoshop, GIMP (data/tiles organization in memory)
  - Knot-DNS (RCU, Copy on write, Cuckoo hashing, )

## The course's background and literature

- Course is based on worldwide recognized book and courses; evaluation Graduate Record Examination – GRE  
Paterson, D., Henessy, J.: Computer Organization and Design, The HW/SW Interface. Elsevier, ISBN: 978-0-12-370606-5
  - John L. Henessy – president of Stanford University, one of founders of MIPS Computer Systems Inc.
  - David A. Patterson – leader of Berkeley RISC project and RAID disks research
- Our experience even includes distributed systems, embedded systems design (of mobile phone like complexity), peripherals design, cooperation with carmakers, medical and robotics systems design

# Topics of the lectures

- Architecture, structure and organization of computers and its subsystems.
- Floating point representation
- Central Processing Unit (CPU)
- Memory
- Pipelined instruction execution
- Input/output subsystem of the computer
- Input/output subsystem (part 2)
- External events processing and protection
- Processors and computers networks
- Parameter passing
- Classic register memory-oriented CISC architecture
- INTEL x86 processor family
- CPU concepts development (RISC/CISC) and examples
- Multi-level computer organization, virtual machines



## Topics of seminars

- 1 - Introduction to the lab
- 2 - Data representation in memory and floating point
- 3 - Processor instruction set and algorithm rewriting
- 4 - Hierarchical concept of memories, cache - part 1
- 5 - Hierarchical concept of memories, cache - part 2
- 6 - Pipeline and gambling
- 7 - Jump prediction, code optimization
- 8 - I / O space mapped to memory and PCI bus
- 9 - HW access from C language on MZ\_APO
- Semestral work

# Classification and Conditions to Pass the Subject

Conditions for assessment:

Category	Points	Required minimum	Remark
4 homeworks	36	12	3 of 4
Activity	8	0	
Team project	24	5	
<b>Sum</b>	<b>60 (68)</b>	30	

Exam:

Category	Points	Required minimum
Written exam part	30	15
Oral exam part	+/- 10	0

Grade	Points range
A	90 and more
B	80 - 89
C	70 - 79
D	60 - 69
E	50 - 59
F	less than 50

# The 1. lecture contents

- Number representation in computers
  - numeral systems
  - integer numbers, unsigned and signed
  - boolean values
- Basic arithmetic operations and their implementation
  - addition, subtraction
  - shift right/left
  - multiplication and division

Motivation: What is the output of next code snippet?

```
int main() {  
    int a = -200;  
    printf("value: %u = %d = %f = %c \n", a, a,  
        *((float*)&a), a);  
  
    return 0;  
}
```

value: 4294967096 = -200 = nan = 8

and memory content is: 0x38 0xff 0xff 0xff  
when run on little endian 32 bit CPU.



## 1<sup>st</sup> lecture

- How they are stored on your computer
  - INTEGER numbers, with or without sign?
- How to perform basic operations
  - Adding, Subtracting,
  - Multiplying

# Non-positional numbers 😊



<http://diameter.si/sciquest/E1.htm>

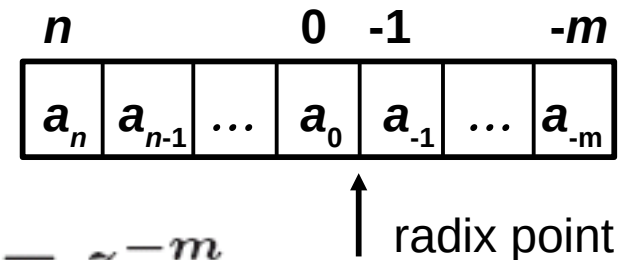


**10, 100, 1000, 10000, 100000, 1 million**

The value is the sum: 1 333 331

## Terminology basics

- Positional (place-value) notation
- Decimal/radix point
- $z$  ... base of numeral system
- smallest representable number  $\epsilon = z^{-m}$
- **Module** =  $Z$  , one increment/unit higher than biggest representable number for given encoding/notation
- **A**, the representable number for given  $n$  and  $m$  selection, where  $k$  is natural number in range  $\langle 0, z^{n+m+1} - 1 \rangle$
- The representation and value



$$0 \leq A = k \cdot \epsilon < Z$$

$$A \sim a_n a_{n-1} \dots a_0, a_1 \dots a_{-m}$$

$$A = a_n z^n + a_{n-1} z^{n-1} + \dots + a_0 + a_1 z^{-1} \dots a_{-m} z^{-m}$$



# Unsigned integers

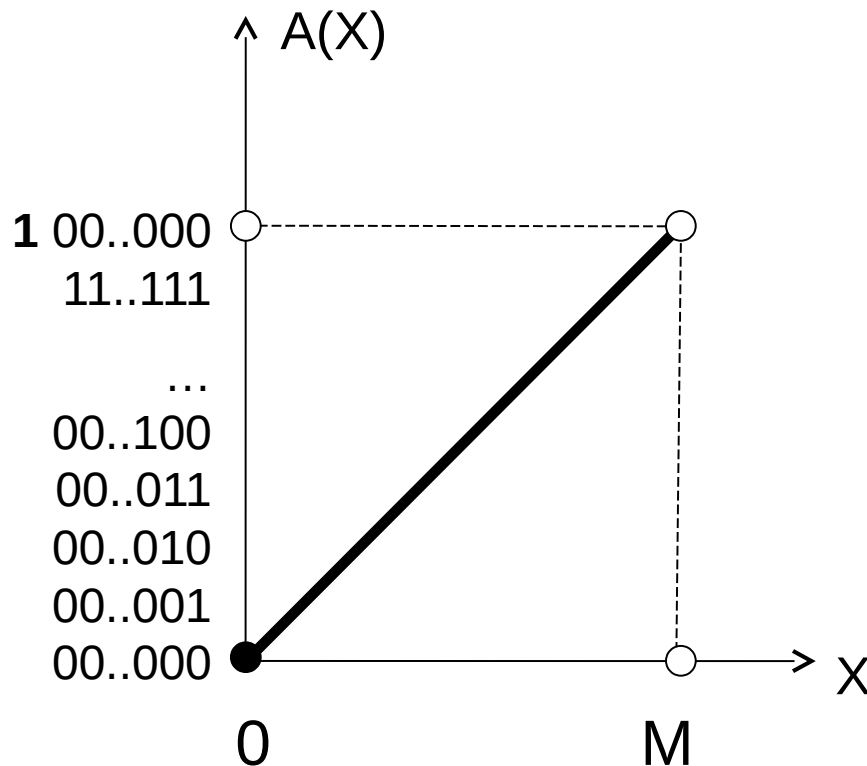
Language C:  
`unsigned int`



## Integer number representation (unsigned, non-negative)

- The most common numeral system base in computers is  $z=2$
- The value of  $a_i$  is in range  $\{0,1,\dots,z-1\}$ , i.e.  $\{0,1\}$  for base 2
- This maps to true/false and unit of information (bit)
- We can represent number  $0 \dots 2^n-1$  when  $n$  bits are used
- Which range can be represented by one byte?  
1B (byte) ... 8 bits,  $2^8 = 256_d$  combinations, values  $0 \dots 255_d = 0b11111111_b$
- Use of multiple consecutive bytes
  - 2B ...  $2^{16} = 65536_d$ ,  $0 \dots 65535_d = 0xFFFF_h$ , (h ... hexadecimal, base 16, a in range 0, ... 9, A, B, C, D, E, F)
  - 4B ...  $2^{32} = 4294967296_d$ ,  $0 \dots 4294967295_d = 0xFFFFFFFF_h$

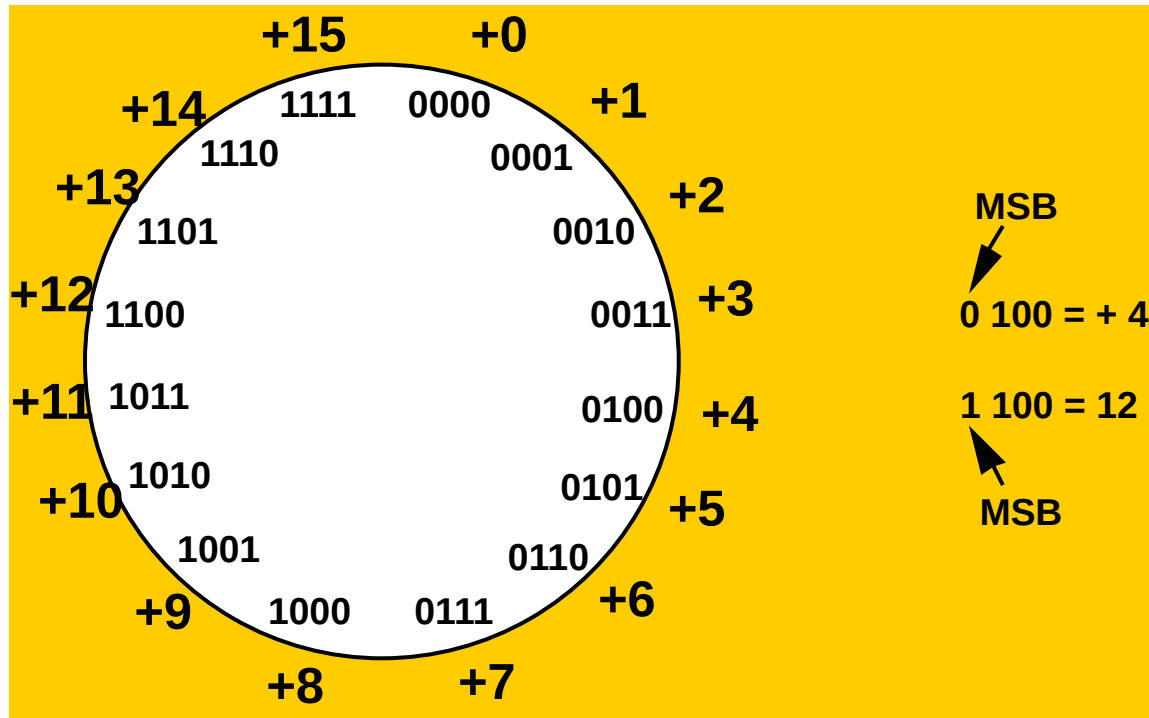
# Unsigned integer



binary value	unsigned int
00000000	$0_{(10)}$
00000001	$1_{(10)}$
⋮	⋮
01111101	$125_{(10)}$
01111110	$126_{(10)}$
01111111	$127_{(10)}$
10000000	$128_{(10)}$
10000001	$129_{(10)}$
10000010	$130_{(10)}$
⋮	⋮
11111101	$253_{(10)}$
11111110	$254_{(10)}$
11111111	$255_{(10)}$

# Unsigned 4-bit numbers

*Assumptions: we'll assume a 4 bit machine word*



■ Cumbersome subtraction



# Signed numbers

Language C:

`int`

`signed int`



# Two's Complement.

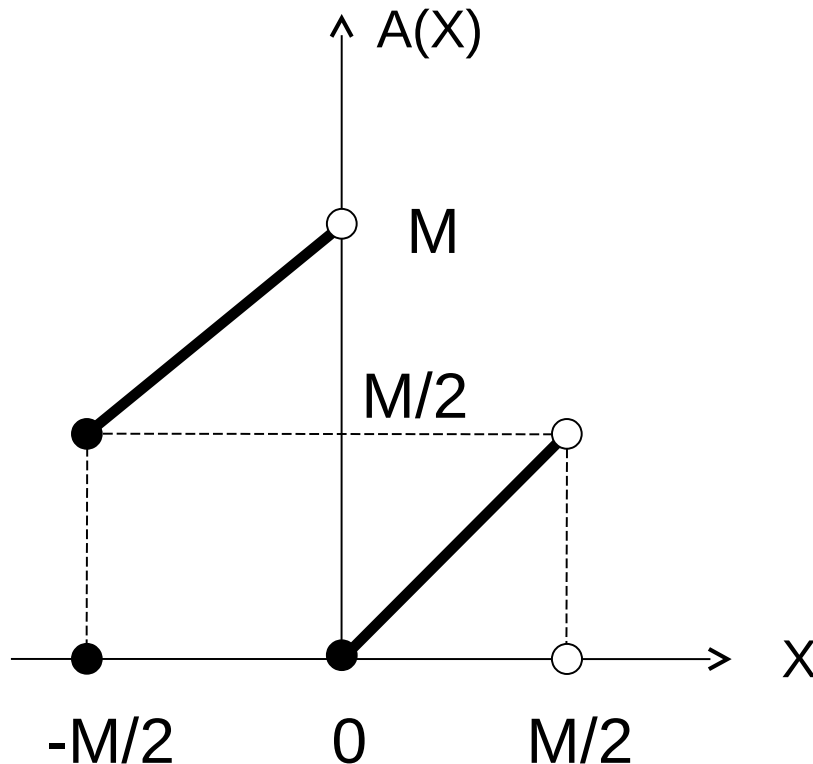
- The most frequent code
- The sum of two opposite numbers with the same absolute value is 00000000H!

Decimal value	4 bit two's complement
6	0110
-6	1010

# Two's Complement

## Dvojkový doplněk – pokračování...

- Pokud **N** bude počet bitů:  
 **$\langle -2^{N-1}, 2^{N-1} - 1 \rangle$**



Binární hodnota	Dvojkový doplněk
00000000	$0_{(10)}$
00000001	$1_{(10)}$
⋮	⋮
01111101	$125_{(10)}$
01111110	$126_{(10)}$
01111111	$127_{(10)}$
10000000	$-128_{(10)}$
10000001	$-127_{(10)}$
10000010	$-126_{(10)}$
⋮	⋮
11111101	$-3_{(10)}$
11111110	$-2_{(10)}$
11111111	$-1_{(10)}$

## Two's complement - examples

- Examples:

- $0_D = 00000000_H,$

- $1_D = 00000001_H,$

- $2_D = 00000002_H,$

- $3_D = 00000003_H,$

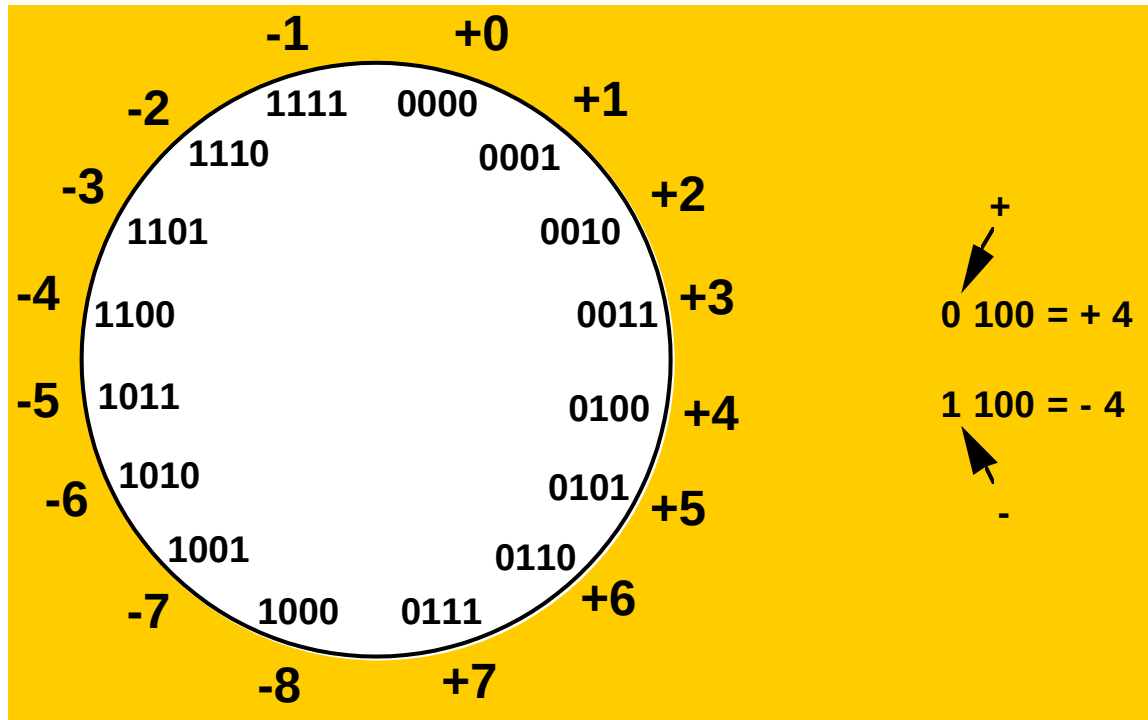
- $-1_D = FFFFFFFF_H,$

- $-2_D = FFFFFFFE_H,$

- $-3_D = FFFFFFFD_H,$

# Number Representations

## **Twos Complement** (In Czech: Druhý doplněk)



Only one representation for 0

One more negative number than positive number

## Two's complement – addition and subtraction

- **Addition**

- $00000000\ 0000\ 0111_B \approx 7_D$     Symbols use:  $0=0_H$ ,  $0=0_B$

- $+ \underline{00000000\ 0000\ 0110_B} \approx 6_D$

- $00000000\ 0000\ 1101_B \approx 13_D$

- **Subtraction** can be realized as addition of negated number

- $00000000\ 0000\ 0111_B \approx 7_D$

- $+ \underline{FFFFFFF\ 1111\ 1010_B} \approx -6_D$

- $00000000\ 0000\ 0001_B \approx 1_D$

- Question for revision: how to obtain negated number in two's complement binary arithmetics?





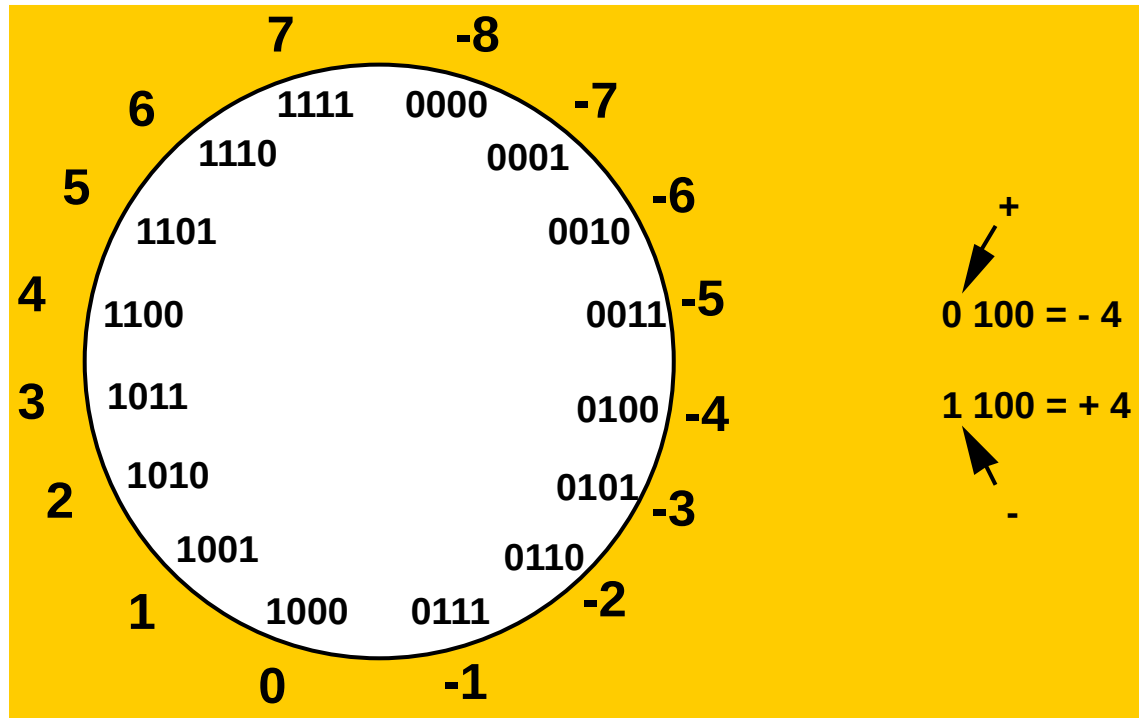
# Other Possibilities

## Integer – biased representation

- Known as excess-K or offset binary as well
- Transform to the representation  $-K \dots 0 \dots 2^n-1-K$   
 $D(A) = A+K$
- Usually  $K=Z/2$
- Advantages
  - Preserves order of original set in mapped set/representation
- Disadvantages
  - Needs adjustment by  $-K$  after addition and  $+K$  after subtraction processed by unsigned arithmetic unit
  - Requires full transformation before and after multiplication

# Number Systems

## Excess-K, offset binary or biased representation

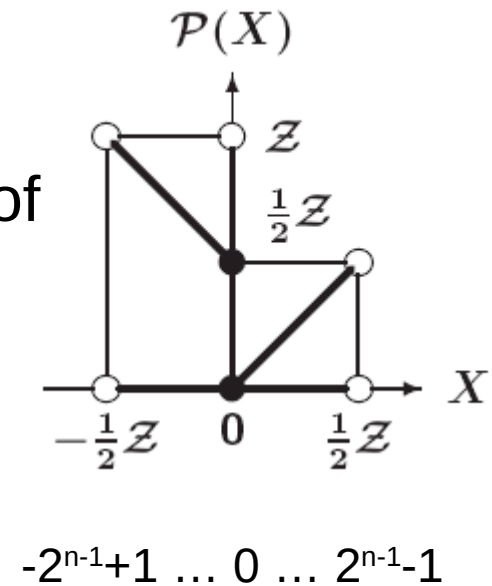


One 0 representation, we can select count of negative numbers -  
*used e.g. for exponents of real numbers..*

Integer arithmetic unit are not designed to calculate with Excess-K numbers  
[Seungryoul Maeng:Digital Systems]

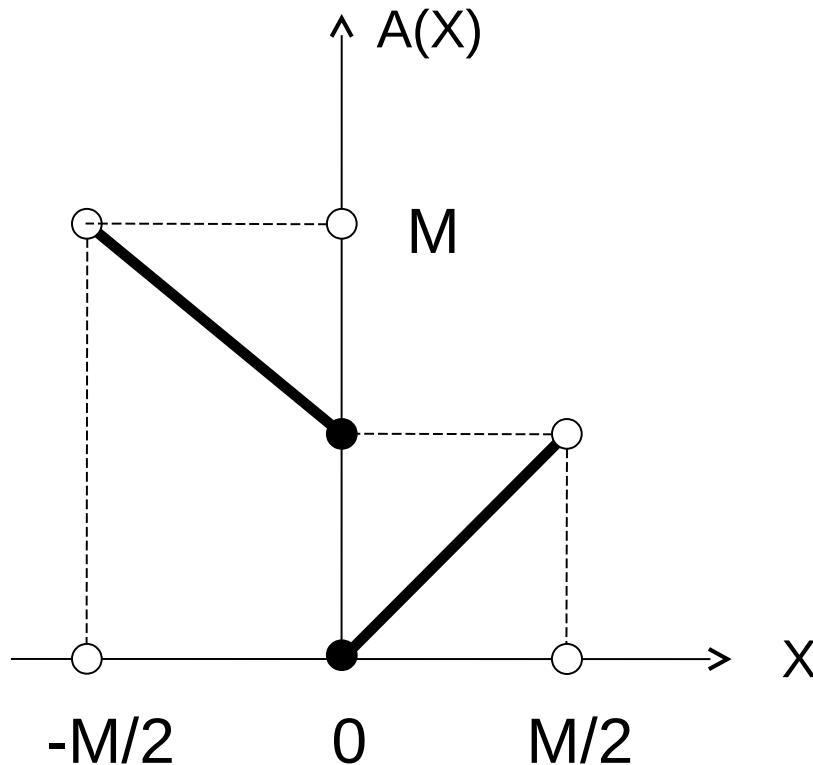
## Integer – sign-magnitude code

- Sign and magnitude of the value (absolute value)
- Natural to humans -1234, 1234
- One (usually most significant – MSB) bit of the memory location is used to represent the sign
- Bit has to be mapped to meaning
- Common use  $0 \approx “+”$ ,  $1 \approx “-”$
- Disadvantages:
  - When location is **k** bits long then only **k-1** bits hold magnitude and each operation has to separate sign and magnitude
  - Two representations of the value 0



# Sign and Magnitude Representation.

$\langle -2^{N-1} - 1, 2^{N-1} - 1 \rangle$

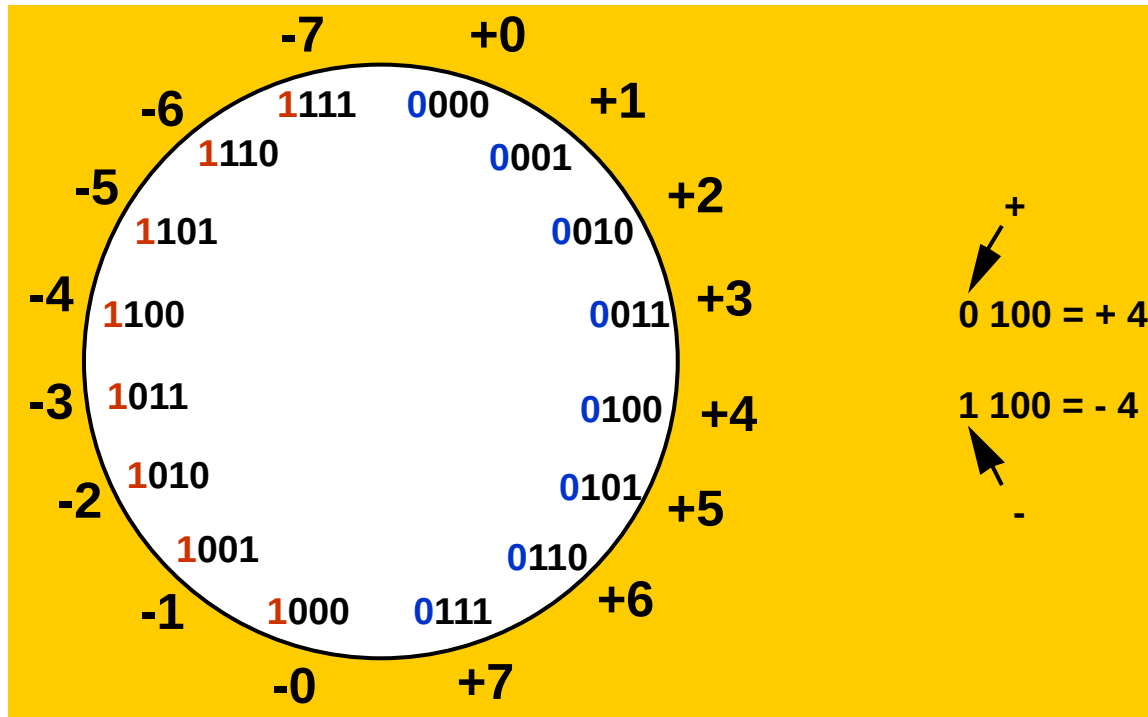


Binary value	Code
<b>00000000</b>	$+0_{(10)}$
<b>00000001</b>	$1_{(10)}$
⋮	⋮
<b>01111101</b>	$125_{(10)}$
<b>01111110</b>	$126_{(10)}$
<b>01111111</b>	$127_{(10)}$
<b>10000000</b>	$-0_{(10)}$
<b>10000001</b>	$-1_{(10)}$
<b>10000010</b>	$-2_{(10)}$
⋮	⋮
<b>11111101</b>	$-125_{(10)}$
<b>11111110</b>	$-126_{(10)}$
<b>11111111</b>	$-127_{(10)}$



# Number Systems

## *Sign and Magnitude Representation*



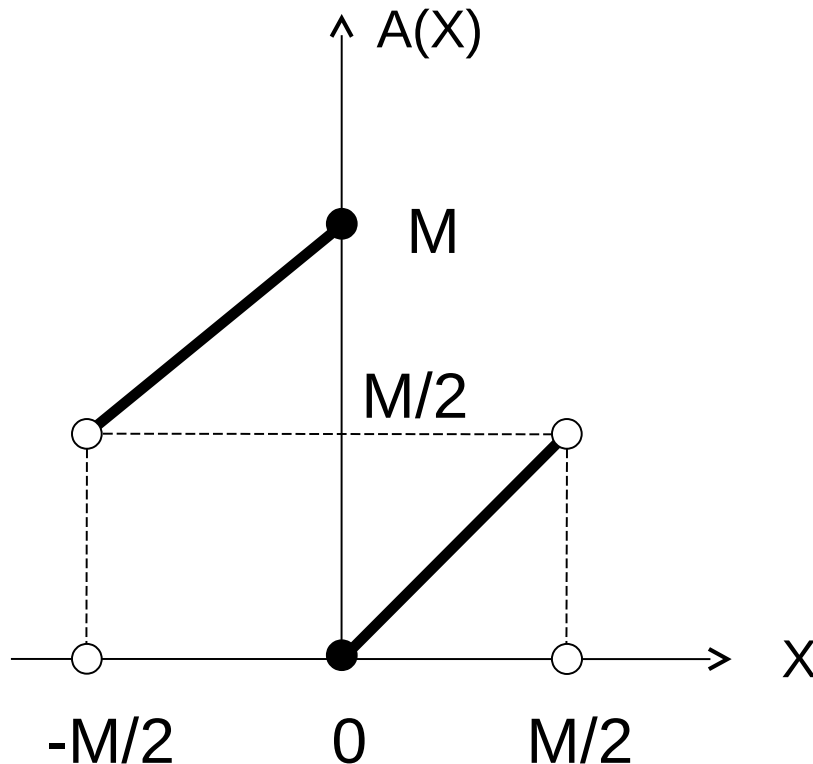
- Cumbersome addition/subtraction
- *Sign+Magnitude usually used only for float point numbers*

## Integers – ones' complement

- Transform to the representation  $-2^{n-1}+1 \dots 0 \dots 2^{n-1}-1$ 
  - $D(A) = A$                       iff  $A \geq 0$
  - $D(A) = Z-1-|A|$             iff  $A < 0$  (i.e. subtract from all ones)
- Advantages
  - Symmetric range
  - Almost continuous, requires hot one addition when sign changes
- Disadvantage
  - Two representations of value 0
  - More complex hardware
- Negate ( $-A$ ) value can be computed by bitwise complement (flipping) of each bit in representation

# Ones Complement

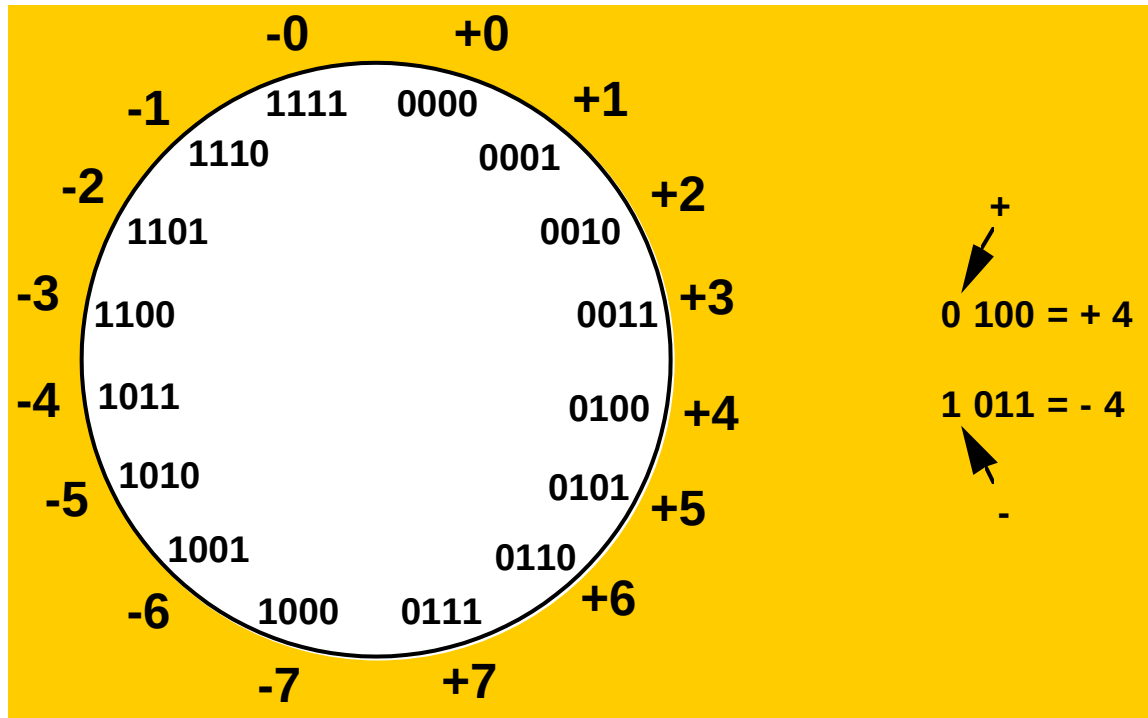
$$\langle -2^{N-1} - 1, 2^{N-1} - 1 \rangle$$



Binary value	Code
<b>00000000</b>	$0_{(10)}$
<b>00000001</b>	$1_{(10)}$
⋮	⋮
<b>01111101</b>	$125_{(10)}$
<b>01111110</b>	$126_{(10)}$
<b>01111111</b>	$127_{(10)}$
<b>10000000</b>	$-127_{(10)}$
<b>10000001</b>	$-126_{(10)}$
<b>10000010</b>	$-125_{(10)}$
⋮	⋮
<b>11111101</b>	$-2_{(10)}$
<b>11111110</b>	$-1_{(10)}$
<b>11111111</b>	$-0_{(10)}$

# Number Systems

## Ones Complement (In Czech: První doplněk)



Still two representations of 0! This causes some problems  
*Some complexities in addition, nowadays nearly not used*



# OPERATION WITH INTEGERS



# Direct realization of adder as logical function

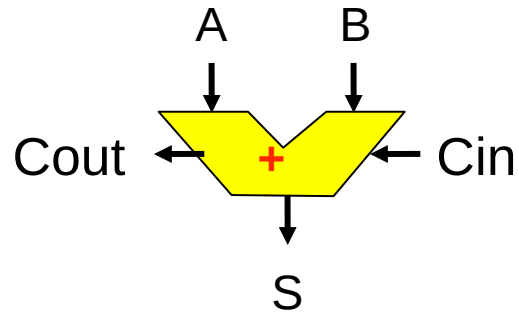
bit width	Number of logic operations for calculating sum
1	3
2	22
3	89
4	272
5	727
6	1567
7	3287
8	7127
9	17623
10	53465
11	115933

Complexity is higher than  $O(2^n)$

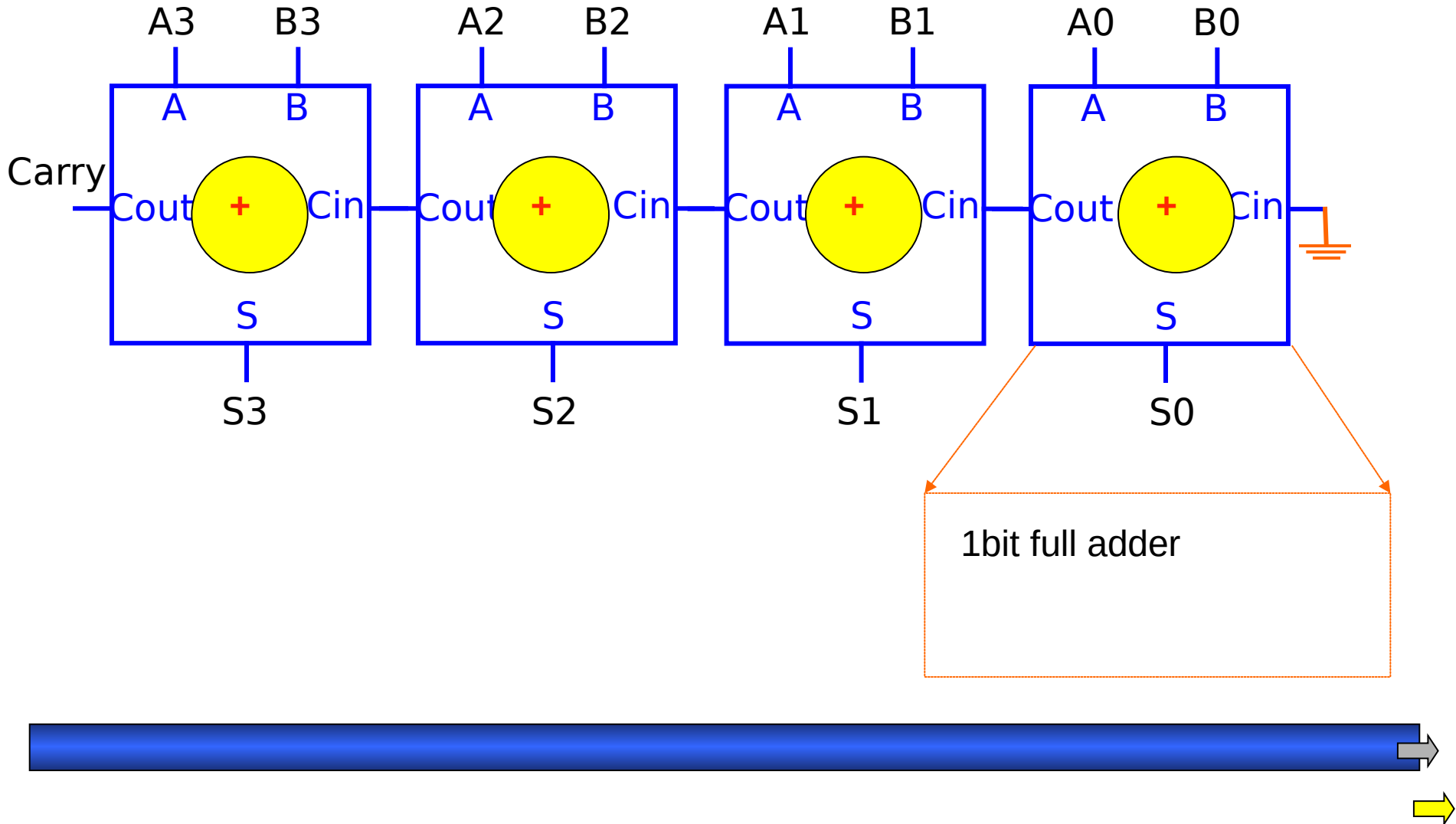
*The calculation was performed by BOOM logic minimizer  
created at the Department of Computer Science CTU-FEE*

# 1bit Full Adder

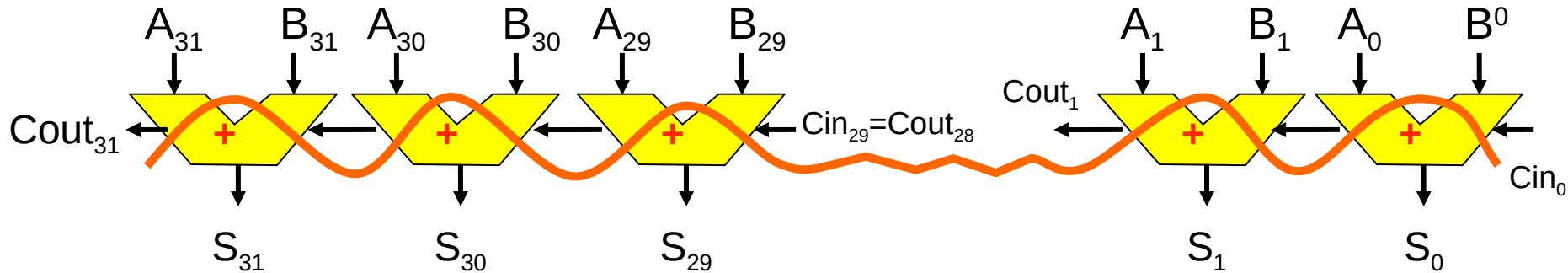
A	0	0	1	1	0	0	1	1
+B	0	1	0	1	0	1	0	1
Sum	00	01	01	10	00	01	01	10
+ Carry-In	0	0	0	0	1	1	1	1
CarryOut Sum	00	01	01	10	01	10	10	11



# Adder



# Simple Adder



## Simplest N-bit adder

we chain 1-bit full adders

"Carry" ripple through their chain

Minimal number of logical elements

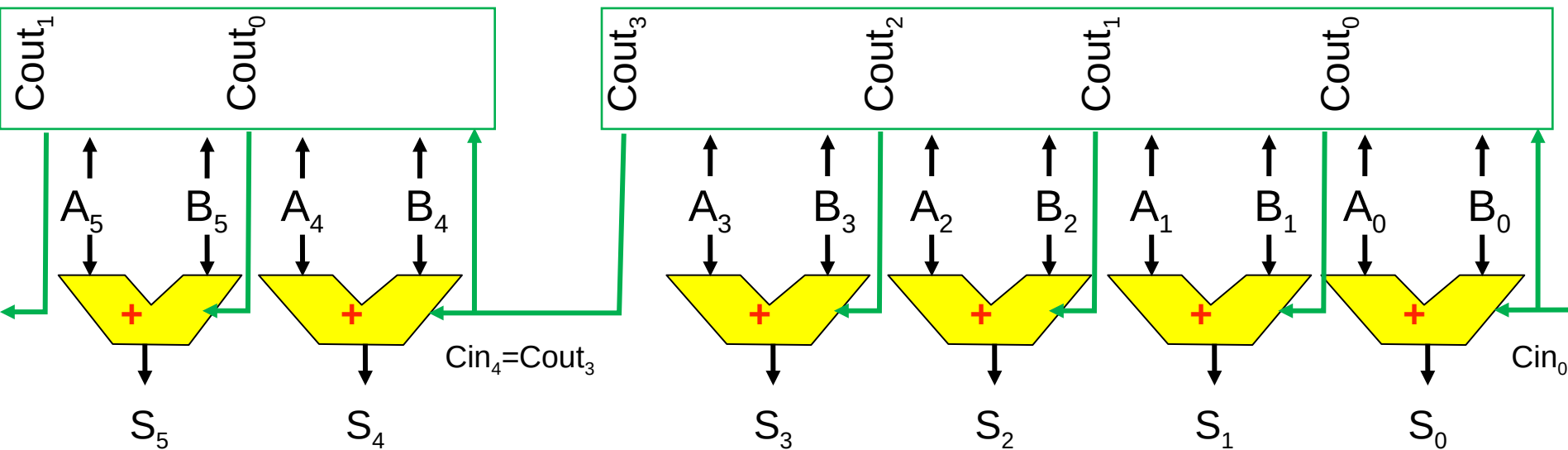
Delay is given by the last Cout -  $2 \cdot (N-1) + 3$  gates of the last adder

=  $(2N+1)$  times propagation delay of 1 gate



# 32bit CLA "carry look-ahead" adder

The carry-lookahead adder calculates one or more carry bits before the sum, which reduces the wait time to calculate the result of the larger value bits



Static "carry look ahead (CLA)" unit for 4 bits





# Increment / Decrement

*Very fast operations that do not need an adder!*

The last bit is always negated, and the previous ones are negated according to the end 1 / 0

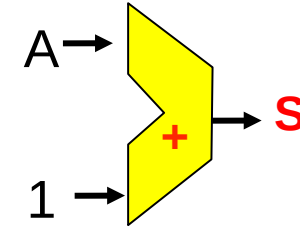
Dec.	Binary 8 4 2 1	+1	Binary 8 4 2 1	-1
0	0000	0001	0000	1111
1	0001	0010	0001	0000
2	0010	0011	0010	0001
3	0011	0100	0011	0010
4	0100	0101	0100	0011
5	0101	0110	0101	0100
6	0110	0111	0110	0101
7	0111	1000	0111	0110
8	1000	1001	1000	0111
9	1001	1010	1001	1000
10	1010	1011	1010	1001
11	1011	1100	1011	1010
12	1100	1101	1100	1011
13	1101	1110	1101	1100
14	1110	1111	1110	1101
15	1111	0000	1111	1110

## Special Case +1/-1

$$S_0 = \text{not } A_0$$

$$S_1 = A_1 \text{ xor } A_0$$

$$S_2 = A_2 \text{ xor } (A_1 \text{ and } A_0)$$

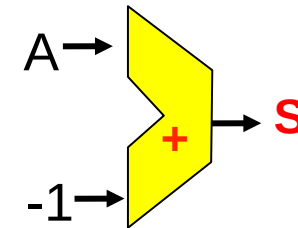


$$\text{Eq: } S_i = A_i \text{ xor } (A_{i-1} \text{ and } A_{i-2} \text{ and } \dots \text{ and } A_1 \text{ and } A_0); i=0..n-1$$

$$S_0 = \text{not } A_0$$

$$S_1 = A_1 \text{ xor } (\text{not } A_0)$$

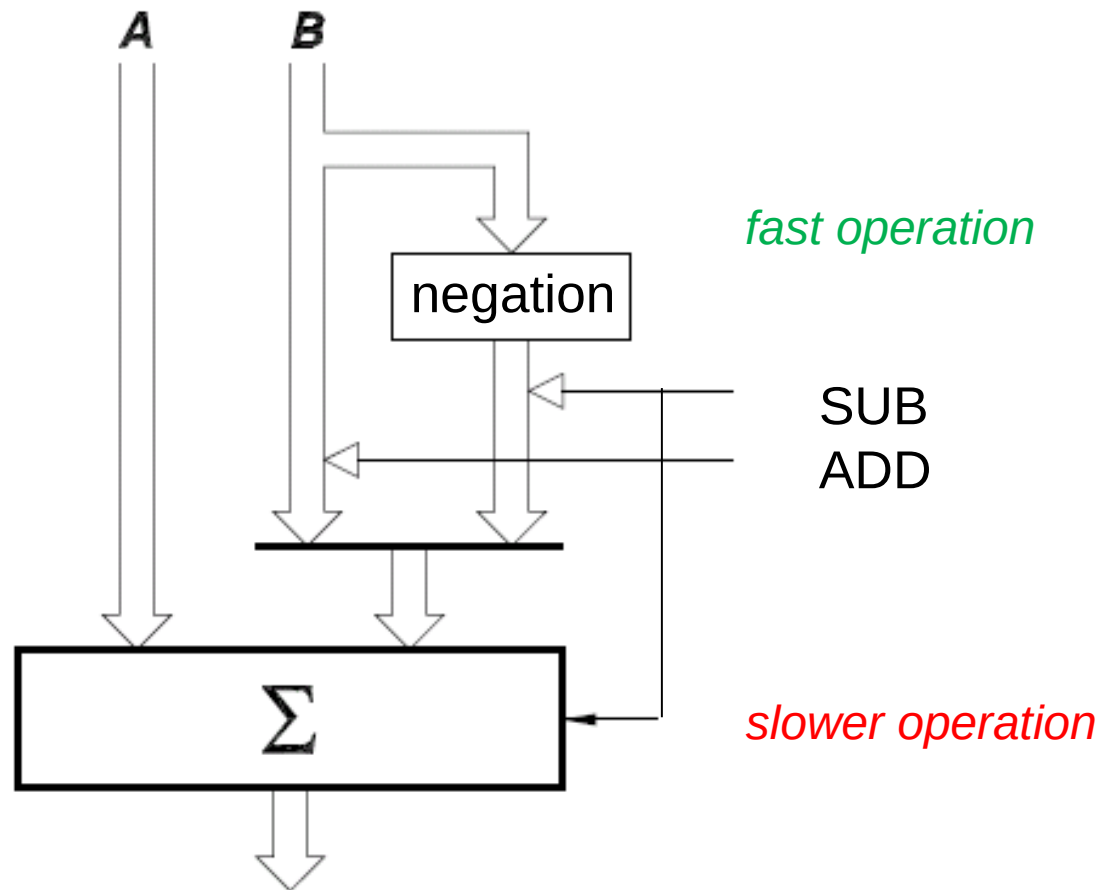
$$S_2 = A_2 \text{ xor } (\text{not } A_1 \text{ and } \text{not } A_0)$$



$$\text{Eq: } S_i = A_i \text{ xor } (\text{not } A_{i-1} \text{ and } \dots \text{ and } \text{not } A_0); i=0..n-1$$

The number of circuits is given by the arithmetic series, with the complexity  $O(n^2)$  where  $n$  is the number of bits. The operation can be performed in parallel for all bits, and for the both +1/-1 operations, we use a circuit that differs only by negations.

# Addition / Subtraction HW

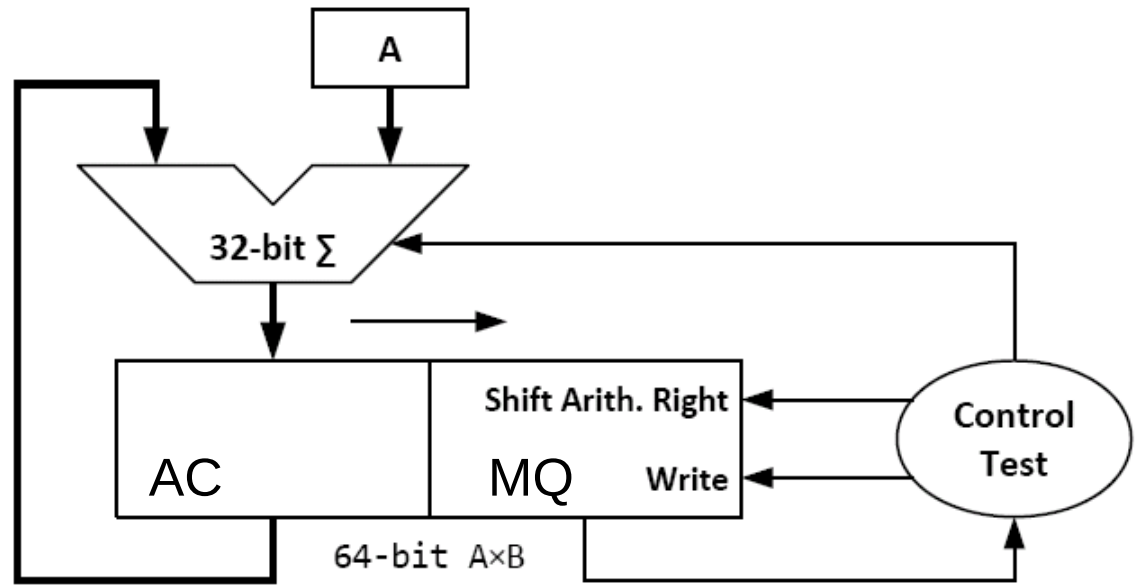


Source: X36JPO, A. Pluháček

## Unsigned binary numbers multiplication

$$\begin{array}{r}
 \text{A} \quad \quad \quad \text{B} \\
 1\ 1\ 0\ 1 \cdot 1\ 0\ 1\ 1 \\
 \hline
 \phantom{1\ 1\ 0\ 1} 0\ 0\ 0\ 0 \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \\
 \phantom{1\ 1\ 0\ 1} 1\ 1\ 0\ 1 \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \\
 \hline
 \phantom{1\ 1\ 0\ 1} 0\ 1\ 1\ 0\ 1 \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \\
 \phantom{1\ 1\ 0\ 1} 1\ 1\ 0\ 1 \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \\
 \hline
 \phantom{1\ 1\ 0\ 1} 1\ 0\ 0\ 1\ 1 \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \\
 \phantom{1\ 1\ 0\ 1} 0\ 0\ 0\ 0 \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \\
 \hline
 \phantom{1\ 1\ 0\ 1} 0\ 1\ 0\ 0\ 1 \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \\
 \phantom{1\ 1\ 0\ 1} 1\ 1\ 0\ 1 \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \\
 \hline
 1\ 0\ 0\ 0\ 1 \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \\
 \phantom{1\ 0\ 0\ 0\ 1} \downarrow \downarrow \downarrow \downarrow \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \\
 \phantom{1\ 0\ 0\ 0\ 1} 1\ 0\ 0\ 0 \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \\
 \phantom{1\ 0\ 0\ 0\ 1} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \phantom{0\ 0\ 0\ 0} \\
 \hline
 \underbrace{\phantom{1\ 0\ 0\ 0}}_{C1} \quad \underbrace{\phantom{0\ 0\ 0\ 0}}_{C0}
 \end{array}$$

# Sequential hardware multiplier (32b case)



The speed of the multiplier is horrible

## Algorithm for Multiplication

```
A = multiplicand;  
MQ = multiplier;  
AC = 0;
```

```
for( int i=1; i <= n; i++) // n – represents number of bits  
{  
if(MQ0 == 1) AC = AC + A; // MQ0 = LSB of MQ
```

```
SR (shift AC MQ by one bit right and insert information about  
carry from the MSB from previous step)  
}  
end.
```

when loop ends AC MQ holds 64-bit result



## Example of the multiply X by Y

Multiplicand  $x=110$  and multiplier  $y=101$ .

<b>i</b>	<b>operation</b>	<b>AC</b>	<b>MQ</b>	<b>A</b>	<b>comment</b>
		000	101	110	initial setup
1	AC = AC+MB	110	101		start of the cycle
	SR	011	010		
2	nothing	011	010		because of $MQ_0 = 0$
	SR	001	101		
3	AC = AC+MB	111	101		
	SR	011	110		end of the cycle

**The whole operation:  $x \times y = 110 \times 101 = 011110$ , (  $6 \times 5 = 30$  )**

# Multiplication in two's complement

Can be implemented, but there is a problem ...

**The intended product is generally not the same as the product of two's numbers!**

Details are already outside the intended APO range.

The best way is the multiplication of their absolute values and decision about its sign.

## Wallace tree based multiplier

$Q = X \cdot Y$ ,  $X$  and  $Y$  are considered as 8bit unsigned numbers

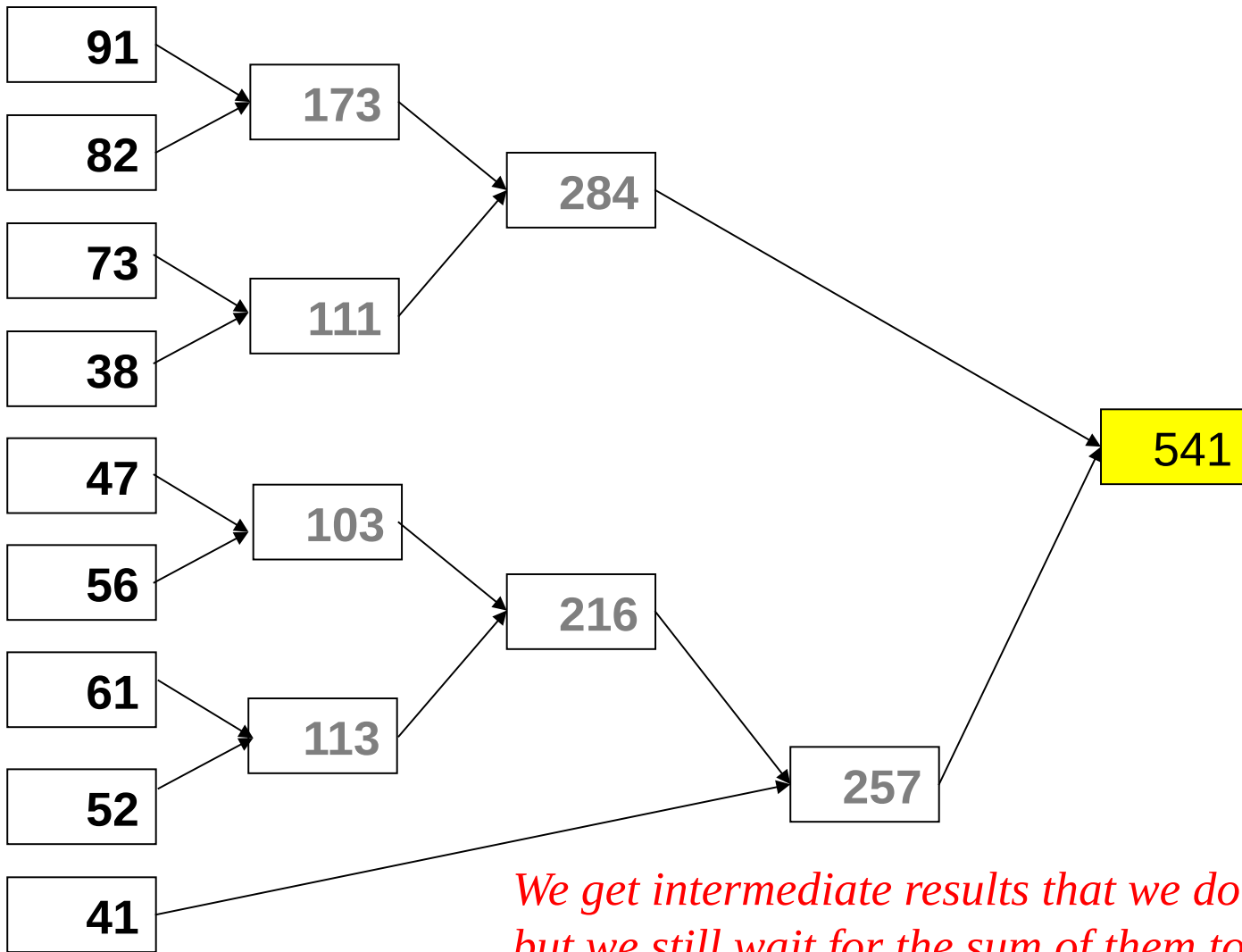
$$(x_7 x_6 x_5 x_4 x_3 x_2 x_1 x_0) \cdot (y_7 y_6 y_5 y_4 y_3 y_2 y_1 y_0) =$$

0	0	0	0	0	0	0	0	$x_7y_0$	$x_6y_0$	$x_5y_0$	$x_4y_0$	$x_3y_0$	$x_2y_0$	$x_1y_0$	$x_0y_0$	P0
0	0	0	0	0	0	0	$x_7y_1$	$x_6y_1$	$x_5y_1$	$x_4y_1$	$x_3y_1$	$x_2y_1$	$x_1y_1$	$x_0y_1$	0	P1
0	0	0	0	0	0	$x_7y_2$	$x_6y_2$	$x_5y_2$	$x_4y_2$	$x_3y_2$	$x_2y_2$	$x_1y_2$	$x_0y_2$	0	0	P2
0	0	0	0	0	$x_7y_3$	$x_6y_3$	$x_5y_3$	$x_4y_3$	$x_3y_3$	$x_2y_3$	$x_1y_3$	$x_0y_3$	0	0	0	P3
0	0	0	0	$x_7y_4$	$x_6y_4$	$x_5y_4$	$x_4y_4$	$x_3y_4$	$x_2y_4$	$x_1y_4$	$x_0y_4$	0	0	0	0	P4
0	0	0	$x_7y_5$	$x_6y_5$	$x_5y_5$	$x_4y_5$	$x_3y_5$	$x_2y_5$	$x_1y_5$	$x_0y_5$	0	0	0	0	0	P5
0	0	$x_7y_6$	$x_6y_6$	$x_5y_6$	$x_4y_6$	$x_3y_6$	$x_2y_6$	$x_1y_6$	$x_0y_6$	0	0	0	0	0	0	P6
0	$x_7y_7$	$x_6y_7$	$x_5y_7$	$x_4y_7$	$x_3y_7$	$x_2y_7$	$x_1y_7$	$x_0y_7$	0	0	0	0	0	0	0	P7
$Q_{15}$	$Q_{14}$	$Q_{13}$	$Q_{12}$	$Q_{11}$	$Q_{10}$	$Q_9$	$Q_8$	$Q_7$	$Q_6$	$Q_5$	$Q_4$	$Q_3$	$Q_2$	$Q_1$	$Q_0$	

The sum of  $P0 + P1 + \dots + P7$  gives result of  $X$  and  $Y$  multiplication.

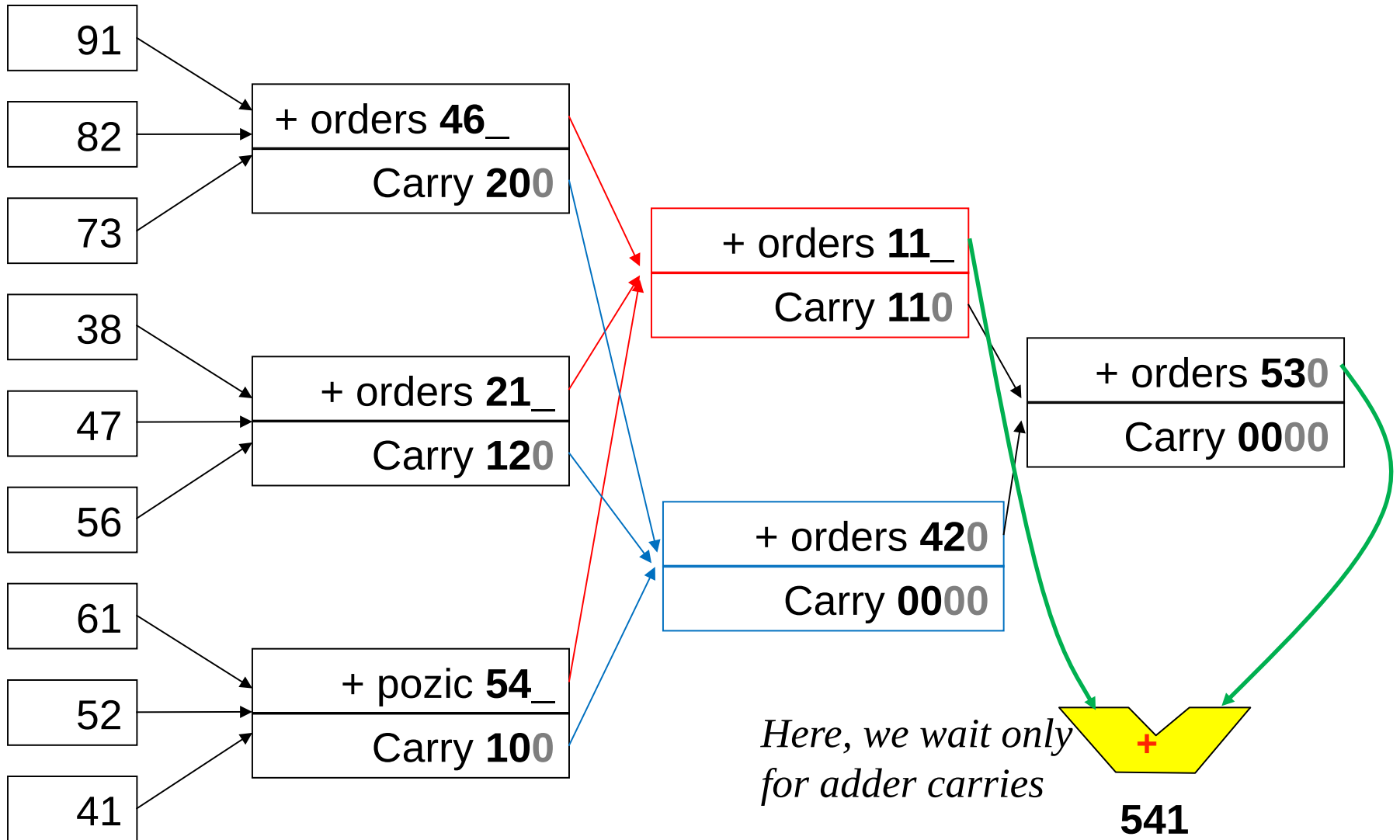
$$Q = X \cdot Y = P0 + P1 + \dots + P7$$

# Parallel adder of 9 numbers



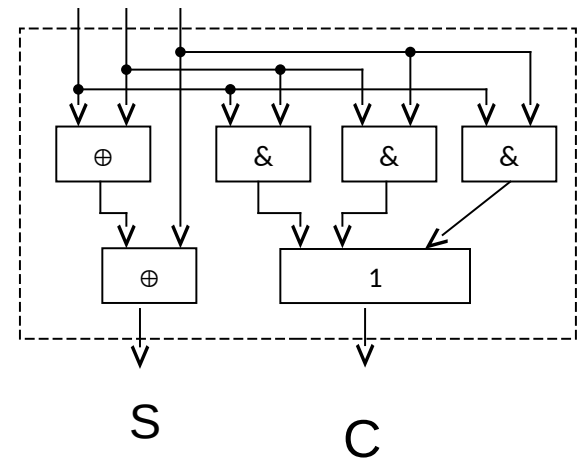
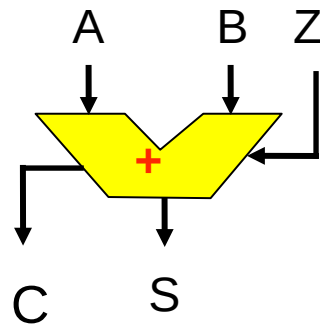
*We get intermediate results that we do not need at all, but we still wait for the sum of them to finish!*

# Decadic Carry-save adder



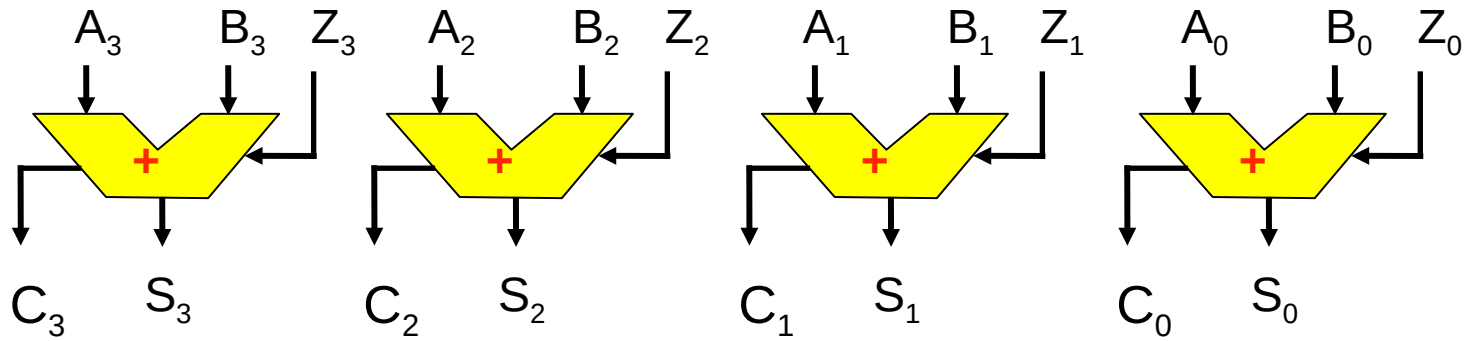
# 1bit Carry Save Adder

A	0	0	1	1	0	0	1	1
+B	0	1	0	1	0	1	0	1
Z=Carry-In	0	0	0	0	1	1	1	1
Sum	0	1	1	0	1	0	0	1
C=Cout	0	0	0	1	0	1	1	1



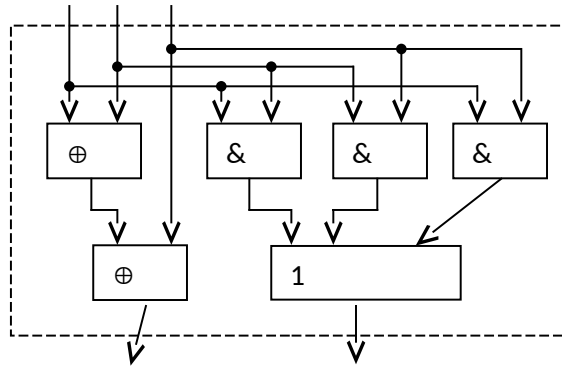


# 3-bit Carry-save adder



# Wallace tree based fast multiplier

The basic element is an CSA circuit (Carry Save Adder)

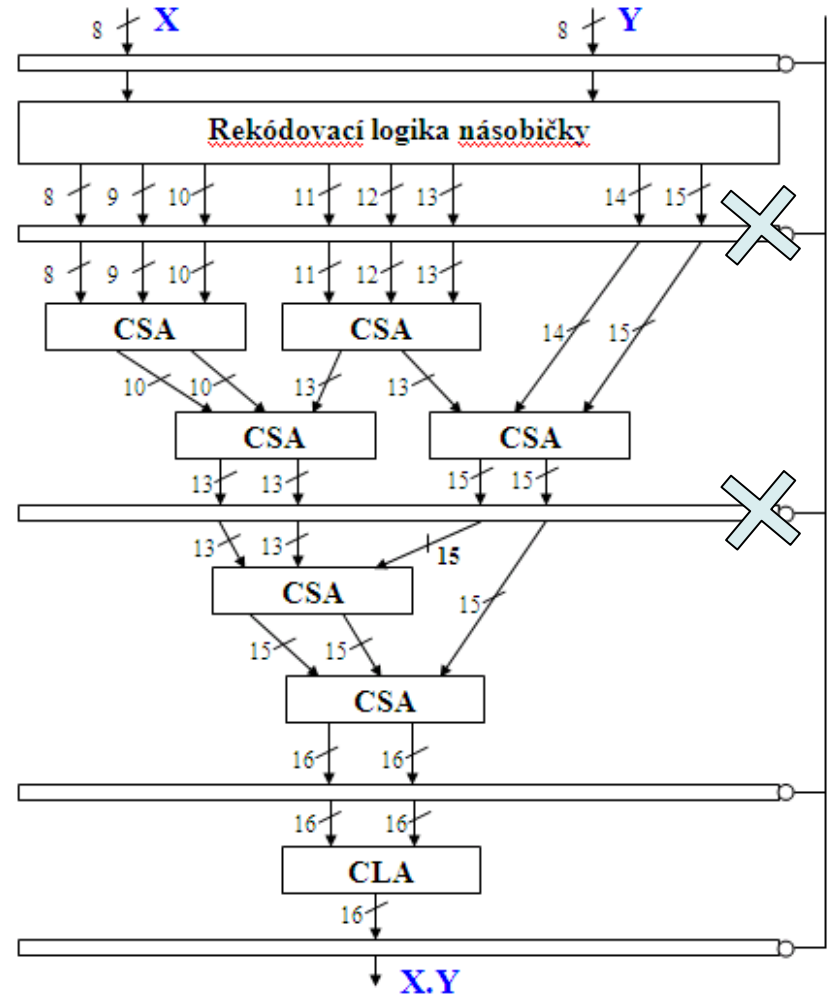


$$S = S^b + C$$

$$S^b_i = x_i \oplus y_i \oplus z_i$$

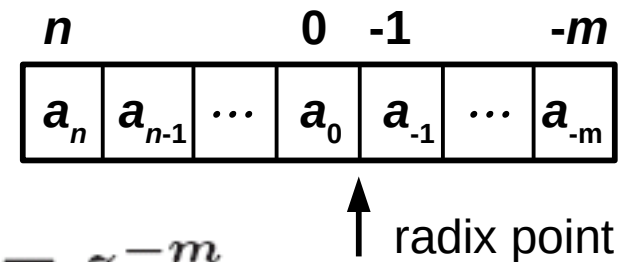
$$C_{i+1} = x_i y_i + y_i z_i +$$

$$z_i x_i$$



# Terminology basics

- Positional (place-value) notation
- Decimal/radix point
- $z$  ... base of numeral system
- smallest representable number  $\epsilon = z^{-m}$
- **Module** =  $Z$  , one increment/unit higher than biggest representable number for given encoding/notation
- **A**, the representable number for given  $n$  and  $m$  selection, where  $k$  is natural number in range  $\langle 0, z^{n+m+1} - 1 \rangle$
- The representation and value



$$0 \leq A = k \cdot \epsilon < Z$$

$$A \sim a_n a_{n-1} \dots a_0, a_1 \dots a_{-m}$$

$$A = a_n z^n + a_{n-1} z^{n-1} + \dots + a_0 + a_1 z^{-1} \dots a_{-m} z^{-m}$$

## Integer number representation (unsigned, non-negative)

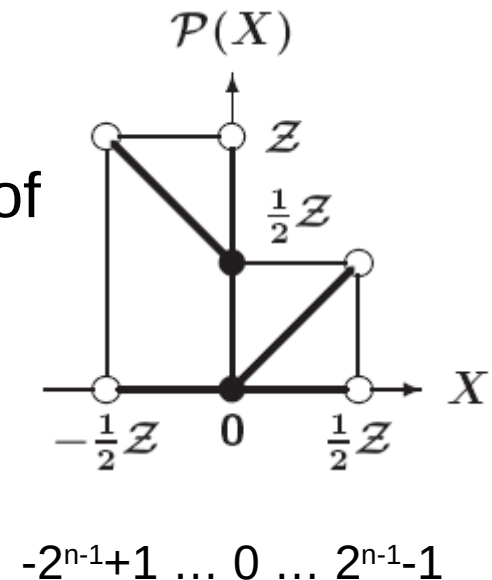
- The most common numeral system base in computers is  $z=2$
- The value of  $a_i$  is in range  $\{0,1,\dots,z-1\}$ , i.e.  $\{0,1\}$  for base 2
- This maps to true/false and unit of information (bit)
- We can represent number  $0 \dots 2^n-1$  when  $n$  bits are used
- Which range can be represented by one byte?
  - 1B (byte) ... 8 bits,  $2^8 = 256_d$  combinations, values  $0 \dots 255_d = 0b11111111_b$
- Use of multiple consecutive bytes
  - 2B ...  $2^{16} = 65536_d$ ,  $0 \dots 65535_d = 0xFFFF_h$ , (h ... hexadecimal, base 16, a in range 0, ... 9, A, B, C, D, E, F)
  - 4B ...  $2^{32} = 4294967296_d$ ,  $0 \dots 4294967295_d = 0xFFFFFFFF_h$

## Signed integer numbers

- Work with negative numbers is required for many applications
- When appropriate representation is used then same hardware (with minor extension) can be used for many operations with signed and unsigned numbers
- Possible representations
  - sign-magnitude code, direct representation, sign bit
  - two's complement
  - ones' complement
  - excess-K, offset binary or biased representation

## Integer – sign-magnitude code

- Sign and magnitude of the value (absolute value)
- Natural to humans -1234, 1234
- One (usually most significant – MSB) bit of the memory location is used to represent the sign
- Bit has to be mapped to meaning
- Common use  $0 \approx "+"$ ,  $1 \approx "-"$
- Disadvantages:
  - When location is **k** bits long then only **k-1** bits hold magnitude and each operation has to separate sign and magnitude
  - Two representations of the value 0



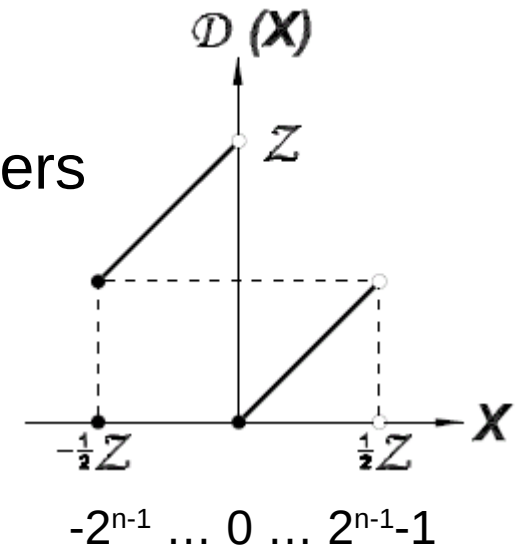
## Integer – two's complement

- Other option is to designate one half of range/combinations for non-negative numbers and other one for positive numbers
- Transform to the representation

$$D(A) = A \quad \text{iff } A \geq 0$$

$$D(A) = Z - |A| \quad \text{iff } A < 0$$

- Advantages
  - Continuous range when cyclic arithmetics is considered
  - Single and one to one mapping of value 0
  - Same HW for signed and unsigned adder
- Disadvantage
  - Asymmetric range  $(-(-1/2Z))$





## Integers – ones' complement

- Transform to the representation  $-2^{n-1}+1 \dots 0 \dots 2^{n-1}-1$ 
  - $D(A) = A$       iff  $A \geq 0$
  - $D(A) = Z-1-|A|$       iff  $A < 0$  (i.e. subtract from all ones)
- Advantages
  - Symmetric range
  - Almost continuous, requires hot one addition when sign changes
- Disadvantage
  - Two representations of value 0
  - More complex hardware
- Negate ( $-A$ ) value can be computed by bitwise complement (flipping) of each bit in representation

## Integer – biased representation

- Known as excess-K or offset binary as well
- Transform to the representation  $-K \dots 0 \dots 2^n-1-K$   
 $D(A) = A+K$
- Usually  $K=Z/2$
- Advantages
  - Preserves order of original set in mapped set/representation
- Disadvantages
  - Needs adjustment by  $-K$  after addition and  $+K$  after subtraction processed by unsigned arithmetic unit
  - Requires full transformation before and after multiplication

## Back to two's complement and the C language

- Two's complement is most used signed integer numbers representation in computers
- Complement arithmetic is often used as its synonym
- “C” programming language speaks about integer numeric type without sign as *unsigned integers* and they are declared in source code as `unsigned int`.
- The numeric type with sign is simply called *integers* and is declared as `signed int`.
- Examples of the values representations when 32 bits are used:
  - $0_D = 00000000_H$ ,
  - $1_D = 00000001_H$ ,  $-1_D = FFFFFFFF_H$ ,
  - $2_D = 00000002_H$ ,  $-2_D = FFFFFFFE_H$ ,
  - $3_D = 00000003_H$ ,  $-3_D = FFFFFFFD_H$ ,
- Considerations about value overflow and underflow from order grit are discussed later.

## Two's complement – addition and subtraction

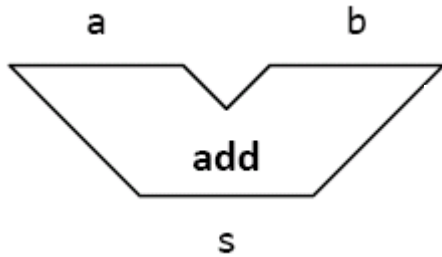
- **Addition**

- $00000000\ 0000\ 0111_B \approx 7_D$     Symbols use:  $0=0_H$ ,  $0=0_B$
- $+ \underline{00000000\ 0000\ 0110_B} \approx 6_D$
- $00000000\ 0000\ 1101_B \approx 13_D$

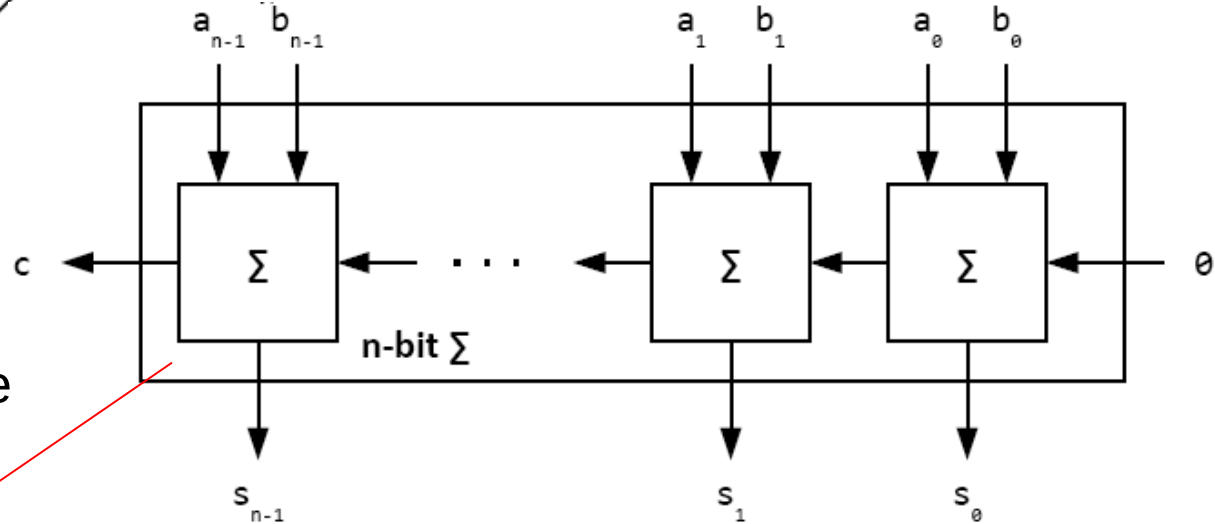
- **Subtraction** can be realized as addition of negated number

- $00000000\ 0000\ 0111_B \approx 7_D$
  - $+ \underline{FFFFFFF\ 1111\ 1010_B} \approx -6_D$
  - $00000000\ 0000\ 0001_B \approx 1_D$
- Question for revision: how to obtain negated number in two's complement binary arithmetics?

# Hardware of ripple-carry adder

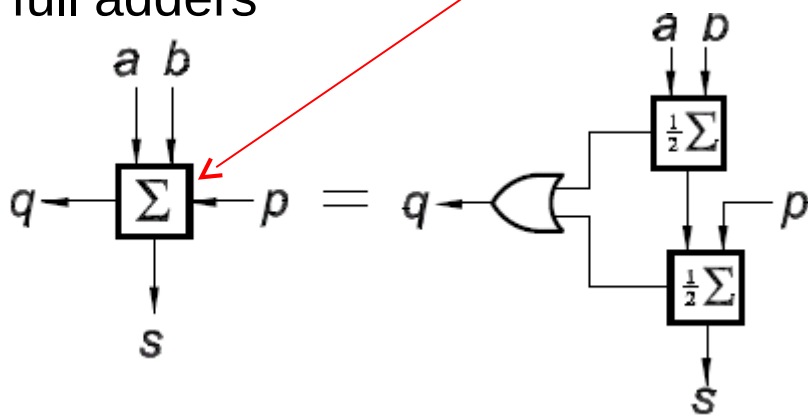


Common symbol for adder

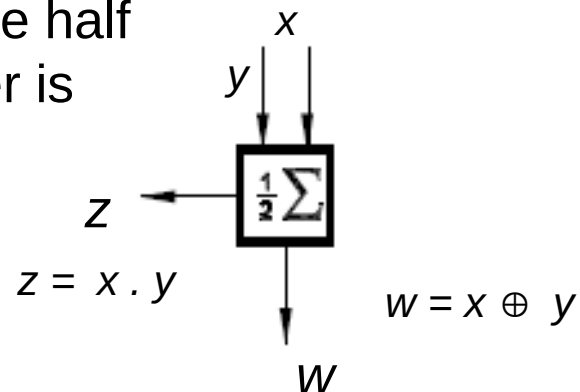


Internal structure

Realized by 1-bit full adders



where half adder is



## Fast parallel adder realization and limits

- The previous, cascade based adder is slow – carry propagation delay
- The parallel adder is combinatorial circuit, it can be realized through sum of minterms (product of sums), two levels of gates (wide number of inputs required)
- But for 64-bit adder  $10^{20}$  gates is required

### Solution #1

- Use of carry-lookahead circuits in adder combined with adders without carry bit

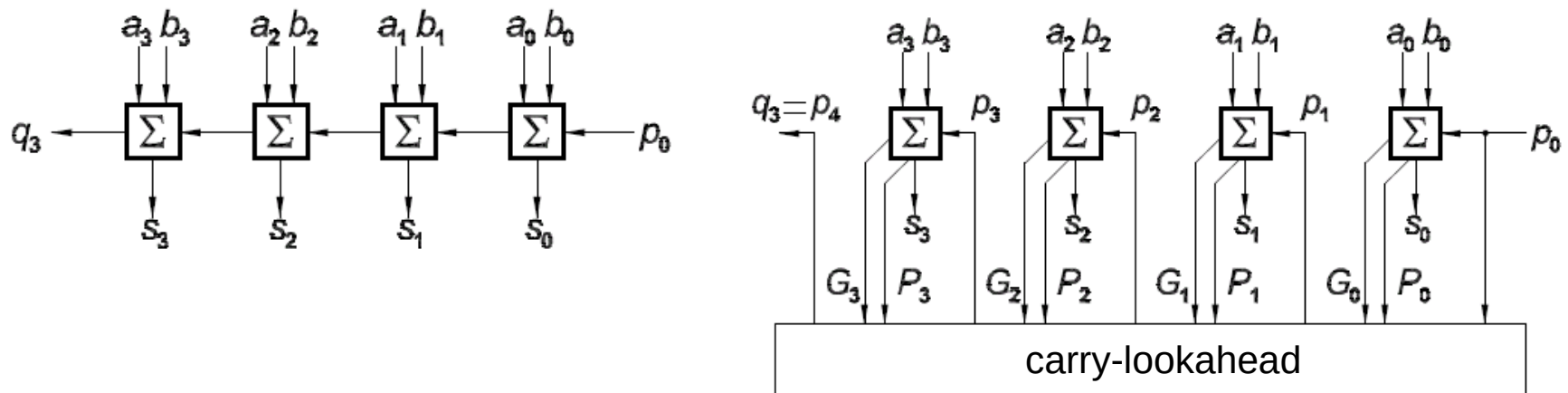
### Solution #2

- Cascade of adders with fraction of the required width

Combination (hierarchy) of #1 and #2 can be used for wider inputs

## Speed of the adder

- Parallel adder is combinational logic/circuit. Is there any reason to speak about its speed? Try to describe!
- Yes, and it is really slow. Why?
- Possible enhancement – adder with carry-lookahead (CLA) logic!



## CLA – carry-lookahead

- Adder combined with CLA provides enough speedup when compared with parallel ripple-carry adder and yet number of additional gates is acceptable
- CLA for 64-bit adder increases hardware price for about 50% but the speed is increased (signal propagation time decreased) 9 times.
- The result is significant speed/price ratio enhancement.



## The basic equations for the CLA logic

- Let:
  - the generation of carry on position (bit) j is defined as:

$$g_j = x_j y_j$$

- the need for carry propagation from previous bit:

$$p_j = x_j \oplus y_j = x_j \bar{y}_j \vee \bar{x}_j y_j$$

- Then:
  - the result of sum for bit j is given by:

$$s_j = c_j (\overline{x_j \oplus y_j}) \vee \bar{c}_j (x_j \oplus y_j) = c_j \bar{p}_j \vee \bar{c}_j p_j = p_j \oplus c_j$$

- and carry to the higher order bit (j+1) is given by:

$$c_{j+1} = x_j y_j \vee (x_j \oplus y_j) c_j = g_j \vee p_j c_j$$

# CLA

The carry can be computed as:

$$c_1 = g_0 \vee p_0 c_0$$

$$c_2 = g_1 \vee p_1 c_1 = g_1 \vee p_1 (g_0 \vee p_0 c_0) = g_1 \vee p_1 g_0 \vee p_1 p_0 c_0$$

$$c_3 = g_2 \vee p_2 c_2 = g_2 \vee p_2 (g_1 \vee p_1 g_0 \vee p_1 p_0 c_0) = g_2 \vee p_2 g_1 \vee p_2 p_1 g_0 \vee p_2 p_1 p_0 c_0$$

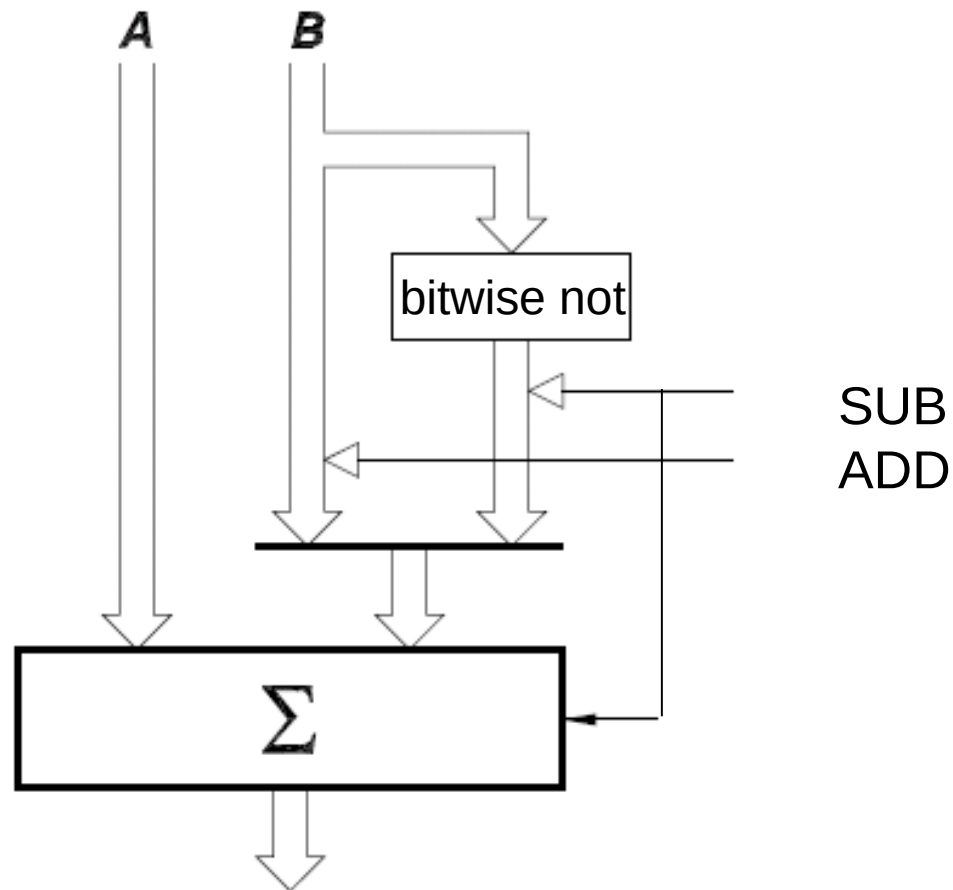
$$c_4 = g_3 \vee p_3 c_3 = \dots = g_3 \vee p_3 g_2 \vee p_3 p_2 g_1 \vee p_3 p_2 p_1 g_0 \vee p_3 p_2 p_1 p_0 c_0$$

$$c_5 = \dots$$

Description of the equation for  $c_3$  as an example:

The carry input for bit 3 is active **when** carry is generated in bit 2 **or** carry propagates condition holds for bit 2 and carry is generated in the bit 1 **or** both bits 2 and 1 propagate carry and carry is generated in bit 0

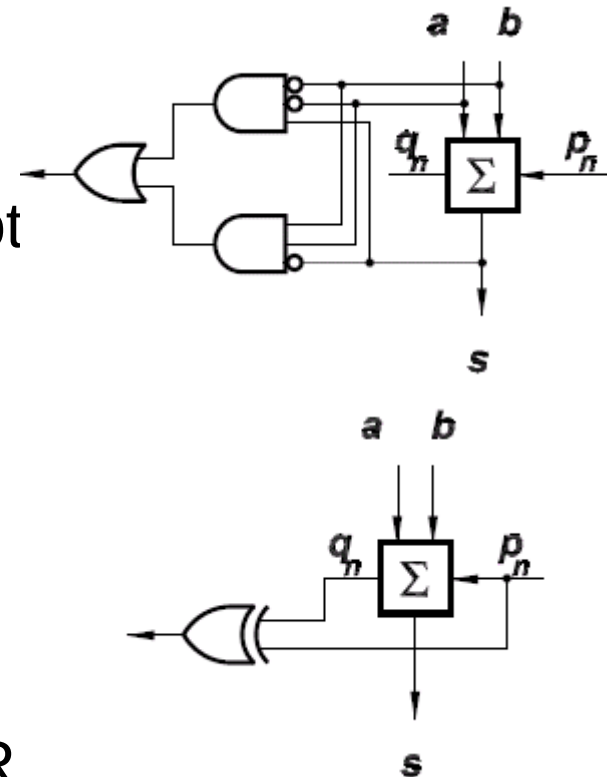
# Arithmetic unit for add/subtract operations



Inspiration: X36JPO, A. Pluháček

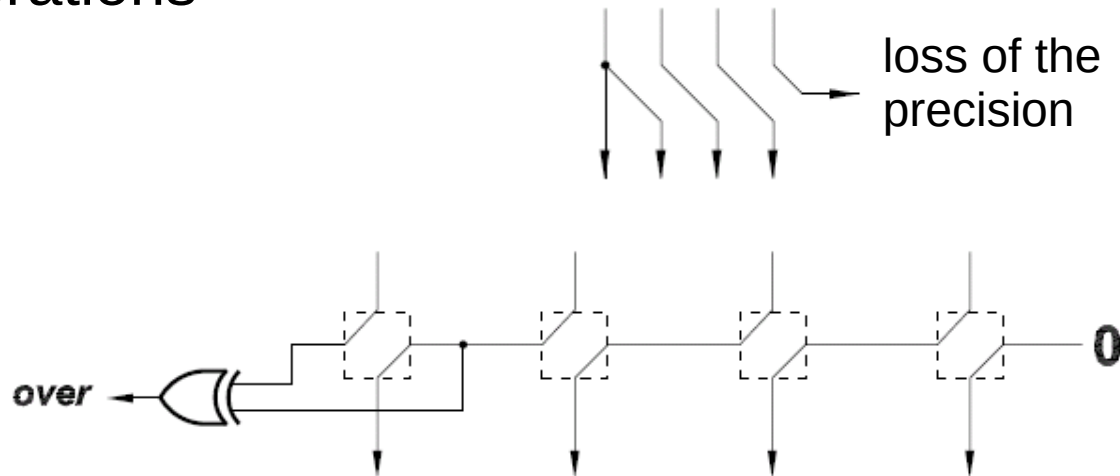
## Arithmetic overflow (underflow)

- Result of the arithmetic operation is incorrect because, it does not fit into selected number of the representation bits (width)
- But for the signed arithmetics, it is not equivalent to the carry from the most significant bit.
- The arithmetic overflow is signaled if result sign is different from operand signs if both operands have same sign
- or can be detected with exclusive-OR of carry to and from the most significant bit



## Arithmetic shift to the left and to the right

- arithmetic shift by one to the left/right is equivalent to signed multiply/divide by 2 (digits movement in positional (place-value) representation)
- Notice difference between arithmetic, logic and cyclic shift operations



- Remark: Barrel shifter can be used for fast variable shifts

## Addition and subtraction for the biased representation

- Short note about other signed number representation

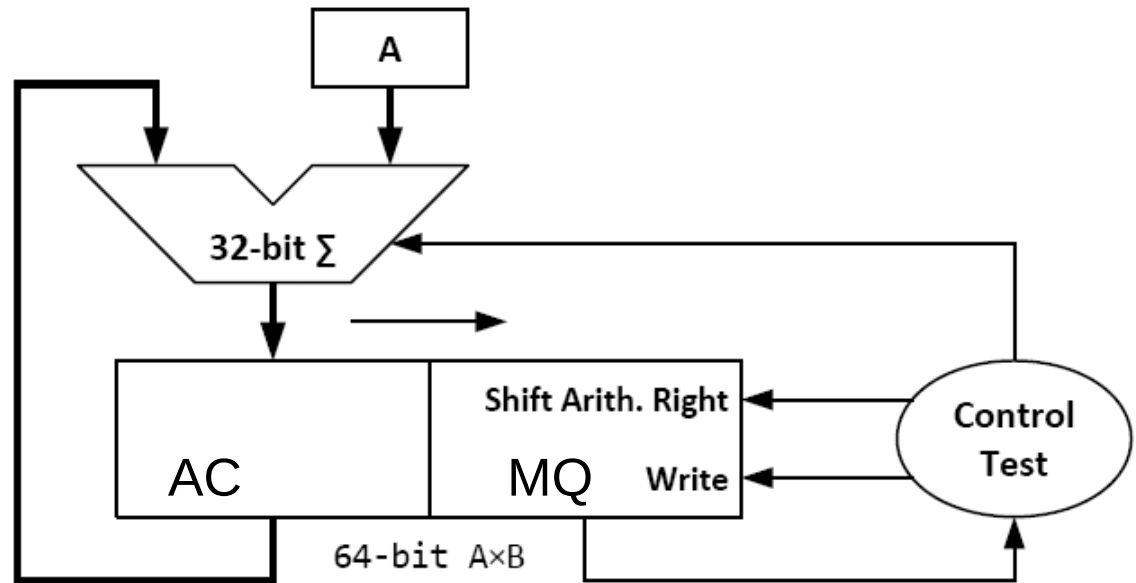
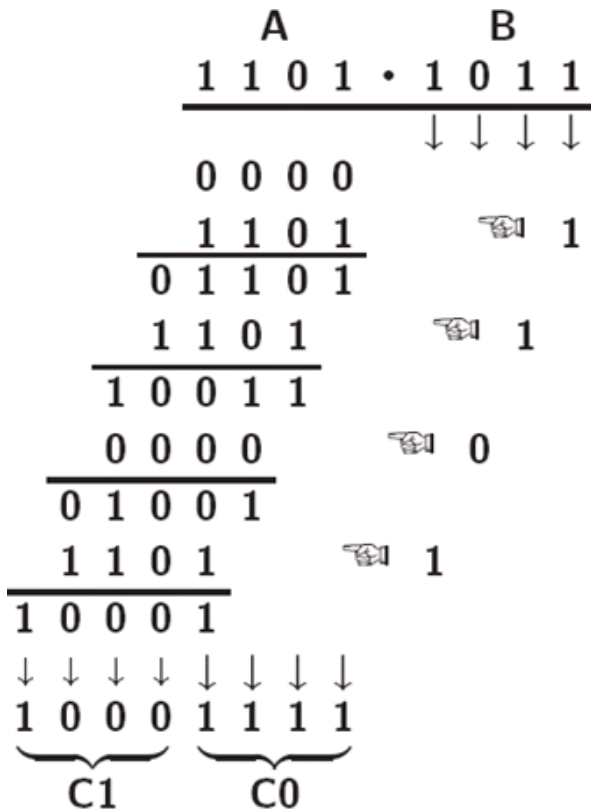
$$\mathcal{A}(A + B) = \mathcal{A}(A) + \mathcal{A}(B) - K$$
$$\mathcal{A}(A - B) = \mathcal{A}(A) - \mathcal{A}(B) + K$$

- Overflow detection
  - for addition:  
same sign of addends and different result sign
  - for subtraction:  
signs of minuend and subtrahend are opposite and sign of the result is opposite to the sign of minuend

# Unsigned binary numbers multiplication

$$\begin{array}{r}
 \begin{array}{cccc}
 & \text{A} & & \text{B} \\
 & 1 & 1 & 0 & 1 & \cdot & 1 & 0 & 1 & 1 \\
 \hline
 & & & & & & \downarrow & \downarrow & \downarrow & \downarrow \\
 & & & 0 & 0 & 0 & 0 & & & \\
 & & & 1 & 1 & 0 & 1 & & \rightarrow & 1 \\
 \hline
 & & & 0 & 1 & 1 & 0 & 1 & & \\
 & & & 1 & 1 & 0 & 1 & & \rightarrow & 1 \\
 \hline
 & & & 1 & 0 & 0 & 1 & 1 & & \\
 & & & 0 & 0 & 0 & 0 & & \rightarrow & 0 \\
 \hline
 & & & 0 & 1 & 0 & 0 & 1 & & \\
 & & & 1 & 1 & 0 & 1 & & \rightarrow & 1 \\
 \hline
 & & & 1 & 0 & 0 & 0 & 1 & & \\
 \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & & \\
 \underbrace{1\ 0\ 0\ 0\ 0}_{\text{C1}} & \underbrace{1\ 1\ 1\ 1}_{\text{C0}}
 \end{array}
 \end{array}$$

# Sequential hardware multiplier (32b case)



The speed of the multiplier is horrible



# Algorithm for multiplication

A = multiplicand;

MQ = multiplier;

AC = 0;

for( int i=1; i <= n; i++) // n – represents number of bits

{

if(MQ<sub>0</sub> == 1) AC = AC + A; // MQ<sub>0</sub> = LSB of MQ

SR (shift AC MQ by one bit right and insert information about carry from the MSB from previous step)

}

end.

when loop ends AC MQ holds 64-bit result

## Example of the multiply X by Y

Multiplicand  $x=110$  and multiplier  $y=101$ .

<b>i</b>	<b>operation</b>	<b>AC</b>	<b>MQ</b>	<b>A</b>	<b>comment</b>
		000	101	110	initial setup
1	AC = AC+MB	110	101		start of the cycle
	SR	011	010		
2	nothing	011	010		because of $MQ_0 = 0$
	SR	001	101		
3	AC = AC+MB	111	101		
	SR	011	110		end of the cycle

**The whole operation:  $x \times y = 110 \times 101 = 011110$ , (  $6 \times 5 = 30$  )**

# Signed multiplication by unsigned HW for two's complement

One possible solution

$$C = A \cdot B$$

Let A and B representations are n bits and result is 2n bits

$$D(C) = D(A) \cdot D(B)$$

$$- (D(B) \ll n) \quad \text{if } A < 0$$

$$- (D(A) \ll n) \quad \text{if } B < 0$$

Consider for negative numbers

$$(2^{2n} + A) \cdot (2^{2n} + B) = 2^{2n} + 2^n A + 2^n B + A \cdot B$$

where  $2^{2n}$  is out of the result representation, next two elements have to be eliminated if input is negative

# Wallace tree based multiplier

$Q=X.Y$ ,  $X$  and  $Y$  are considered as 8bit unsigned numbers

$$(x_7 x_6 x_5 x_4 x_3 x_2 x_1 x_0) \cdot (y_7 y_6 y_5 y_4 y_3 y_2 y_1 y_0) =$$

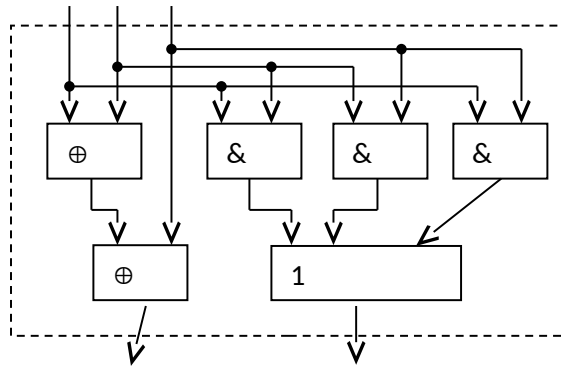
0	0	0	0	0	0	0	0	$x_7y_0$	$x_6y_0$	$x_5y_0$	$x_4y_0$	$x_3y_0$	$x_2y_0$	$x_1y_0$	$x_0y_0$	P0
0	0	0	0	0	0	0	$x_7y_1$	$x_6y_1$	$x_5y_1$	$x_4y_1$	$x_3y_1$	$x_2y_1$	$x_1y_1$	$x_0y_1$	0	P1
0	0	0	0	0	0	$x_7y_2$	$x_6y_2$	$x_5y_2$	$x_4y_2$	$x_3y_2$	$x_2y_2$	$x_1y_2$	$x_0y_2$	0	0	P2
0	0	0	0	0	$x_7y_3$	$x_6y_3$	$x_5y_3$	$x_4y_3$	$x_3y_3$	$x_2y_3$	$x_1y_3$	$x_0y_3$	0	0	0	P3
0	0	0	0	$x_7y_4$	$x_6y_4$	$x_5y_4$	$x_4y_4$	$x_3y_4$	$x_2y_4$	$x_1y_4$	$x_0y_4$	0	0	0	0	P4
0	0	0	$x_7y_5$	$x_6y_5$	$x_5y_5$	$x_4y_5$	$x_3y_5$	$x_2y_5$	$x_1y_5$	$x_0y_5$	0	0	0	0	0	P5
0	0	$x_7y_6$	$x_6y_6$	$x_5y_6$	$x_4y_6$	$x_3y_6$	$x_2y_6$	$x_1y_6$	$x_0y_6$	0	0	0	0	0	0	P6
0	$x_7y_7$	$x_6y_7$	$x_5y_7$	$x_4y_7$	$x_3y_7$	$x_2y_7$	$x_1y_7$	$x_0y_7$	0	0	0	0	0	0	0	P7
$Q_{15}$	$Q_{14}$	$Q_{13}$	$Q_{12}$	$Q_{11}$	$Q_{10}$	$Q_9$	$Q_8$	$Q_7$	$Q_6$	$Q_5$	$Q_4$	$Q_3$	$Q_2$	$Q_1$	$Q_0$	

The sum of  $P0+P1+\dots+P7$  gives result of  $X$  and  $Y$  multiplication.

$$Q = X.Y = P0 + P1 + \dots + P7$$

# Wallace tree based fast multiplier

The basic element is an CSA circuit (Carry Save Adder)

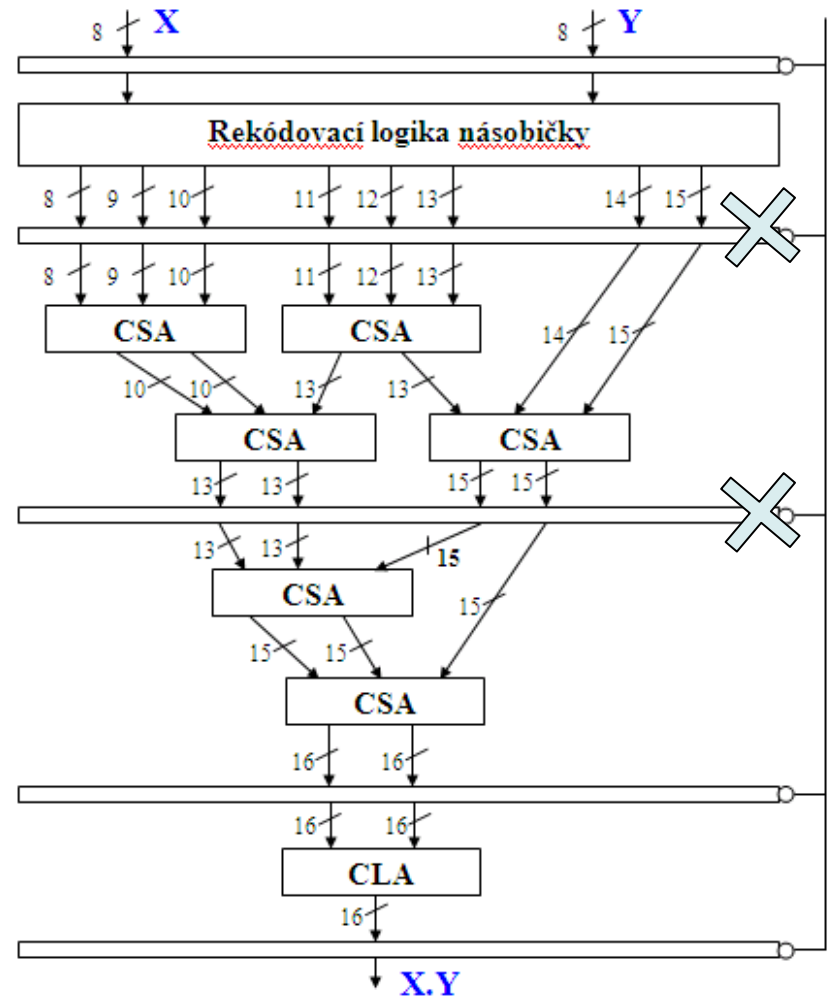


$$S = S^b + C$$

$$S^b_i = x_i \oplus y_i \oplus z_i$$

$$C_{i+1} = x_i y_i + y_i z_i +$$

$$z_i x_i$$



# Hardware divider

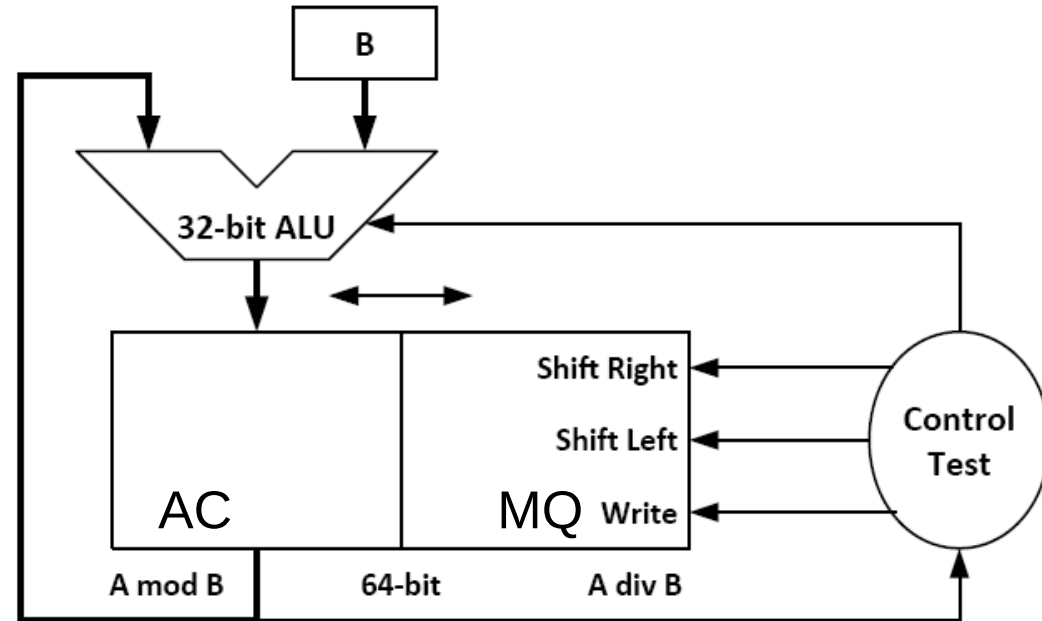
$$\boxed{111 : 011}$$

	0 0 0 1 1 1	:	0 0 1 1	
⊖	1 1 0 0	:		negate hot one
	1	:		
	0 1 1 1 0	:		- ⇒ 0
	↓ ↓ ↓ ↓	:		
	1 1 0 1	:		
⊕	0 0 1 1	:		
	1 0 0 0 0 1	:		+ ⇒ 1
	↓ ↓ ↓ ↓	:		
	0 0 0 1	:		
⊖	1 1 0 0	:		
	1	:		
	0 1 1 1 0	:		- ⇒ 0
⊖	0 0 1 1	:		return
	1 0 0 0 1	:		
	0 0 1	—	reminder	0 1 0 — quotient

# Hardware divider logic (32b case)

$111 : 011$  dividend = quotient  $\times$  divisor + remainder

⊖	0 0 0 1 1 1	:	0 0 1 1	
	1 1 0 0	:	:	negate
	1	:	:	hot one
	0 1 1 1 0	:	:	- $\Rightarrow$ 0
	↓ ↓ ↓ ↓	:	:	
	1 1 0 1	:	:	
⊕	0 0 1 1	:	:	
	1 0 0 0 0 1	:	:	+ $\Rightarrow$ 1
	↓ ↓ ↓ ↓	:	:	
	0 0 0 1	:	:	
⊖	1 1 0 0	:	:	
	1	:	:	
	0 1 1 1 0	:	:	- $\Rightarrow$ 0
	↓ ↓ ↓ ↓	:	:	
	0 0 1 1	:	:	return
⊖	1 0 0 0 1	:	:	
	0 0 1	-	-	remainder
	0 1 0	-	-	quotient



## Algorithm of the sequential division

MQ = dividend;

B = divisor; (Condition: divisor is not 0!)

AC = 0;

```
for( int i=1; i <= n; i++)    {
    SL (shift AC MQ by one bit to the left, the LSB bit is kept on zero)
    if(AC >= B) {
        AC = AC - B;
        MQ0 = 1;    // the LSB of the MQ register is set to 1
    }
}
```

→ Value of MQ register represents quotient and AC remainder



## Example of X/Y division

Dividend  $x=1010$  and divisor  $y=0011$

i	operation	AC	MQ	B	comment
		0000	1010	0011	initial setup
1	SL	0001	0100		
	nothing	0001	0100		the if condition not true
2	SL	0010	1000		
		0010	1000		the if condition not true
3	SL	0101	0000		$r \geq y$
	<b>AC = AC - B; MQ<sub>0</sub> = 1;</b>	0010	0001		
4	SL	0100	0010		$r \geq y$
	<b>AC = AC - B; MQ<sub>0</sub> = 1;</b>	0001	0011		end of the cycle

**$x : y = 1010 : 0011 = 0011$  reminder 0001, ( $10 : 3 = 3$  reminder 1)**

## Higher dynamic range for numbers (REAL/float)

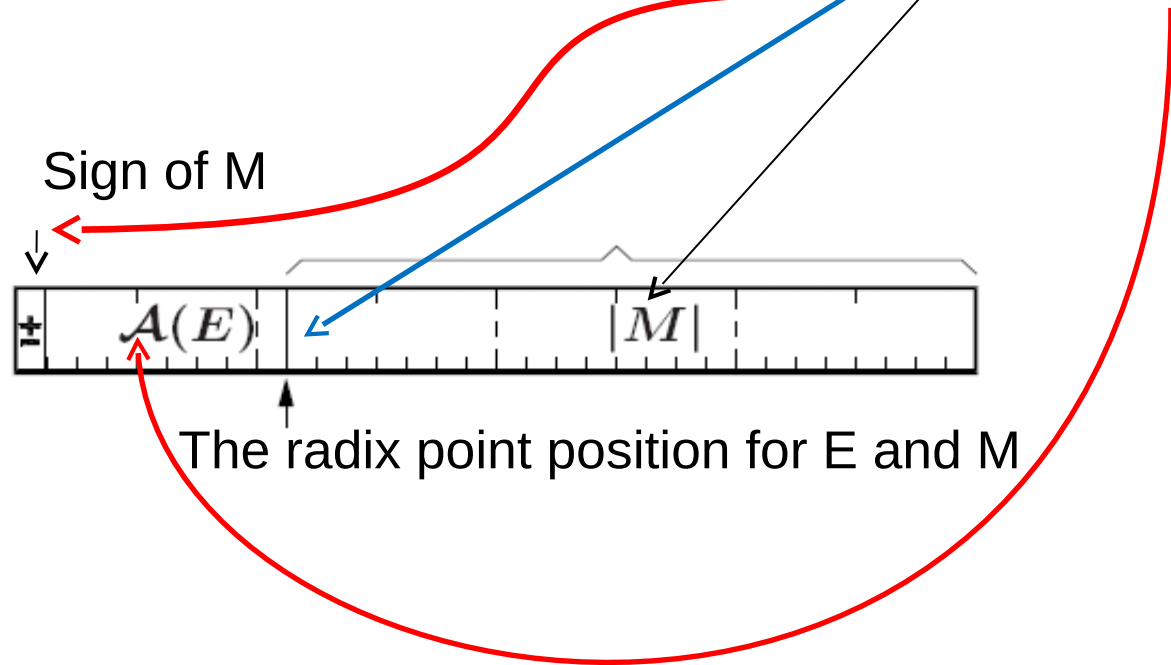
- Scientific notation, semilogarithmic, floating point
  - The value is represented by:
    - EXPONENT (E) – represents scale for given value
    - MANTISSA (M) – represents value in that scale
    - the sign(s) are usually separated as well
- Normalized notation
  - The exponent and mantissa are adjusted such way, that mantissa is held in some standard range.  $\langle 0.5, 1 \rangle$  or  $\langle 1, 2 \rangle$  for considered base  $z=2$
  - Generally: the first digit is non-zero or mantissa range is  $\langle 1, z \rangle$

## Standardized format for REAL type numbers

- Standard IEEE-754 defines next REAL representation and precision
  - single-precision – in the C language declared as `float`
  - double-precision – C language `double`

# Examples of (de)normalized numbers in base 10 and 2

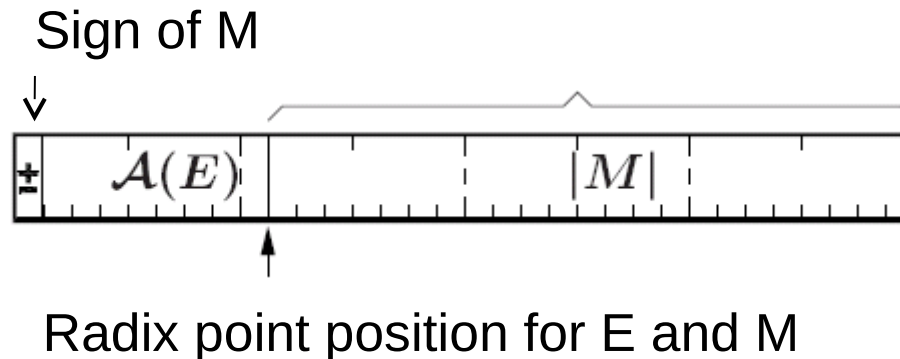
- $-2.34 \times 10^{56}$  ← normalized
  - $+0.002 \times 10^{-4}$  ← not normalized
  - $+987.02 \times 10^9$  ← not normalized
- binary
- $\pm 1.xxxxxxx_2 \times 2^{yyyy}$



# The representation/encoding of floating point number

- Mantissa encoded as the sign and absolute value (magnitude) – equivalent to the direct representation
- Exponent encoded in biased representation ( $K=127$  for single precision)
- The implicit leading one can be omitted due to normalization of  $m \in \langle 1, 2 \rangle$  – 23+1 implicit bit for single

$$X = -1^s 2^{A(E)-127} m \quad \text{where } m \in \langle 1, 2 \rangle$$
$$m = 1 + 2^{-23} M$$

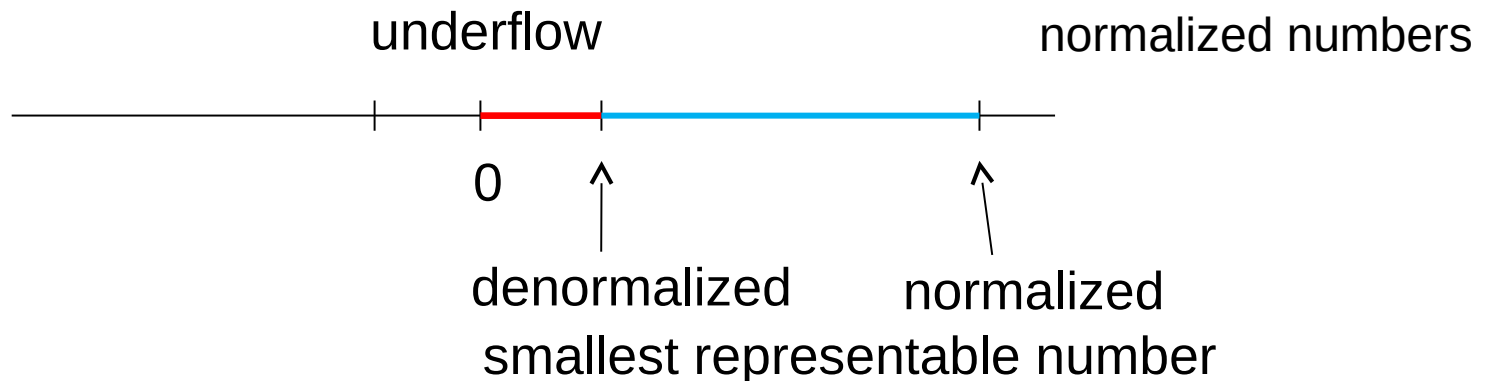


## Implied (hidden) leading 1 bit

- Most significant bit of the mantissa is one for each normalized number and it is not stored in the representation for the normalized numbers
- If exponent representation is zero then encoded value is zero or denormalized number which requires to store most significant bit
- Denormalized numbers allow to keep resolution in the range from the smallest normalized number to zero

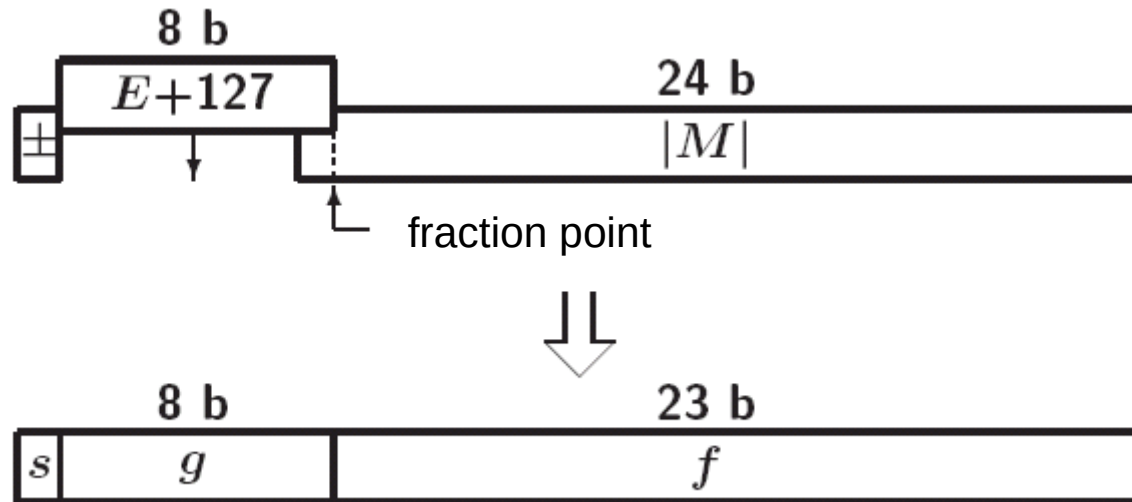
# Underflow/lost of the precision for IEEE-754 representation

- The case where stored number value is not zero but it is smaller than smallest number which can be represented in the normalized form
- The direct underflow to the zero can be prevented by extension of the representation range by denormalized numbers



# ANSI/IEEE Std 754-1985 – 32b a 64b formats

ANSI/IEEE Std 754-1985 — single precision format — 32b



ANSI/IEEE Std 754-1985 — double precision format — 64b

$g \dots 11b$

$f \dots 52b$



# Representation of the fundamental values

## Zero

Positive zero	<b>0 00000000 00000000000000000000000000000000</b>	<b>+0.0</b>
Negative zero	<b>1 00000000 00000000000000000000000000000000</b>	<b>-0.0</b>

## Infinity

Positive infinity	<b>0 11111111 00000000000000000000000000000000</b>	<b>+Inf</b>
Negative infinity	<b>1 11111111 00000000000000000000000000000000</b>	<b>-Inf</b>

## Representation corner values

Smallest normalized	<b>* 00000001 00000000000000000000000000000000</b>	<b><math>\pm 2^{(1-127)}</math> <math>\pm 1.1755 \cdot 10^{-38}</math></b>
Biggest denormalized	<b>* 00000000 11111111111111111111111111111111</b>	<b><math>\pm (1 - 2^{-23}) 2^{(1-126)}</math></b>
Smallest denormalized	<b>* 00000000 00000000000000000000000000000001</b>	<b><math>\pm 2^{-23} 2^{-126}</math> <math>\pm 1.4013 \cdot 10^{-45}</math></b>
Max. value	<b>0 11111110 11111111111111111111111111111111</b>	<b><math>(2 - 2^{-23}) 2^{(127)}</math> <b><math>+3.4028 \cdot 10^{+38}</math></b></b>

## Not a number (NaN)

- All ones in the exponent
- Mantissa not equal to the zero
- Used, where no other value fits (i.e.  $+\text{Inf} + -\text{Inf}$ ,  $0/0$ )
- Compare to  $(X+ +\text{Inf})$  where  $+\text{Inf}$  is sane result

# IEEE-754 special values summary

sign bit	Exponent representation	Mantissa	Represented value/meaning
0	$0 < e < 255$	any value	normalized positive number
1	$0 < e < 255$	any value	normalized negative number
0	0	$> 0$	denormalized positive number
1	0	$> 0$	denormalized negative number
0	0	0	positive zero
1	0	0	negative zero
0	255	0	positive infinity
1	255	0	negative infinity
0	255	$\neq 0$	NaN – does not represent a number
1	255	$\neq 0$	NaN – does not represent a number

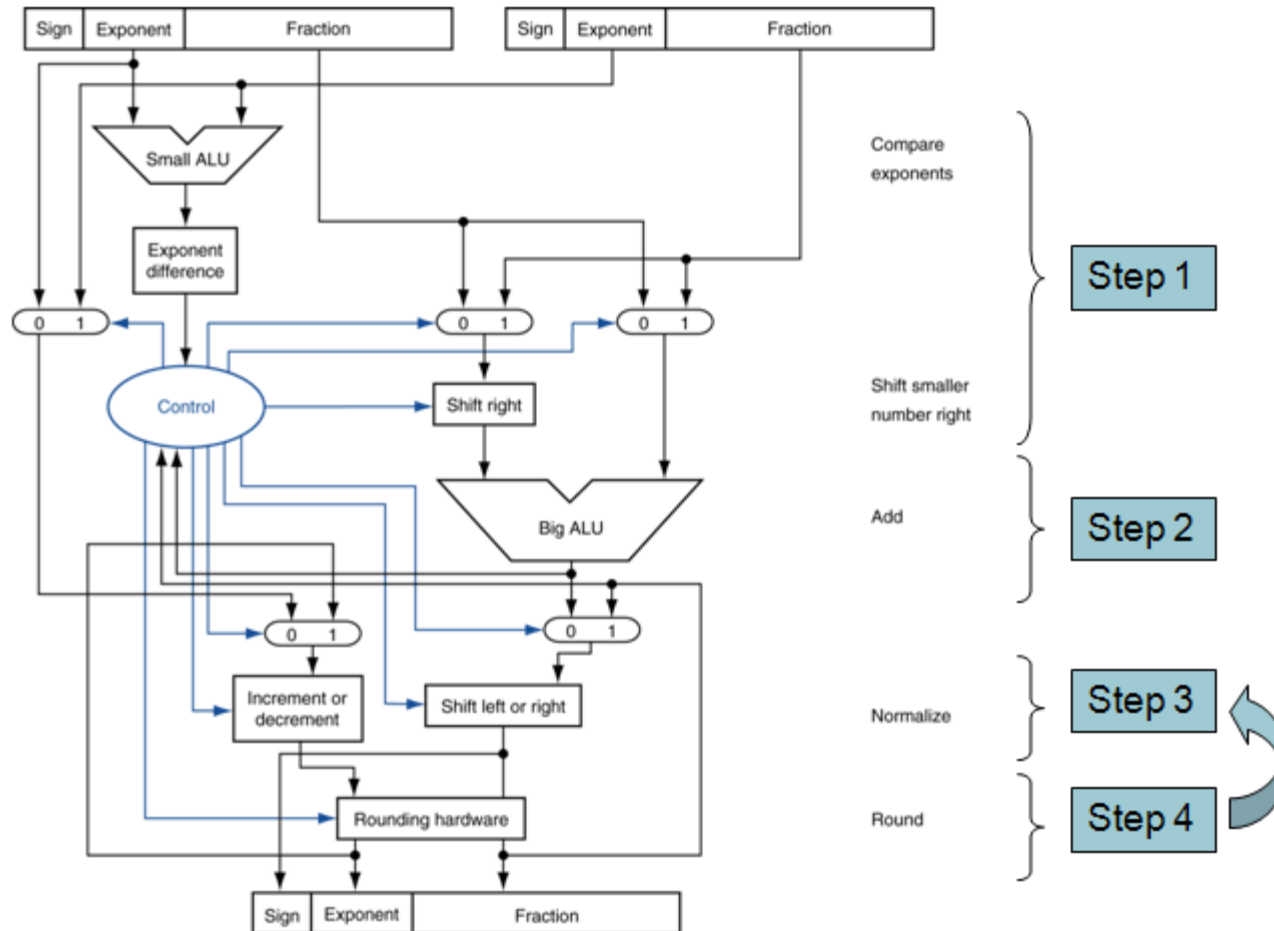
## Comparison

- Comparison of the two IEEE-754 encoded numbers requires to solve signs separately but then it can be processed by unsigned ALU unit on the representations
$$A \geq B \iff A - B \geq 0 \iff D(A) - D(B) \geq 0$$
- This is advantage of the selected encoding and reason why sign is not placed at start of the mantissa

## Addition of floating point numbers

- The number with bigger exponent value is selected
- Mantissa of the number with smaller exponent is shifted right – the mantissas are then expressed at same scale
- The signs are analyzed and mantissas are added (same sign) or subtracted (smaller number from bigger)
- The resulting mantissa is shifted right (max by one) if addition overflows or shifted left after subtraction until all leading zeros are eliminated
- The resulting exponent is adjusted according to the shift
- Result is normalized after these steps
- The special cases and processing is required if inputs are not regular normalized numbers or result does not fit into normalized representation

# Hardware of the floating point adder



## Multiplication of floating point numbers

- Exponents are added and signs xor-ed
- Mantissas are multiplied
- Result can require normalization
  - max 2 bits right for normalized numbers
- The result is rounded
  
- Hardware for multiplier is of the same or even lower complexity as the adder hardware – only adder part is replaced by unsigned multiplier

# Floating point arithmetic operations overview

**Addition:**  $A \cdot z^a, B \cdot z^b, b < a$  unify exponents  
 $B \cdot z^b = (B \cdot z^{b-a}) \cdot z^{b-(b-a)}$  by shift of mantissa  
 $A \cdot z^a + B \cdot z^b = [A + (B \cdot z^{b-a})] \cdot z^a$  sum + normalization

**Subtraction:** unification of exponents, subtraction and normalization

**Multiplication:**  $A \cdot z^a \cdot B \cdot z^b = A \cdot B \cdot z^{a+b}$   
 $A \cdot B$  - normalize if required  
 $A \cdot B \cdot z^{a+b} = A \cdot B \cdot z \cdot z^{a+b-1}$  - by left shift

**Division:**  $A \cdot z^a / B \cdot z^b = A/B \cdot z^{a-b}$   
 $A/B$  - normalize if required  
 $A/B \cdot z^{a-b} = A/B \cdot z \cdot z^{a-b+1}$  - by right shift