

# Low Power SoC Design

Christian Piguet, CSEM, Neuchâtel, Switzerland

Topical Workshop on Electronics for Particle Physics

TWEPP-09, Paris

21-25 September, 2009

# Introduction

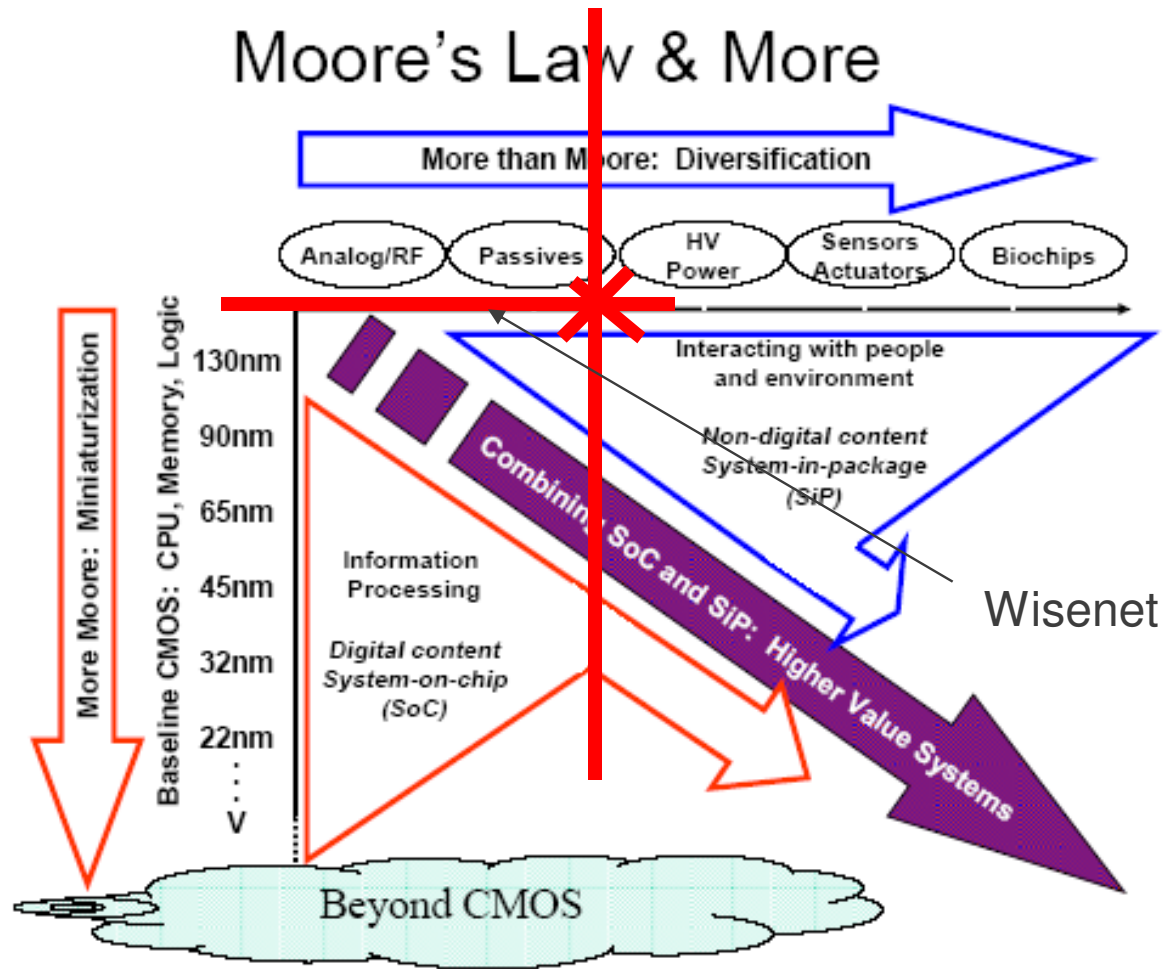
---

- The design of System-on-Chip (SoC) in very deep submicron technologies becomes a very complex task that has to bridge very high level system description (complexity) to low level consideration (technology defaults and variations).
- This talk will describe some of these low level main issues, such as dynamic and static power consumption, temperature, technology variations, interconnect, DFM, reliability and yield, and their impact on high-level design, such as the design of multi-Vdd, fault-tolerant, redundant or adaptive chip architectures.
- Some very low power System-on-Chip (SoC) will be presented in three domains: wireless sensor networks, vision sensors and mobile TV.

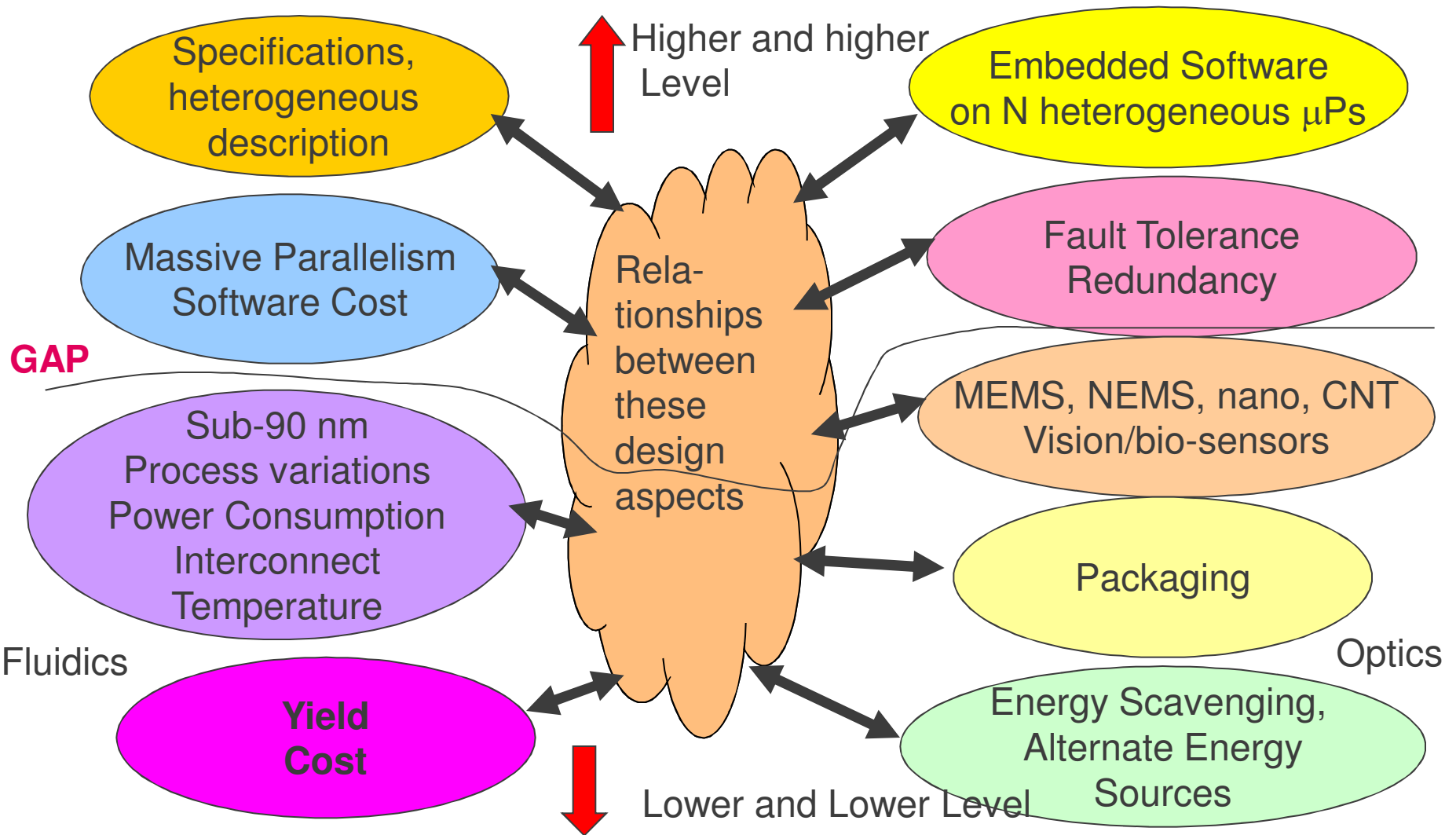
# 1. What is different today?

- More Moore
- More than Moore
- Beyond CMOS
  
- And it is mandatory to simulate and verify everything

It is the famous diagram that everybody knows quite well  
 What does it mean, practically?

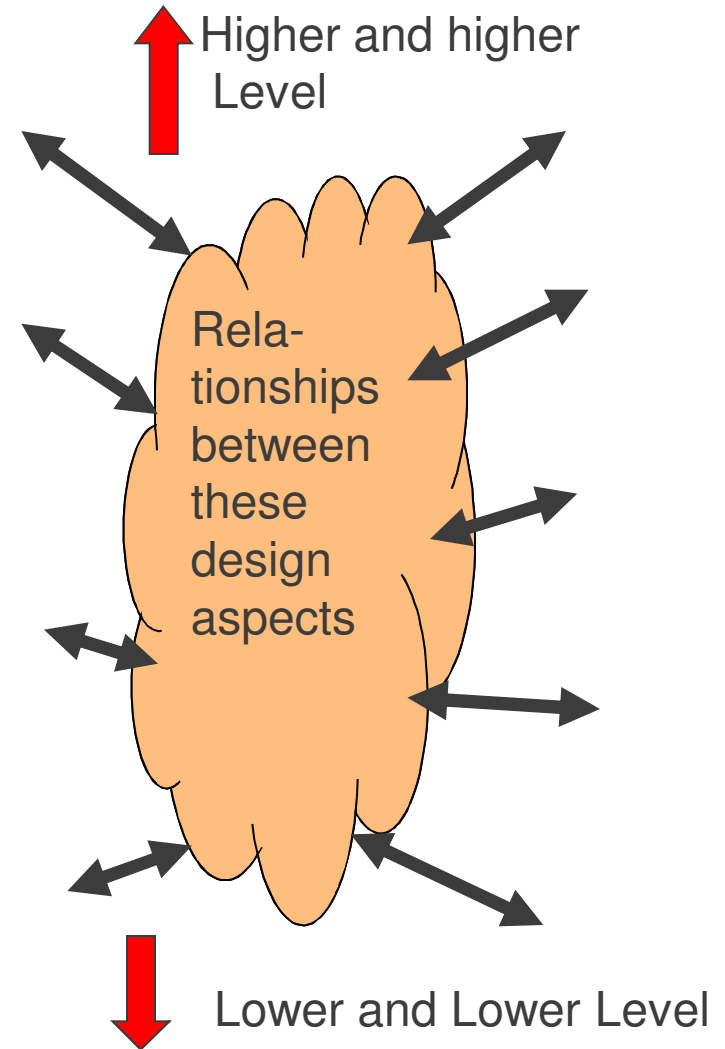


# A very long list of problems... very complex...



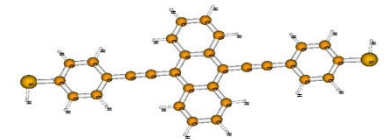
# Why it is so complex...

- The relationships between these design aspects are very complex
- **It is extremely interdisciplinary**
- We are going higher and higher, it is ARTEMIS platform
- We are going lower and lower, it is ENIAC platform
- With a huge gap between the two, that is larger and larger!!



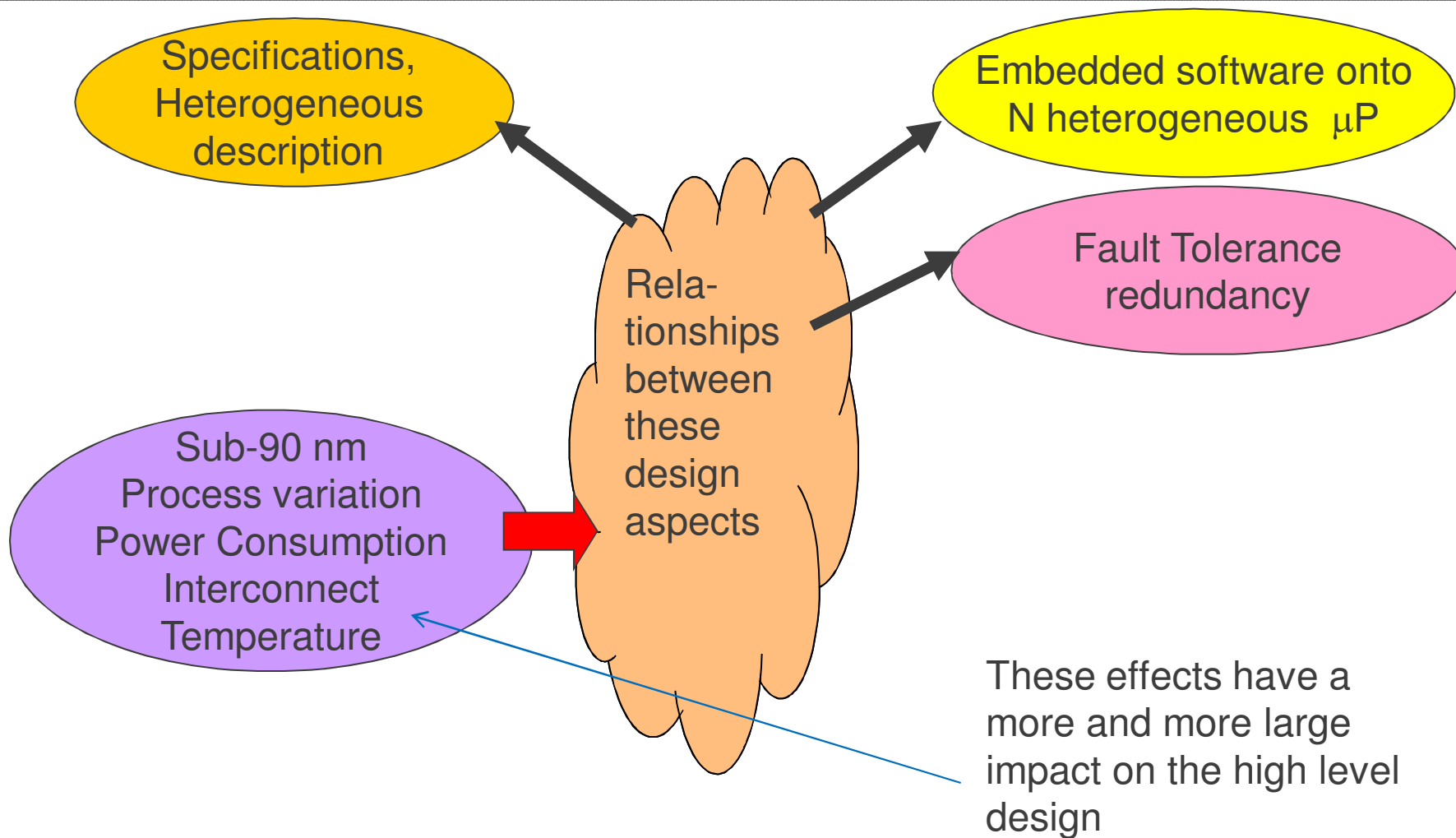
## 2. Interdependency from low level towards high level

- The problems that we have at low level:
  - Dynamic power consumption, current peak
  - Static power consumption, various types of leakage
  - Temperature
  - Technology variations
  - Interconnect delays
  - Reliability and Yield
  - CMOS « end of scaling » around 11 nm, 2015? « beyond CMOS » (CNT, molecular)



- And we have to shift all these effects to high level to be capable of designing heterogeneous SoC that take into account these effects. **What is the impact on the design of architectures?**

# To shift these sub 65 nm effects to high level design



To go up

Artemis

New Disruptive Architectures for Heterogeneous Systems  
MultiCores, Huge Embedded Software, Synchronization

They touch each other,  
But is it still for a long time?

?

It is very difficult to shift these  
information to high level and to  
take them into account for a  
high level synthesis

Power, Leakage, Peak Current, Reliability, Yield,  
Temperature, Complexity, Fault-Tolerance, Memory  
dominated, DRAM, ... new devices, new materials,  
Cost

Eniac



European Nanoelectronics  
Initiative Advisory Council

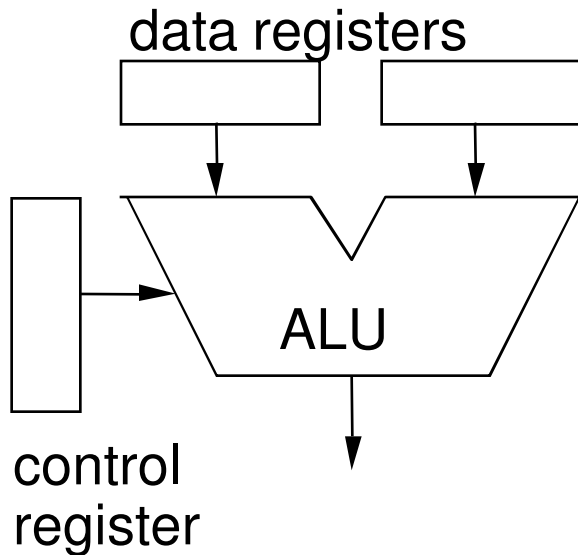


# Example 1: Dynamic Power

---

- Many techniques have been proposed (and some are widely used today) for reducing dynamic power
- Gated Clock, Logic parallelization, asynchronous, adiabatic, bus encoding, standard cell libraries, complex gate decomposition, transistor sizing
- Gated clock is widely used (to cut the clock when the unit is idle)
- Parallelism has a strong impact on high level design. Working with many parallel cores or execution units at low supply voltage is always beneficial for the dynamic power.
- However, it is another story for leakage due the significant increase in terms of number of transistors.

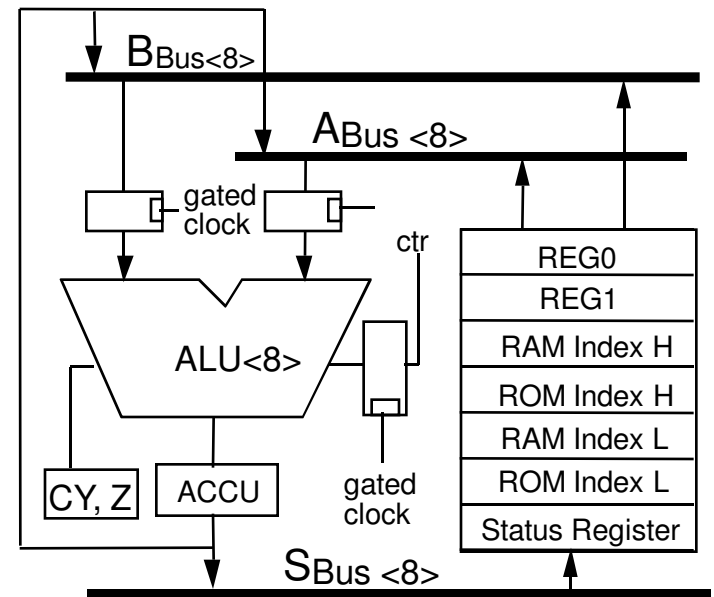
# Gated Clock in DFF-based Designs



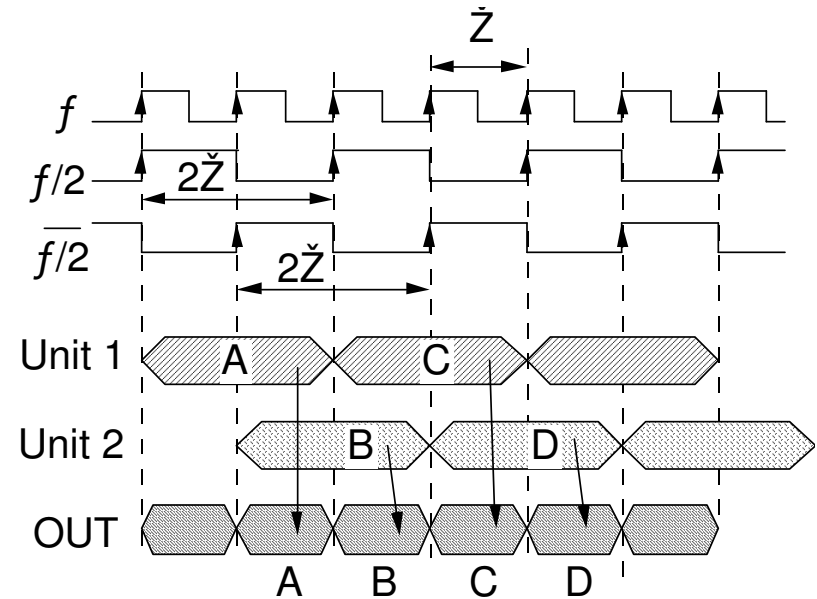
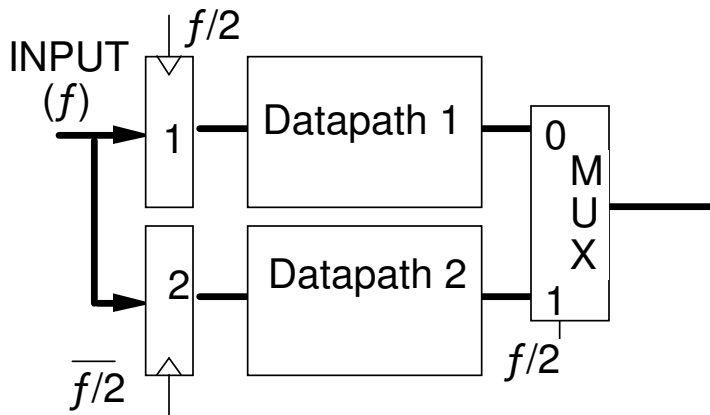
To minimize the activity of a combinational circuit (ALU), registers are located at the inputs of the ALU. They are loaded at the same time --> very few transitions in the ALU

These registers are at the same time pipeline registers (a pipeline for free !)

The pipeline mechanism does not result in a more complex architecture, but reduces the power



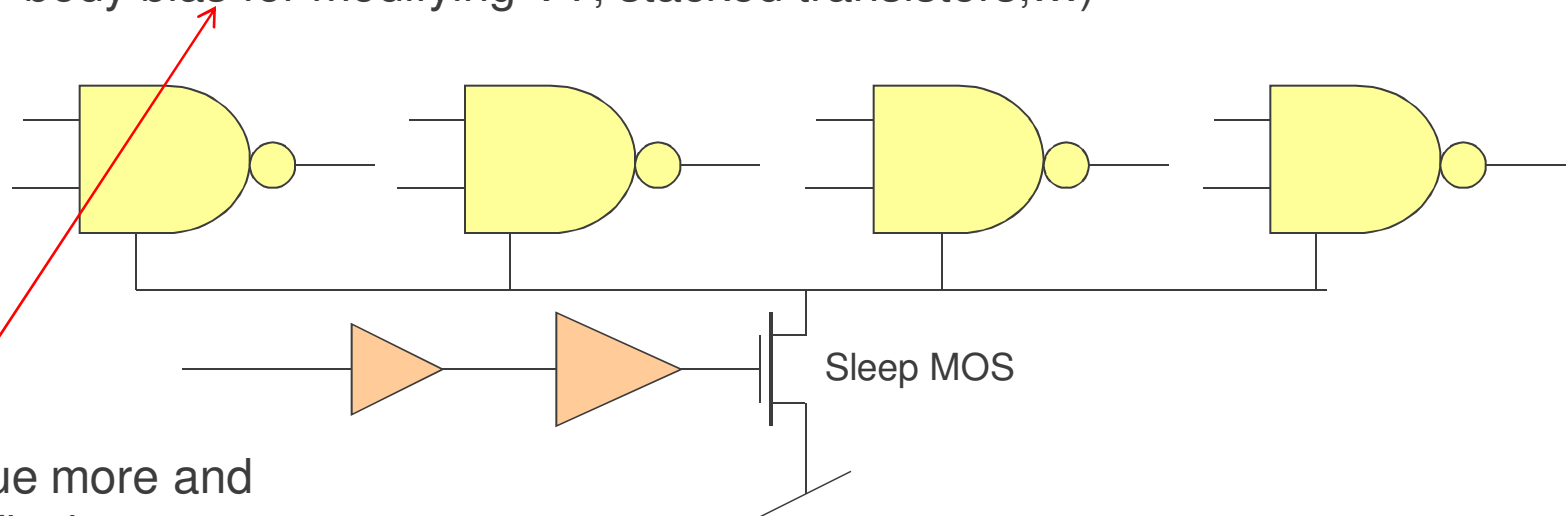
# Logic Parallelization



$$P = M \cdot C \cdot f/M \cdot V_{dd}^2 = C \cdot f \cdot V_{dd}^2$$

## Example 2: Leakage, impact on architectures

- Leakage current (MOS, CNT, nanowires,..) when they are off, leakage exponential increase with  $V_T$  decrease, 50% of total power is leakage
- There are many techniques at low or at circuit level for reducing leakage. A well-known and used technique is to have sleep transistors to cut the supply voltage for idle blocks, but other techniques are also available (several  $V_T$ , body bias for modifying  $V_T$ , stacked transistors,...)



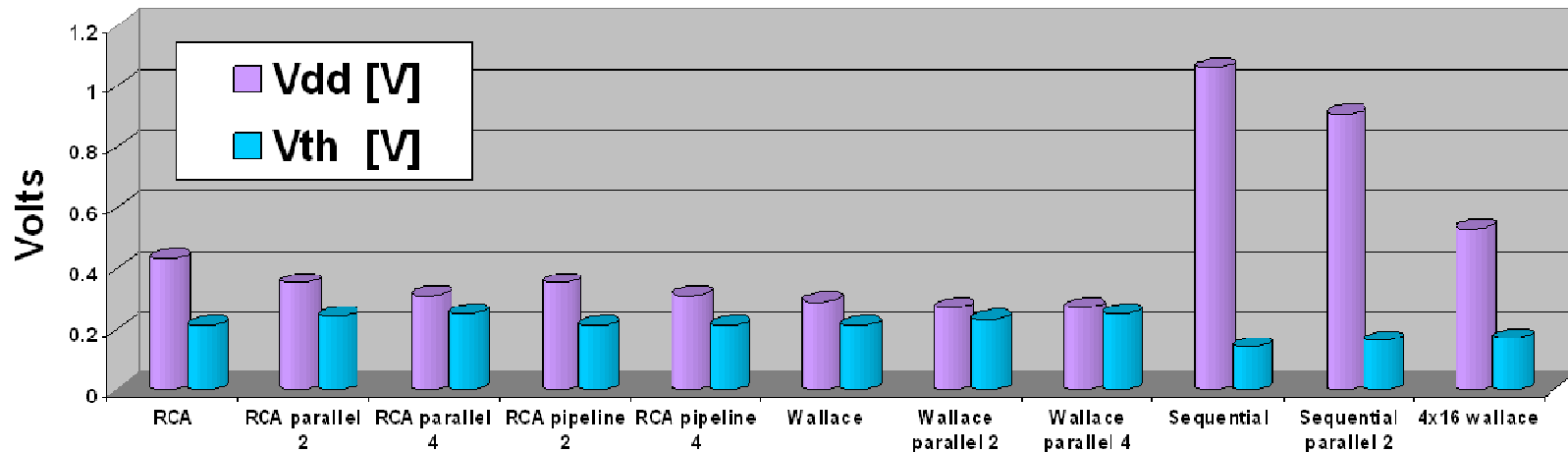
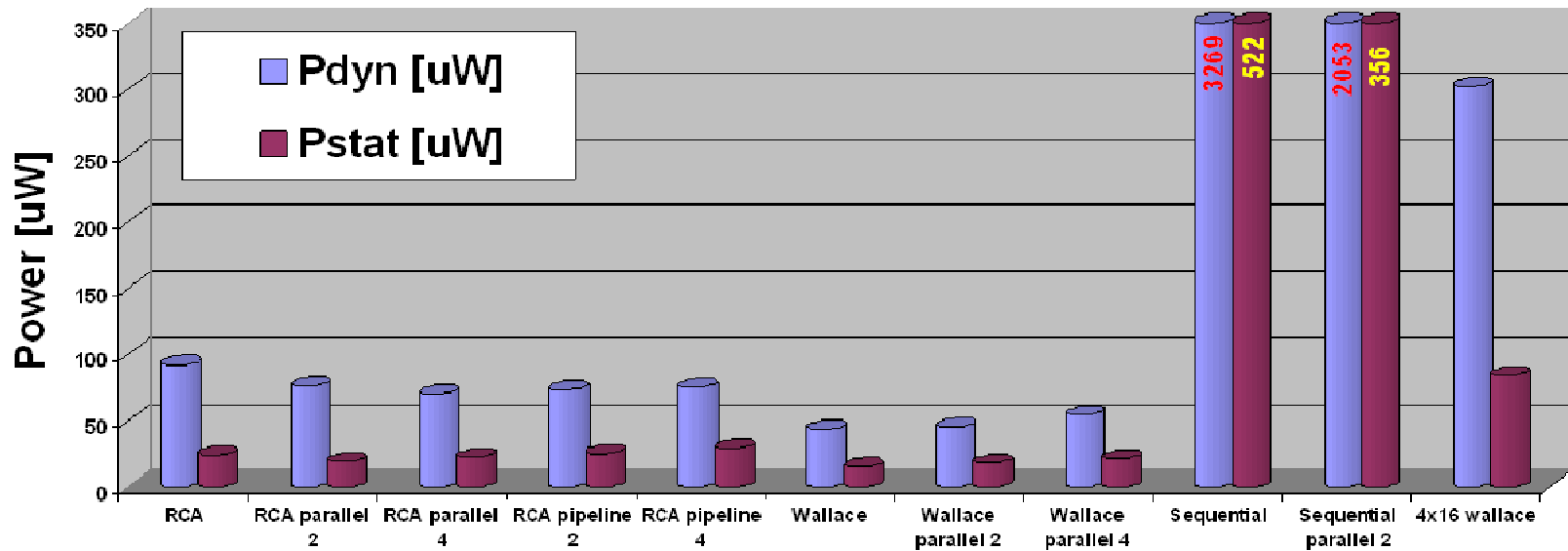
Technique more and more difficult to use

# Leakage Reduction at Architectural Level

---

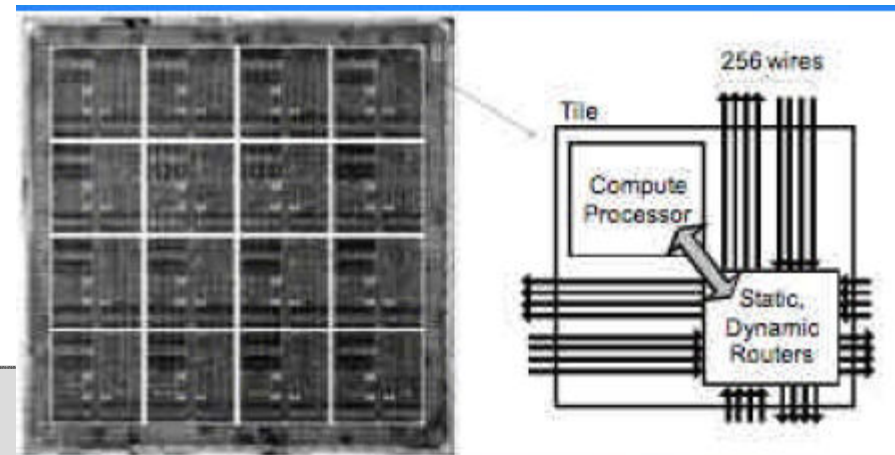
- Is it possible to find an architecture, for the same logic function, that provides a lower leakage, or a lower total power (dynamic + static) at the same speed?
- One searches for an optimal total power considering that the two parameters  $V_{dd}$  and  $V_T$  are free
- Consequently, a too sequential architecture will present a high  $V_{dd}$  and a small  $V_T$  to reach the required speed (and a very large dynamic as well as a very large static power).
- Similarly, if the selected architecture presents a too large parallelism, the number of transistors will so large that leakage and total power will be too large
- Example with a  $16 \times 16$  bit multiplier

# 11 architectures of multipliers



## Example 3: Interconnect delays

- Quite obvious: for every technology reduction factor  $S$ , wire delay is increased by a factor  $S^2$ !!
- It is a severe problem for busses, but it is an extremely dramatic problem for clock distribution.
- Consequently, the influence on architectures is large: 1) everything could be clockless or asynchronous, or 2) GALS
- Any architecture becomes an array of  $N \times N$  zones (isochronous), so to multicores, massive parallelism, synchronization problems
- NoC: Network-on-Chip



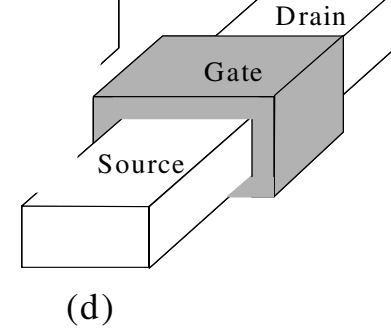
## Example 4: Process Variations

---

- On the same die, technology variations from transistor to transistor (systematic, random), oxide thickness, W and L variations, doping, temperature and Vdd variations, even soft errors)
- VT Variations, factor 1.5 on delays and 20 on leakage
- Big impact on yield
- Which impact on architectures?
- For instance, to go for Multi-core, big parallelism. One can work at lower frequency for the same computation throughput. Consequently, the processor cores (at lower frequencies) are less sensitive to process variations on delay.
- Better to work at high Vdd (at low Vdd, 0.5 Volt, very sensitive to variations)
- Architectures with Fault-Tolerance, spatial or timing redundancy



# What about ABB (body bias) technique



- For over-100nm technologies, Adaptive Body Biasing (ABB) is a good technique for compensating the variations. Both forward body biasing (FBB) and reverse body biasing (RBB) can be used in sub- $V_T$  regime.
- However, body factor is almost zero in emerging Multi-Gate devices which are promising candidate for future electronics.
- For instance, Double-gated FinFET, tri-gated, nanowire body, and gate-all-around (GAA) MOSFETs
- In multi-gate devices, body factor is much smaller than in single-gate devices because of the enhanced coupling between gate and channel and because the lateral gates shield the device from the electric field from the back-gate
- Measurements in [Singh] show that in GAA devices body factor is exactly zero.

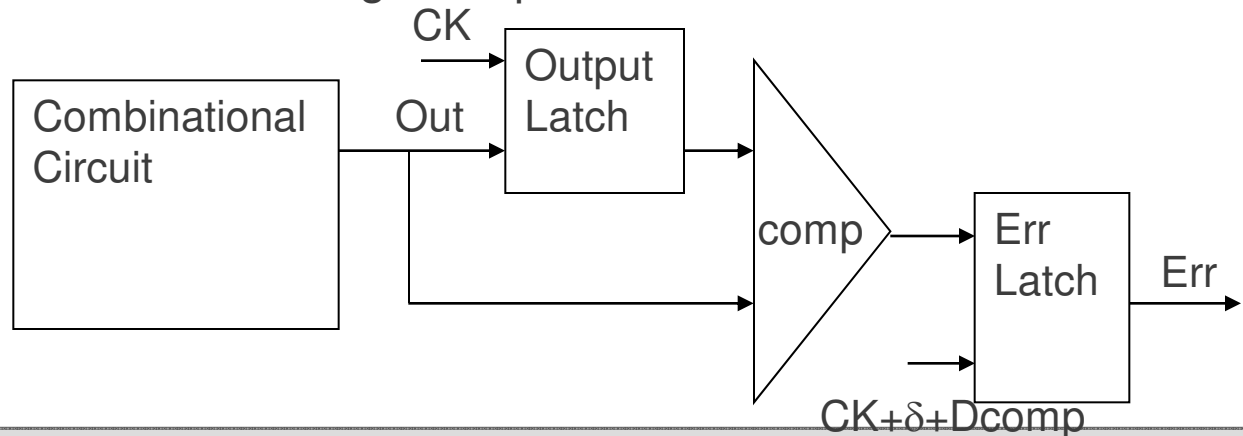
# Number of cells in Library

- Resistance to PV is better with long critical paths, as the technology variations are better compensated with a large number of cells
- For the same logic function, a way to have more cells in a given critical path is to provide a library with few simple cells
- Carnegie Mellon: “It can be shown that with a small set of Boolean functions ... (and careful selection of lithography friendly patterns)...we mitigate PV”
- Block architecture: for an full adder, which is the best architecture (ripple carry, carry look-ahead) and VDD for reducing the effect of PV?
- A ripple carry adder at 500 mV provides same speed and same power than a carry look-ahead adder at 400 mV with 2 times less sensitivity to PV
- Using low-power slow circuits in higher VDD voltage is better than using high-power fast circuits in lower VDD!

# Spatial or timing redundancy

- A system is not composed of reliable units, one has to consider that every unit could fail. However, the system could not fail!!
- A possible architecture is to use massive parallelism while presenting redundant units that could take over the work of faulty units.
- One can have spatial redundancy (very expensive) or timing redundancy (quite expensive in terms of throughput)
- And big problem, what about speed and power??
- To compare the result of a given operation at 2 different time frames

- RAZOR



# PCMOS

- PCMOS, Probabilistic CMOS, each gate has a probability of failure.
- This work characterizes an explicit relationship between the probability  $p$  with which the CMOS switch computes correctly, and its associated energy consumed by each switching step across technology generations
- Each basic logic gate (NOT, NAND, NOR) has a given probability to provide a correct result for a given input
- For instance, a truth table indicates that for input 100 (correct output is “0”), probability for the output to be “1” is  $\frac{1}{4}$  while probability for the output to be “0” is  $\frac{3}{4}$
- Using such basic gates to synthesize more complex functions (adder, flip-flops, etc..), there are many different schematics that perform the same function: so which schematic will provide the better probability to have a correct adder “sum” and “carry” (or “Q” for the flip-flop)

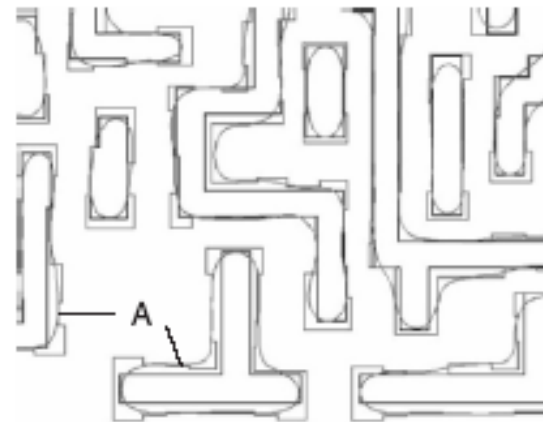
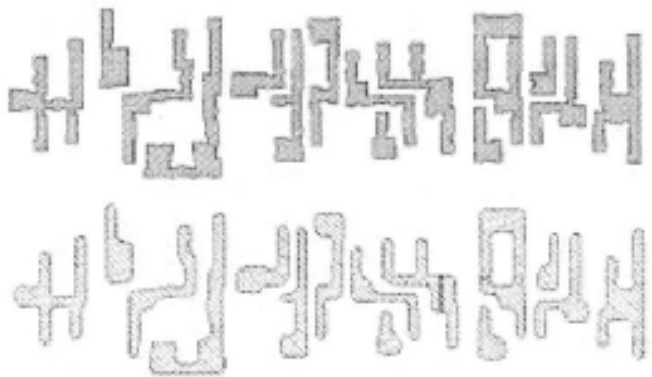
# Subthreshold Logic

---

- Logic circuits based on transistors operated in weak inversion (also called subthreshold)
- This technique has been revived recently and applied to complete subsystems operated below 200 mV.
- It has been demonstrated that minimal energy circuits are circuits operated in subthreshold regime with  $V_{dd}$  reduced to under  $V_T$ , resulting in lower frequencies and larger clock period.
- So dynamic power is reduced, static power is decreased, but the static energy is increased as more time is required to execute the logic function. So there is an optimum in energy.
- This optimal energy is also depending on logic depth and activity factor. The minimal  $V_{dd}$  (and minimal energy) is smaller for small logical depth and for large activity factors.

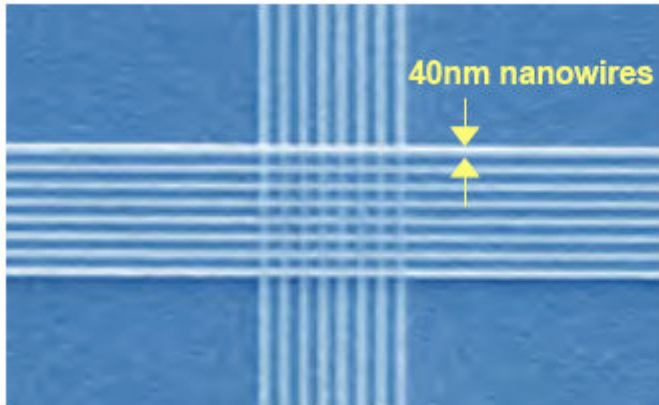
## Example 5: Yield and DFM

- DFM (Design For Manufacturing)
- Mask smallest dimensions are well below the lithographic light wavelengths.
- One has a lot of strange effects, bad line extension, missing small geometries, etc...
- To limit those variations by using regular layouts such as PLA or ROM for combinational circuits



# Impact of DFM on architectures

- Architectures based on regular blocks at layout level
- ROM, PLA, gate matrix, SoC fully dominated by memories

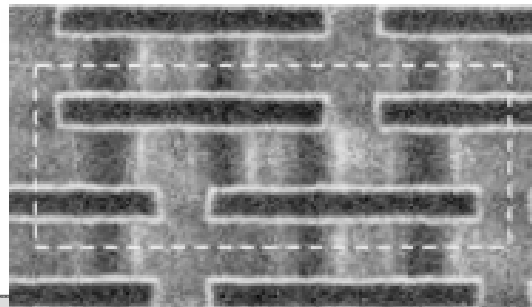


8 x 8 crossbar  
switch built with  
40nm nanowires

HP Labs – 2005



Gate  
Matrix  
CSEM  
1985



SRAM INTEL  
45 nm

## Example 6: Alternative Energy Sources

---

- Many alternative and very diverse energy sources, with «energy scavenging»
  - batteries, fuel cells
  - Scavenging from environment, vibrations, thermoelectricity, solar cells, piezo, human energy sources, walking, shoes, mechanical watches
  - One has to generate inside the SoC multiples supply voltages with very diverse peak currents (some  $\mu\text{A}$ , some mA, up to 10 or 100 mA)
- « Power Management » circuits become very complicated (DC-DC, regulators).
- On top of this, one requires to add DVS and DVFS (Dynamic Voltage Frequency Scaling)



# Impact of energy sources on architectures

---

- One has to manage these energy sources as well as DVFS, idle modes, recharge modes, etc...
- This can be performed by the Operating System, it is however quite complex
- This part of the embedded software has to interact with the application embedded software, that increases the overall complexity.

# Conclusion Up to this Point: Increased Complexity

---

- More and more low level effects that have to be taken into account
- The impacts of these low level effects on to the high level synthesis process of SoCs is more and more difficult to understand and to take them into account
- Sure that only the low level effects have been presented here, but there are also effects at high level that have to be taken at low level:
  - Architecture for executing efficiently a given language
  - Asynchronous architecture requiring special standard cell library
  - Parallelizing compiler onto N processors, which are the constraints on to the processor architecture?

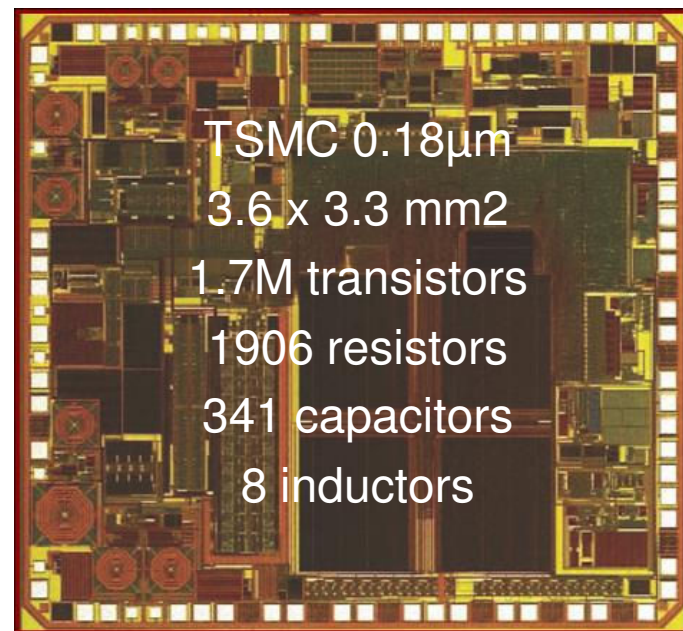
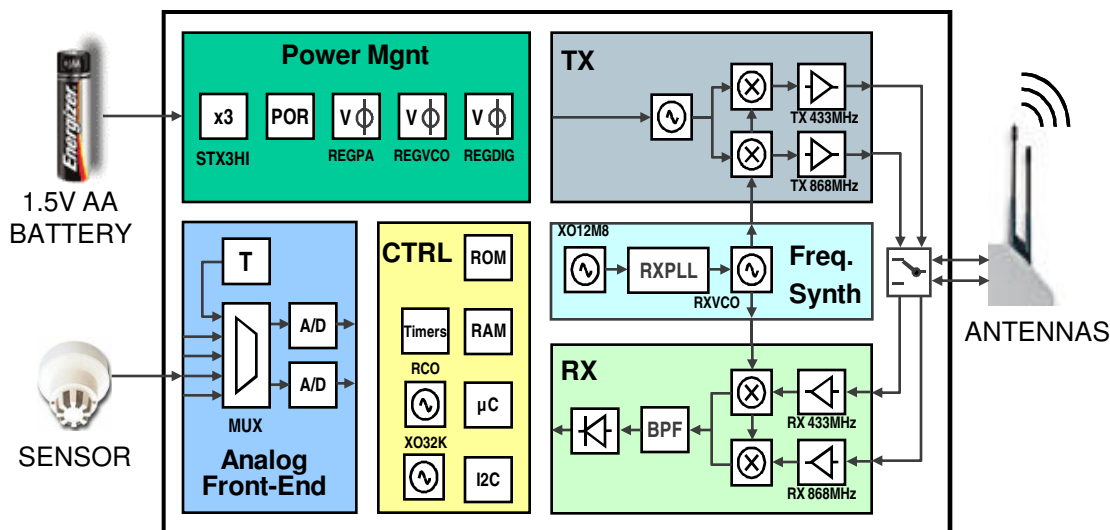
## 3. Heterogeneous SoC Examples

---

- Wireless Sensor Network (called Wisenet)
- Vision sensor-based SoC (called icycam)
- OFDM Mobile TV chip
- Icyflex-based Radio SoC (called icycom)

## 3.1 WiseNET SoC – first CSEM SoC (2004)

- 433MHz / 868MHz Rx & Tx
- CoolRISC  $\mu$ C with low leakage SRAM
- Analog sensor interface
  - 10-bit ADC
  - $4\mu\text{A}$   $\Delta$ - $\Sigma$  ADC with  $10\mu\text{V}$  resolution
- Power Management with step-up converters
- Low voltage SoC:  $0.9\text{V}$ - $1.5\text{V}$
- $I_{\text{RX}}=2\text{mA}$ ,  $I_{\text{TX}}=30\text{mA}$  @ 10dBm
- $I_{\text{avg}}=25\mu\text{A}$  @ 0.1%-1.0% duty cycle
- 12.5-100kb/s in FSK, 2kb/s in OOK



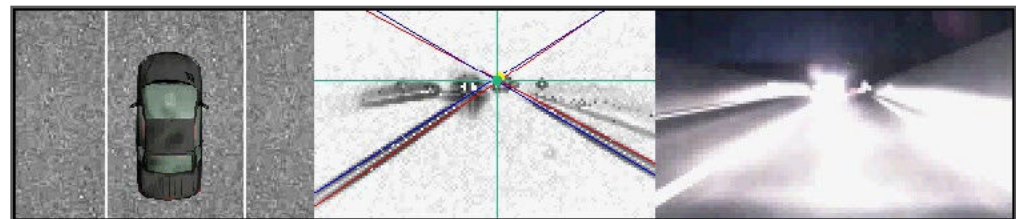
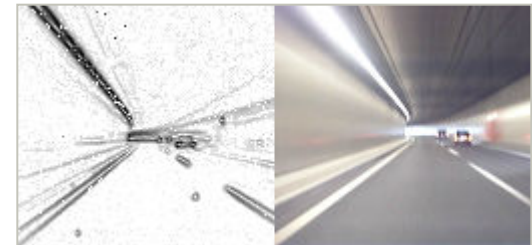
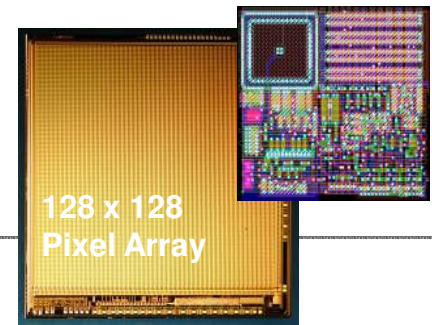
## 3.2 Vision sensors

Extract **image features** from **contrast** magnitude and orientation and output **only** pixels with **pertinent** information

- Much **less data** to process
- **Real-time** operation, low-power and low-cost
- **Independent** of the **illumination** conditions
- Great **stability** (constancy) of representation
- Easy **motion** detection
- High dynamic range (**120dB**)

### Automotive applications:

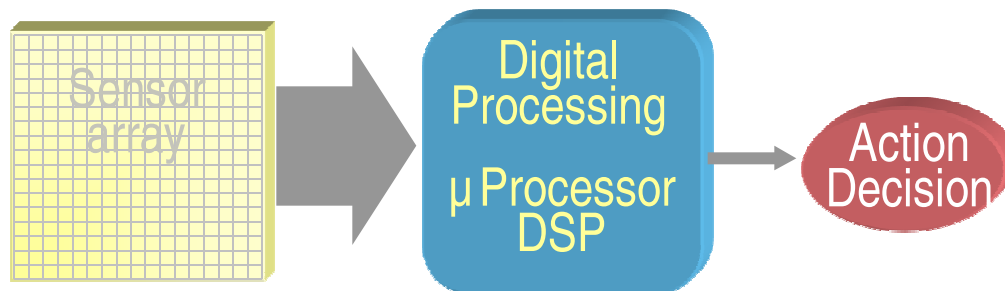
- **Lane departure warning**
- **Collision warning**
- **Driver surveillance**



# Sensory Information Processing: Vision Sensor

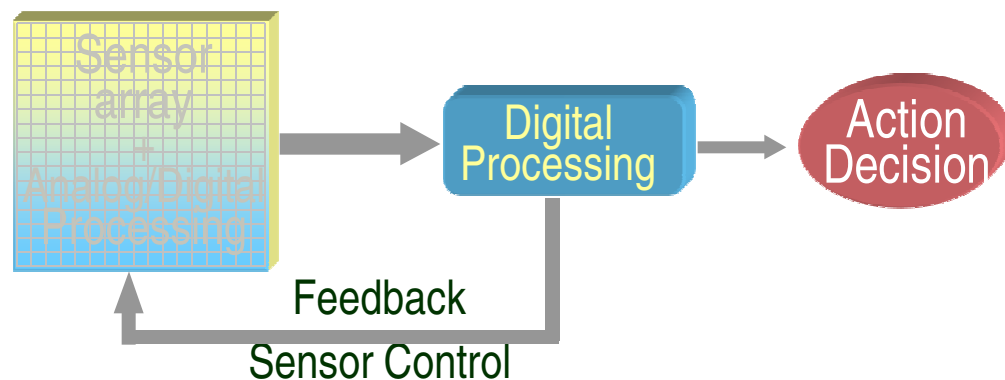
- Classical Approach

- Imaging
- ADC
- « Number crunching »



- Our Approach

- On-Chip Image Processing
- light digital post-processing



- Extraction & Processing of only sensor informations that are relevant for a given task

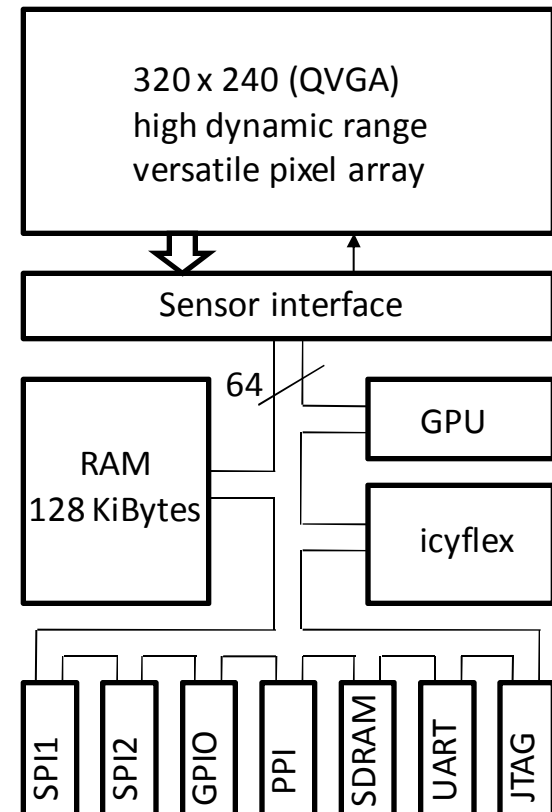
# icycam, a SoC for Vision Applications

- Icycam is a circuit combining on the same chip a **32-bit icyflex processor operated at 50 MHz**, and a high dynamic range versatile pixel array, integrated on a 0.18  $\mu\text{m}$  optical process.

-The pixel array has a resolution of 320 by 240 pixels (QVGA CMOS), with a **pixel pitch of 14  $\mu\text{m}$** .

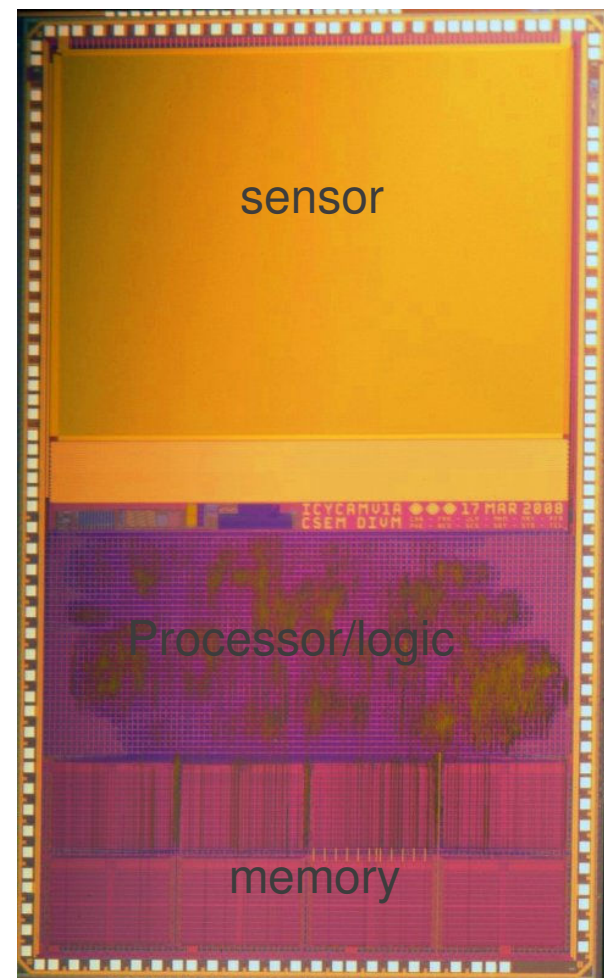
- Able to extract on the fly the local contrast magnitude (relative change of illumination between neighbour pixels) and direction when data are transferred from the pixel array to the memory

-DMA, SPI, PPI, GPIO, UART, SDRAM, JTAG



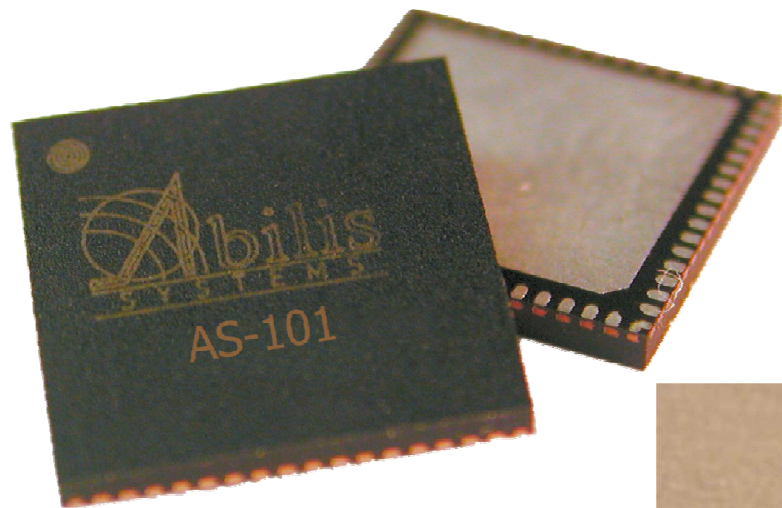
# SoC with local pixel processing and icyflex

- Data transfer between the pixel array and memory or peripherals is performed by group of 4 (10 bits per pixel) or 8 (8 bits per pixel) pixels in parallel at system clock rate.
- These image data can be processed with the icyflex's Data Processing Unit (DPU) which has been complemented with a Graphical Processing Unit (GPU) tailored for vision algorithms, able to perform simple arithmetical operations on 8- or 16-bit data grouped in a 64-bit word.
- Internal SRAM being size consuming, the internal data and program memory space is 128 KiBytes. This memory range can be extended with an external SDRAM up to 32 MiBytes.
- Tower Semiconductor, 180 nm, 43 mm<sup>2</sup>

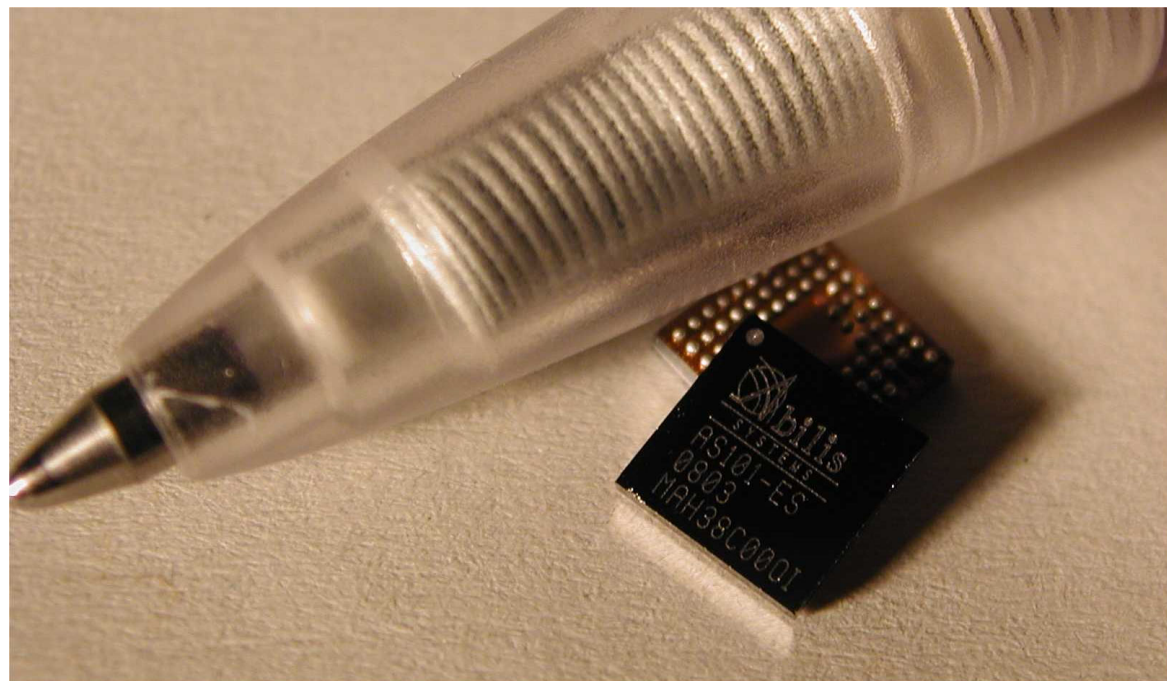




## 3.3 Mobile TV Chip (DVB-T/H) by Abilis



CSEM DSP MACGIC  
License



# MACGIC test Chip

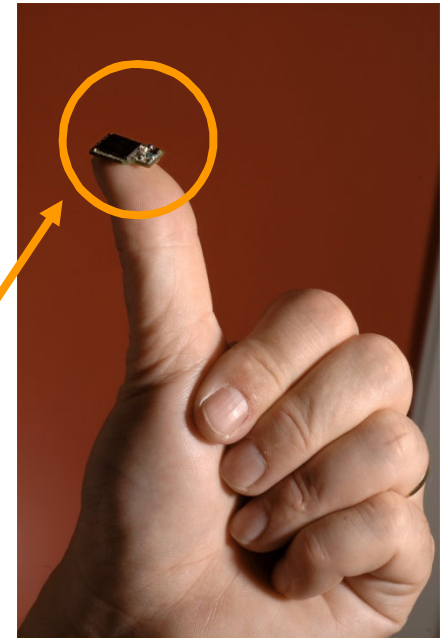
- Technology: TSMC 0.18  $\mu\text{m}$
- Integration date: March 2005
- Size: 3.6 mm by 5.0 mm = 18 mm<sup>2</sup>
- Gates: 200'000 (excluding the RAMs)
- Standard cell lib: CSEM's CSL 6.0 for low power consumption
- Memory: 96 KBytes of CSEM's low power RAM
- Signal pads: 50 (HDU ports are not available as pads)
- Voltage RAM/core: 0.9V to 1.8V
- Voltage pads: VDDcore to 3.3V
- Frequency: 7 MHz @0.7V; 35MHz @1V; 80 MHz @1.8V
- Power consumption: 72  $\mu\text{W}/\text{MHz}$  @0.7V; 150  $\mu\text{W}/\text{MHz}$  @ 1V

# The Abilis chip contains three MACGIC cores

- **Abilis:** To become the world leading supplier of semiconductor solutions of multimode, digital TV receiver and broadband wireless connectivity for mobile terminals

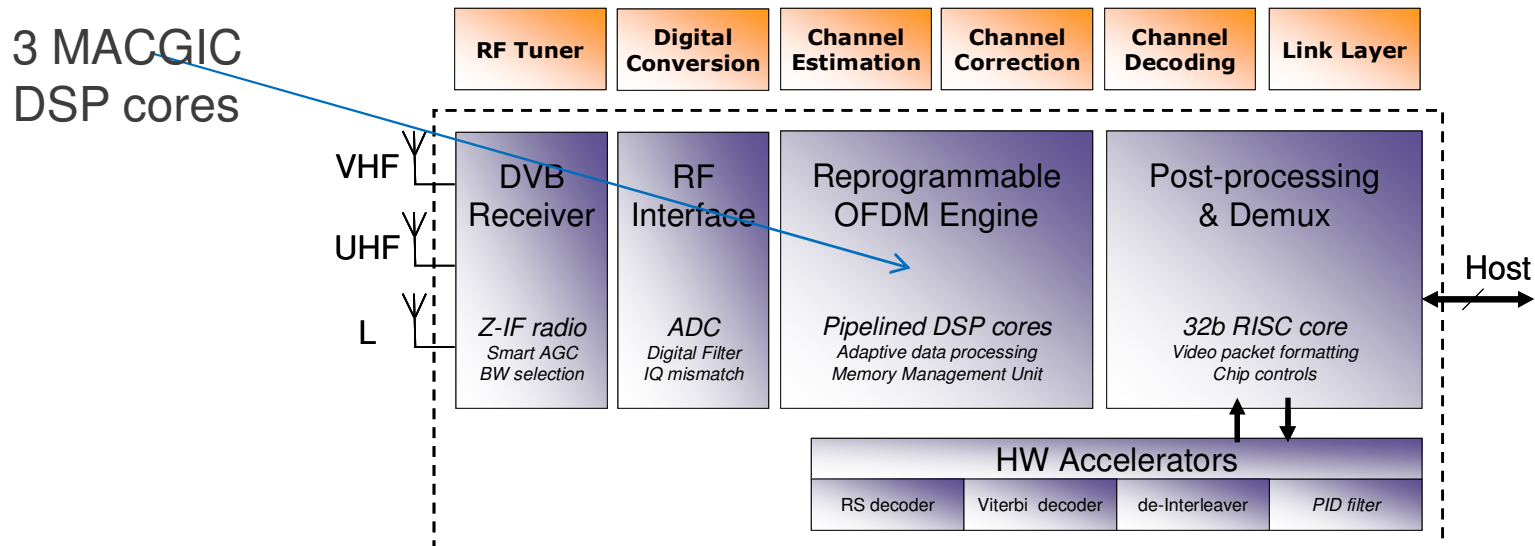


myTV



# Abilis Technology

- Low Power Software Define Radio (OFDM Engine)
- Multi modes System Partitioning
- Quad-band, Low Power RF
- IBM CMOS 90nm technology



## 3.4 icycom: a System-on-Chip for RF applications

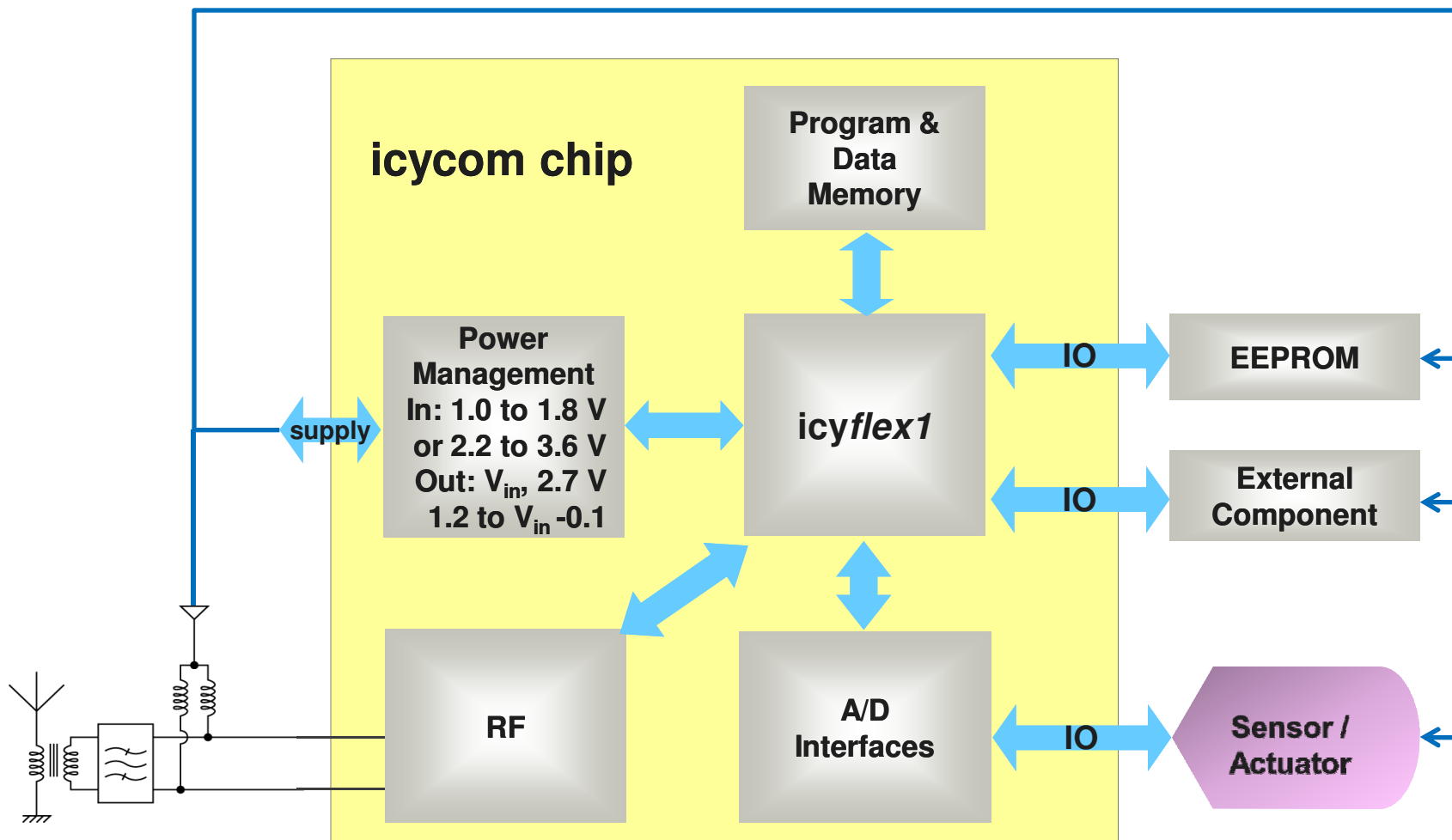
- **icyflex1** runs at up to 3.2 MHz, Avg dyn power: 120  $\mu$ W/MHz @ 1.0 V
- RF: 865 ~ 915 MHz, FSK (incl. MSK, GFSK), 4FSK, OOK, OQPSK
  - TX: 10 dBm
  - RX: -105 dBm at 200 kb/s (BER =  $10^{-3}$ )
- Power management
  - Power supplies for external devices, use of single alkaline or lithium cells
  - Low power modes: multiple standby modes
- 10 bit ADC
- SRAM: 96 KiBytes (with BIST)
- DMA, RTC, Timers, Watchdog, I2C, SPI, I2S, GPIO, UART, JTAG
- TSMC, 180 nm, generic

## Comparison of Starcore, CoolFlux, Macgic and icyflex

Features	Starcore	Macgic	icyflex	CoolFlux
Bits per Instruction	128-bit	32-bit	32-bit	32-bit
Data Word width	16-bit	32/24-bit	32-bit	24-bit
Number of MAC	4	4	2	2
Memory Transfer	8	8	4	2
Operations per cycle	32	32	16	8
Number of equivalent NAND gates	600k	150k	110k	45k
Clock cycles for FFT 256	** 1'614	1'410	2'700	* 5'500
Average Power per MHz @ 1V	* 350 $\mu$ W	170 $\mu$ W	*120 $\mu$ W	* 75 $\mu$ W
Power per MHz @ 1V for FFT	* 600 $\mu$ W	300 $\mu$ W	*215 $\mu$ W	* 130 $\mu$ W
Normalized energy for FFT @ 1V	2.3	1	1.4	1.7

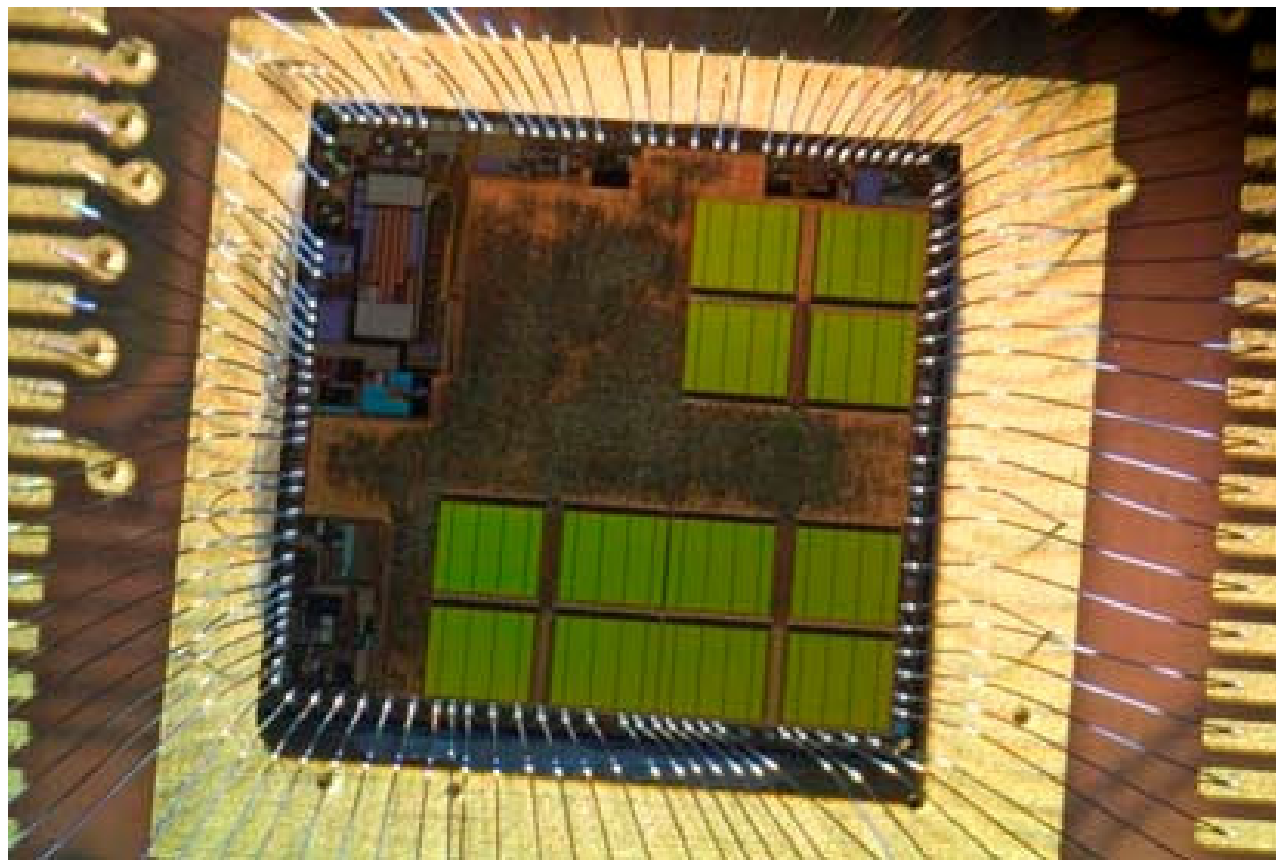
\*\*single precision \*estimated

# icycom: a System-on-Chip for RF applications (cont'd)



## Icycom Chip (2009)

- The application context is miniature wireless sensor network nodes with years of autonomy.





## 4. « Disruptive » Architectures and Systems??

---

- One can look at various Roadmaps
- The end of CMOS «scaling» is predicted around 11 nanometers, around 2013 and 2017
- After 2017, we should move to « Beyond CMOS »
- However, today, there is no clear alternating route to replace CMOS
- CNT, nanowires, molecular switches etc... it is not so clear for architectures and systems requiring billions of switches and how to interconnect them with billions of wires
- Nevertheless, there is an interesting approach in hybrids CMOS and nano, it will be heterogeneous...
- With these nano-elements, one has sometimes the same problems at low level (leakage, process variations), but we could also imagine or hope that some of these effects would disappear!!

## However, one can take some risk by having new ideas..

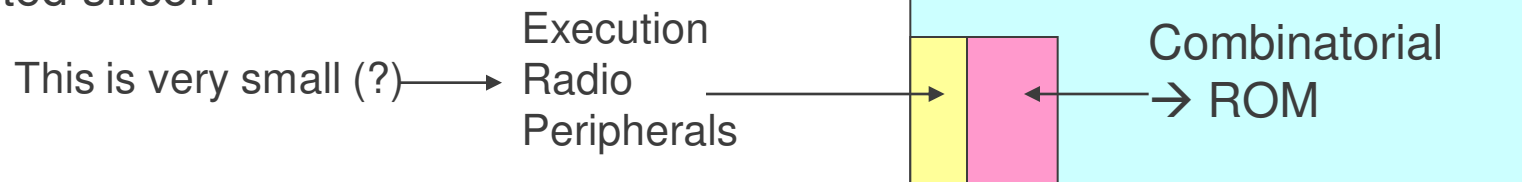
- **A single universal SoC or MPSoC platform:** everybody has to use the same hardware, consequently, design is fully concentrated on embedded software.
  - Very expensive to develop, about 100 M€, and one could ask if is reasonable for applications sensitive to power consumption and to some other performances...

- **A SoC or MPSoC dominated by memories**

Memories are automatically generated, so the hardware part to design is very small.

But one re-generates only the used memories...

Less wasted silicon



# And still new Architectures...

---

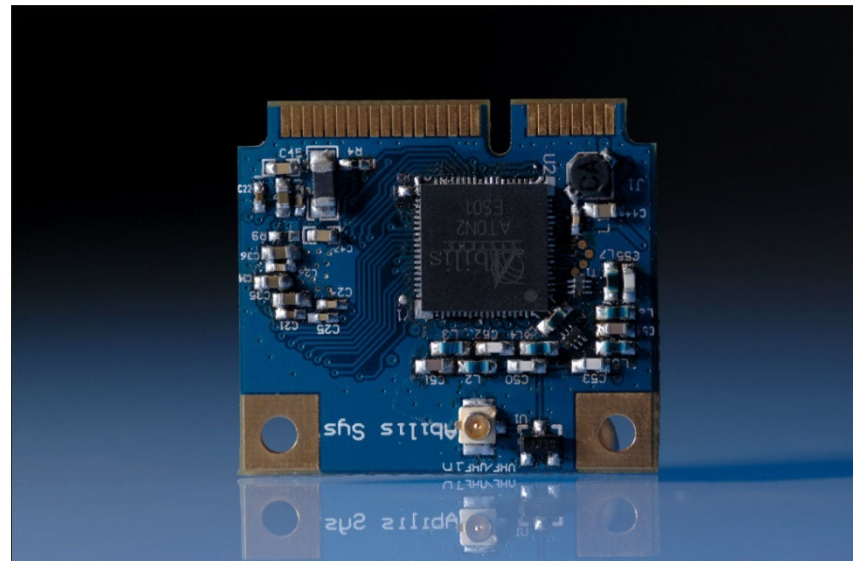
- **SoC or MPSoC with 1'000 parallel processors** « A View of Berkeley »
  - Not the same than multicore (2 à 32)
  - Small logic blocks of 50K gates, and... a lot of memory...
- **Architecture with nano-elements**
  - Completely different, bottom-up design methodology (not top-down)
  - Sure very and very regular circuits and layouts
  - Applications which will be completely different than Pentium
- **Other idea?**

## 5. Conclusion

---

- Diagnostic is clear:
  - Complexity increases, interdisciplinary too
  - More and more interactions between all design levels
  - We are going higher and higher (ARTEMIS) but also lower and lower (ENIAC), resulting in a gap which increases
  - We see design teams that are also more and more heterogeneous
- What we have to do is also more difficult to define:
  - Focused research, yes, but primarily interdisciplinary research
  - To talk to many people, to understand more and more people
  - To go to conferences...

# Thank you for your attention.



Abilis SoC and printed board

# References (I)

- [www.csem.ch](http://www.csem.ch)
- J. Rabaey, “Managing Power Dissipation in the Generation-after-Next Wireless Systems”, FTFC’99, June 1999, Paris, France
- E. Vittoz, “Weak Inversion for Ultimate Low-Power Logic”, Chapter 16 in « Low-Power Electronics Design”, CRC Press, November 2004, edited by Christian Piguet.
- S. Hanson, B. Zhai, D. Blaauw, D. Sylvester, A. Bryant, X. Wang, “Energy Optimality and Variability in Subthreshold Design”, International Symposium on Low Power Electronics and Design, 2006 , pp. 363-365 .
- K. Roy, S. Mukhopadhyay, H. Mahmoodi-Meimand, “Leakage Current Mechanisms and Leakage Reduction Techniques in Deep-Submicrometer CMOS Circuits”, Proc. of the IEEE, Vol. 91, No.2, 2003, pp305-327.
- Christian Schuster, Jean-Luc Nagel, Christian Piguet, Pierre-André Farine, « Leakage reduction at the architectural level and its application to 16 bit multiplier architectures”, PATMOS ’04, Santorini Island, Greece, September 15-17, 2004
- C. Schuster, C. Piguet, J-L. Nagel, P-A. Farine, «An Architecture Design Methodology for Minimal Total Power Consumption at Fixed Vdd and Vth”, Journal of Low-Power Electronics (JOLPE), Vol. 1, No. 1, April 2005, pp. 1-8.
- C. Schuster, J-L. Nagel, C. Piguet, P-A. Farine, “Architectural and Technology Influence on the Optimal Total Power Consumption”, DATE 2006, Munchen, March 6-10, 2006
- N. Singh, et al., “High-performance fully depleted silicon nanowire (diameter  $\leq 5$  nm) gate-all-around CMOS devices,” IEEE Electron Device Lett., vol. 27, no. 5, pp. 383–386, May 2006.

## References (II)

- Zhai, B., Blaauw, D.; Sylvester, D. & Flautner, K. "Theoretical and Practical Limits of Dynamic Voltage Scaling", *Design Automation Conference, 2004* , pp. 868-873.
- Hanson, S.; Zhai, B.; Blaauw, D.; Sylvester, D.; Bryant, A. & Wang, "Energy Optimality and Variability in Subthreshold Design", *International Symposium on Low Power Electronics and Design, ISLPED 2006*, pp. 363-365.
- B. Zhai, et al., "A 2.60pJ/Inst Subthreshold Sensor Processor for Optimal Energy Efficiency," *VLSI Ckts Symp, 2006*.
- Joyce Kwong, et. al., "A 65nm Sub-Vt Microcontroller with Integrated SRAM and Switched-Capacitor DC-DC Converter," *ISSCC'08*, pp. 318-319, 2008.
- C. Piguët, G. Berweiler, C. Voirol, E. Dijkstra, J. Rijmenants, R. Zinszner, M. Stauffer, M. Joss, "ALADDIN : A CMOS Gate-Matrix Layout System", *Proc. of ISCAS 88*, p. 2427, Espoo, Helsinki, Finland, 1988.
- M. Haykel Ben Jamaa, Kirsten E. Moselund, David Atienza, Didier Bouvet, Adrian M. Ionescu, Yusuf Leblebici, and Giovanni De Micheli, "Fault-Tolerant Multi-Level Logic Decoder for Nanoscale Crossbar Memory Arrays", *Proc. ICCAD'07*, pp. 765-772
- Christian C. Enz et al. « WiseNET: An Ultralow-Power Wireless Sensor Network Solution », *IEEE Computer*, August 2004, pp. 62-70.
- C. Arm, S. Gyger, J.-M. Masgonty, M. Morgan, J.-L. Nagel, C. Piguët, F. Rampogna, P. Volet, "Low-Power 32-bit Dual-MAC 120  $\mu$ W/MHz 1.0 V icyflex DSP/MCU Core", *ESSCIRC 2008*, Sept. 2008, Edinburgh, Scotland, U.K.
- C. Arm, J.-M. Masgonty, M. Morgan, C. Piguët, P.-D. Pfister, F. Rampogna, P. Volet; "Low-Power Quad MAC 170  $\mu$ W/MHz 1.0 V MACGIC DSP Core", *ESSCIRC 2006*, Sept. 19-22. 2006, Montreux, Switzerland
- <http://www.abiliss.com>