

## Systèmes multiprocesseurs embarqués adaptatifs

Gilles SASSATELLI

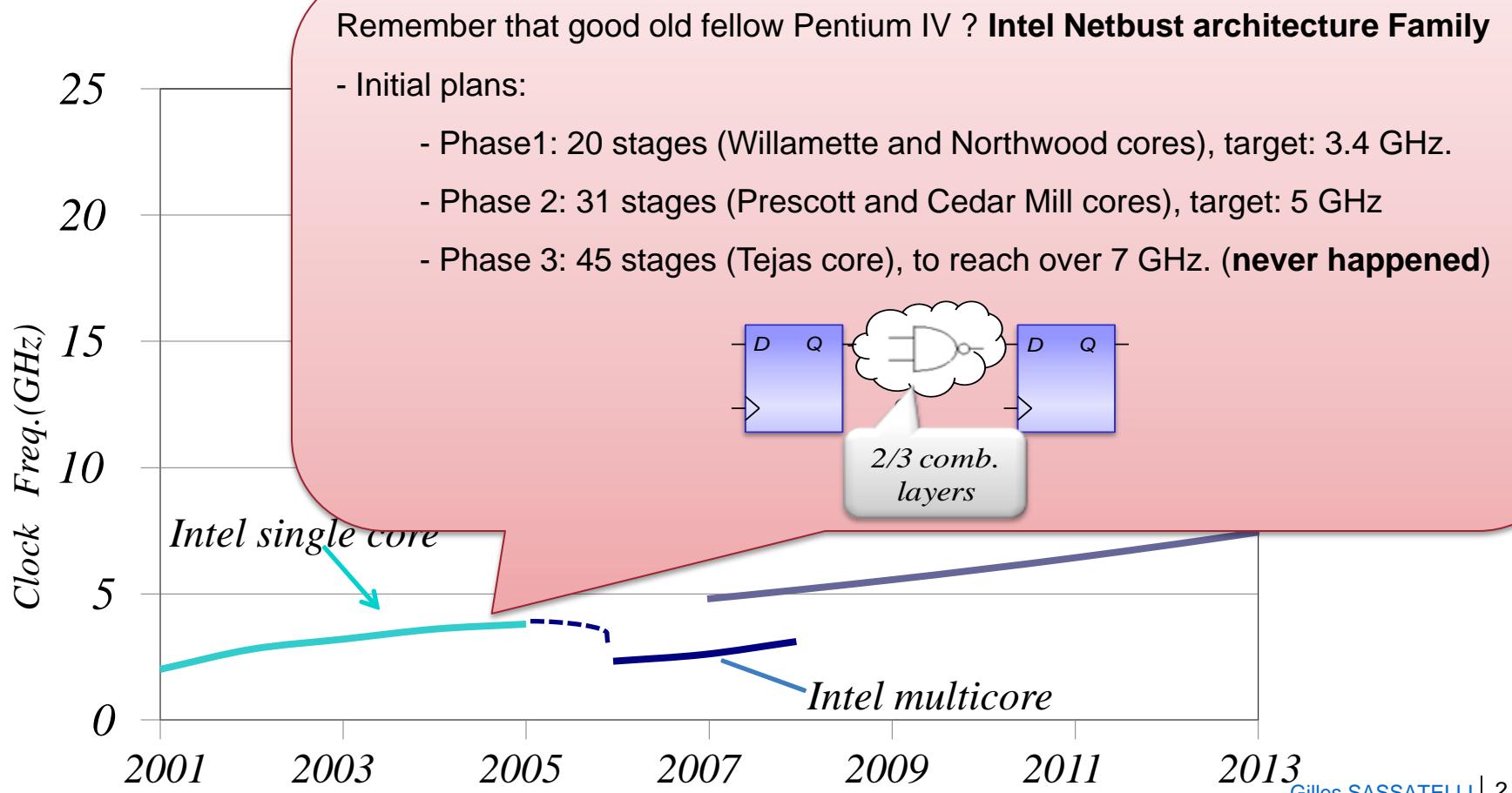
SEmba 2010

Autrans, October 18, 2010

# Further increasing performance?

## ■ Shrinking is not allowing this anymore

- Shedding light from General Purpose computing...



# Performance... not the only challenge

## ■ Hitting ILP wall, the Power wall, what are the options?

- ILP reached a plateau, frequency cannot be scaled further up...
- ...how do we utilize the additional transistors Moore's delivers every 18 month?

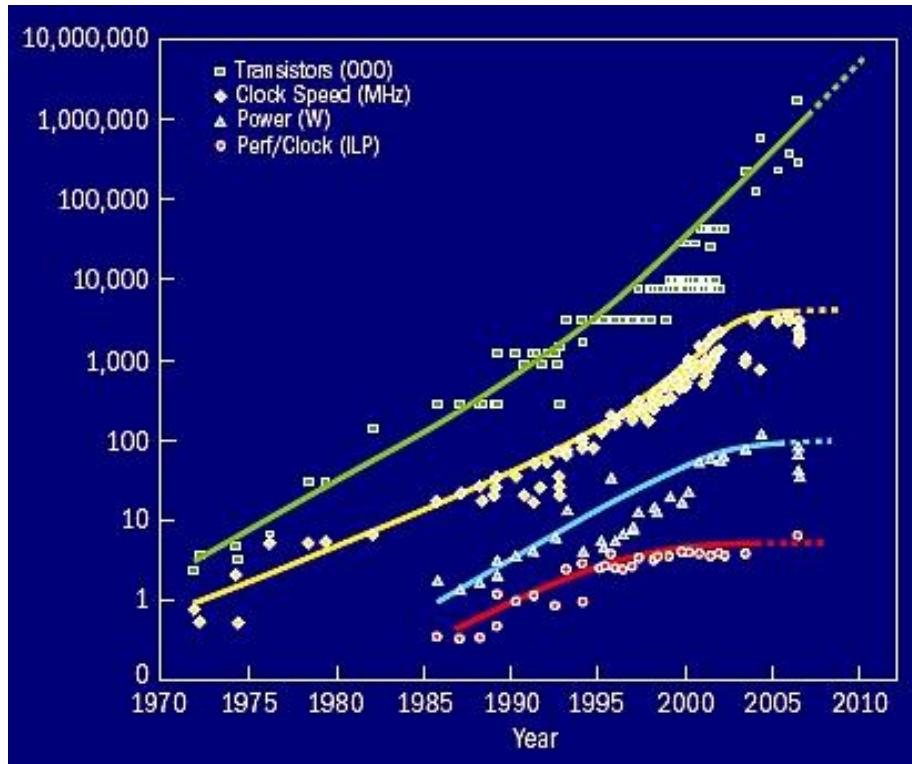
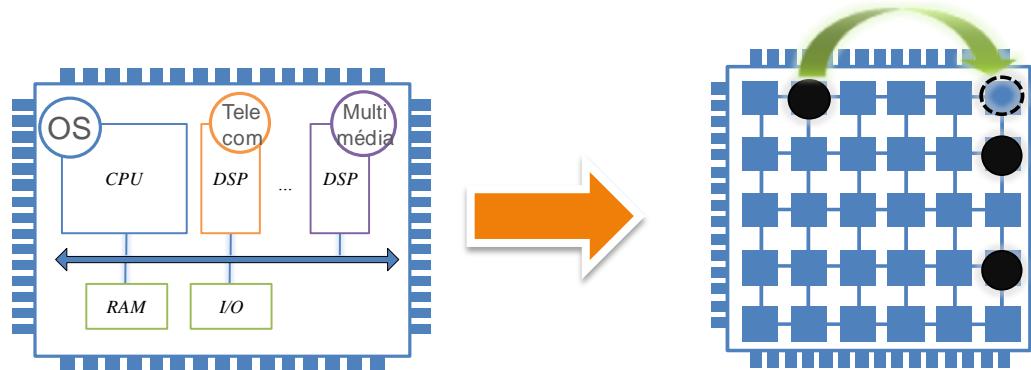


Illustration: A. Tovey Source: D. Patterson, UC-Berkeley

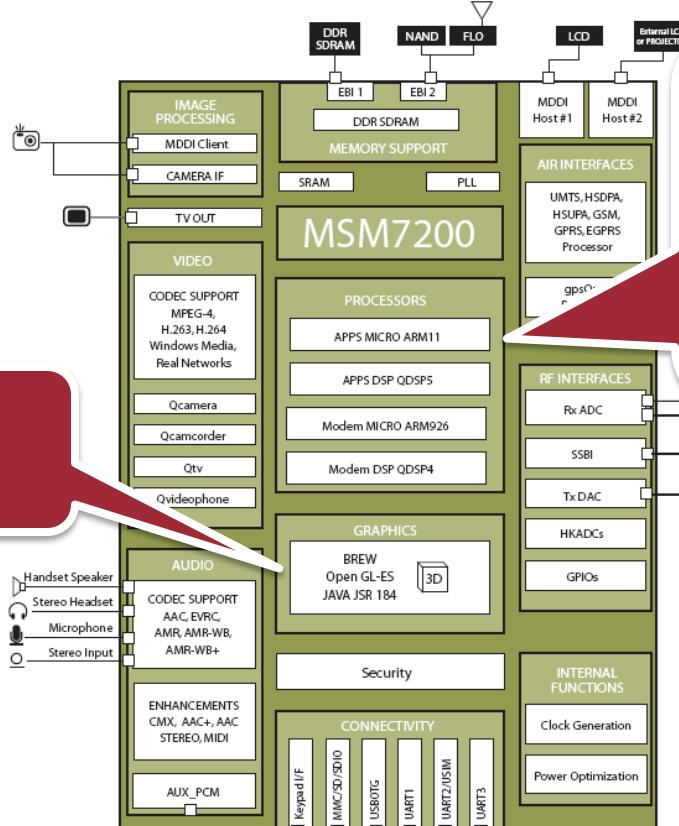
# 0- PRINCIPLES: Homogeneous, Distributed & Adaptive



# Shared memory / message passing

## Example in the Embedded domain: Qualcomm MSM7200

- MPSoC + specialized ISPs + HW accelerators (1, 2, 3)
- 4 cores heterogeneous / shared memory design + a number of accelerators / IFs
- Packs up each and every possible accelerator,
- Customer pays royalties for using optional features (3D acceleration, USB host, etc.)



2D/3D, Java Accelerators

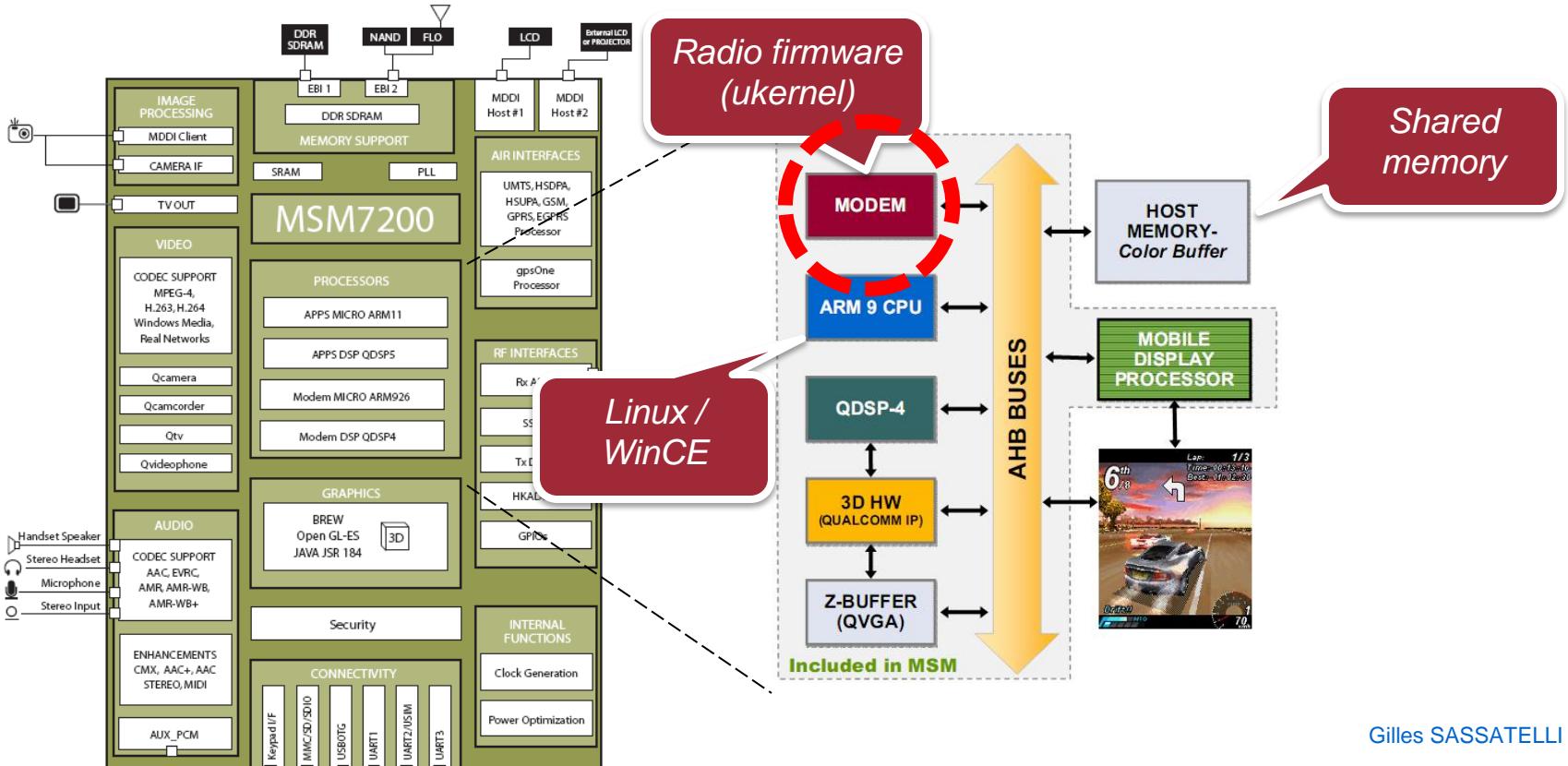
4 differentiated CPUs:

- ARM 11 (Application proc.)
- ARM 9 (modem)
- 2 DSPs (Audio + Modem)

# Shared memory / message passing

## Example in the Embedded domain: Qualcomm MSM7200

- MPSoC + specialized ISPs + HW accelerators (1, 2, 3)
- 4 cores heterogeneous / shared memory design + a number of accelerators / IFs
- Packs up each and every possible accelerator,
- Customer pays royalties for using optional features (3D acceleration, USB host, etc.)

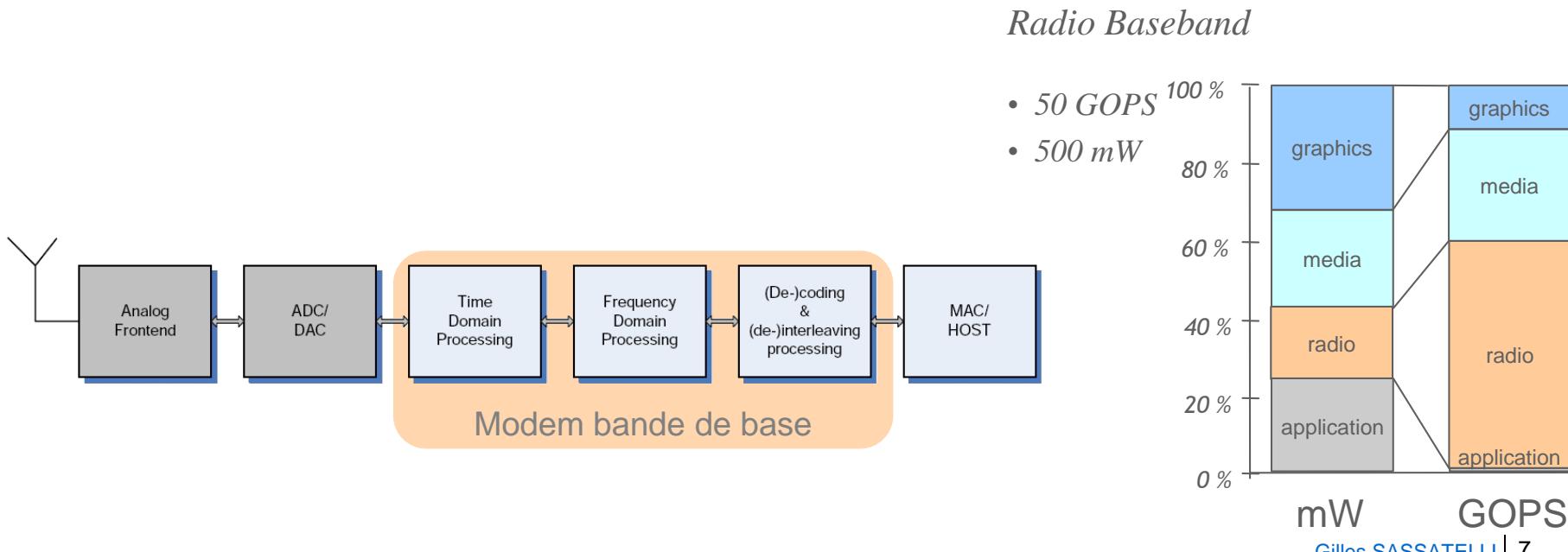


# Scaling performance: the modem

LTE: Long Term Evolution

## ■ Baseband complexity from GSM to 4G

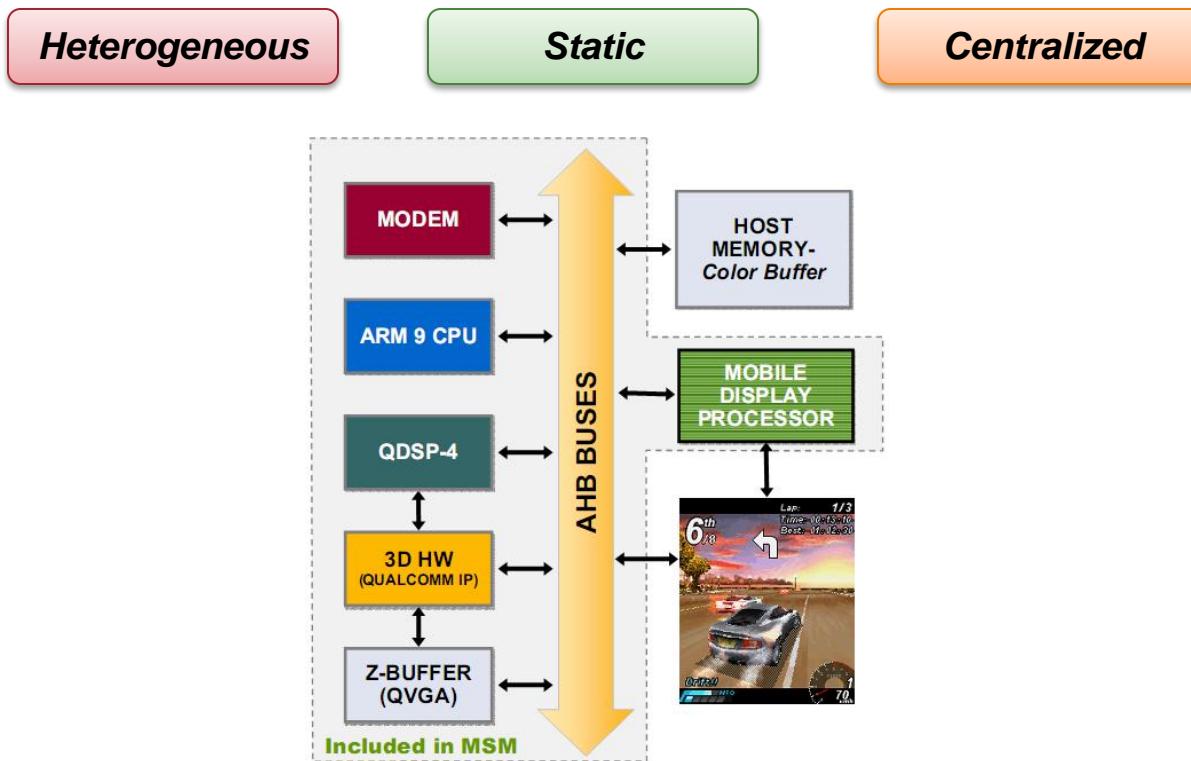
Standard	2G GSM	2.5G GPRS	3G UMTS	3.5G HSPA	4G LTE
Débit (Mb/s)	0.01	0.1	1	10	100
Complexité (MOPS)	5	50	500	5000	50000
Consommation (mW)	100	200	300	400	500
Noeud Technologique	130nm	90nm	65nm	45nm	32nm



# Scaling – what's at stake?

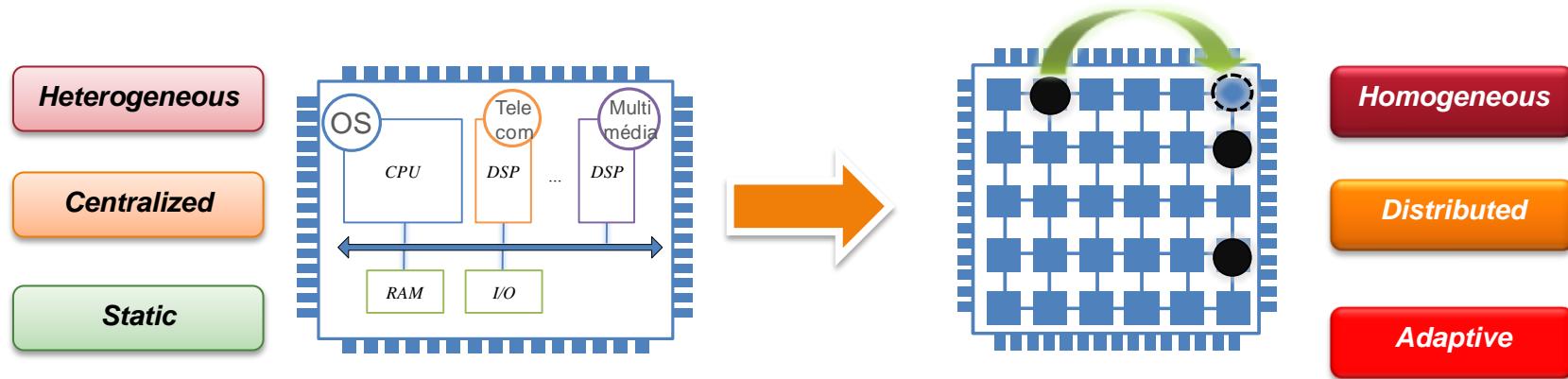
## ■ Does this template scale?

- Heterogeneity makes it difficult to increase processor count (TTM, progr.)
- And yields to decreased overall silicon usage (ex: 2D/3D seldom used)



# Scaling?

- 180 U-turn: Homogeneous, distributed & adaptive systems



# 1- AN EXPLORATORY MPSoC: HS-SCALE



**HS-Scale Tracer - LIRMM**

**Processor 0:**  
30051040 ns Nor of Instructions: 0  
30052680 ns Allocated space (task  
stack) 8384  
30343940 ns Task Object =  
30624980 ns Task Header =  
30625920 ns Task Header =  
30903400 ns Task Header =  
00000000  
30904520 ns 30974800 ns

**Processor 2:**  
out frozen  
13225600 ns Sent WD life message  
to NPU id1  
132901200 ns We do not receive  
WD message from the neighbour since  
9 ticks -> PB -> Create Diagnostic  
Neighbour thread  
133514840 ns Create new task  
(name: OSDiagnosticNeighbour), Start  
Address: 100e694  
133007760 ns DiagnosticThread ->

**Processor 1:**  
30122720 ns Nor of Instructions: 0  
30123160 ns Allocated space (task  
stack) 8384  
30413320 ns Task Object =  
30695160 ns Task Object =  
30696200 ns 30973760 ns Task Header =  
00000000  
30974800 ns 30974800 ns

**Processor 3:**  
I=5 -> tcip found a listening socket...  
64007320 ns Received Watchdog  
UDP Packet 00030000, ticks:7,  
ticksLastWD:4, delta:3  
64306520 ns Sent WD life message  
to NPU id1  
64306520 ns --> HW Inter->-Intf1  
I=5 -> tcip found a listening socket...  
99736360 ns Received Watchdog  
UDP Packet 00030000, ticks:10,  
ticksLastWD:7, delta:3  
100035560 ns Sent WD life message

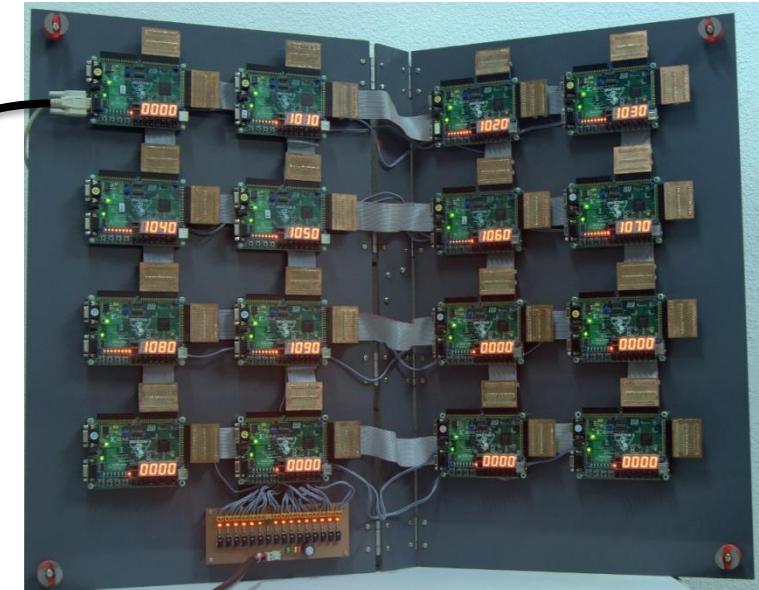
**HouseMade Debugger**

**Processor 0 Processor 1 Processor 2 Processor 3**

Value (hex)	0	Program Counter	0	Opcode
Reg 0	0	PC Next	0	Exception PC
Reg 1	0	pfr	0	Target
Reg 2	0	OP	0	
Reg 3	0	RS	0	
Reg 4	0	RT	0	
Reg 5	0	RD	0	
Reg 6	0	RE	0	
Reg 7	0	RNC	0	
Reg 8	0	Imm	0	
Reg 9	0	Status	0	
Reg 10	0	Exception Id	0	
	SLL: r1 r2 r3< t1>-rc (NOP)	INFO Status	0	
	0 0 r2 0	WakeUp	0	

Current Time: 40 ns  
cycles

Hardware Break Point:  
 Enable  
 Processor 0  
 Hex  
 Dec  
 PC = 0 Run for Run to 100000000  
 ns cycles



## Concepts:

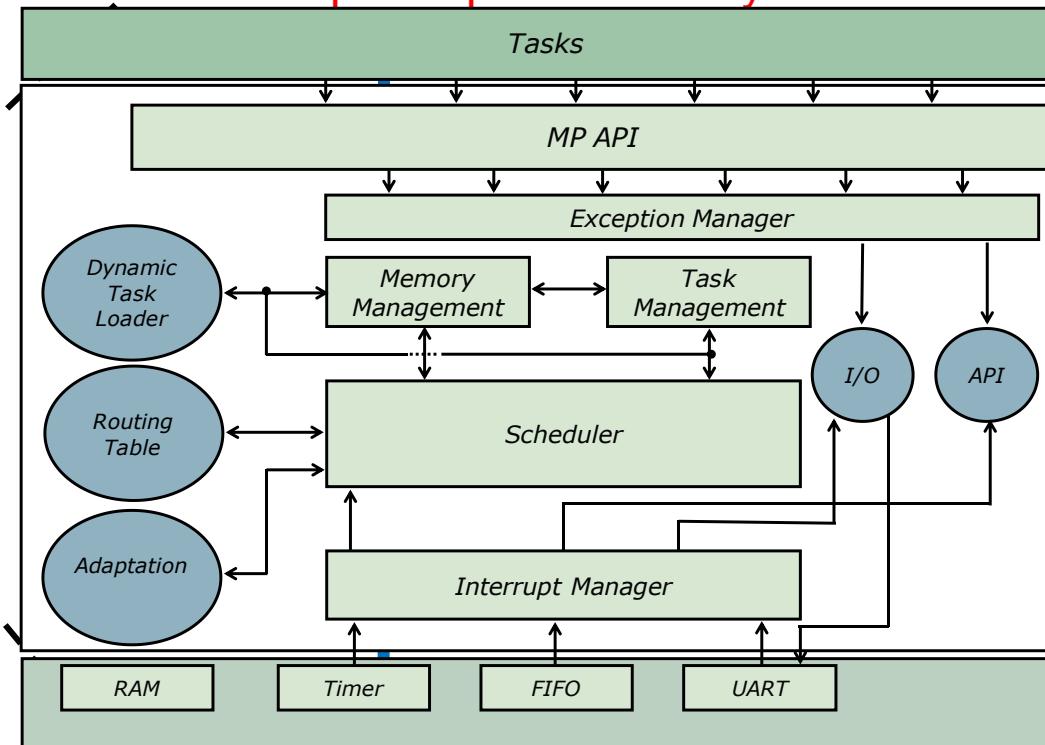
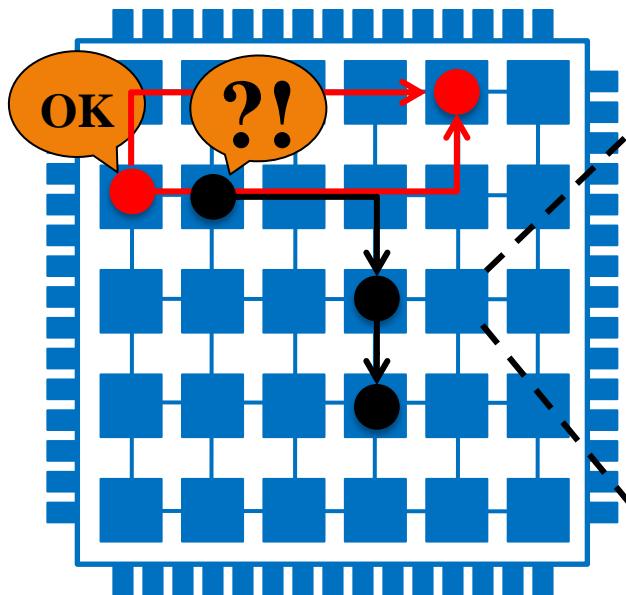
- A grid of tiny networked (Hermes XY NoC) PEs
- PE made of simplest RISC + RAM +  $\mu$ K + NI → Private memories, message passing
- Distributed decision-making, in SW

*Homogeneous*

*Distributed*

*Adaptive*

20-30kB preemptive uK w. dynamic loader

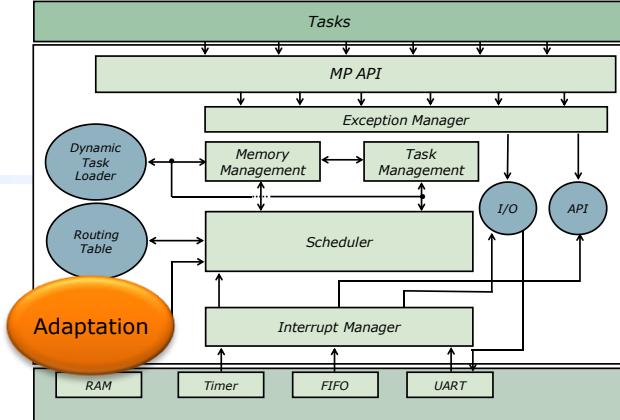
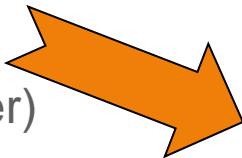


# Adaptation

*Adaptive*

## ■ What are the main levers?

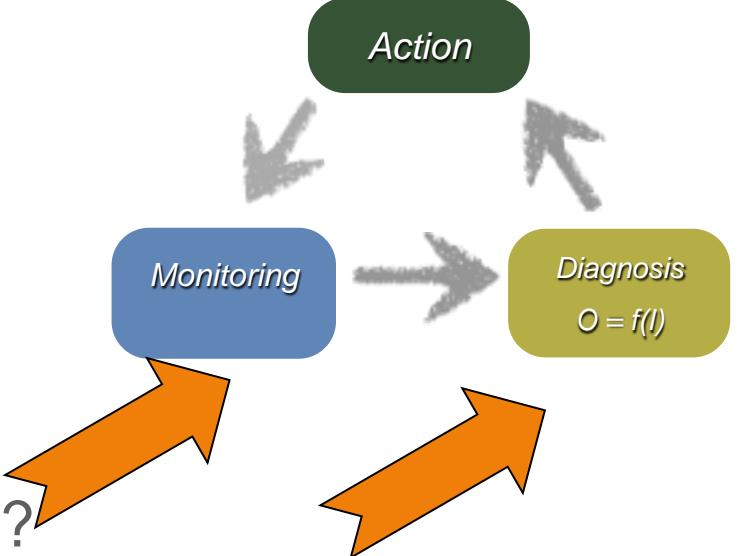
- Voltage/Frequency scaling (perf/power)
- Route reconfiguration
- Rescheduling
- Task migration
- Algorithmic reconfigurations
- ...



*Action*

## ■ And what about decision making?

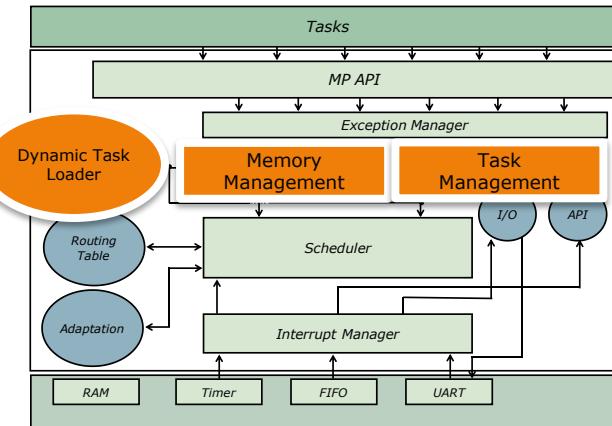
- Scenario-driven
- monitoring-driven
- Iterative
- Heuristic
- ...



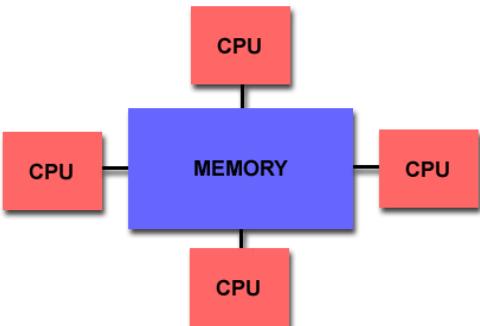
# Task migration

## Task migration incurs a cost..

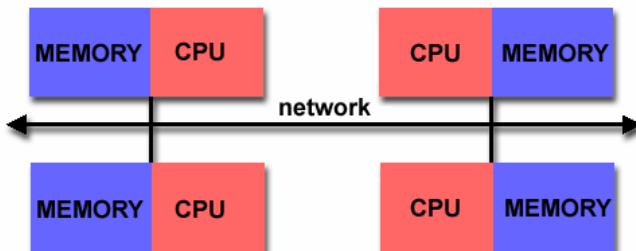
- Contrarily to shared / coherent memory MPSoCs, task code and context have to be migrated physically..
- Usually requires MMU HW / dynamic loader in SW



Shared memory

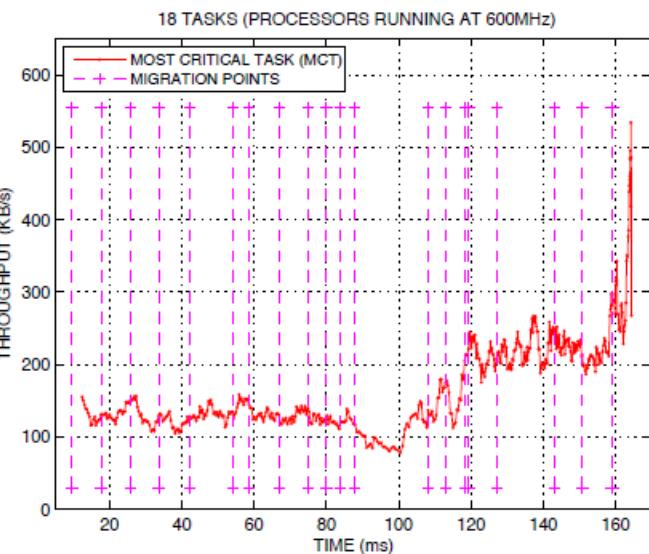
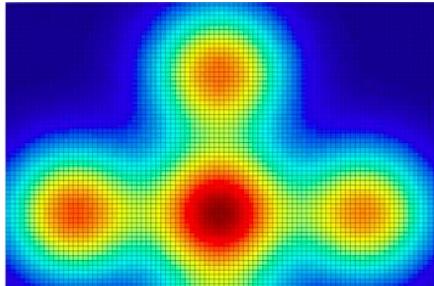


Private memories

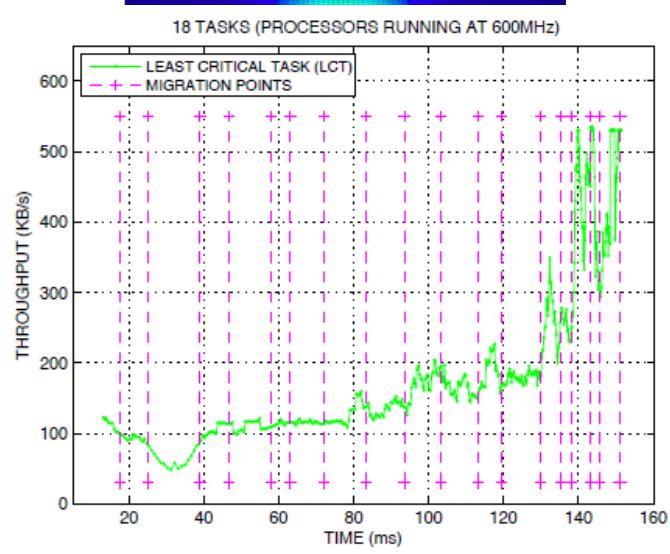
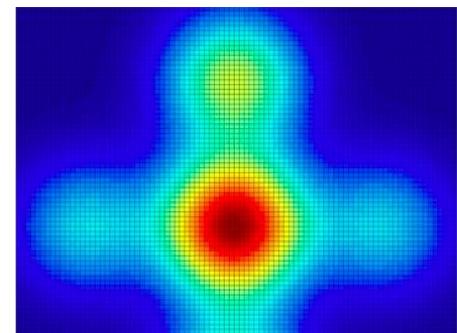


## ■ Neighbor load balancing « hot / cold potatoe »

*Most Critical Task*



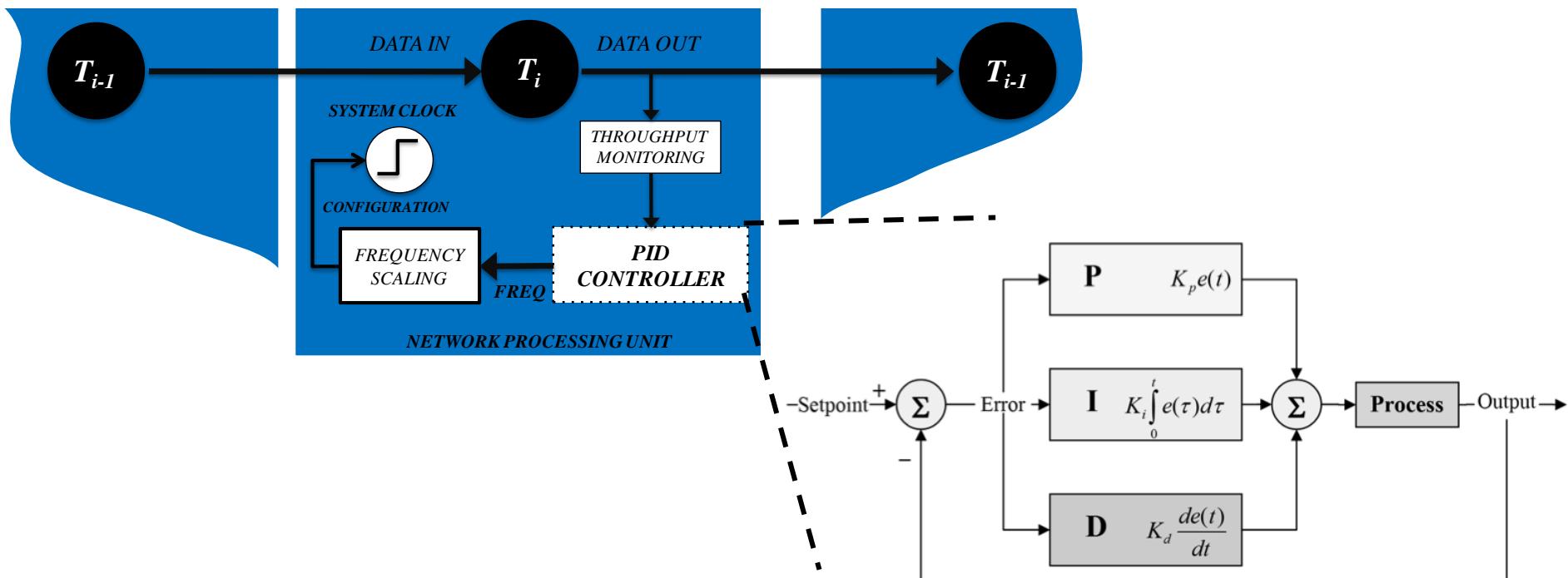
*Least Critical Task*



# Example decision making strategy

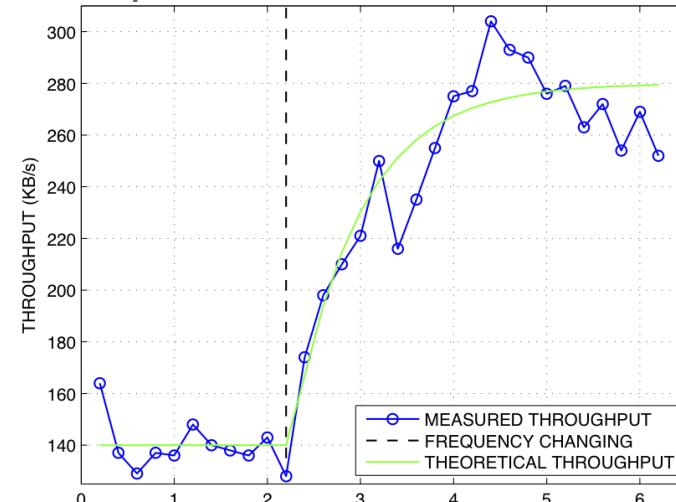
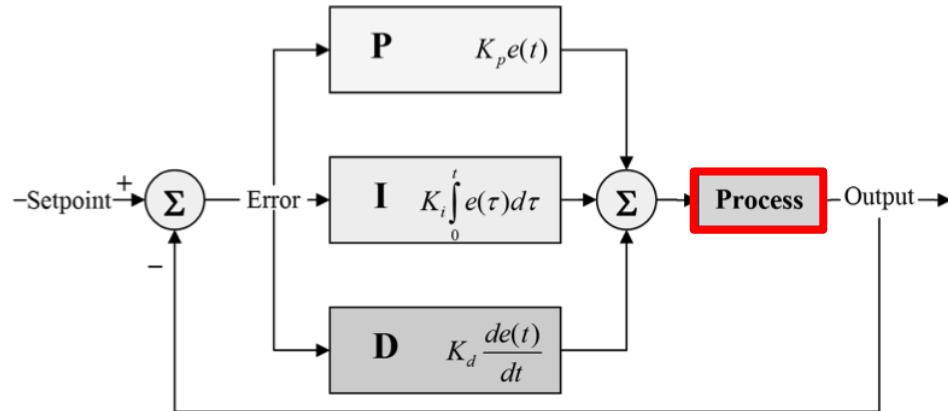
## ■ Use of well-known feedback control-loop techniques..

- Frequency scaling for ‘boosting’ performance whenever necessary (task migration) and revert to the lowest possible otherwise (saving power)
- All of this at task-level, overlaid network of PID controllers



## Strategy:

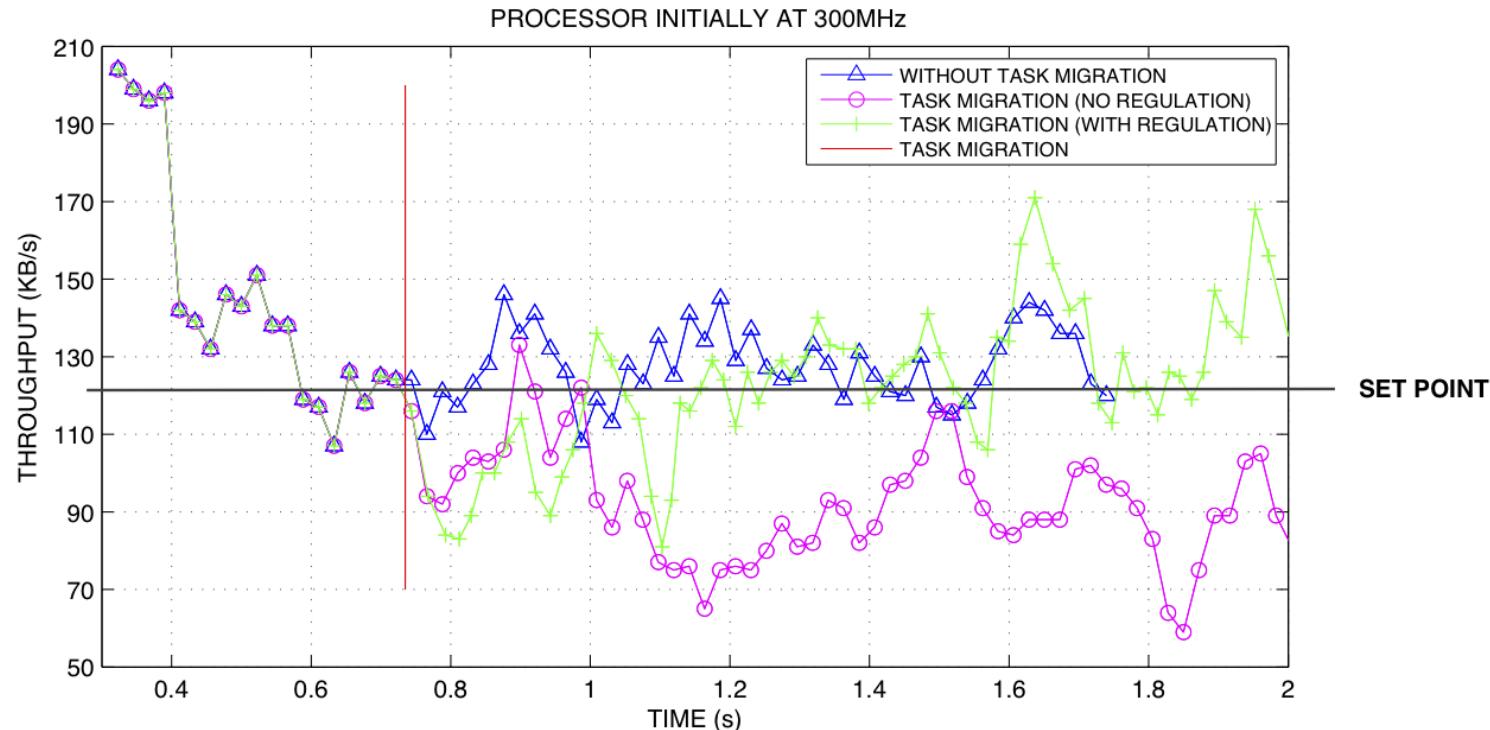
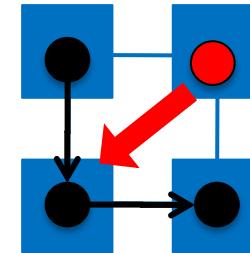
- 1- model the process when exposed to a step



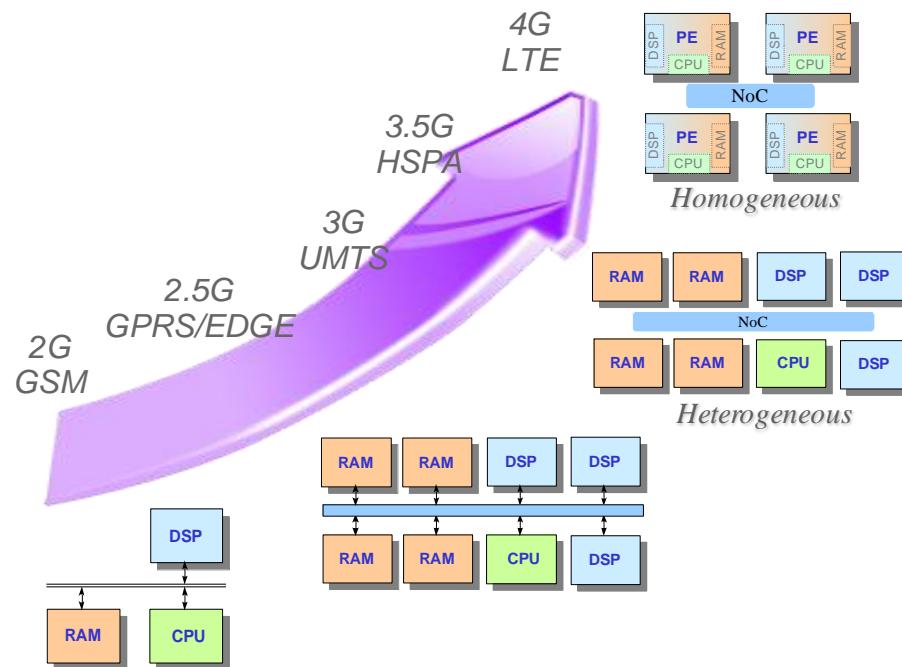
# Use of PID controllers

- And then see how the system behaves when undergoing perturbations

- Audio/video decoder (MJPEG / ADPCM): 7 tasks
- Here an incoming task (task migration)



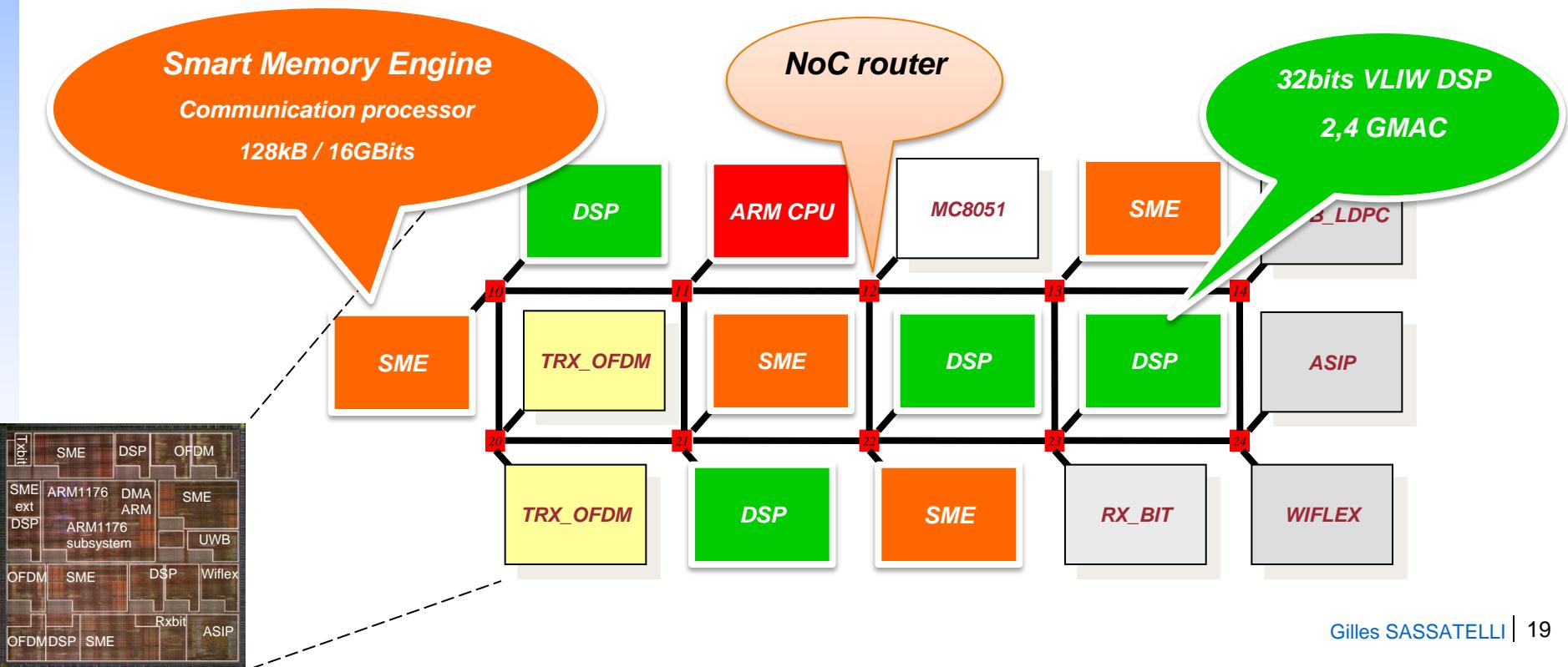
## II- CASE STUDY: GENEPY, a MPSoC for baseband processing



# State of the art

## MAGALI cea leti

- Tuned for 3GPP-LTE telecom, 65nm ST CMOS process
- Powerful energy efficient solution

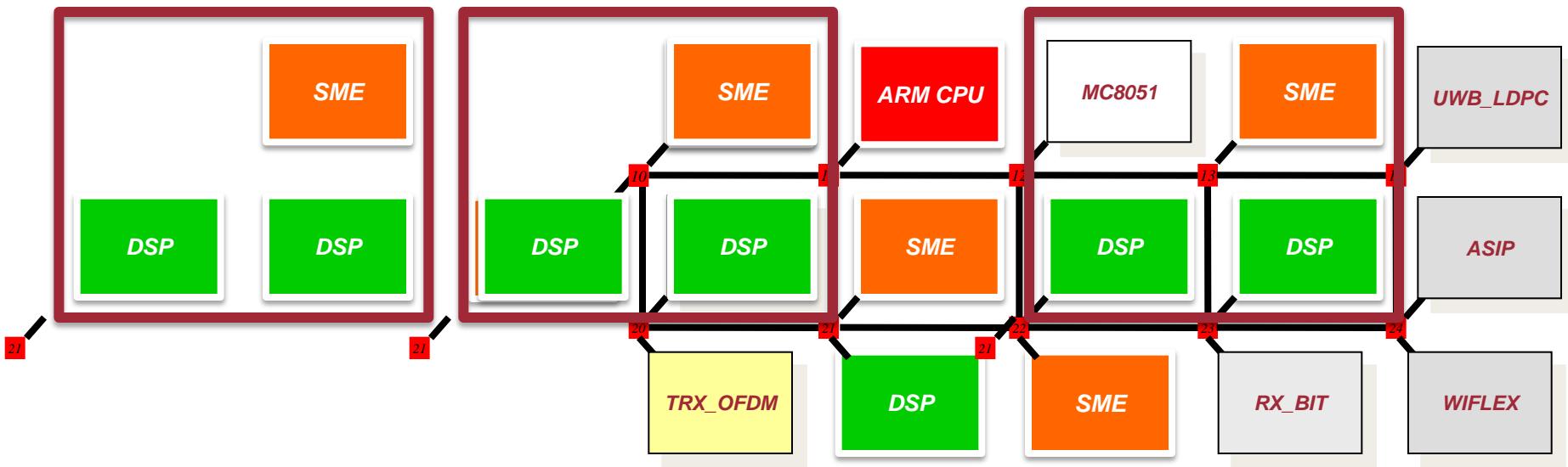


# Going homogeneous

## Keep the ingredients, change the recipe:

Homogeneous

- Define a tile that embeds 2 DSPs and a SME
- Keep the ARM CPU for control purposes
- And instantiate as many tiles as necessary



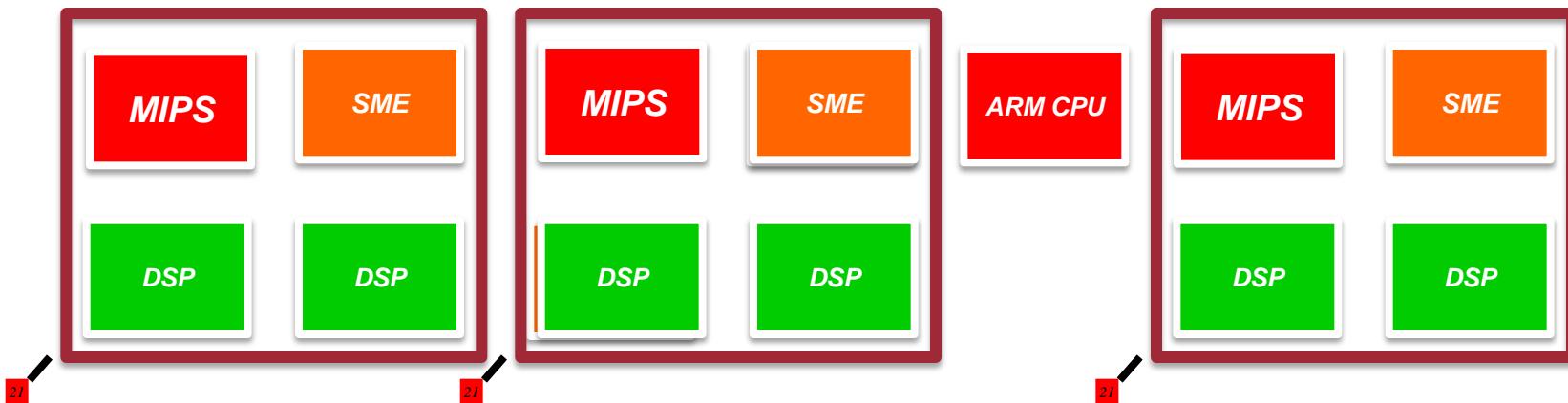
# Going homogeneous, distributed control

## Distributed control?

- Go a little further, remove the ARM and use a local control CPU

Homogeneous

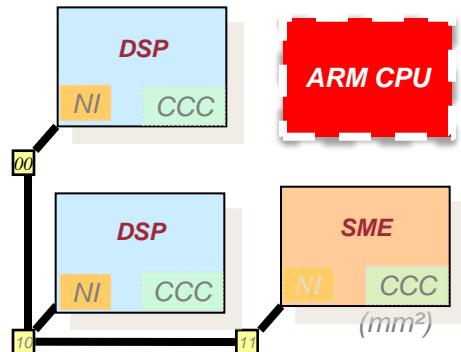
Distributed



# Area

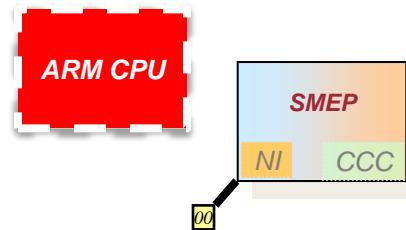
## ■ ST65nm - 400MHz

*Centralized control*



TOTAL : 2,768

*Semi-distributed control*



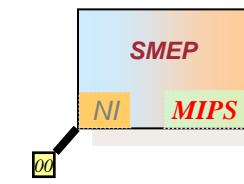
TOTAL : 2,303

*Homogeneous*

*Homogeneous*

*Distributed*

*Fully distributed control*



TOTAL : 2,392

(mm<sup>2</sup>)

(mm<sup>2</sup>)

(mm<sup>2</sup>)

Routeur : 0,159

Routeur : 0,159

2 DSPs : 0,792

2 DSPs : 0,816

SME : 1,226

SME : 1,216

Mover : 0,033

Mover : 0,033

NI + CCC : 0,126

NI + MIPS : 0,201

TOTAL : 2,303

TOTAL : 2,392

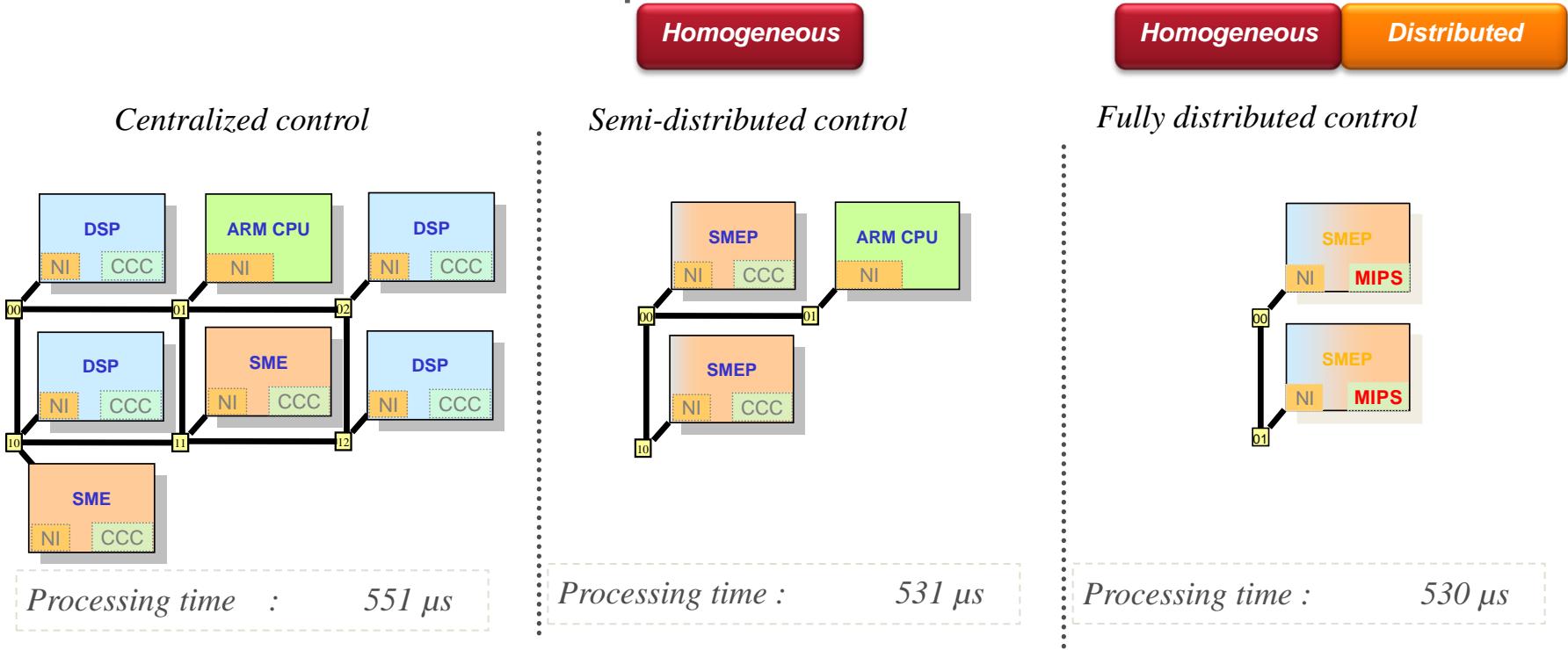
- 17%

- 14%



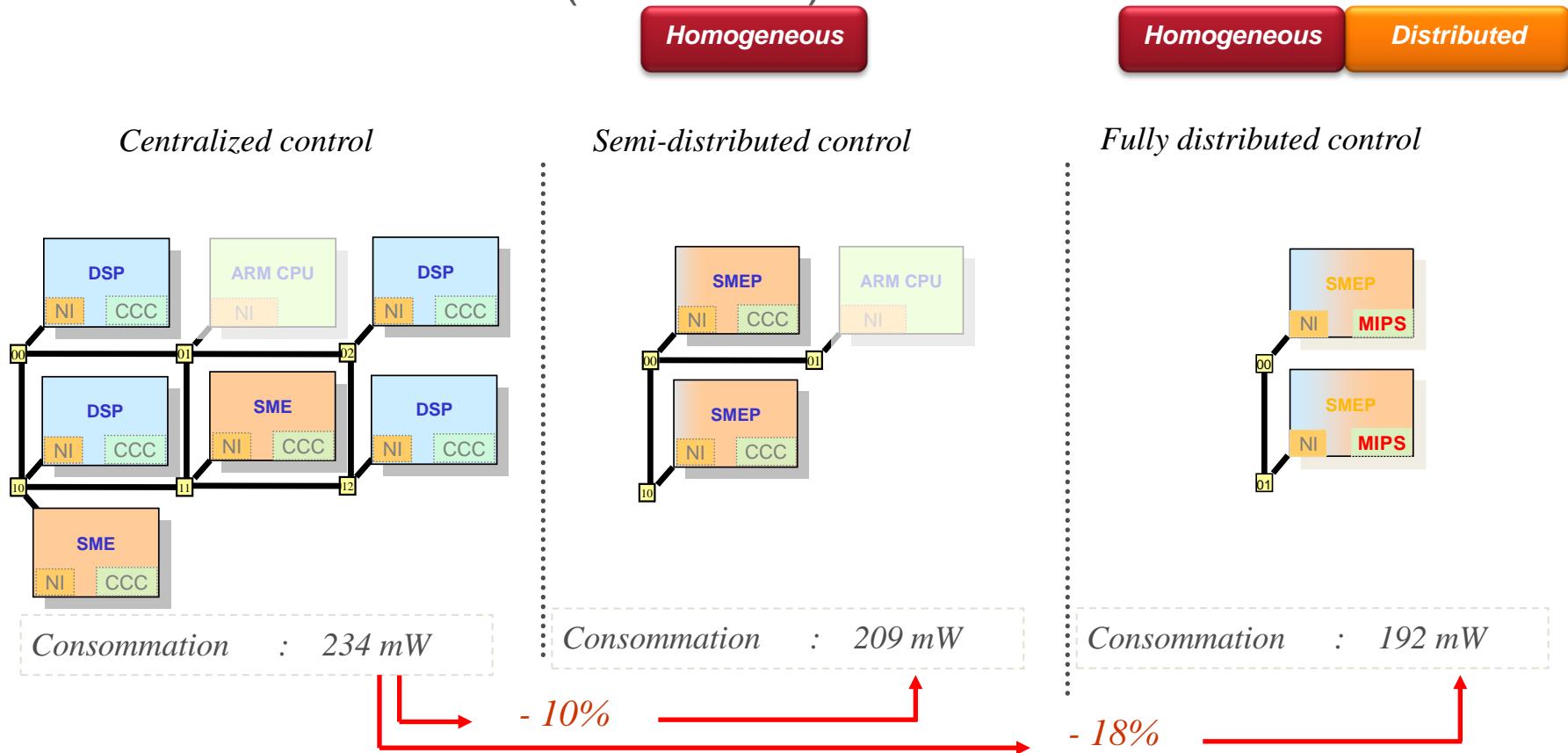
# Performance

## ■ 400MHz – 1 frame computation



# Power

## Post P&R extraction (w/o ARM)



# Task mapping

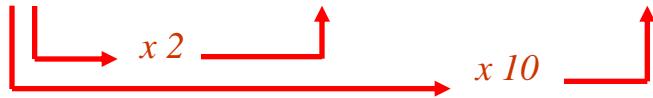
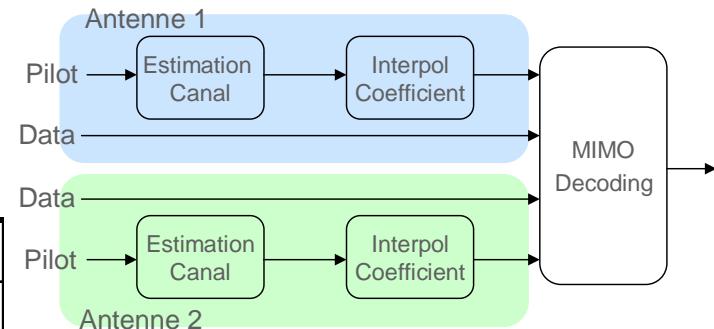
**Homogeneous**

**Distributed**

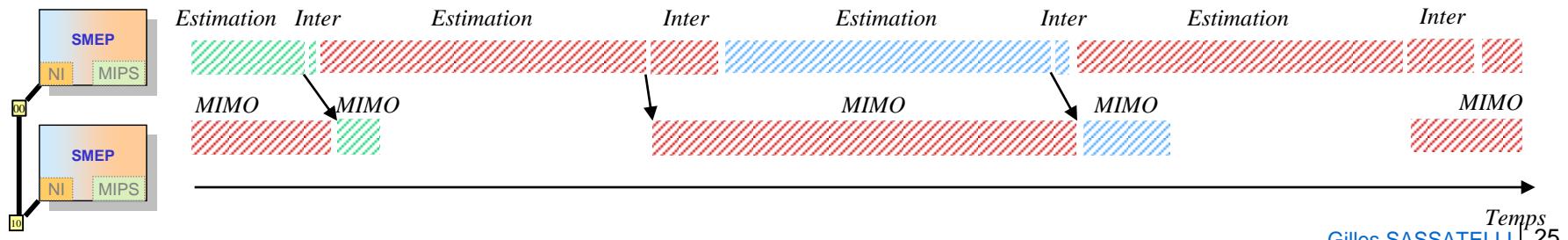
- 3GPP-LTE features adaptive modulation

- Frames of multiple modes have to be supported,
- With the associated processing overhead:

	Mode 1	Mode 3	Mode 5
Estimation	91000 cycles	91000 cycles	91000 cycles
Interpolation	1750 cycles	3500 cycles	17500 cycles
MIMO decoding	12000 cycles	24000 cycles	120000 cycles

## Static mapping



# And adaptive mapping

Homogeneous

Distributed

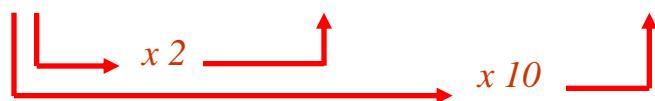
Adaptive

- 3GPP-LTE features adaptive modulation

- Frames of multiple modes have to be supported,
- With the associated processing overhead:

Up to another 7% power savings

	Mode 1	Mode 3	Mode 5
Estimation	91000 cycles	91000 cycles	91000 cycles
Interpolation	1750 cycles	3500 cycles	17500 cycles
MIMO decoding	12000 cycles	24000 cycles	120000 cycles

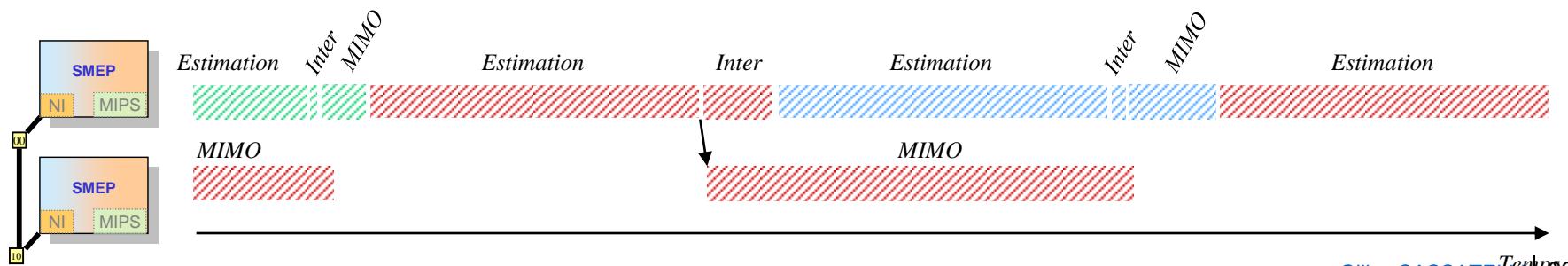


Adaptive mapping of

MIMO decoding :

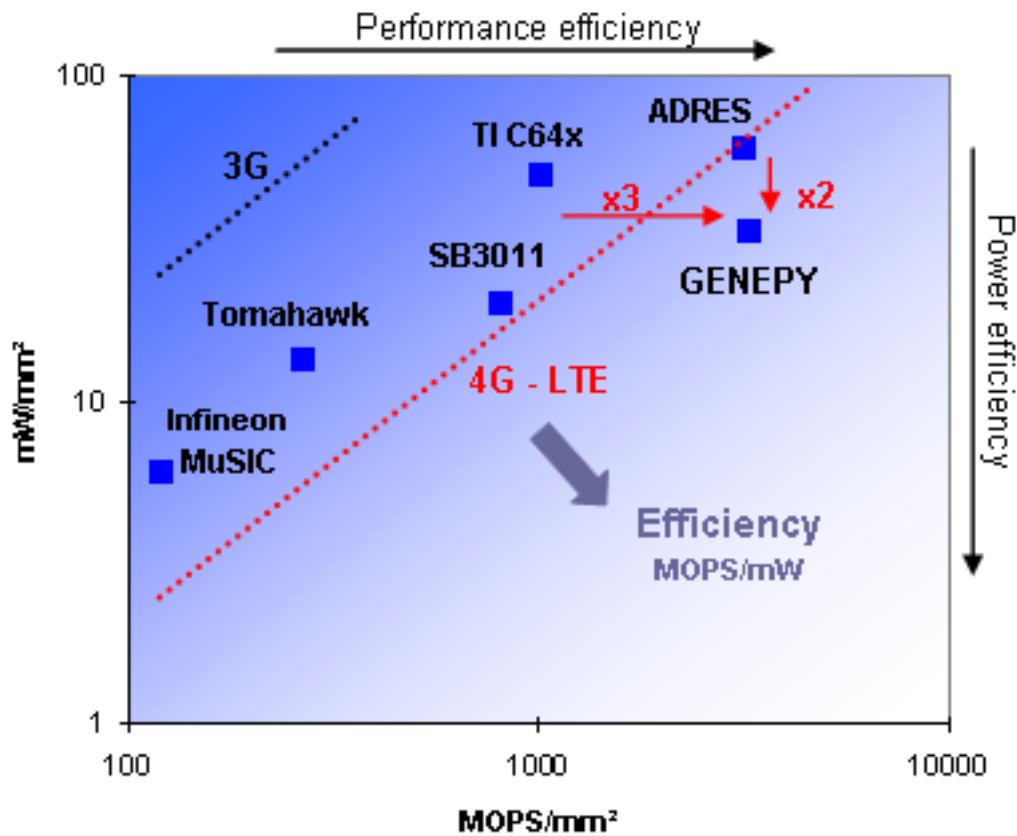
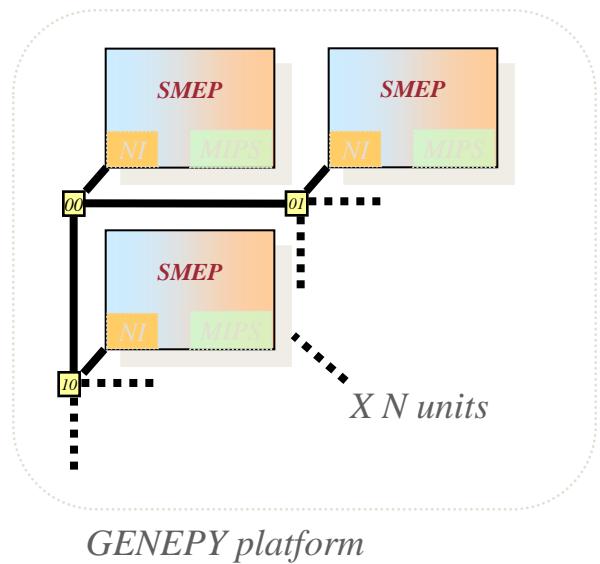
- Mode 1,3 : SMEP 1
- Mode 5 : SMEP 2

### Adaptive mapping



# Vs. State of the art

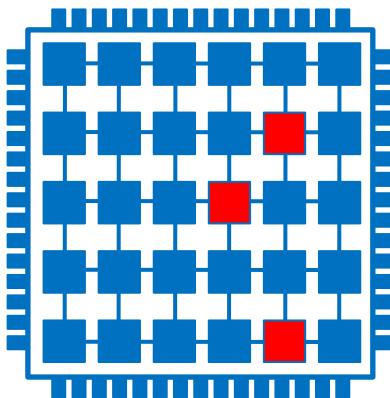
- Compared to other MIMO capable devices



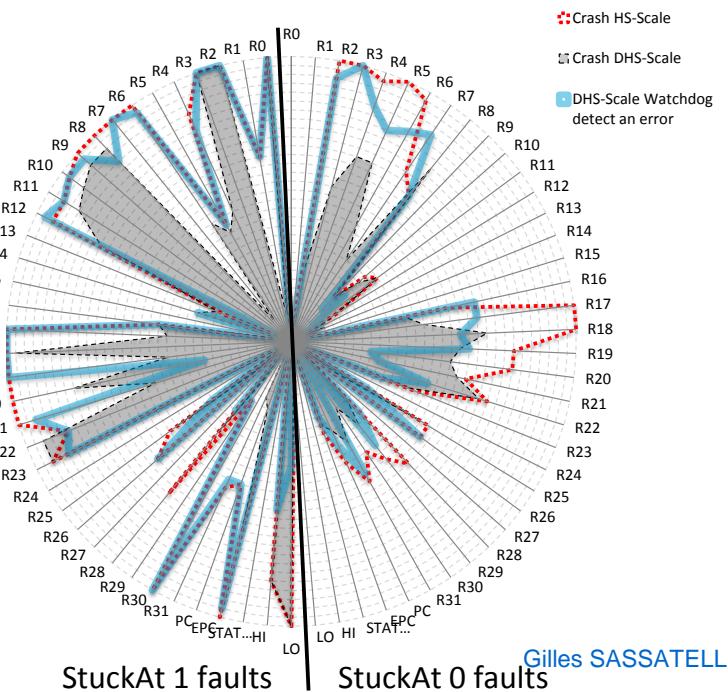
Results normalized against 65nm – 1.2V :

Infineon MuSiC, Tomahawk, SandBlaster SB3011, TI C64x, IMEC ADRES

### III- PUTTING HOMOGENEITY TO USE: DScale, a reliable exploration platform

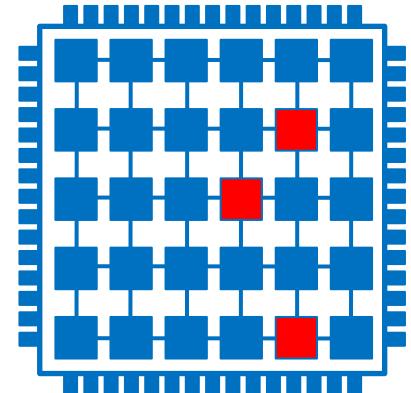


Crash repartition due to StuckAt faults on NPU(1,1)

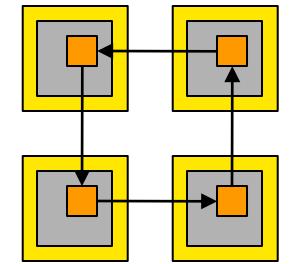


# Last one: improving reliability

- Is there a means to exploit the natural redundancy of such systems?
  
- Targeted Applications
  - Typical non life critical (smartphones, Pay TV, etc.)
  
- Performance constraints
  - Low area / power overhead
  - Transient loss of availability / functionality / data acceptable
  
- Target crash faults



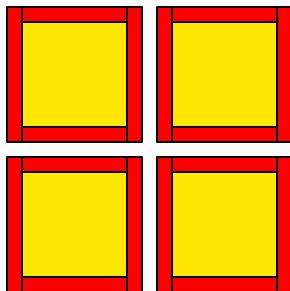
# Implementing recovery



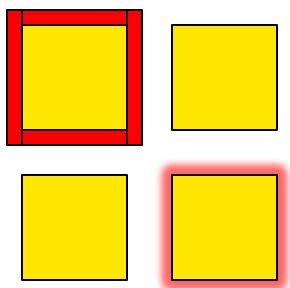
*Fault Detection*



*System-wide RESET*



*Error diagnostic upon reboot*

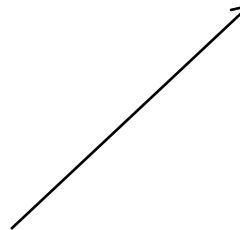
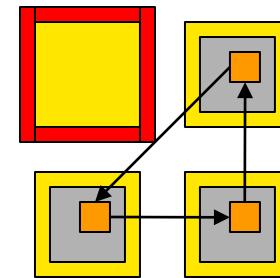


*Local defect isolation*

*or*

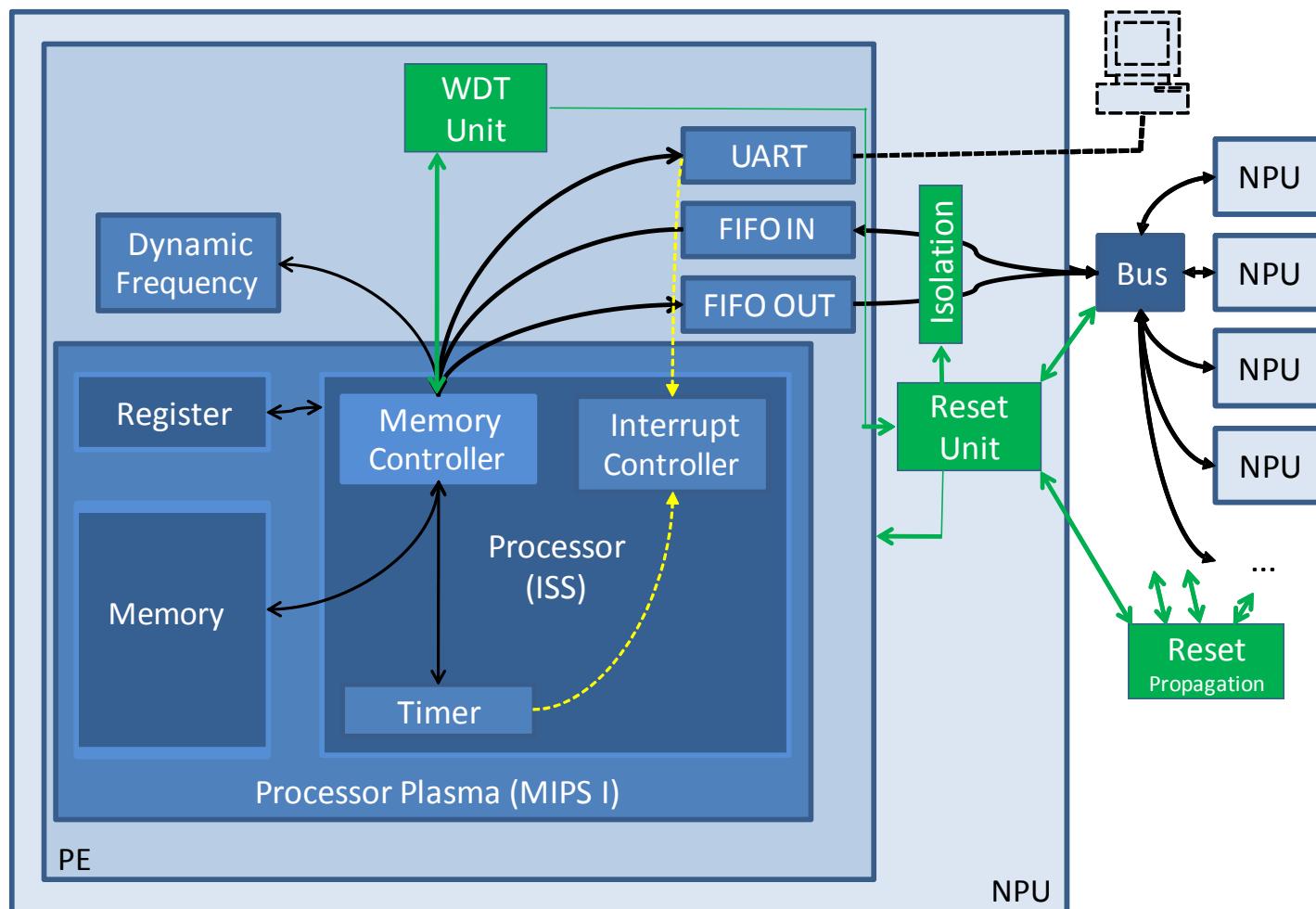
*Local process compensation*

*Application remapping*



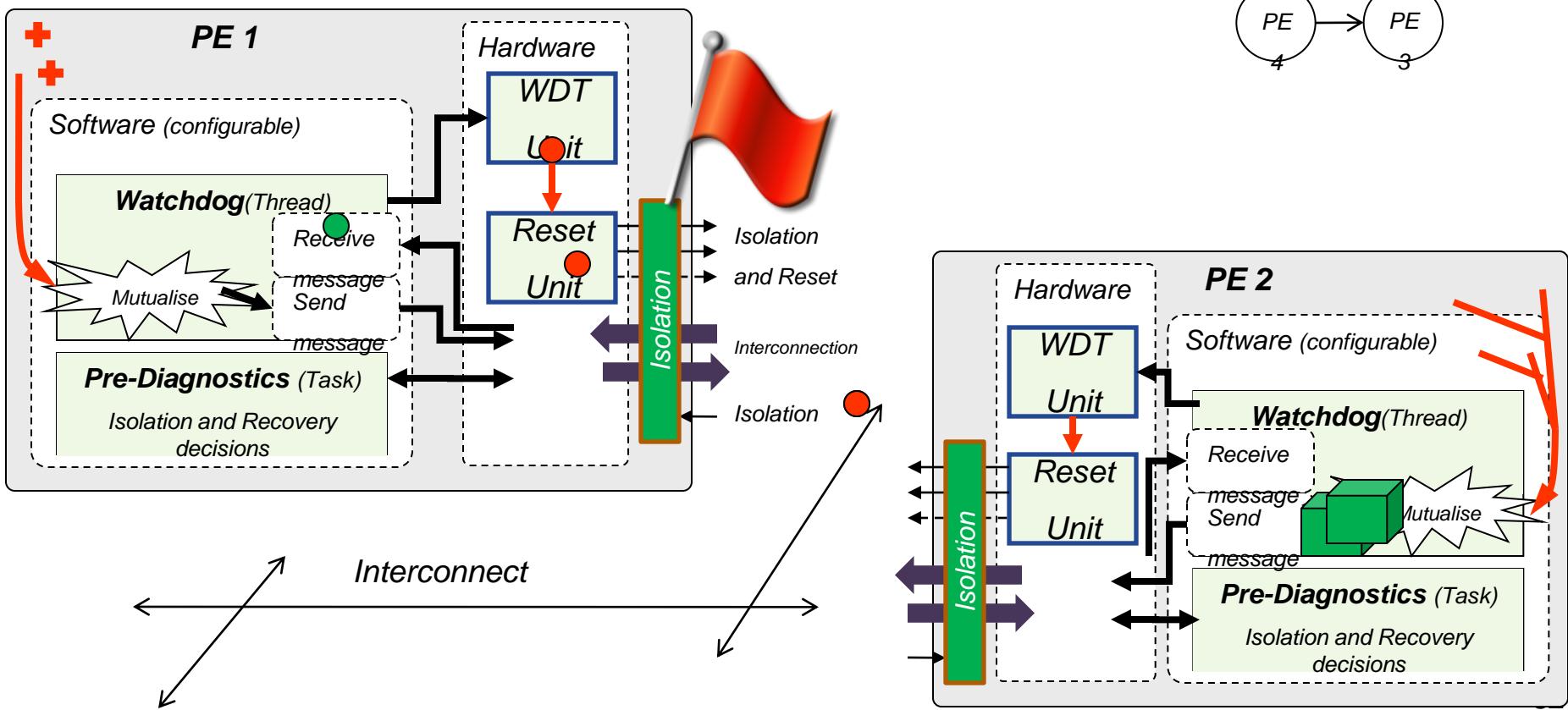
# Implementing recovery

## ■ Additional hardware



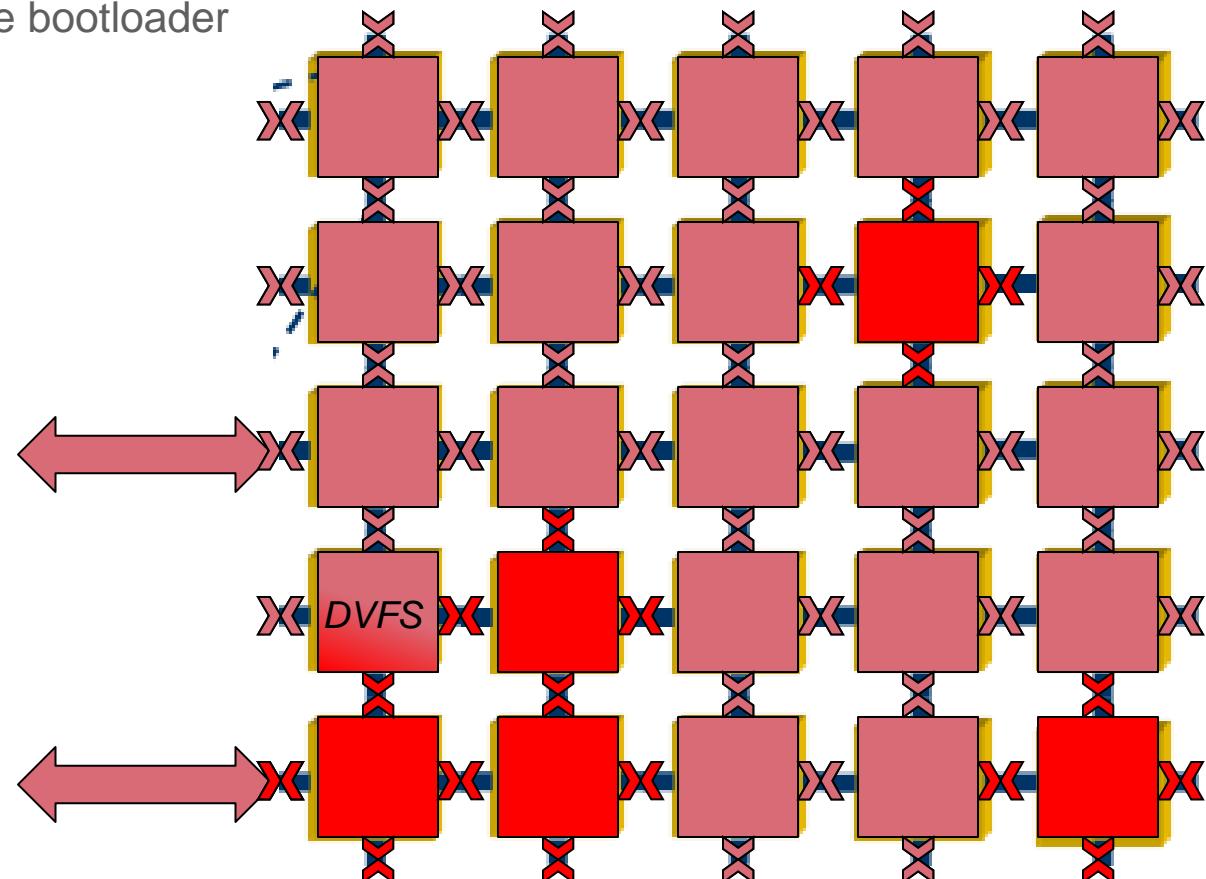
# Implementing recovery

## Functionning



# Implementing recovery

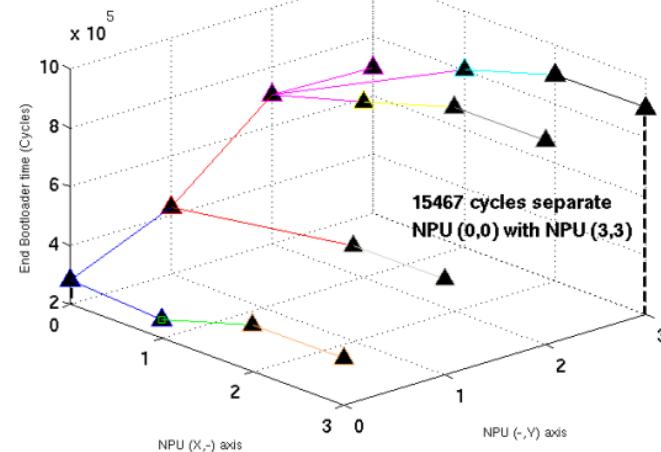
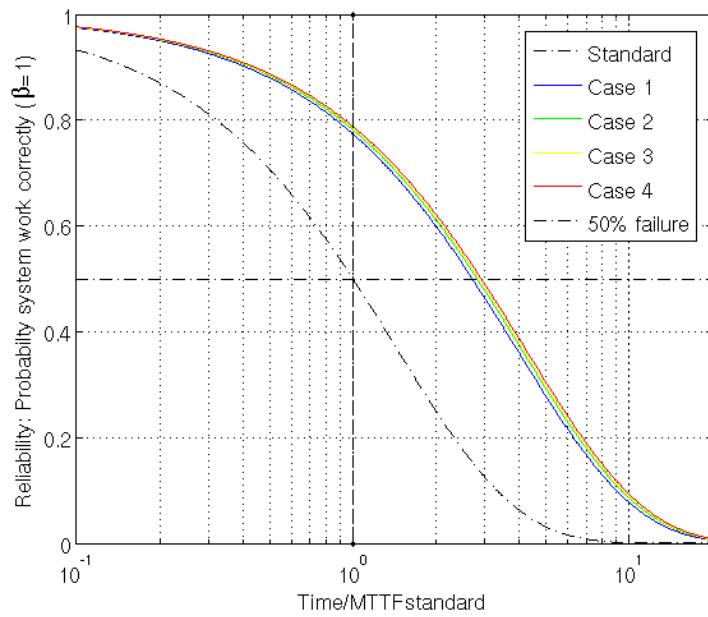
- Depending on error type and frequency, the distributed boot loader can:
  - Do fine tuning (DVFS)
  - Maintain the defective element isolated
  - Performed by the bootloader



# How reliable?

## Accepting a performance overhead

- Less PEs, others take over computation
- MTTF improvements limited by #interface PEs (case 1-4) and interconnect



# Conclusion

- Shifting part of the application handling at run-time (vs. design-time) is a promising alternative: Adaptation
- Combined with homogeneity, a means to fill the performance gap, to an extent. Further, allows for facilitated design, reliability strategies
- Finally control has to be distributed (hierarchical) so that performance scales