

# Exploiting heterogeneity for greater energy efficiency of SoCs

**From modeling to power management**

*FETCH 2020*

*Sébastien Bilavarn*

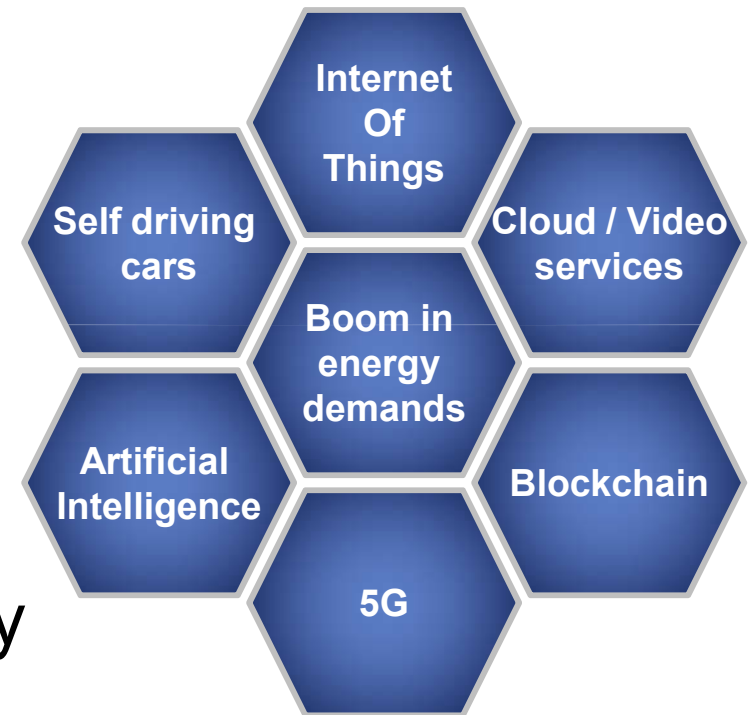
- Heterogeneity and energy efficiency
- Problem decomposition
  - **System level modeling and design**
  - **Power and resource management**
  - **Elements of energy efficiency improvements**
- Application study
- Conclusion and perspectives

# Technological challenges

- Driven by critical energy efficiency requirements

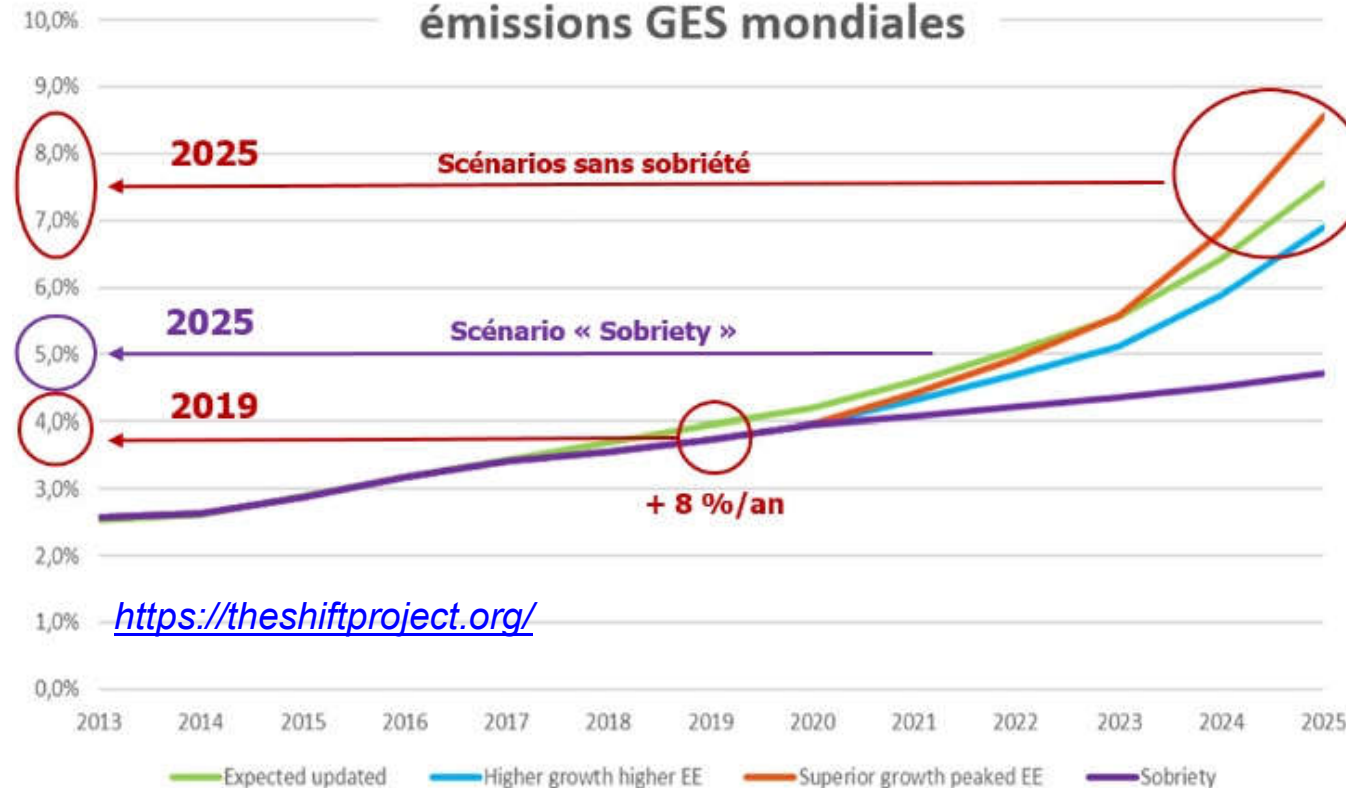
- Rise of connected devices
- More communications (5G, ...)
- Bigger infrastructure development (cloud, datacenters, VoD, ...)
- Demanding application domains (blockchain, AI, automotive, ...)

→ Boom in energy demands for many application domains



# Technological challenges

## Part du Numérique dans les émissions GES mondiales



Accelerated global warming:  
0,2 °C since 2011 – 2015

In 2019, almost 4% of world  
greenhouse gas (GES)  
emissions are due to the  
production and usage of  
digital technology.

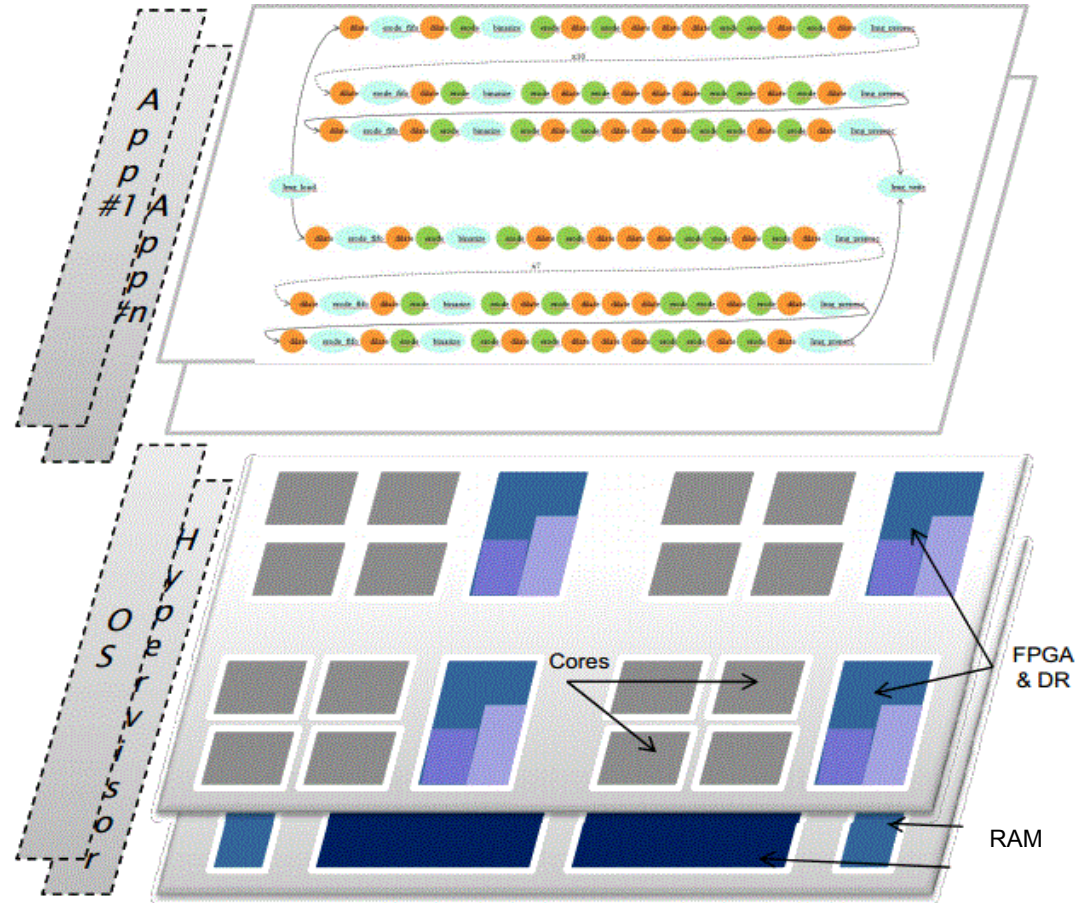
This is even more than the  
usual 2% attributed to civil air  
transport.

With a continuous increase of around 8% per year, this share could double in 2025 and reach 8% – the current level of emissions for cars and two-wheeled vehicles presently

→ Significant energy efficiency improvements are needed (>> x10)

# Problem statement

- General SoC / application mapping overview
  - Heterogeneity
    - big.LITTLE, DynamIQ
  - Complexity
  - Abstraction
- Determining factors of efficiency
  - System properly designed
    - System level modeling and design
  - Manage execution of the system
    - Scheduling, power and resource management



# System level modeling and design

- Power modeling at system level
  - Homogeneous multiprocessor platforms
- Allowing higher degrees of abstraction
  - Load balancing leads to characteristic power values
  - Promote homogeneous multithreading
    - Simple power estimation

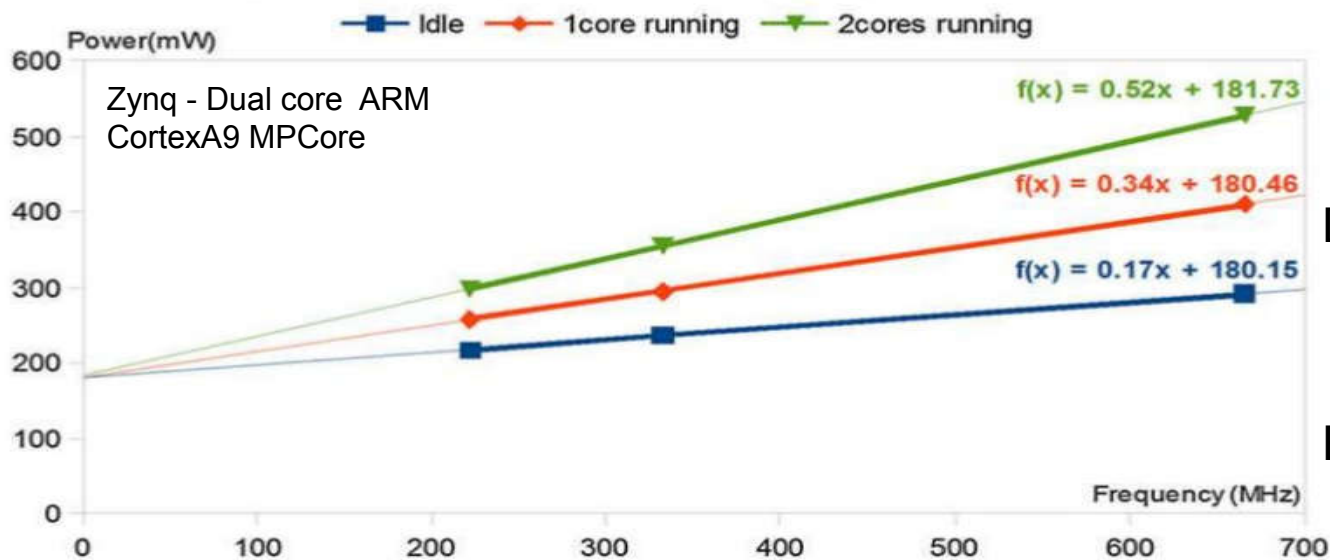
Parallel H.264 decoder

Quad core ARM11 MPCore	8 threads	4 threads	2 threads	1 thread	0 thread
<b>1 core</b>	633 mW	634 mW	636 mW	636 mW	361 mW
<b>2 cores</b>	885 mW	878 mW	880 mW	637 mW	364 mW
<b>3 cores</b>	1034 mW	932 mW	882 mW	637 mW	363 mW
<b>4 cores</b>	1256 mW	1187 mW	867 mW	642 mW	355 mW

ANR PHERMA (2007 – 2010)

# System level modeling and design

- Power modeling at system level
  - Homogeneous multiprocessor platforms / DVFS
- Allowing higher degrees of abstraction
  - $P = f(N_{\text{cores}}, F)$



Reliable energy estimation

- At function level
- From limited, meaningful parameters

Design space exploration

- Function mapping
- Parallelism

Extensions

- Multi/many core
- Core heterogeneity

→ Realistic power analysis of full application mapping on multiprocessor



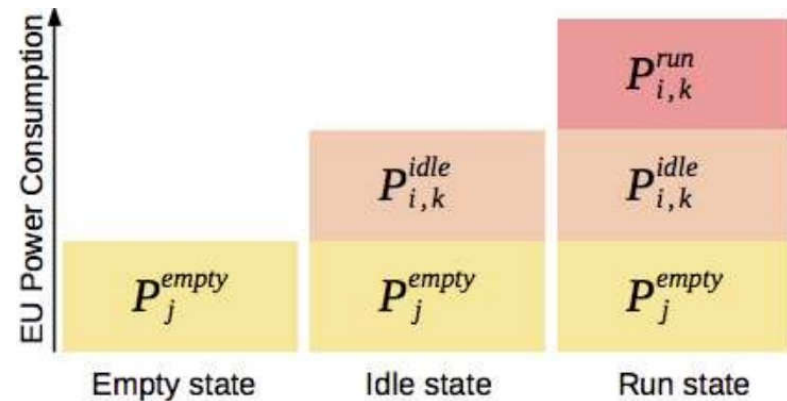
# System level modeling and design

## ■ Power modeling at system level

- Heterogeneous multiprocessor platforms – Reconfigurable acceleration

## ■ FPGA energy model

- At function level (HLS)



- Allow full Hw/Sw functions and global application power characterization

Function (i)	Execution Unit (j)	$T_{ij}, ms$	$P_{i,j}^{idle}/P_{i,j}^{run}, mW$	$N_{cell}; N_{bram}; N_{dsp}$
<i>Img_load</i> (i = 1)				
<i>dilate</i> (i = 2)	Core (j = [1; 2])	17.5	...	...
	RR (j ≥ 3)	4.3	38/63	2718; 0; 0
<i>erode_fifo</i> (i = 3)				

Sw: 7,21 mJ  
Hw: 0,472 mJ } 15X

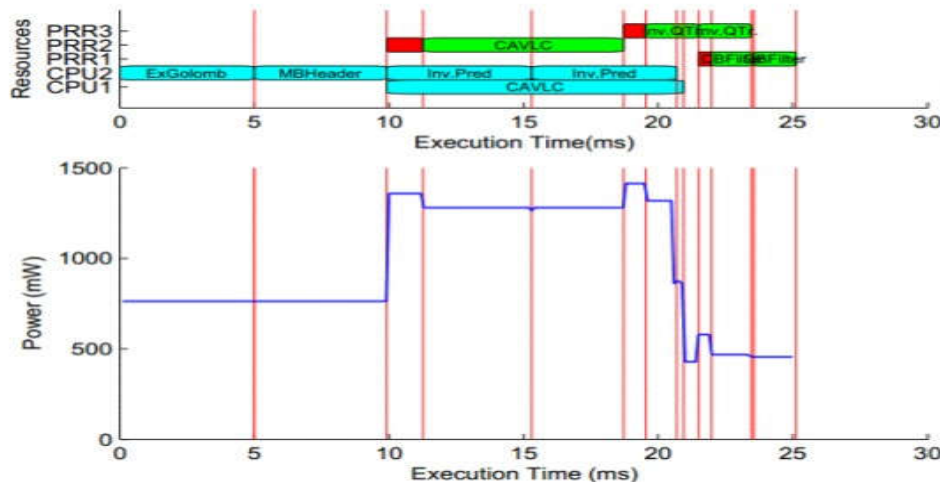
→ Realistic analysis of full application mapping on multiprocessor / Hw accel



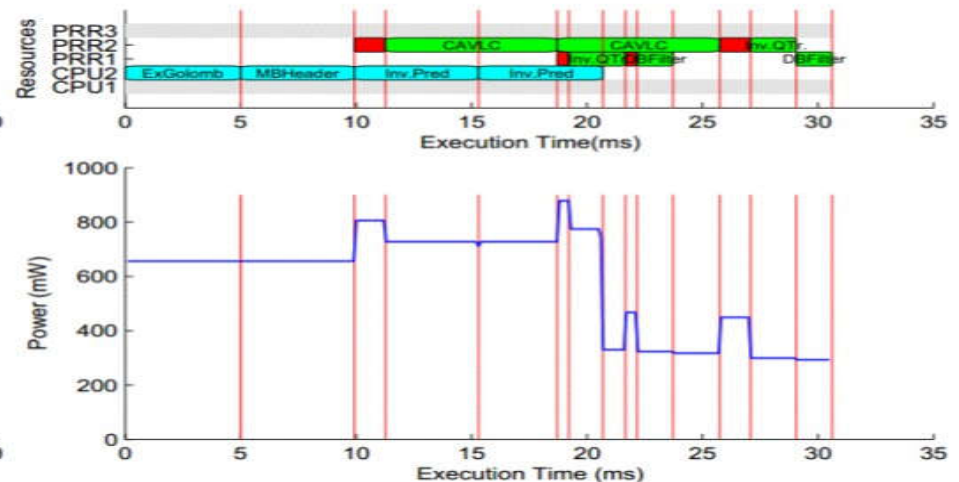
# System level modeling and design

- Power modeling at system level
  - Heterogeneous platforms – Dynamic reconfiguration (DPR)
- Reconfiguration control modeling
- Higher design space complexity
  - FPGA partitioning (Partial Reconfigurable Regions, PRR)
  - Scheduling

*ANR Open-PEOPLE (2009 – 2012)*



Best Performance(BP) solution: 2 cores, 3 PRRs

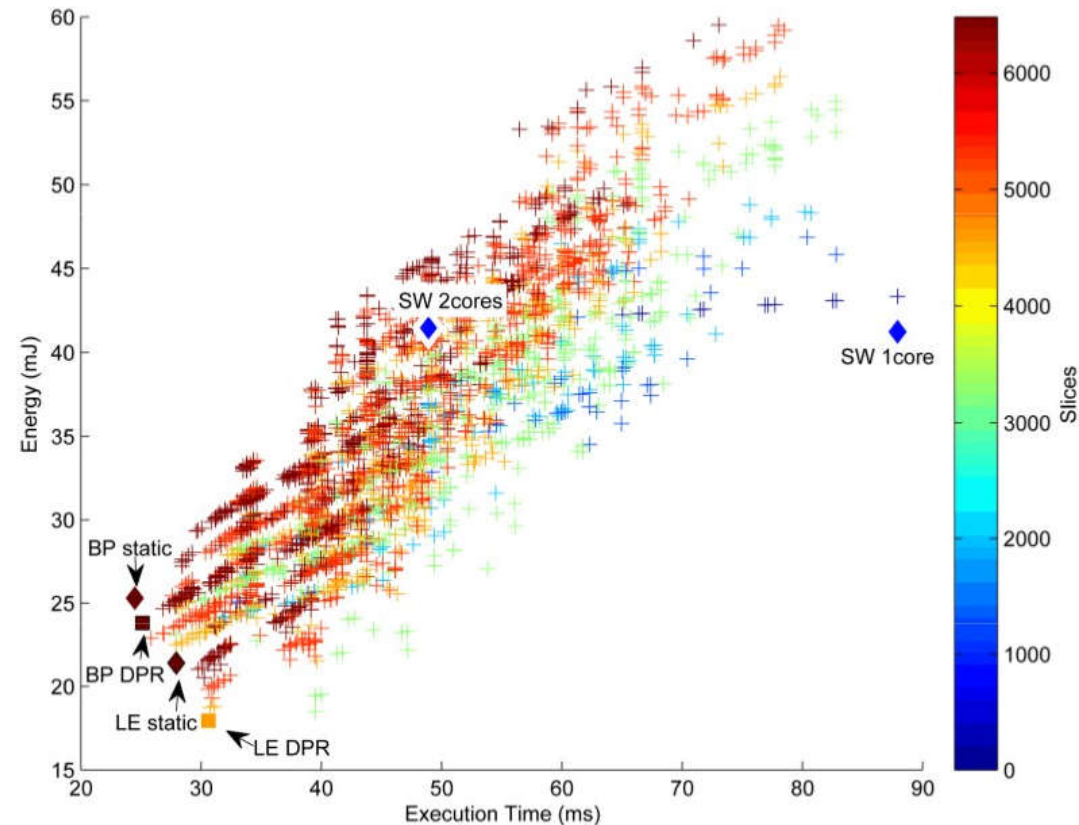


Low Energy (LE) solution: 1 core, 2 PRRs

→ Realistic analysis of full application mapping on multiprocessor / DPR accel

# System level modeling and design

- System level design
    - Heterogeneous platforms – DPR
  - Improve system level design and mapping
    - More focus on energy
      - Low complexity power estimation
    - Design space exploration
      - Methodic approach
      - FPGA partitioning and scheduling (DPR)
- Further savings from fine management of power at runtime

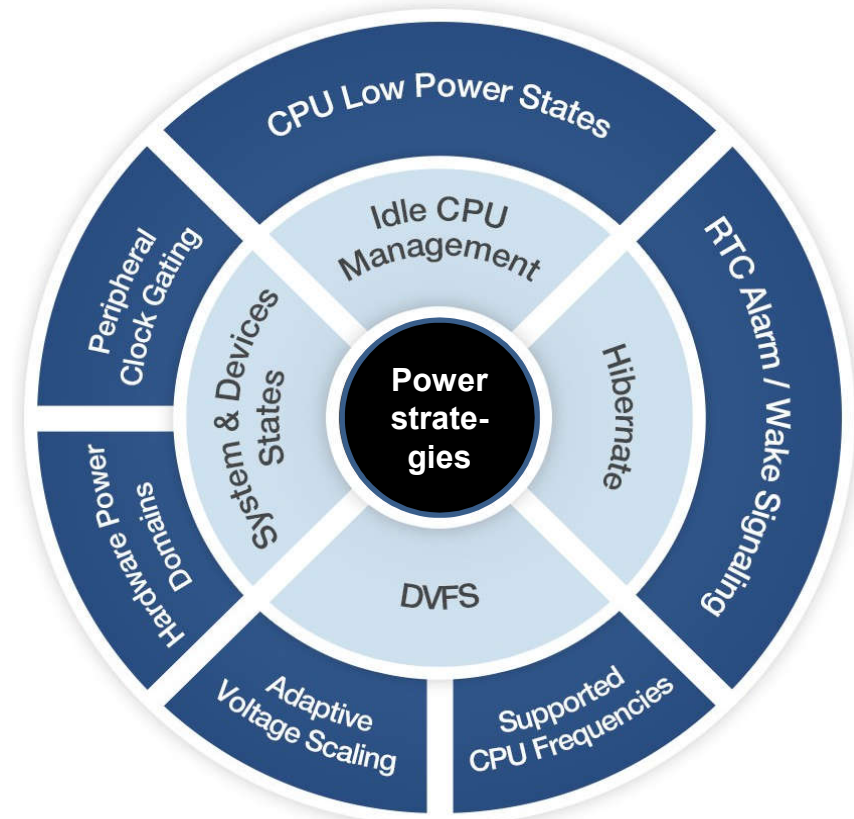


Exhaustive exploration:  
10 functions / 6 Hw functions  
> 1M solutions explored in a few seconds

# Power and resource management

- Power management and power strategies
  - Workload based (general purpose)
    - If the system is busy, run at maximum frequency
    - When the system is less busy, decrease frequency
- For video processing, this strategy leads to maximum power consumption
- There's room for more energy savings

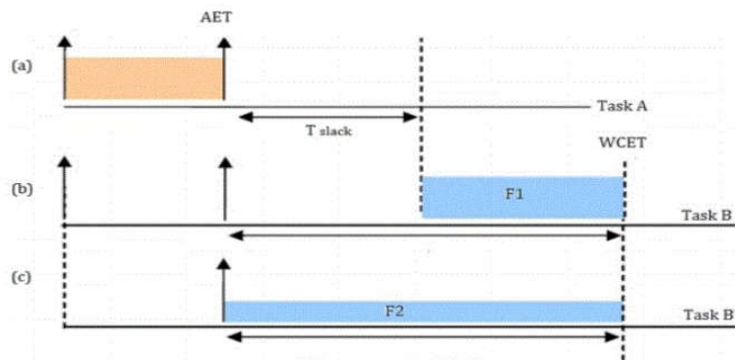
SoC power consumption



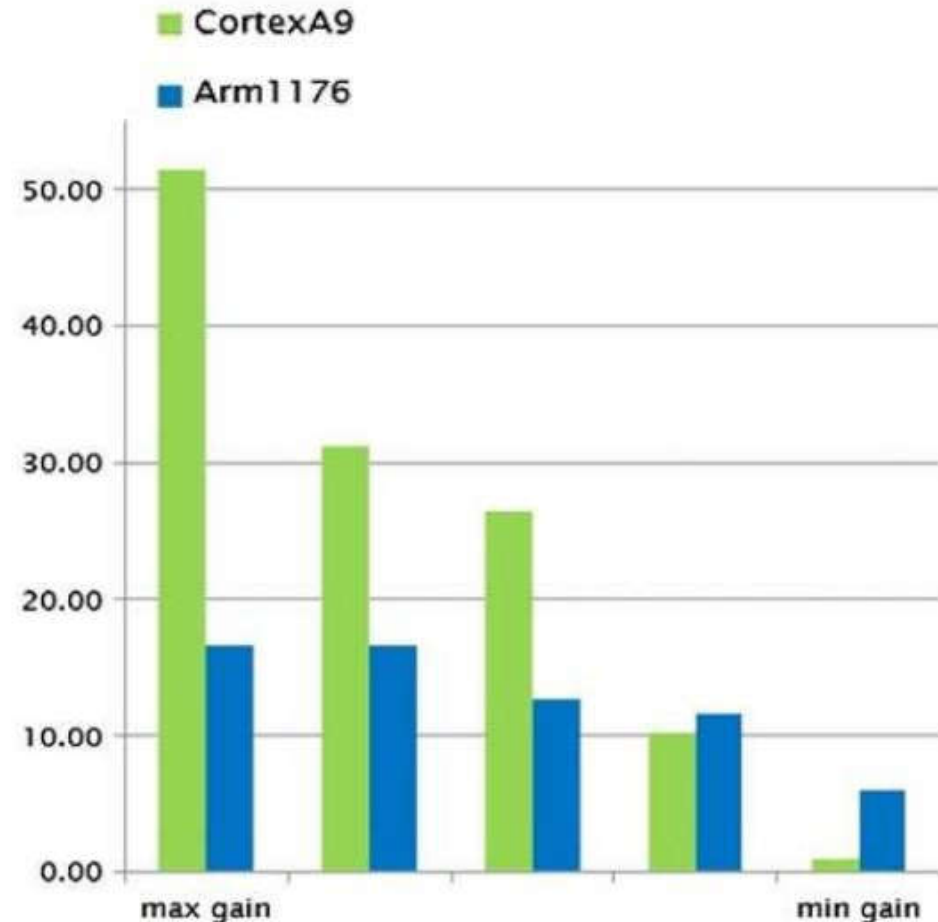
→ Specialized power strategies

# Power and resource management

- Specialized power strategies
  - Deadline (real-time) scheduling / DVFS
- Energy gains
  - Application dependent (slack)
  - Platform dependent (operating points, silicon technology)
    - System level analysis of power strategies
      - e.g. DVFS or idle states ?



COMCAS (CATRENE 2009 – 2012)



# Power and resource management

- Specialized power strategies
  - Energy Aware Heterogeneous Scheduler (EAHS)
  - Homogeneous multicore (without DVFS) + DPR
- Earliest Deadline First (EDF)
- Releasing real-time constraint to promote energy efficiency
  - A tuning parameter allows to choose different power / performance tradeoffs
- At each scheduling event, a cost function evaluates all possible implementations for ready tasks
  - Decisions rely on efficient decision support and estimators
- Energy gains
  - 44% against Sw execution
  - H.264 decoder / Zynq: 18 functions / 12 Hw functions

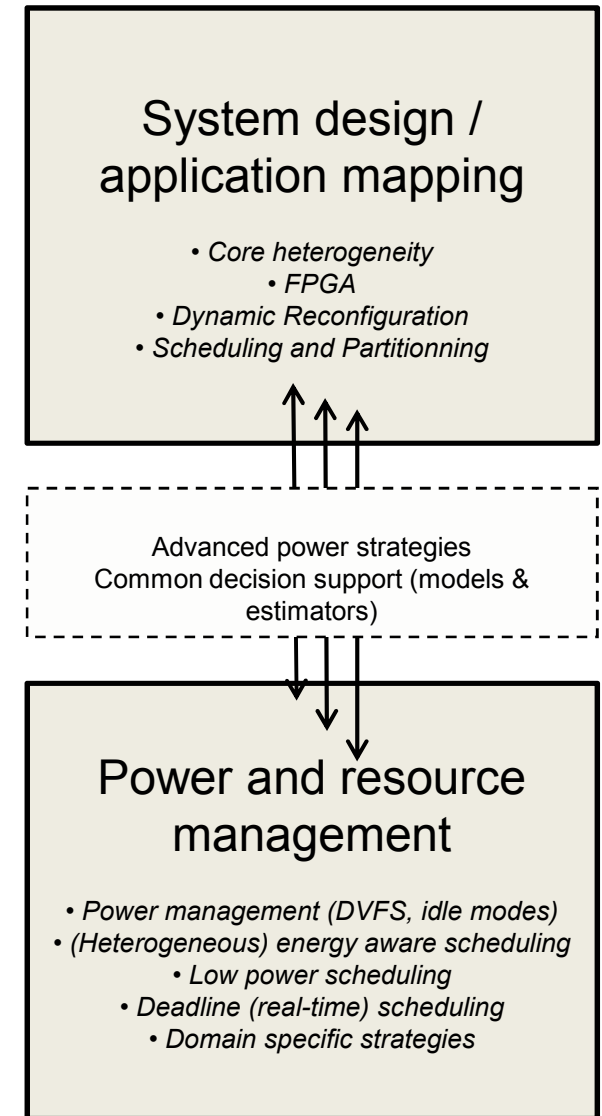
# Power and resource management

- Further saving potential lies in specialized low power scheduling
  - Up to 50%, probably more when addressing heterogeneity
  - Use specialized power strategies based on application domain properties (e.g. real-time, video processing, ...)
  - Include power management analysis at early development stages to check the efficiency of power strategies
    - e.g. DVFS, idle states, DPR, combination of strategies
    - What is the best option for my application ?
  - Complex to implement (kernel OS based)  
→ Provide implementation support

# Energy efficiency improvements

- FoTRReSS: an integrated development methodology
  - Centered on energy efficiency
  - Improving global cohesion
    - Allow a better use (and full design) of advanced power strategies
    - Include power strategies in exploration to improve the relevance of solutions defined
    - Common decision support for exploration and strategies (same estimators for coherence of decisions)
  - Allow full prototyping of solutions (application and strategies)
    - Middleware
    - Development of advanced power strategies (in userpace Linux, patented technique)

<http://fortress-toolbox.unice.fr/doku.php>



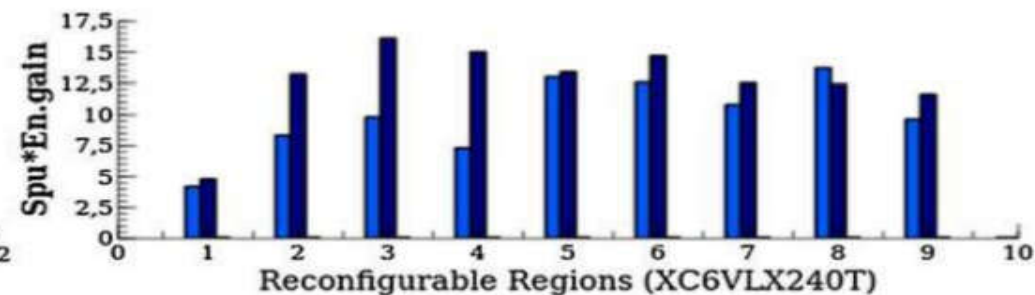
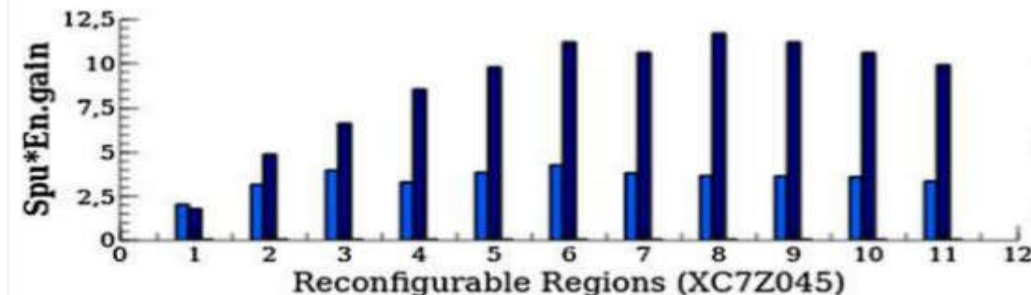


# Energy efficiency improvements

## ■ Application: License Plate Recognition

- Gains from heter. only **x4 – x7**
- Gains from integr. méth. **x0.9 – x3**
- Total **x4 – x20**

■ EDF  
■ EAHS



- Gains from EAHS heterogeneous scheduler
    - **65%** (Zynq) and **15%** (Virtex6/MB) in average
  - Variations occur and depend on various parameters
    - Parallelism, partitionning, size, silicon technology, cores, ...
- Necessity of a global and methodic approach to fully exploit the very high level of heterogeneity

BENEFIC (CATRENE 2013 – 2016)

# Conclusion and perspectives

- Contributions to the exploitation of heterogeneity
  - Energy modeling at higher abstraction levels
  - Power and resource management: EA(H)S, ...
  - The significant increase of complexity involves major design challenges and needs comprehensive solutions to make the best out of heterogeneous resources, dynamic management and proper use of different technologies
- Lessons
  - The prevailing approach still addresses the different problems, technologies and techniques with high separation of concerns
  - A strong potential lies in the investigation of more global approaches, but this requires ambitious research on the long run
- Perspectives
  - Confrontation with other application domains (Exascale, IA, aut. systems)
  - Specialized strategies (AI, VoD, Exascale...)
  - Technologies (DPR + DVFS, NV-RAM, Coarse Grained Reconf...)