# Technology showstoppers for the next generation of devices (*)

Prof. William Fornaciari

Politecnico di Milano – Dipartimento di Elettronica e Informazione

william.fornaciari@polimi.it

www.elet.polimi.it/~fornacia

**Milano, June 2009**

**(*) Most of the information contained in this presentation are courtesy of STMicroelectronics: Talk of R.Zafalon @ DATE'09**

# Outline

➢ Market Application rush

➢ CMOS Roadmap: 3 main showstoppers

➢ Why bothering for low power systems?

  – Technology Scaling, Trends & Roadmap

  – Leakage Aware design strategies

  – Cost of heat removal: packaging and reliability

  – Memory architectures

  – Increased market share of mobile electronics

➢ European proposal: ST Computing Platform's Roadmap

  – Platform 2012

➢ Not Only Mobile…!

➢ Conclusion

# 30 Years of Electronics Industry CAGR
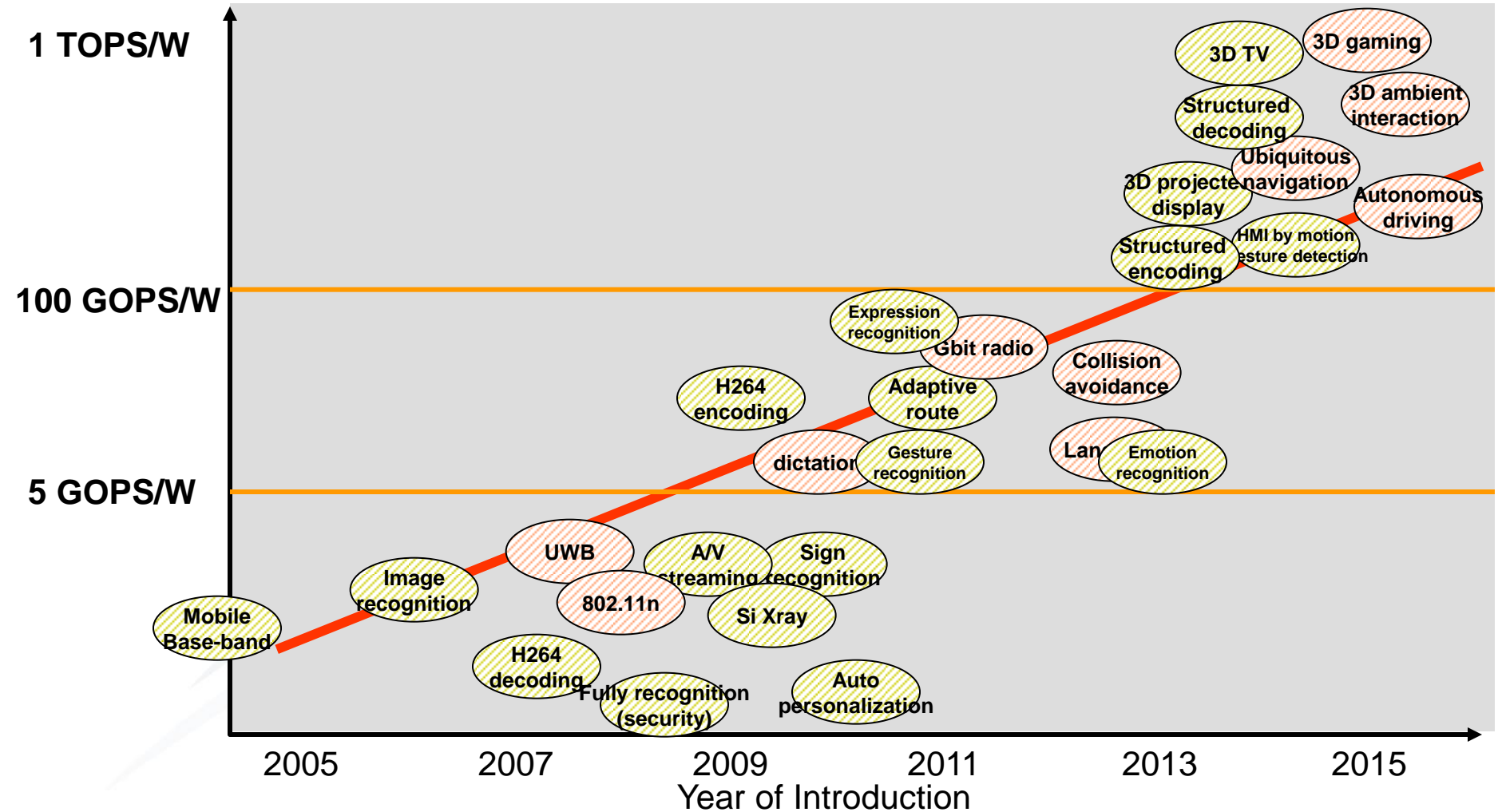
**Semic. Capex: 17%**

**Semic. Market: 15%**

**Electronic Systems: 8 %**

**WW GDP: 3,4%**

# Market Application rush

# CMOS Roadmap: 3 main showstoppers

**Pat Gelsinger, CTO Intel Corp.**

**Quote from DAC'04 Keynote:**

*Power is the only limiter !!*

**CMOS Roadmap: 3 main showstoppers:**

1. **Subthreshold Leakage Current ( $I_{off}$ )**

2. **Huge Process Variation Spread**

3. **Interconnect Performance and Signal Integrity**

# Why bothering for low power systems?

➤ Practical market issue:

  ✓ Increasing market share of mobile, asking for longer cruising life

  ✓ Limitations of battery technology

➤ Economic issue:

  ✓ Reducing packaging costs and achieving energy savings

➤ Technology issue:

  ✓ Enabling the realization of high-density chips (heat poses severe constraints to reliability)
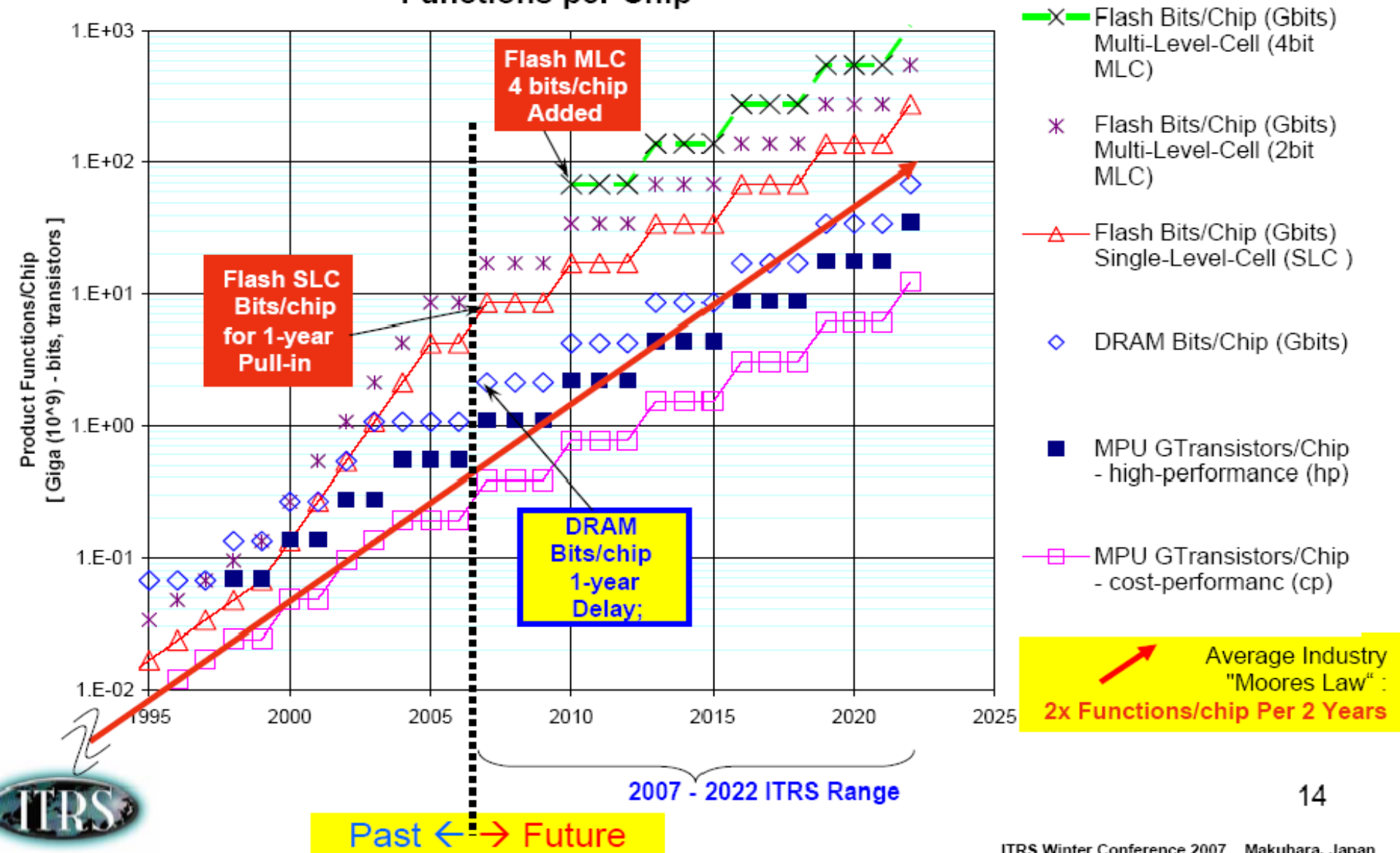
# CMOS at core of chip making still for many years

- The theoretical limit for transistor gate length on silicon is around 1.5nm.

  - ✓ Today's 65nm CMOS process has a gate length of 42nm: i.e **28X larger** than the theoretical limit!

  - ✓ In 32nm, the gate length is 21nm

    i.e. **14X above limit**

- The gate delay determines the fundamental speed of the logic. The theoretical limit is 0.04ps

  - ✓ Today's 65nm logic NAND2 is ~1ps, i.e. **24X slower**!

- Transistor density, i.e. the number of device which can be squeezed into a chip, reaches the limit around 1.8 billion Tx per cm².

  - ✓ Today's 65nm CMOS device is **7.5X larger**! (i.e. 750Kgate/mm² = 2.4M Tx/mm² = 240M Tx/cm²)

- Performance as measured by clock speed, fell off Moore's Law during the last decade, thanks to Multi Processors computing architectures.

# ITRS Roadmap 2007 vs Moore's law



2007 ITRS Product Technology Trends - Functions per Chip

[DRAM and Flash Updated]

14

ITRS Winter Conference 2007   Makuhara, Japan
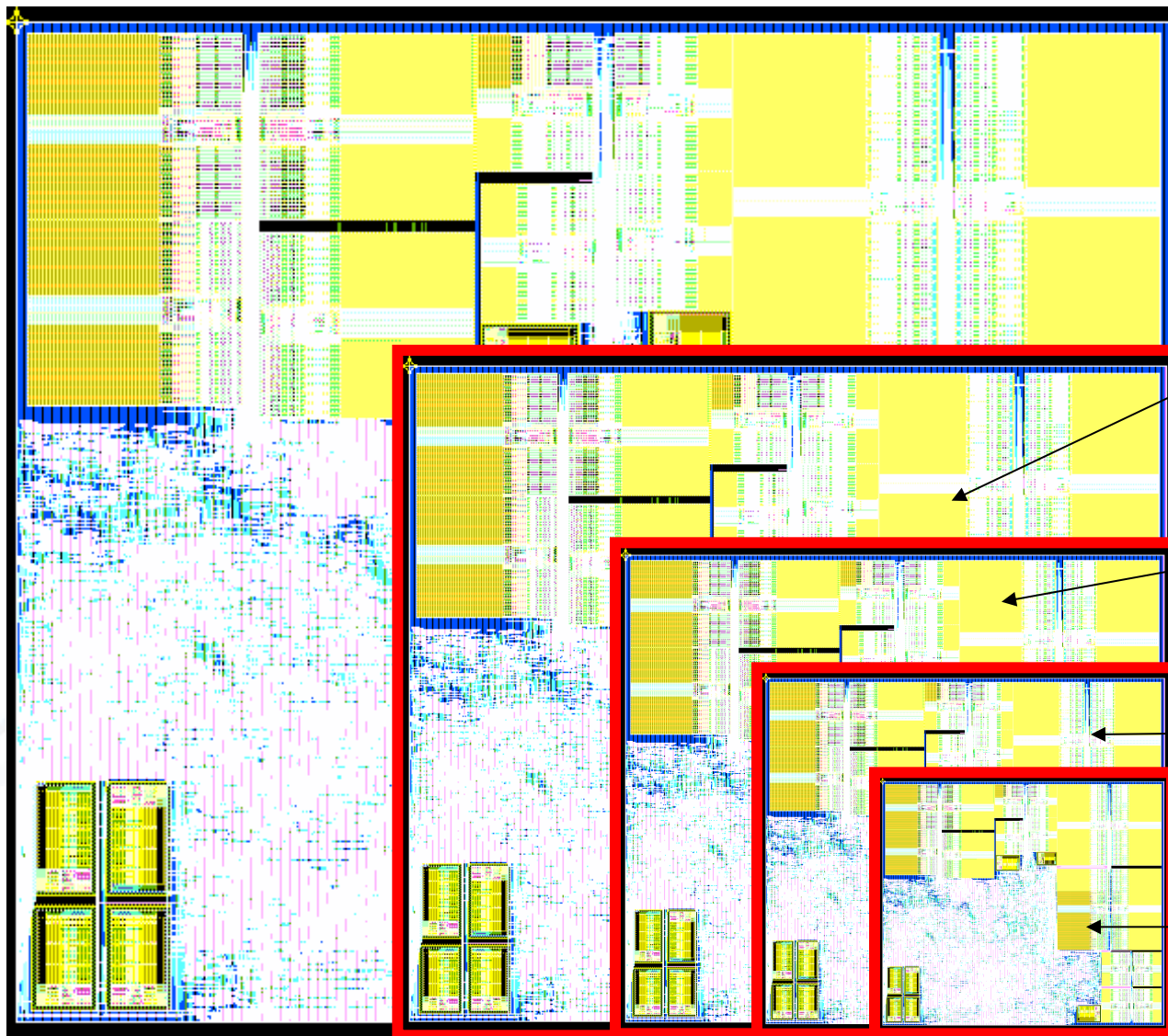
# Squeezing costs of computing cores



**ARM 9**

180 nm
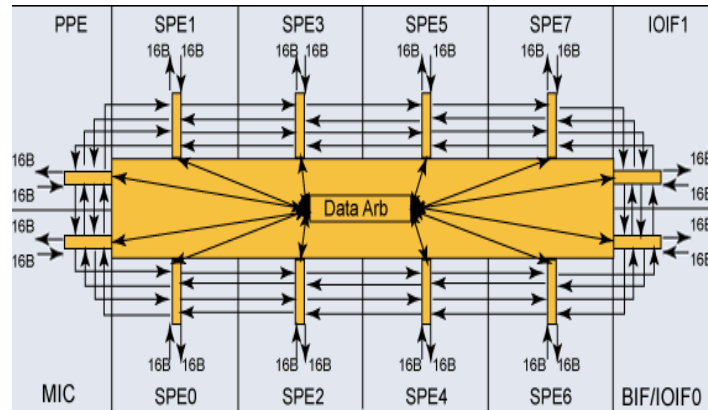11.8 mm2

130 nm,
5.2 mm2

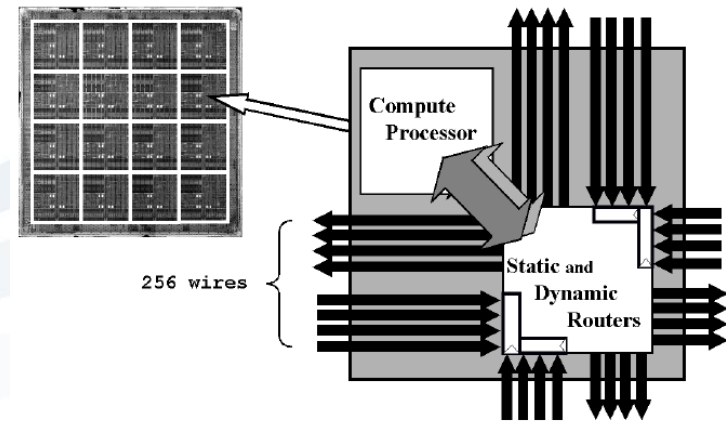90 nm,
2.6 mm2

65 nm
1.4 mm2

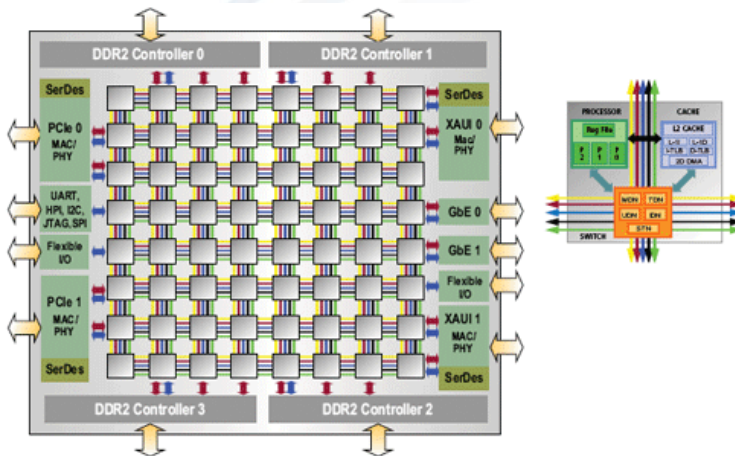45 nm
0.75 mm2
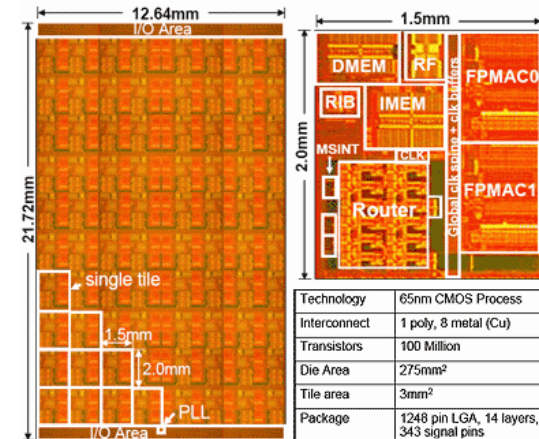
# From multi-core to many-core

## IBM - Cell/B.E.
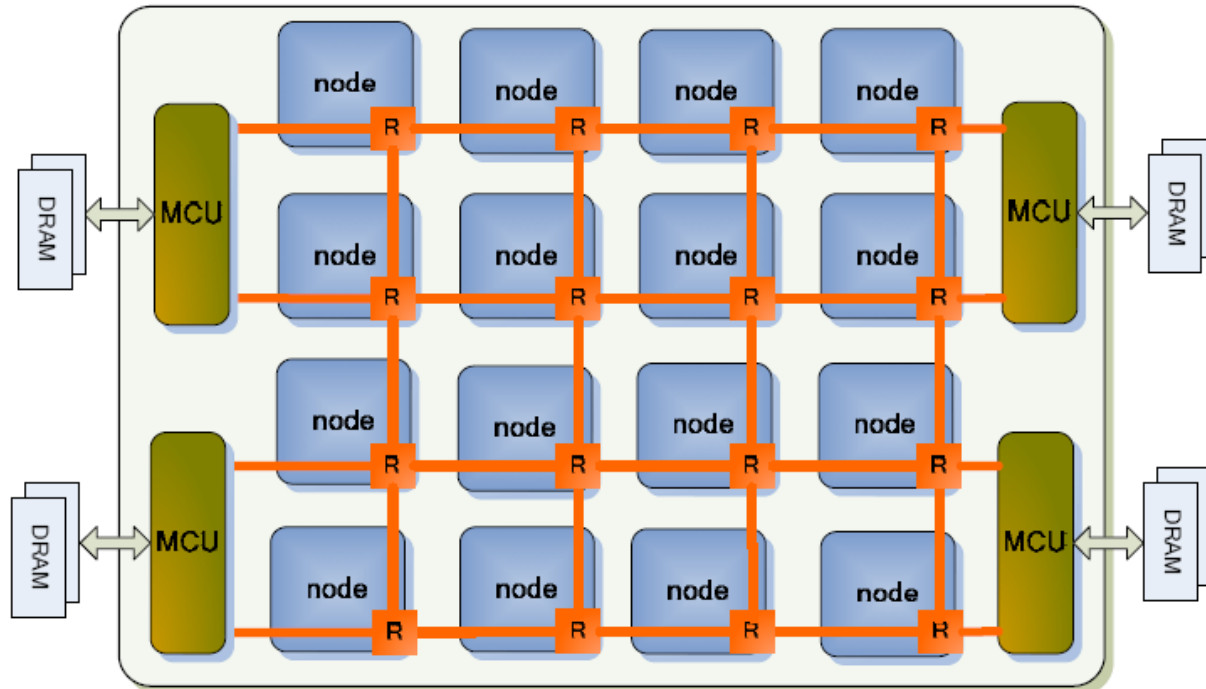


## MIT - RAW



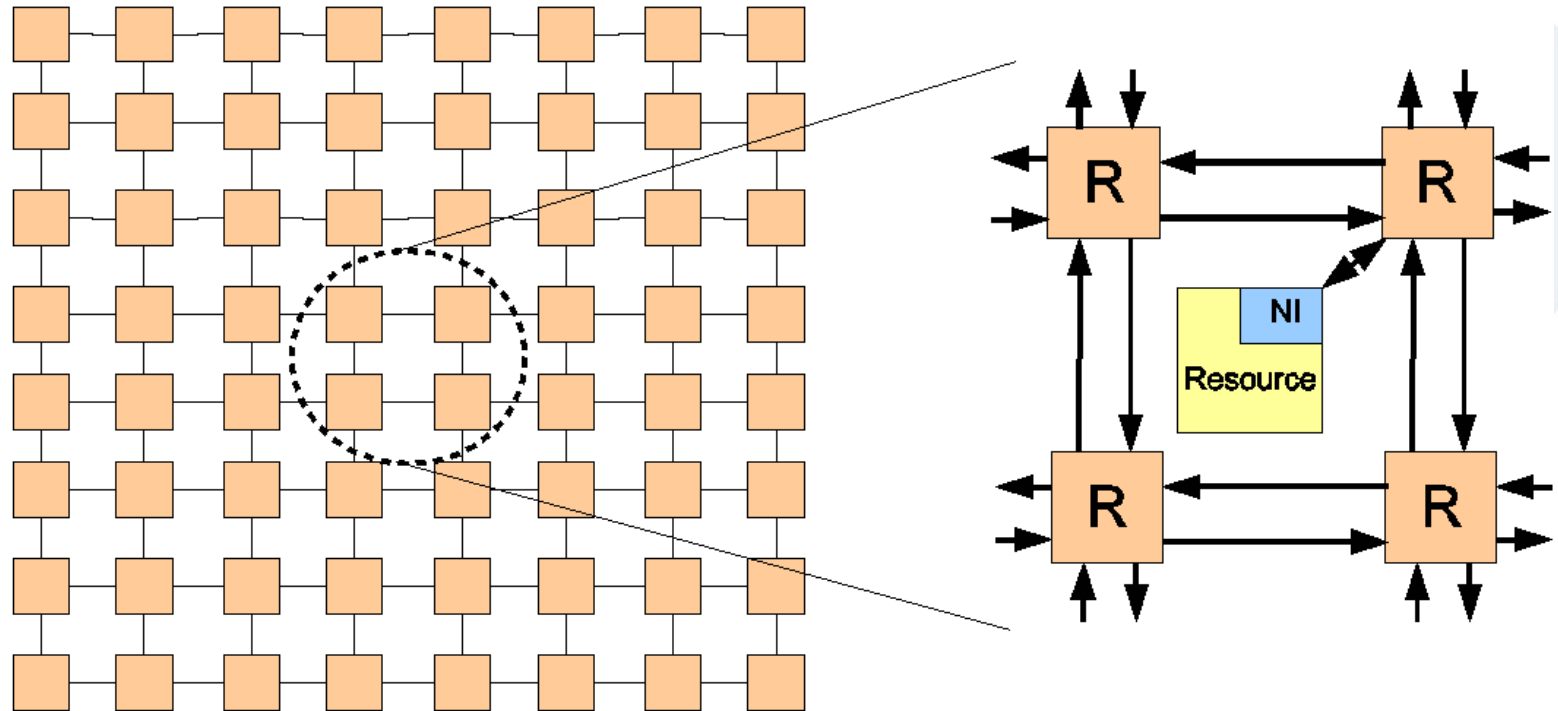## Tilera - TILE64



## Intel -Terascale

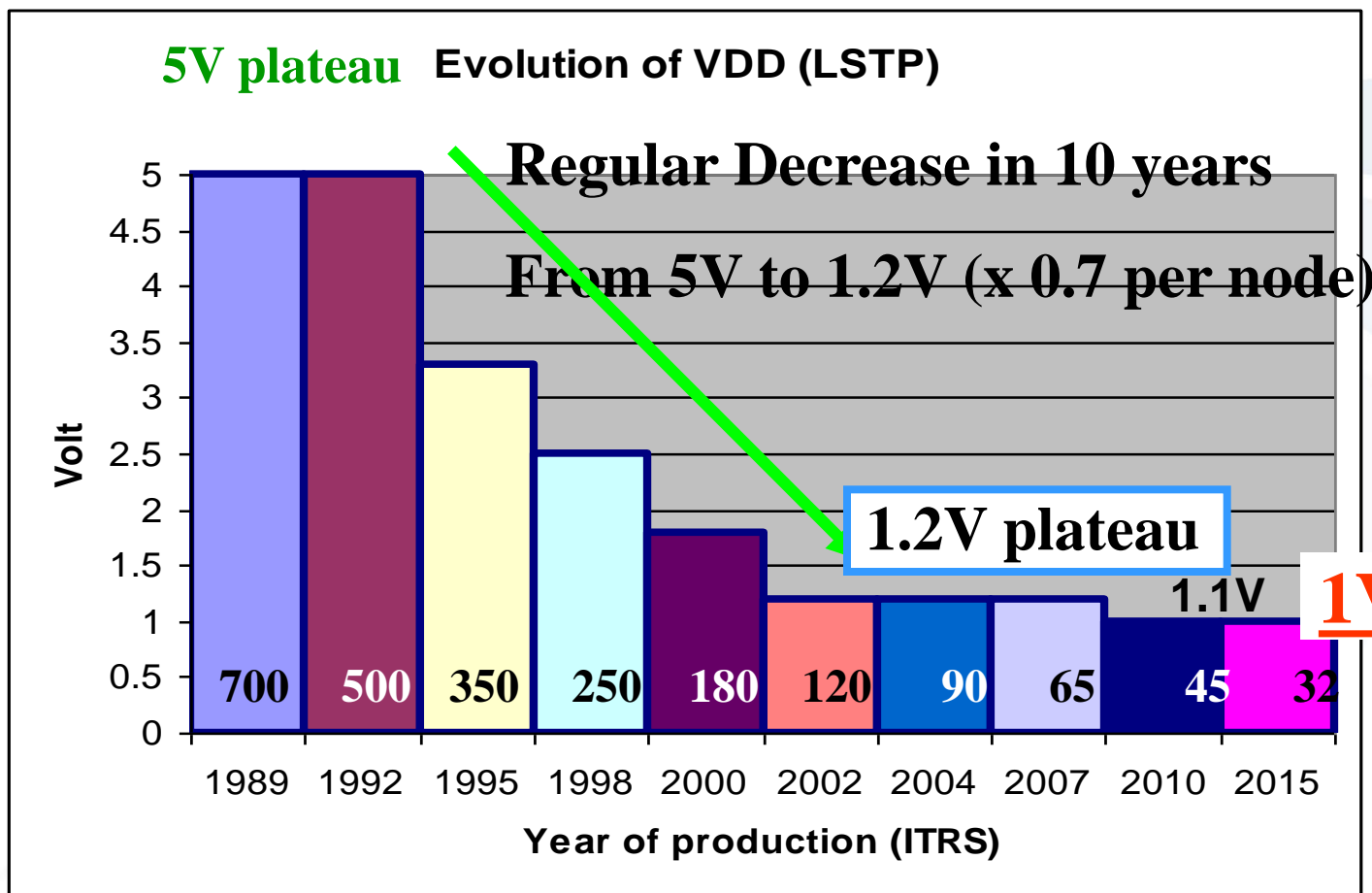# Towards 1B Multi-core architectures

➢ **Multi-core architecture: a tiled homogeneous multi-core architecture for general embedded purpose** (Godson-T)

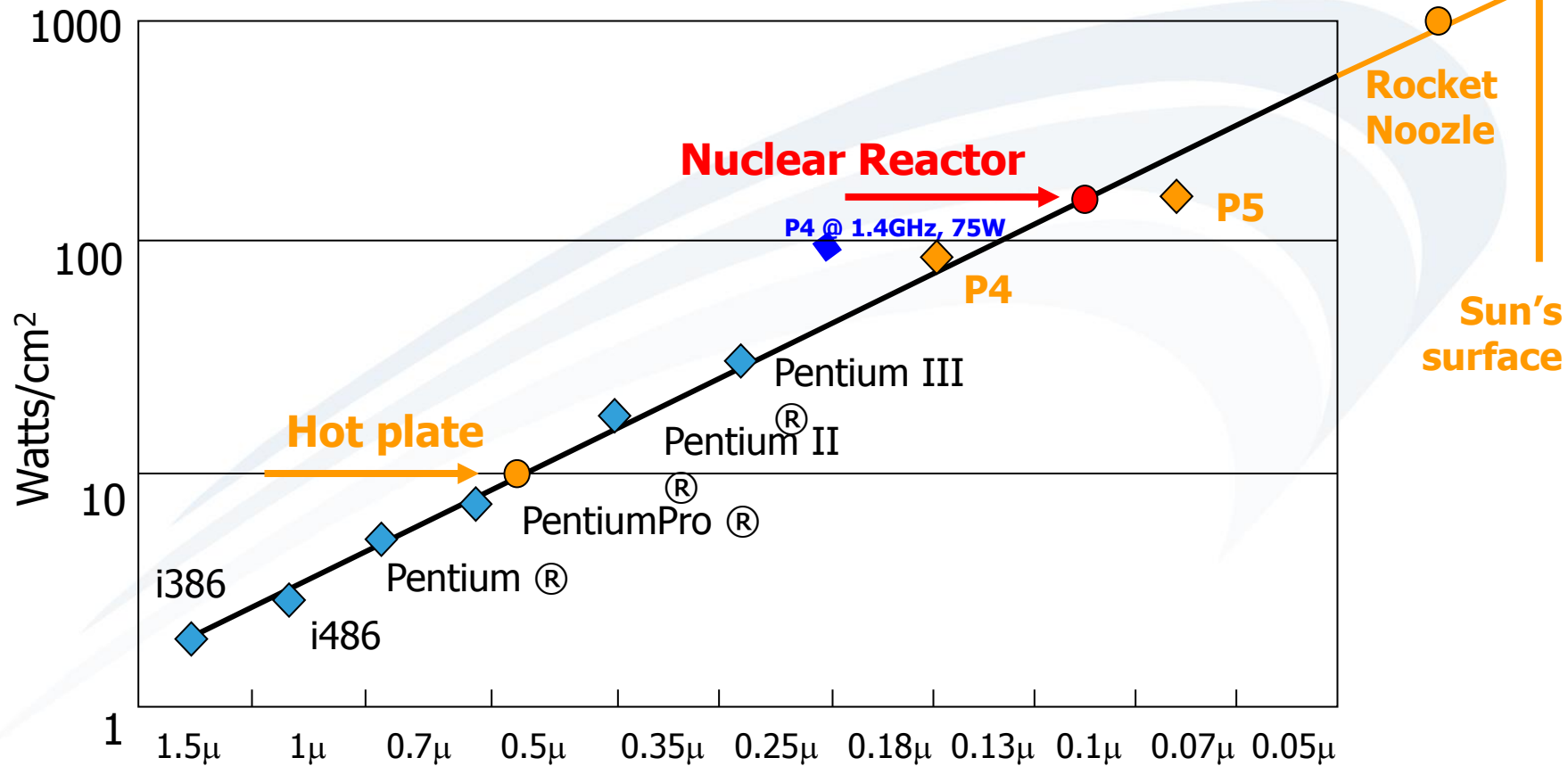# Many-core computing fabric template

# VDD (no more) scaling is increasing the «power crisis»



**5V plateau**

**Evolution of VDD (LSTP)**

**Regular Decrease in 10 years**

**From 5V to 1.2V (x 0.7 per node)**

**1.2V plateau**

1.1V

**1V plateau?**

Volt (y-axis): 5, 4.5, 4, 3.5, 3, 2.5, 2, 1.5, 1, 0.5, 0

Bar values: 700, 500, 350, 250, 180, 120, 90, 65, 45, 32

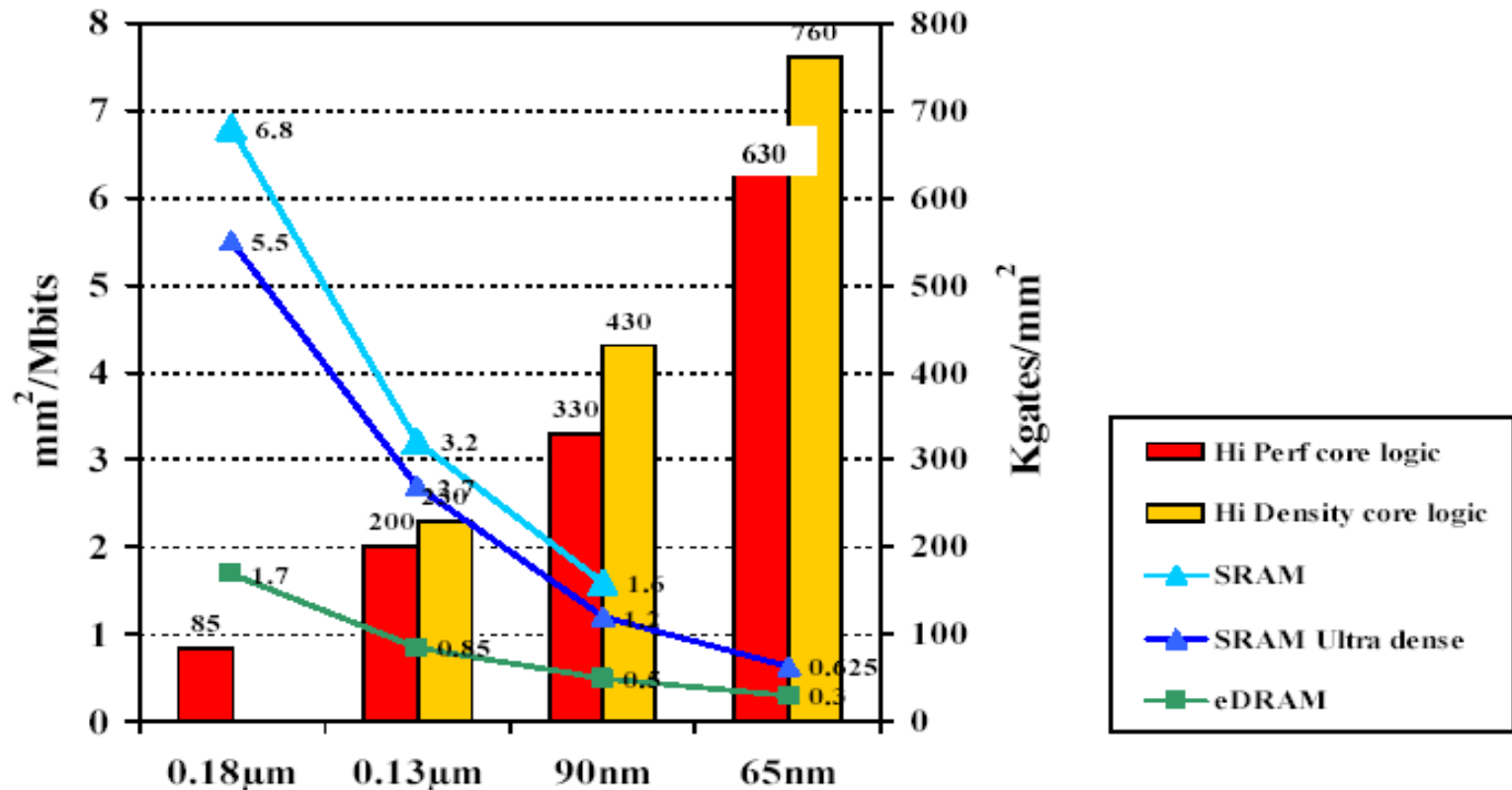Year of production (ITRS): 1989, 1992, 1995, 1998, 2000, 2002, 2004, 2007, 2010, 2015

# Power Trend for microprocessors

> Power density in Intel's microprocessors:

# CMOS Logic Tech Overview



Source: STMicroelectronics

# 90/65/45nm Speed vs Leakage



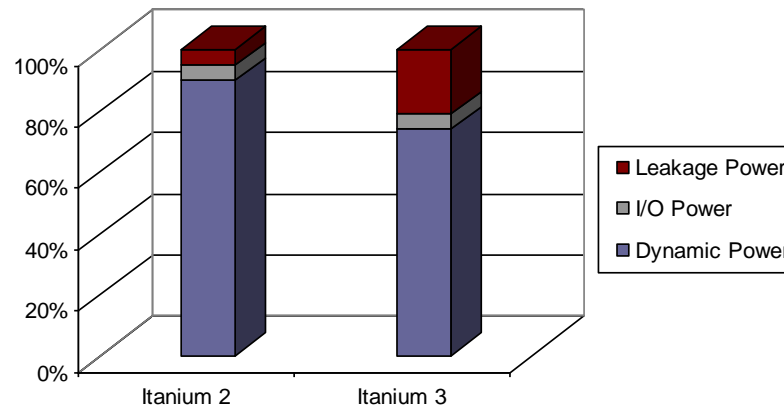Source: STMicroelectronics

# Technology Scaling

➢Increasing contribution of leakage power:
✓Example: ASICs [source: STMicroelectronics]

**Power Density (Watts/cm$^2$)**

Chart: Stacked bar chart with y-axis 0 to 150 (25 increments). X-axis: 250nm, 180nm, 130nm, 90nm, 65nm. Legend: Leakage Power, Dynamic Power.

✓Example: Microprocessors [source: Intel].

**Itanium 2:**
**180nm, 1.5V, 1.0GHz, 221MTx (core+cache)**

**Itanium 3:**
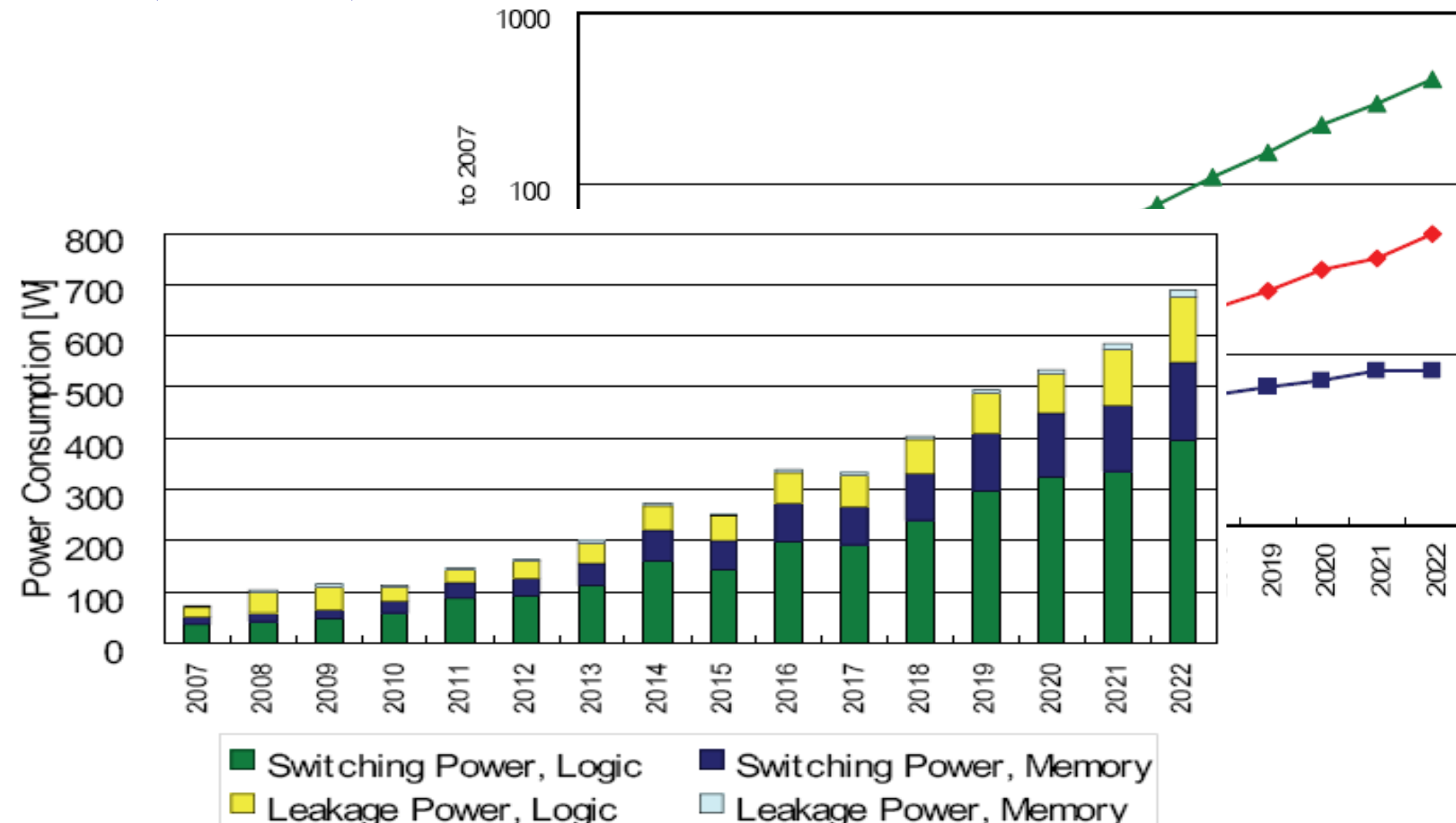**130nm, 1.3V, 1.5GHz, 410MTx (core+cache)**

Chart: Stacked bar chart, y-axis 0% to 100% (20% increments). X-axis: Itanium 2, Itanium 3. Legend: Leakage Power, I/O Power, Dynamic Power.

# SoC Requirements for MP platforms (1)

➢ Processing performance is expected to grow more than 200x in the next 15 years.



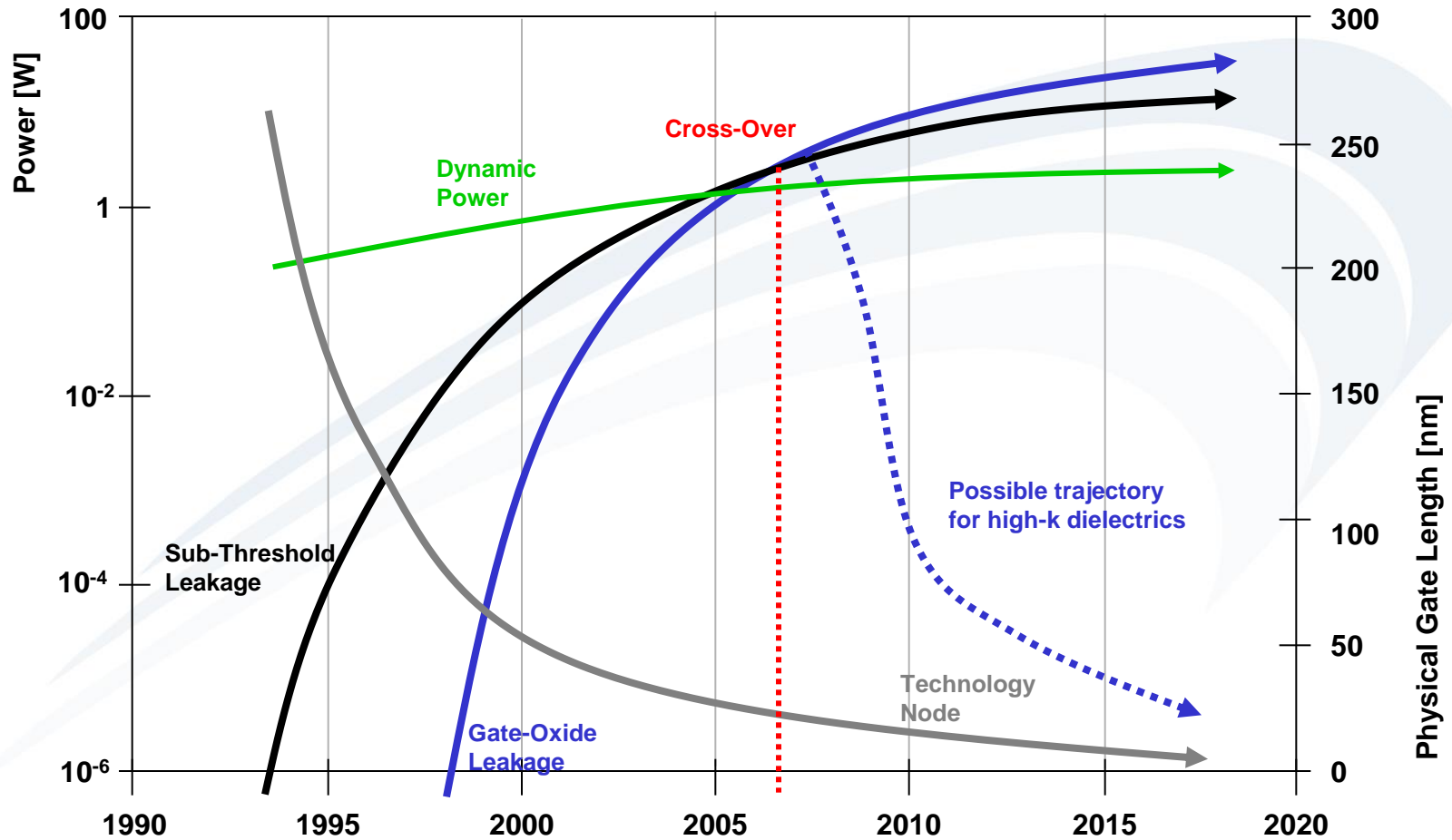| | | |
|---|---|---|
| Processing Performance (Normalized to 2007, Right Y Axis) | Switching plus Leakage Power (Normalized to 2007, Left Y Axis) | Design Effort, Power (Normalized to 2007, Left Y Axis) |

➢ # PE per chip; Processing Performance; ND2's max switching frequency
(normalized to 2007)

# Dynamic vs. Leakage Power



Source: ITRS Roadmap

# Semiconductor's Challenge



## Moore's Law at Work!

# Leakage crisis: Is it a technology issue only?

- ➢ **Trends**:

    - ✓ nominal Vdd getting stable around 1V

    - ✓ MOS's Vth linearly scales to keep costant speed

    - ✓ But… leakage grows exponentially with Vth reduction !!

    - ✓ sub-threshold current from 100 to 1000 pA/um

    - ✓ gate leakage to become larger that sub-threshold

    - ✓ total static power from 21E-12 to 60E-12 W/Transistor

- ➢ SOI has major disadvantages w.r.t. sub-threshold reduction!

# "Leakage Aware" design strategy includes

## A. Gate/Circuit-level techniques

Use of multiple $V_{th}$
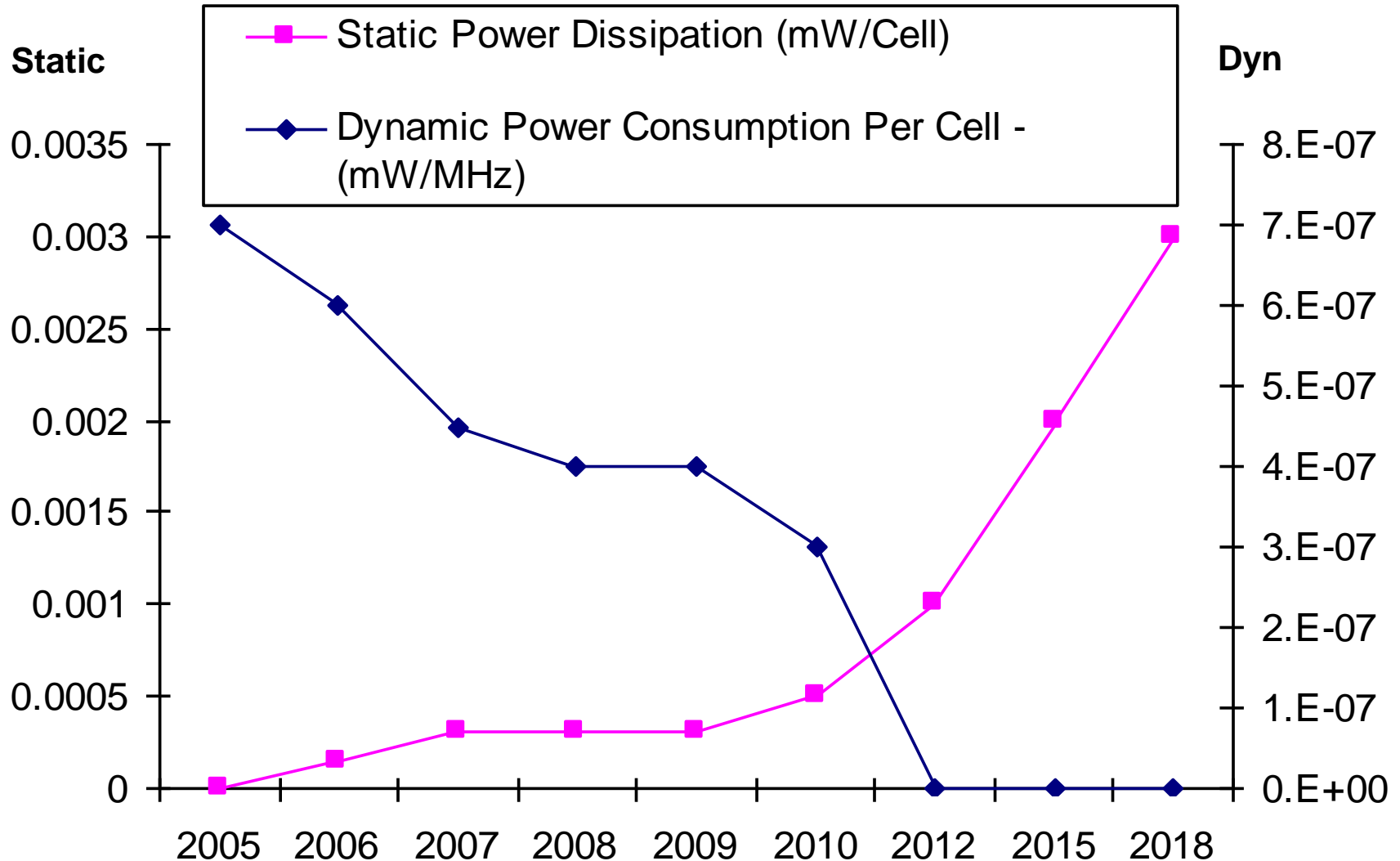
- Dual-$V_{th}$ design.
- Mixed-$V_{th}$ (MVT) CMOS design.
- MTCMOS.

✓ Sleep transistor insertion/Voltage islands

✓ State retention FFs

✓ **Techniques for memory circuits**

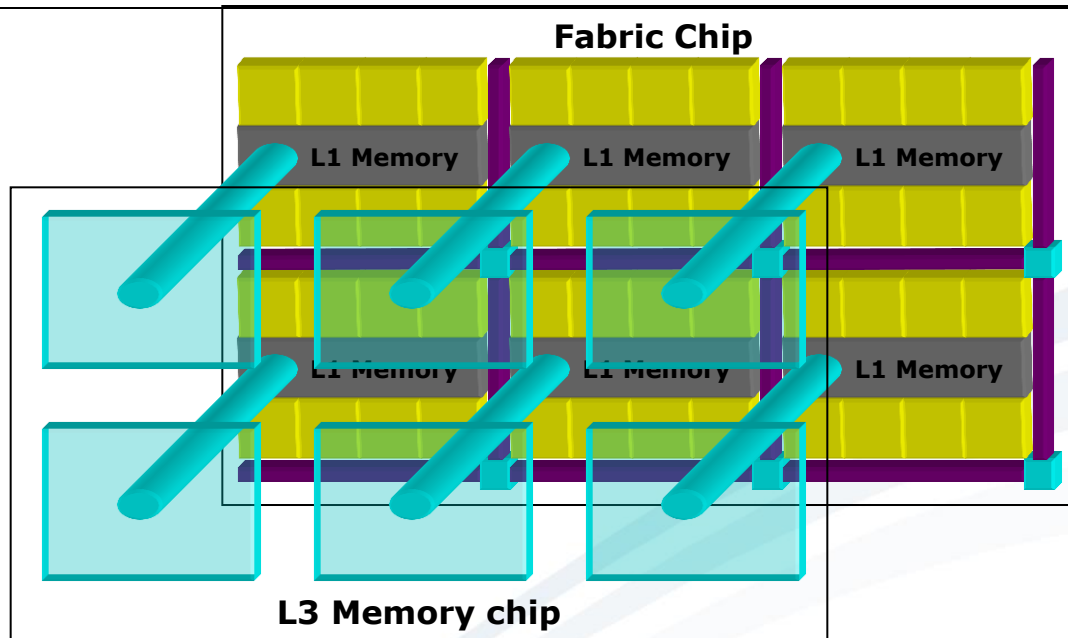Cell state (stored value) determines exactly which transistors "leak"

✓ **State-preserving** techniques:

- Only suitable choice for non-cache memories (e.g., scratchpad).

✓ **State-destroying** techniques:

- Suitable for caches (can invalidate values).

## C. Architectural techniques

✓ Adaptive Body Biasing (ABB).

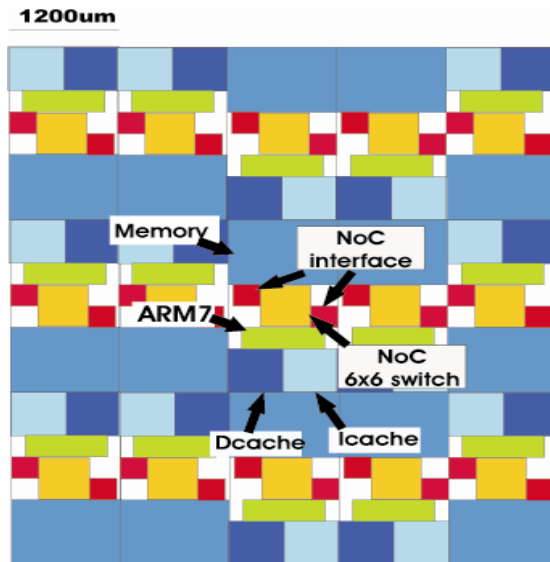✓ Adaptive Voltage Scaling (AVS).

✓ $V_{th}$ hopping.

✓ Multiple $V_{BB}$

# Memory Driver
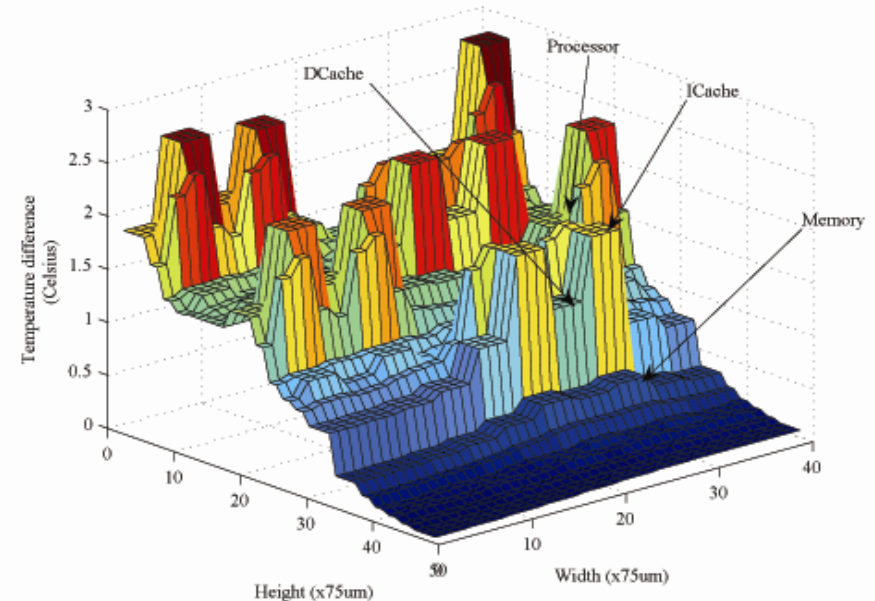
# Technological opportunities. 3D Architecture

**Fabric Chip**

**L1 Memory**   **L1 Memory**   **L1 Memory**

**L1 Memory**   **L1 Memory**   **L1 Memory**

**L3 Memory chip**

- L2 memory
  - ✓ Can be reduced to minimum size
  - ✓ > 50% more processing per mm²

- L3 memory
  - ✓ Separate chip
  - ✓ Multiple high bandwidth direct connections
  - ✓ Connection roadmap
    - – Bumping
    - – Through Si Via
    - – RF pad

- **Huge memory bandwidth available**
- **Drastic reduction of power budget**
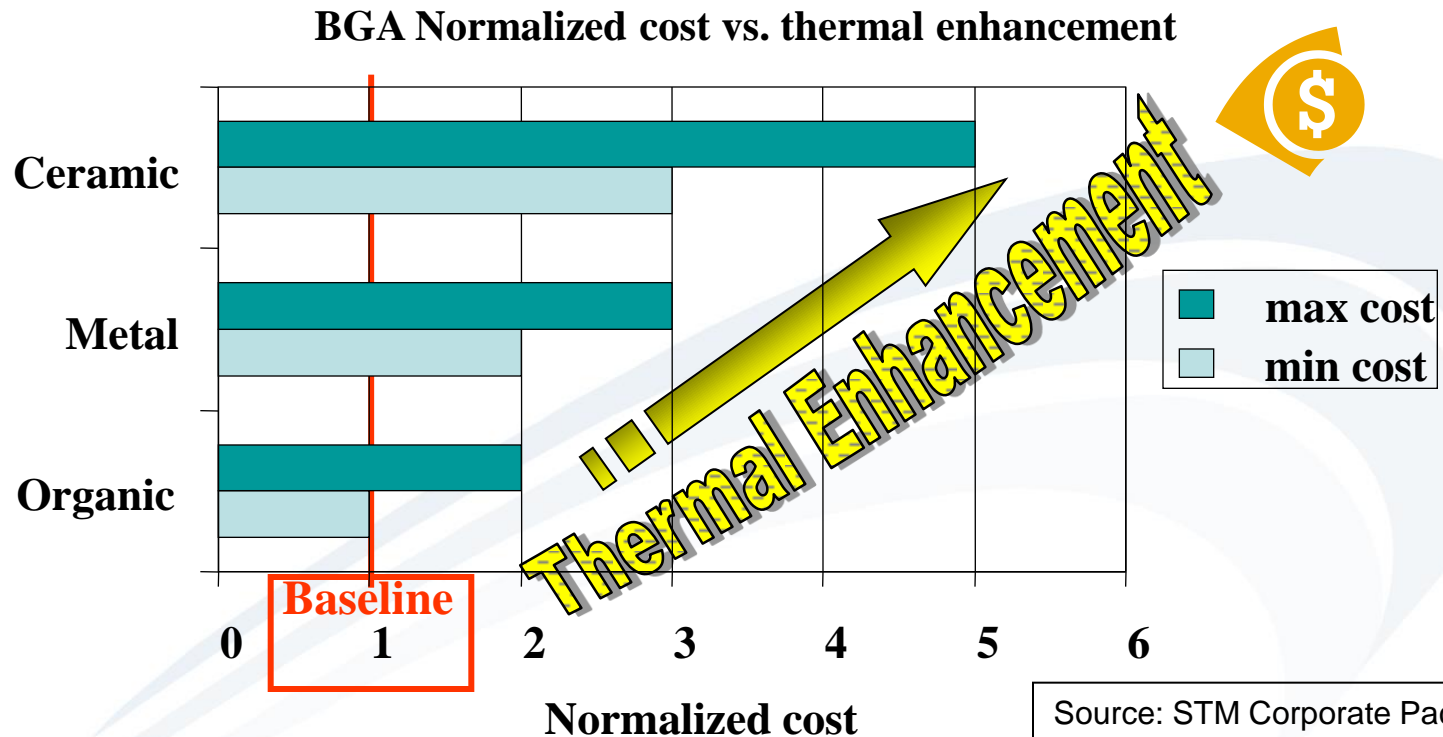
**Chip floorplan**

**Steady state temperature**

Some hot spots in steady state:

§ Silicon is a good thermal conductor (only 4x worse than Cu)

and temperature gradients are likely to occur on large dies

§ Lower power density than on a high performance CPU
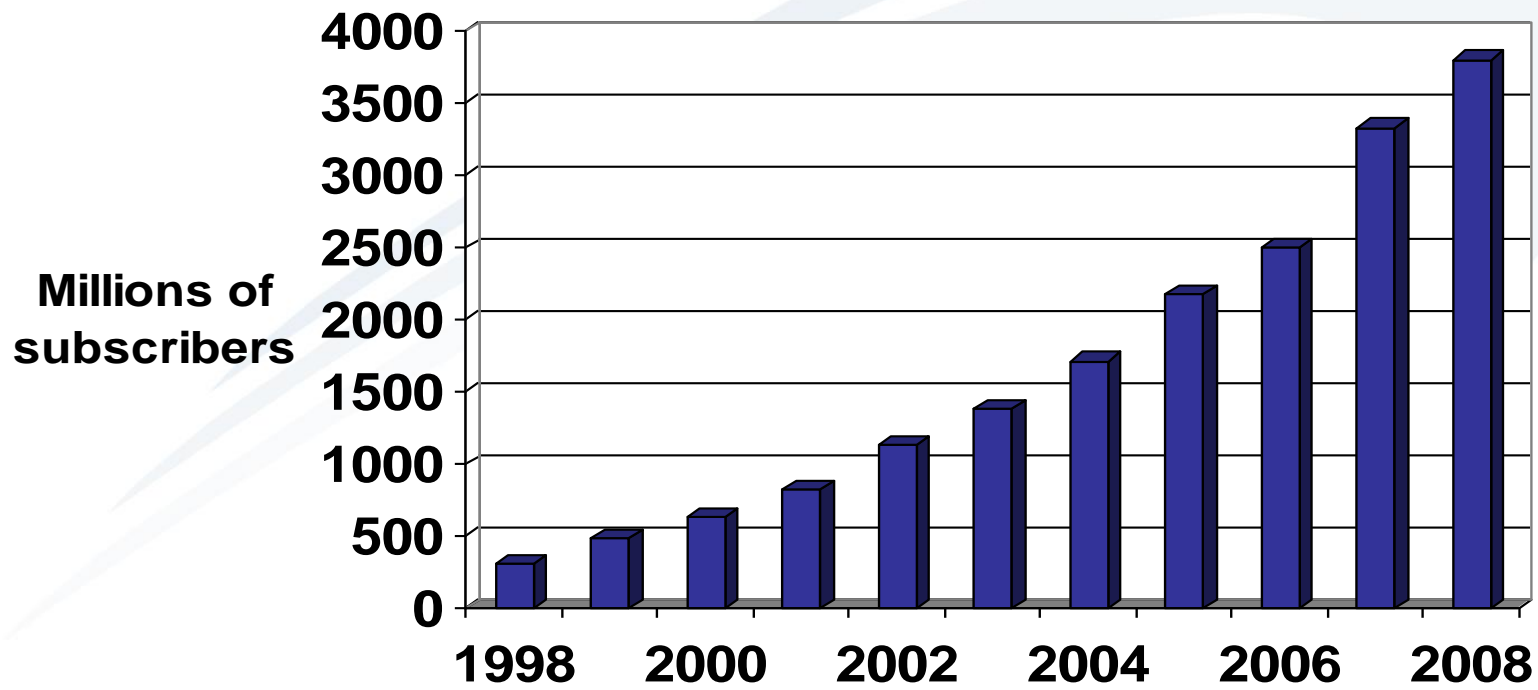
(lower frequency and less complex HW)

# Thermal Management Challenge



**BGA Normalized cost vs. thermal enhancement**

Source: STM Corporate Packaging

➢ **BGA package rough** (Cost-performance ÷ High-performance)

    ✓ max power density = 50÷60 W/cm2

    ✓ Cost per pin = 0.25÷1.1 ¢/pin   (~ 90 pins/cm2)

    ✓ Max pincount = 500÷2500+

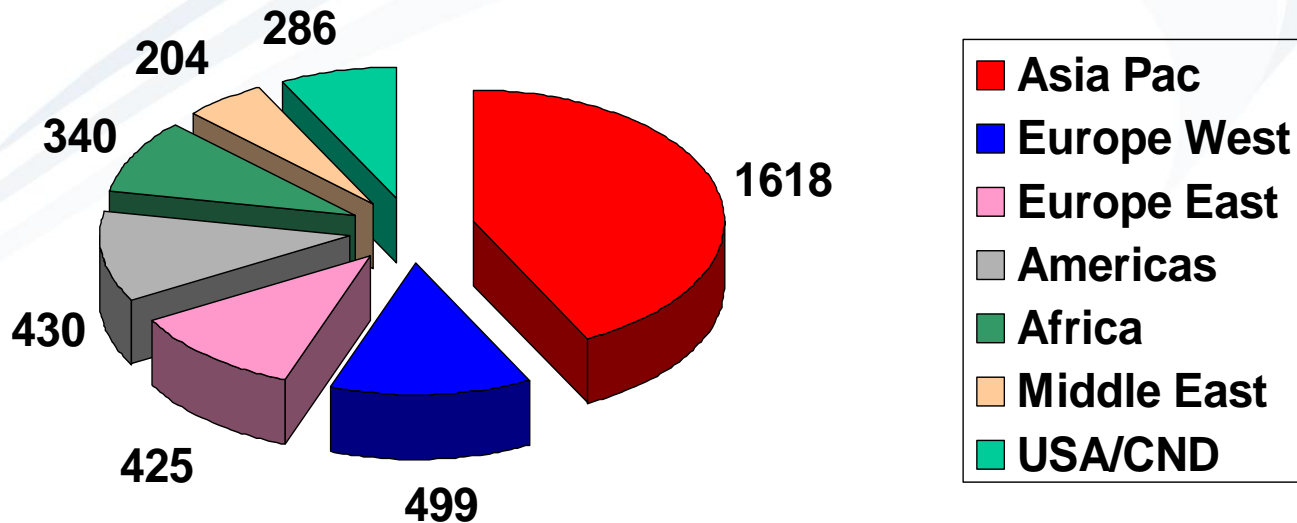# Increased share of Mobile Phone Subscribers

- Cellular Phones: GSM+CDMA
  - ✓ The fastest growing communication technology of all time.
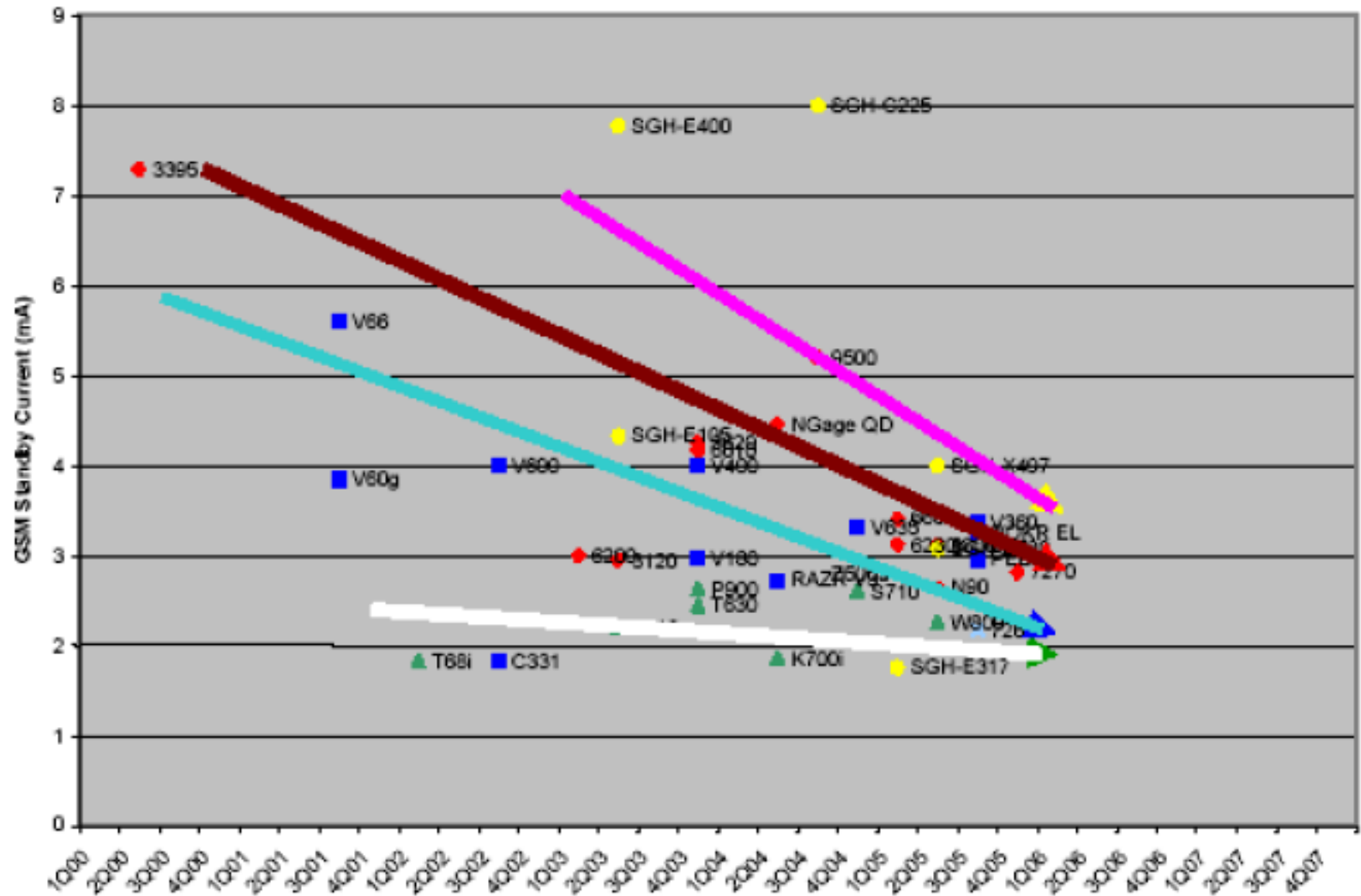- The billionth subcriber user was connected in Q1 2002

# Mobile Phones Regional Split at Q4-2008

➢ 3,804 M subscribers as of Q4-2008

➢ Mobile Broadband Network (HSPA) subscribers has reached 58 M from 11 M on 2007 (i.e. 4 M/Month growth rate).
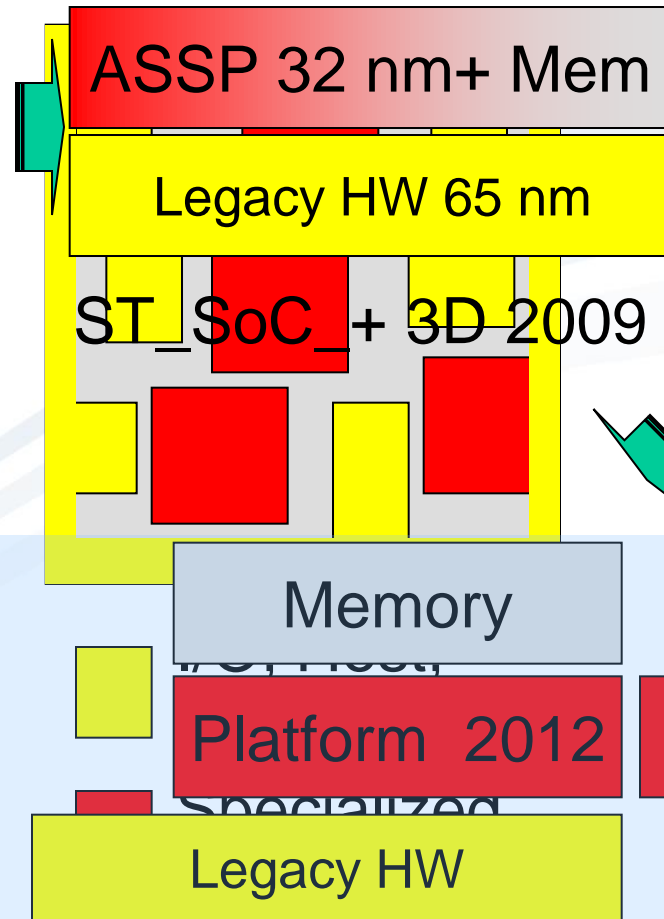
## GSM Regional Statistics Q4-2008



286
204
340
1618
430
425
499

**Legend:**
- Asia Pac
- Europe West
- Europe East
- Americas
- Africa
- Middle East
- USA/CND

# Cellular Phone's standby current

# ST's New SoC Roadmap

ST_SoC 65 nm
Today

Memory

ASSP 32 nm+ Mem

Platform 2012

Legacy HW 65 nm

Legacy HW 65 nm

ST_SoC_+ 3D 2009

Flex_ST_SoC

Memory

Memory

Platform 2012

Platform 2012

Legacy HW

Legacy HW

ST Small SoC

ST Big SoC

# ST's goal

*Replace Heterogeneous mixed HW/SW specialized sub-systems by a single scalable and programmable computing fabric while solving manufacturability issues*
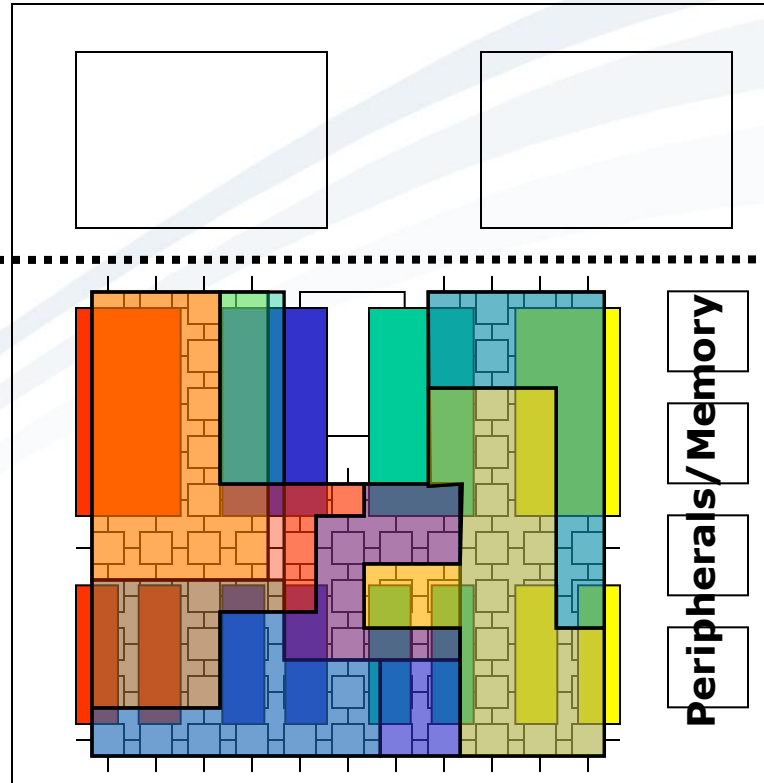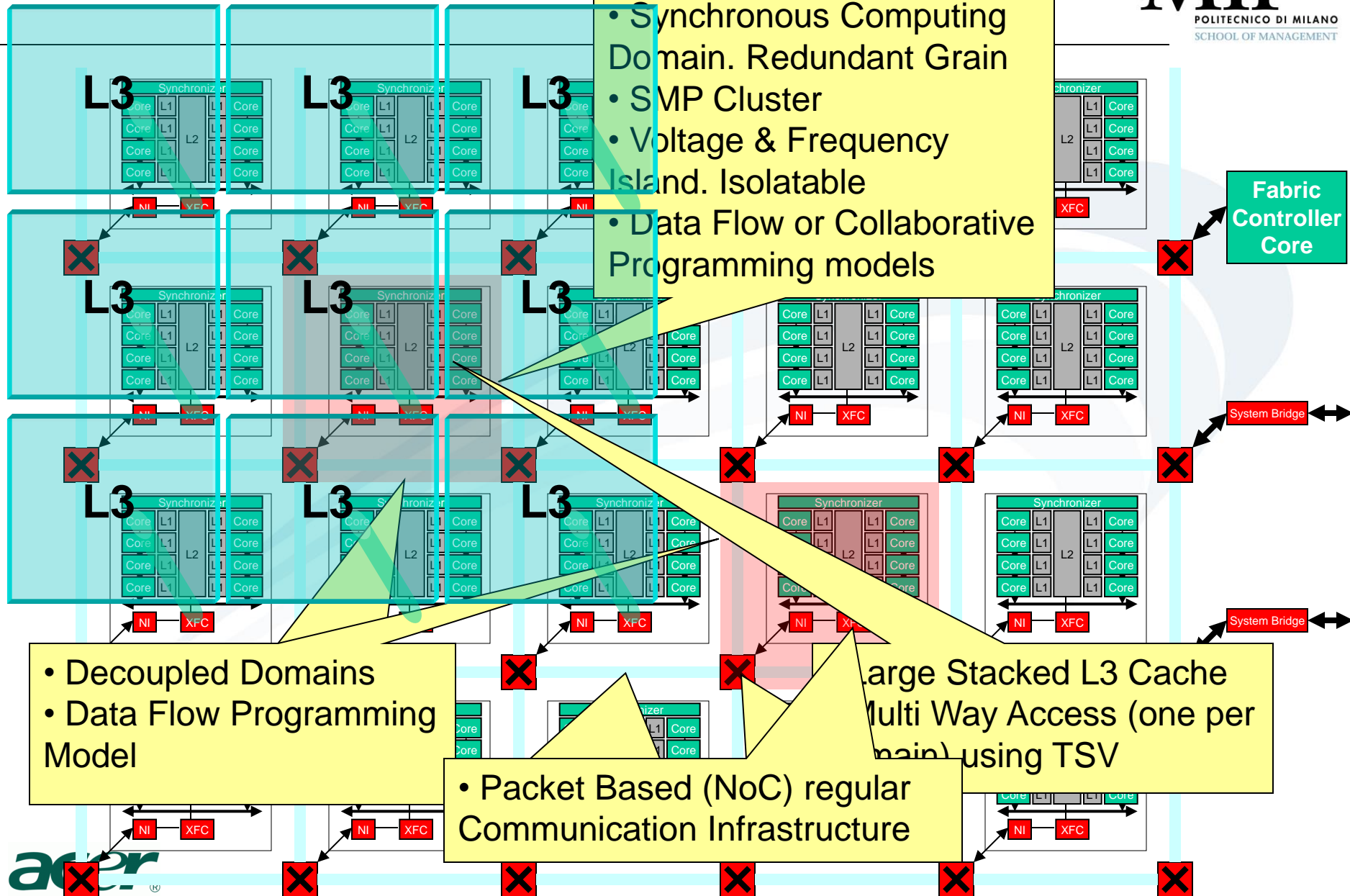
**GP Core, Host**

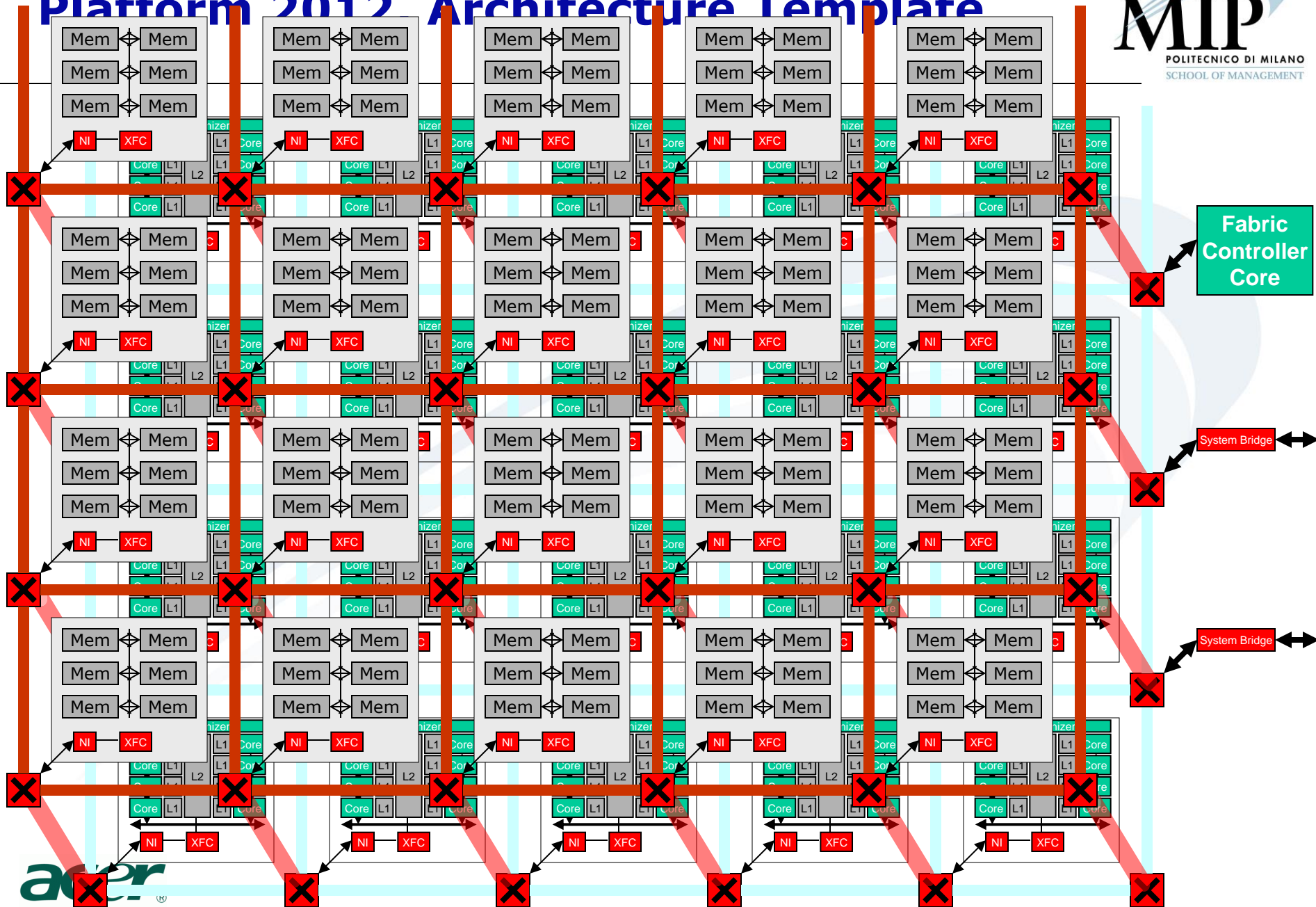**Application Oriented Cores**

**High Density Programmable Fabric**

**Specialized HW**

**Peripherals/Memory**

Physical Sub Systems identical to logical sub systems

Logical Suschystems are not really logical!

# Platform 2012: Architecture Template



- Synchronous Computing Domain. Redundant Grain
- SMP Cluster
- Voltage & Frequency Island. Isolatable
- Data Flow or Collaborative Programming models

- Decoupled Domains
- Data Flow Programming Model

- Large Stacked L3 Cache Multi Way Access (one per domain) using TSV

- Packet Based (NoC) regular Communication Infrastructure

Fabric Controller Core

System Bridge

System Bridge

# Platform 2012: Tool stack, A Constructive Approach

# Platform 2012: SW Stack

**MIP** POLITECNICO DI MILANO

## System

| Run time | IT, I/O mngt | System Monitoring |
|---|---|---|

System Hardware Dependant Software

## Application Execution Engine

• Focus on Dynamic execution for 2nd silicon demonstrator at this level

| | Global Memory Mgt | Inter Communication |
|---|---|---|

Tolerance | Variability | Monitoring

## Fabric

### Fabric Quality of Services (FQoS Middleware)

Fabric Hardware Dependant Software

## Application Execution Engine

• Focus on Dynamic execution for 1st silicon demonstrator at this level

| | Memory Mgt | Communication |
|---|---|---|

Energy | Tolerance | Variability | Monitoring

## Cluster
acer®

### Cluster Quality of Services (CQoS Middleware)

Cluster  Hardware Dependant Software

# Milestones

| | Local Cluster | Fabric 1 | Local Cluster + 3D | Fabric 2 |
|---|---|---|---|---|
| Application | 720p, 250 Gops | 1080p 1 Tops | 1080p 500 Gops | 1080p 2 Tops |
| Programming Model | Streaming, SMP | Streaming, SMP | Streaming SMP Client-server | Streaming SMP Client-server |
| Tools & SW | Local Dynamic | Global Static Local Dynamic | Local Dynamic | Local & Global Dynamic |
| Variability | Local | Local + Global | Local | Local + Global |
| Redundancy | On/Off | On/Off | On/Off | Dynamic |
| 3D | V0 if available *(Not Mandatory)* | V0 | V1 | V1 |

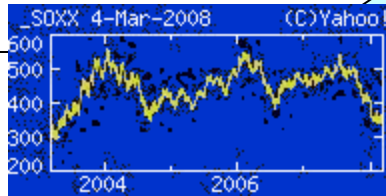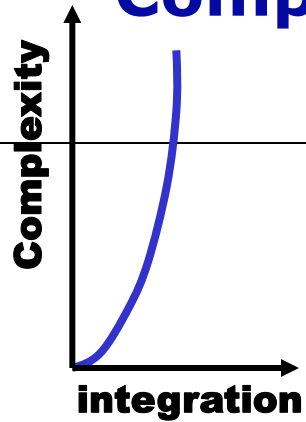Techno Validation
Reusable as an IP

Integrated into product SoC

# … and … not Only Mobile!

- 20% of electrical energy consumed in Amsterdam is used for Telecom

- In the US, Internet is responsible for 9% of the electrical energy consumed nation-wide

  - ✓ This grows to 13% with all computer applications

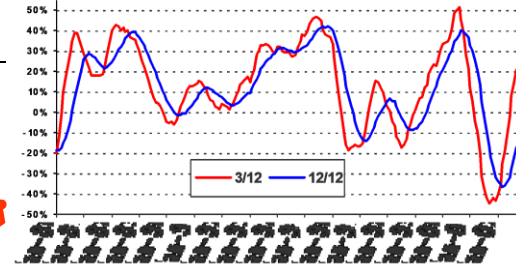- Transfering 2 MBytes of data through the internet consumes the energy of 1 pound of coal (1 pound=0.453 Kg)

Source: 2000 CO2 conference, Amsterdam, NL

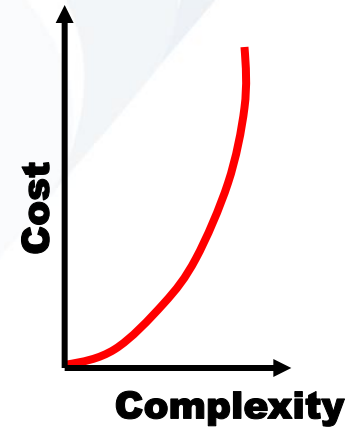# Complexity goes non-linear



**SOXX Finance**

**Markets**

**LITHO & DFM**

**NON LINEAR THINKING IS REQUIRED!!**

**LINEAR APPROACH TO** *non linear* **PROBLEMS!**

**SW complexity**

**IC Design verification**

acer

# Conclusion

- ➢ Semiconductor market is still CMOS dominated:
  - ✓ Switching and leakage power.
- ➢ Leakage will become dominant for technology nodes below 65nm.
  - ✓ Leakage power optimization must be addressed from both technology and design points of view.
- ➢ Multi Processing Platform:
  - ✓ Not yet supported by "production grade" SW tool chain
- ➢ Higher-level approaches are still in their infancy:
  - ✓ Results are promising.

- ➢ The Embedded System's industry calls for a REVOLUTION!

# Industry's Needs

- **Ultra low power systems**
- **Ultra low power cognitive radio**
- **Improve SW productivity**
- **Micro-Nano systems**

  **System In Package**

- **System On Wafer**

  **Many Cores**