4004

1971

8086

1978

PowerPC 601

1992

Pentium 4 Prescott
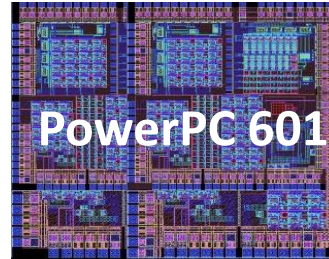
2004

# 45-year CPU Evolution: 1 Law -2 Equations

Daniel Etiemble

LRI – Université Paris Sud

Xeon X7560

2010

Power9

2017

Nvidia Pascal

2016

# Exponential Evolution

- *Are there some fundamental rules?*
  - **Moore's Law**

  - Program execution time                    **Hennessy-Patterson**

$$T_{ex} = IC * (CPIcpu + CPImem) * T_c = \frac{IC}{IPC * F}$$

Instruction count        Avg Cycles/Instruction        Cycle time
                         Computing – Waiting data

  - CMOS power dissipation

$$P_d = V_{dd} * I_{leakage} + \alpha * \sum C_i * V_{dd}^2 * F$$

# Moore's law



**N increases**

**Nb transistors/chip doubles every N months (12/18/24)**

Daniel Etiemble

# Technological nodes

**CMOS technological nodes (nm)**



- From one node to the next one
  - **(rough approximation)**
  - Gate delay /1.4  =>  Higher clock frequencies
  - Increase of transistor count per area unit.

- Consequences
  - From multi-chips to one chip
  - Larger on-chip memories
  - More functionalities per chip

But MISMATCHES



Cache hierarchies

# Evolution of cache hierarchies



**386**

**Pentium**

**Pentium 4**

**Power6**

*Reducing CPImem*
**Latency & Bandwidth**

**Power 8**

**NUCA**

# Improving Performance $\quad T_{ex} = \dfrac{IC}{IPC * \textcolor{red}{F}}$

## Increase $\textcolor{red}{F}$ (Technological nodes)

- Gate delay /1.4  => Higher clock frequencies
- Increased Performance

## Increase F

- Increased Power and Power density

$$P_d = V_{dd} * I_{leakage} + \alpha * \sum C_i * V_{dd}^2 * F$$

**HEAT WALL**



2017 : F in the 3-4 GHz range
Except water cooled IBM z14 CPU (5.2 GHz)

# Improving Performance $T_{ex} = \dfrac{IC}{\textbf{IPC} * \textbf{F}}$

**ILP in mono-processor CPU**

- Scalar CPU (IPC <1)
  - Pipeline and superpipeline
- Superscalar CPU (IPC >1)
  - In-order CPUs
  - **Out-of-order CPU**
- VLIW CPU (IPC >1)

- Intrinsic limits of ILP in a sequential program & " HW diminishing return".
  - Larger buffer to extract µops to launch, but launching width remains constant !



Legend: ROB, µop/cycle, Reservation stations, Int. Physical Registers

**ROB Renaming**     **Register renaming**

3 to 4 µops/clock

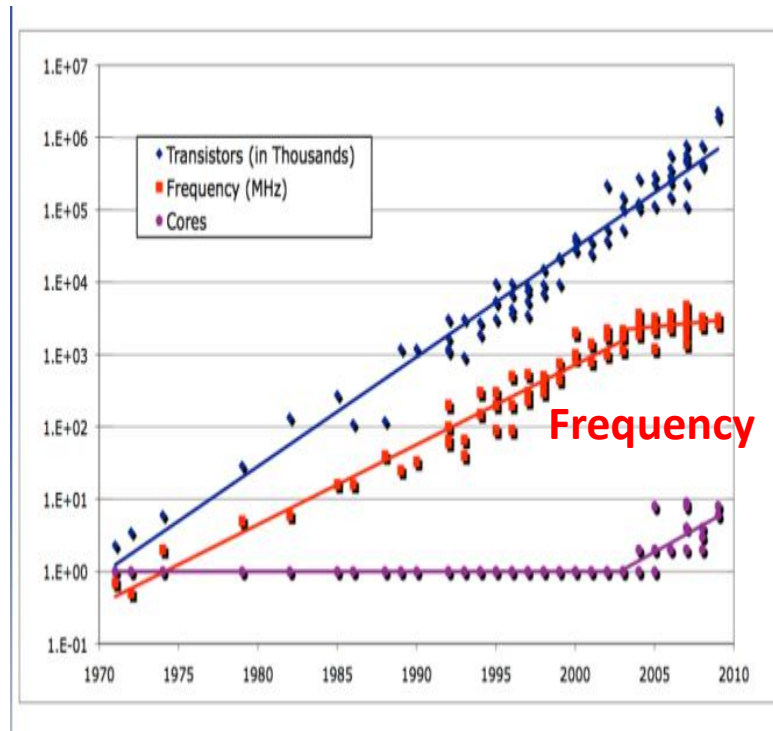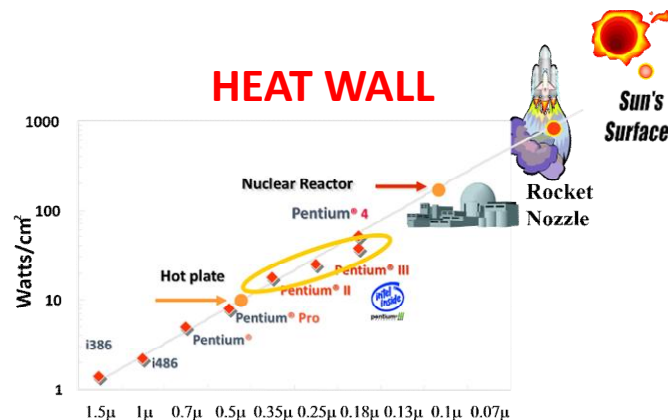| | | | |
|---|---|---|---|
| 256 | | | |
| 64 | | | |
| 16 | | | |
| 4 | | | |
| 1 | **1995** | **2000** | **2011** | **2013** |
| | PENTIUM PRO | PENTIUM4 | SANDY BRIDGE | HASWELL |

# Improving Performance $T_{ex} = \dfrac{IC}{IPC * F}$

## Data Parallelism in mono-processor CPU

CPU

SIMD
1 instruction with several data

SSE2/3/4 – Neon – Altivec
AVX – AVX2 – AVX 512…

| X3 | X2 | X1 | X0 | source 1/dest. |

| Y3 | Y2 | Y1 | Y0 | source 2 |

op   op   op   op

| X3 op Y3 | X2 op Y2 | X1 op Y1 | X0 op Y0 | source 1/dest. |

### Instruction Decoder and Warp Scheduler

GPU

SIMT
1 instruction for several threads

CUDA Core — Dispatch Port — Operand Collector — FP Unit — INT Unit — Result Queue

registers

thread

# CPU + GPU

- Two different programming models
- Two chips or one chip?
  – CPU + GPU or APU



| 2014 | Tesla K40 + CPU | Nvidia Tegra K1 |
|---|---|---|
| Single Precision Peak | 4.2 TeraFlops | 326 GFlops |
| Single Precision SGEMM | 3.8 TeraFlops | 290 GFlops |
| Memory | 12GB @ 288GB/s | 2GB @ 14.9GB/s |
| Power (CPU + GPU) | ~ 385Watt | <11Watts |
| Performance Per Watt | 10SP GFlops Per Watt | 26SP GFlops Per Watt |

# Improving Performance $T_{ex} = \dfrac{IC}{IPC * F}$

**Data or/and Thread Parallelism in Parallel Architectures**

**Dispatch IC among several processors (or cores)**
- Except for simple cases or embarrasingly parallel applications, the dispatch depends on the architecture, the programming model, Amdahl law, etc.
- In some architectures, communication times must be included.

OpenMP
Pthreads

MPI

| CPU | CPU | CPU | CPU |

Bus or Interconnection

Memory

Interconnection

| CPU | CPU | CPU | CPU |

Memory Memory Memory Memory

Multiprocessors        Multi-computers

Switching from sequential to parallel programming
- Limited to servers and super-computers during the « free lunch » period

# Free lunch… (by Intel)
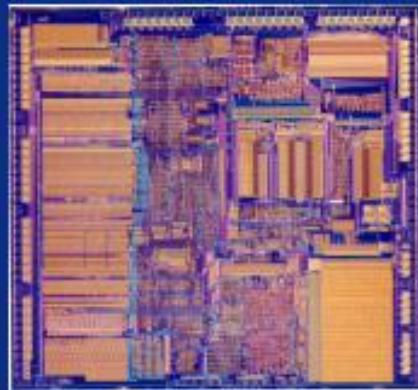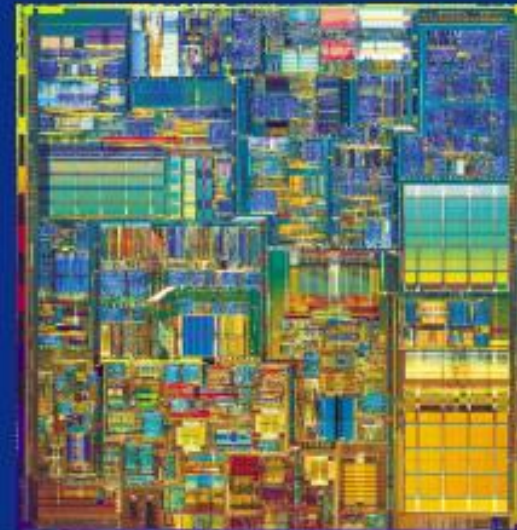
**F**

**IPC**

**IC (SIMD)**



## Scaling at its best

**386 Processor**

**Pentium® 4 Processor**

May 1986
@16 MHz core
275,000 1.5μ transistors
~1.2 SPECint2000

17 Years
200x
200x/11x
1000x

August 27, 2003
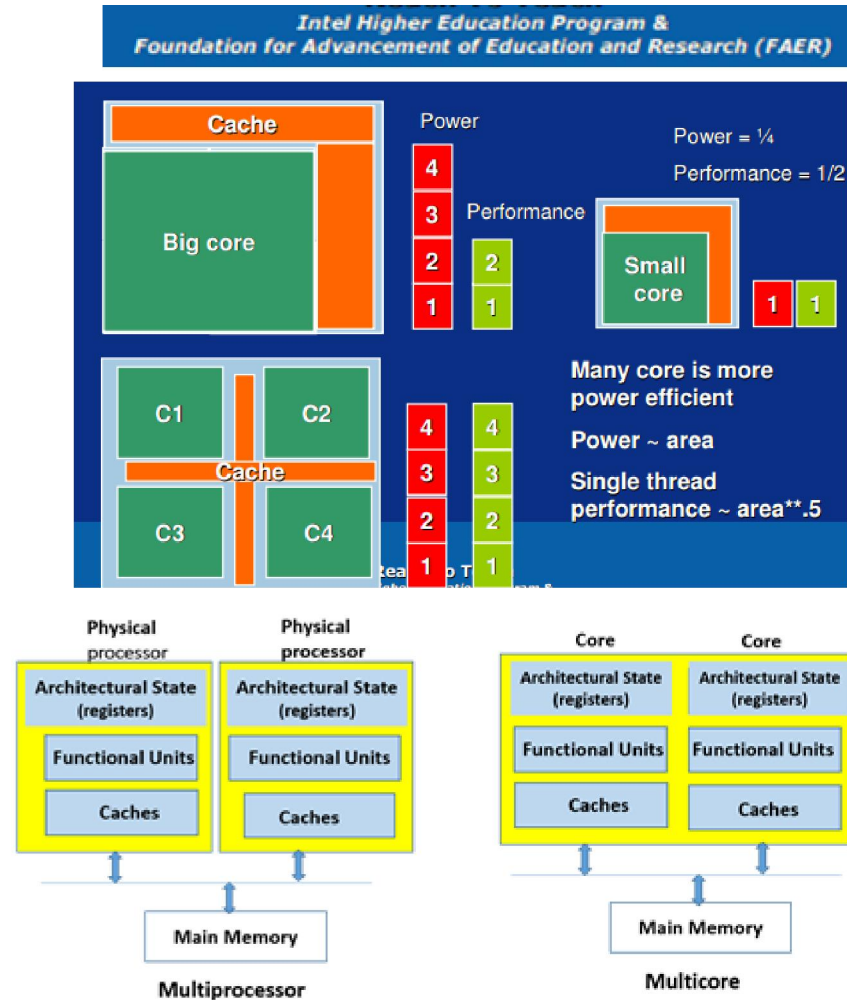@3.2 GHz core
55 Million 0.13μ transistors
1249 SPECint2000

**Reach To Teach**
*Intel Higher Education Program &*
*Foundation for Advancement of Education and Research (FAER)*

9

Daniel Etiemble

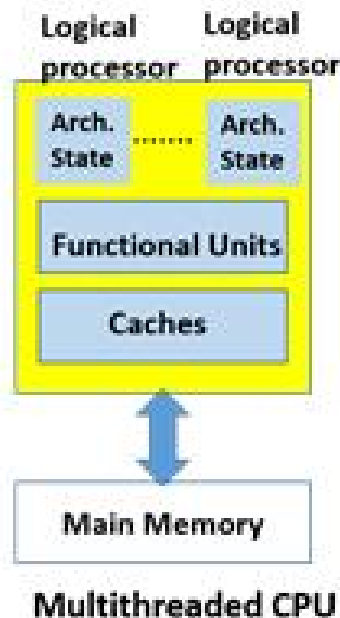# From mono-processors to multi-cores

- Performance – Power trade-off

- Intrinsic limitations of mono-processors, even multithreaded

- Multi-processors to multi-cores

- Clusters of multi-cores.

# Multithreaded CPUs

$$T_{ex} = \frac{IC}{IPC * F}$$

- Sequential programs
- Several programs (multiprogramming): $IC = \sum IC_i$
- Several threads (TLP) : $IC = \sum IC_i$

Fine grain multithreading
- Switching in one clock from one thread to next one on pipeline hazards
- Sun Niagara, Oracle servers

Simultaneous multithreading
- Issuing instructions from different threads at each clock
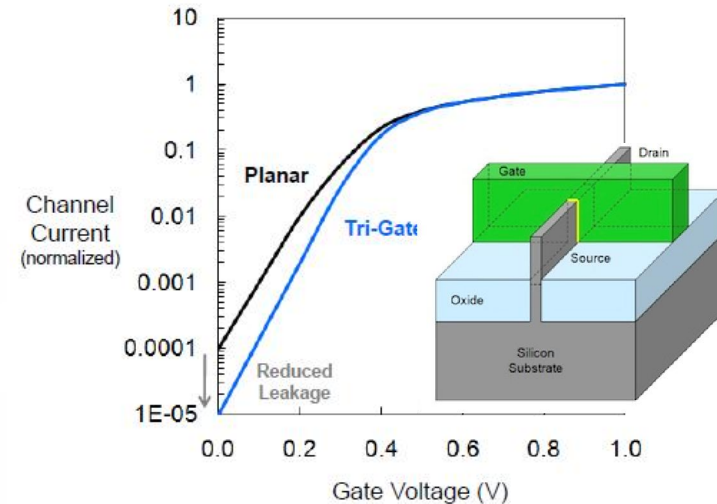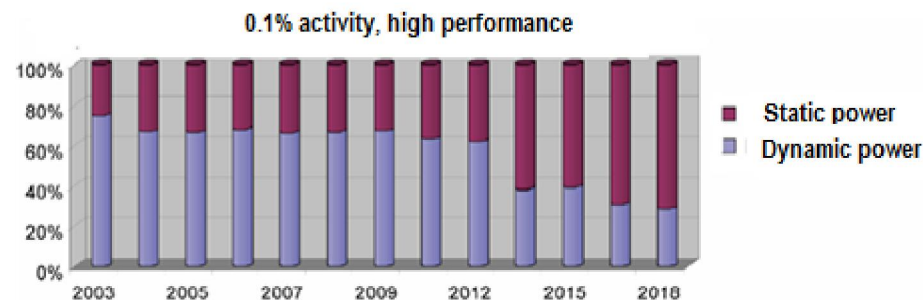- Intel Hyperthreading (2), IBM Power (2 to 8)

**Multithreading reduces $CPI_{Mem}$**

**Multi-cores with multithreaded cores**

Logical processor    Logical processor

| Arch. State | ...... | Arch. State |

Functional Units

Caches

Main Memory

**Multithreaded CPU**

# Reducing static power dissipation

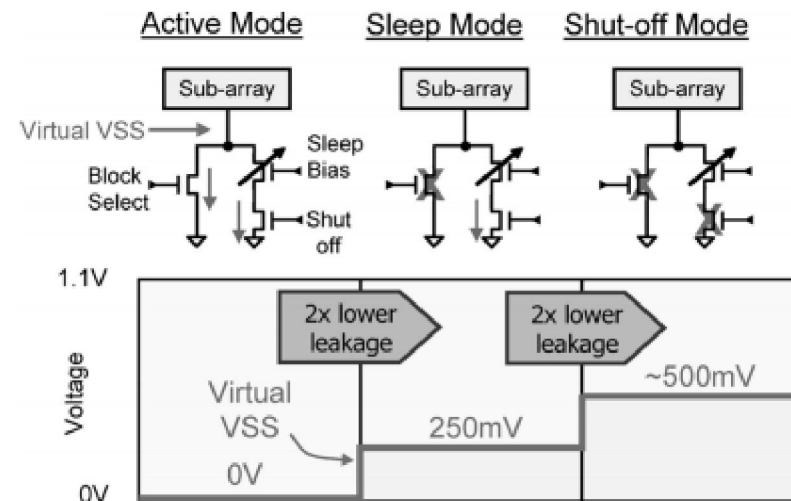$$P_{d\ static} = V_{dd} * I_{leakage}$$





**TECHNOLOGY**
- Ex: Intel Tri-gate

**CIRCUITERY**
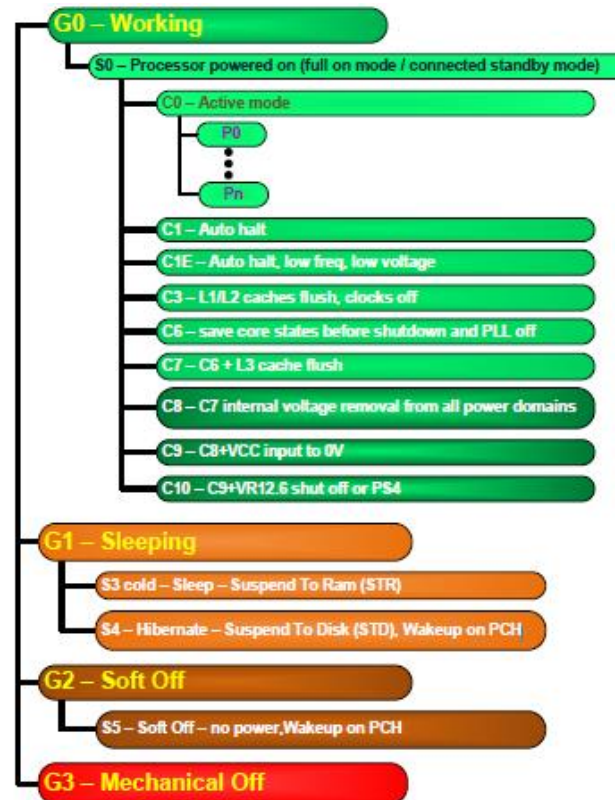- Ex: Virtual Vss
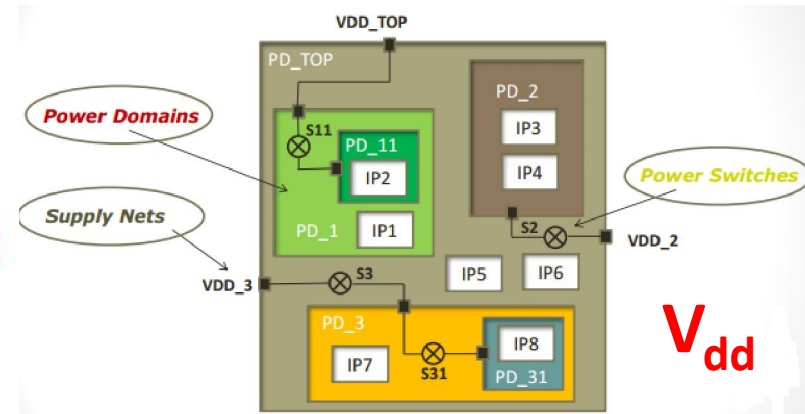65-nm Xeon CPU L3 cache

# Reducing dynamic power dissipation

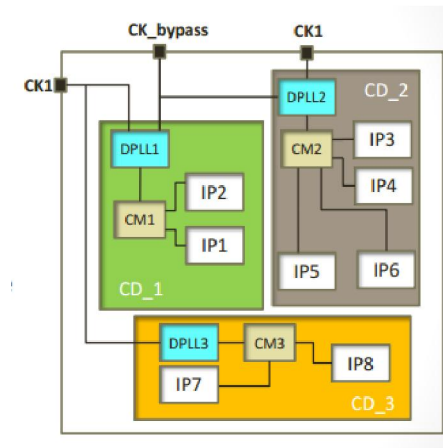$$P_{d\,dynamic} = \alpha * \sum C_i * V_{dd}^2 * F$$

$\alpha$

Several Operating
Modes per blocks:
Ex: 5th generation
Intel Cores



**Power domains**

$V_{dd}$



**Clock domains**

$F$

# What about future?

- Moore's law: towards fundamental limits
- Execution time:
  - F: significant changes are doubtful
  - IPC:
    - Limits of ILP in cores
    - New PIM architectures ? (Data access issues: "Memory Wall")
  - IC continues to decrease
    - more job per instruction
      - SIMD width (256-512-1024…) and data size (16-bit FP)
      - New 2D instructions: Tensor cores (Nvidia), Matrix Multiplication Unit (Google TPU)…
    - More cores
      - From multi-cores to many-cores. Exponential increase of core number???
- Power dissipation
  - As long as CMOS will be used…

# Concluding remarks

- **Main trends (**not details) of CPU evolution can be explained by
  - Moore's law
  - $T_{ex} = IC * (CPICPU + CPIMem) * T_c = \dfrac{IC}{IPC * F}$
  - $P_d = V_{dd} * I_{leakage} + \alpha * \sum C_i * V_{dd}^2 * F$

- Valid for software programmable processors as long as CMOS technology will be used.

- Mixed HW-SW architectures (FPGA) are more complex to modelize.

- Evolution of CPU architectures is driven by new specific applications (AI, IoT, ...).

- Business... (as usual). Ex: "proprietary" versus "open-source" ISAs.