# Parallel Programming and Heterogeneous Computing
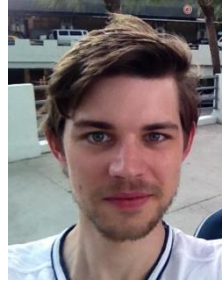
Introduction

Max Plauth, Sven Köhler, Felix Eberhardt, Lukas Wenzel and Andreas Polze

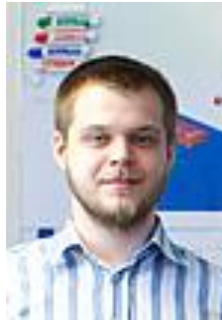Operating Systems and Middleware Group

# Teaching Team (PhD Students)

- **Max Plauth**
  - GPU Accelerators
  - Virtualization of Accelerators

- **Felix Eberhardt**
  - NUMA Topologies
  - Dynamic Scale-Up/ Down

- **Sven Köhler**
  - On-Core Accelerators (SIMD, Crypto)
  - Programming Models for Accelerators

- **Lukas Wenzel**
  - Computer Architecture
  - FPGA Accelerators

**ParProg 2020 Introduction**

Chart **2**

# Course Topics

A. The Parallelization Problem

- Power wall, memory wall, Moore's law
- Terminology and metrics

B. Shared Memory Parallelism

- Theory of concurrency, hardware today and in the past
- Programming models, optimization, profiling

C. Heterogeneous Computing

- On-Chip Accelerators (e.g. SIMD, special purpose accelerators, etc.)
- External Accelerators (e.g. GPUs, FPGAs, etc.)

D. Shared Nothing Parallelism

- Theory of concurrency, hardware today and in the past
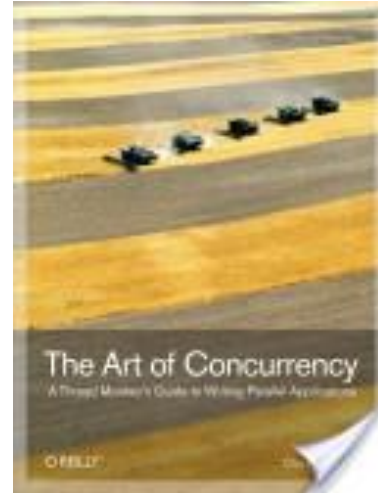- Programming models, optimization, profiling

**ParProg 2020
Introduction**

Chart **3**

# Course Design

*This is a course about concepts, not a programming tutorial!*

- Weekly lectures
- Practical assignments
- Oral exam ~30 minutes

- Literature list on course home page
- Based on a course by Prof. Dr. Peter Tröger



*The Art of Concurrency:
A Thread Monkey's Guide to
Writing Parallel Applications*

Clay Breshears
O'Reilly Media, Inc.
2009

# Course Design
# Lectures

- Lectures are given in flipped classroom style
  - 75% time: Videos delivering lecture content
  - 25% time: Online meetings for recap sessions, questions and discussions
- Videos are posted every Wednesday on the website
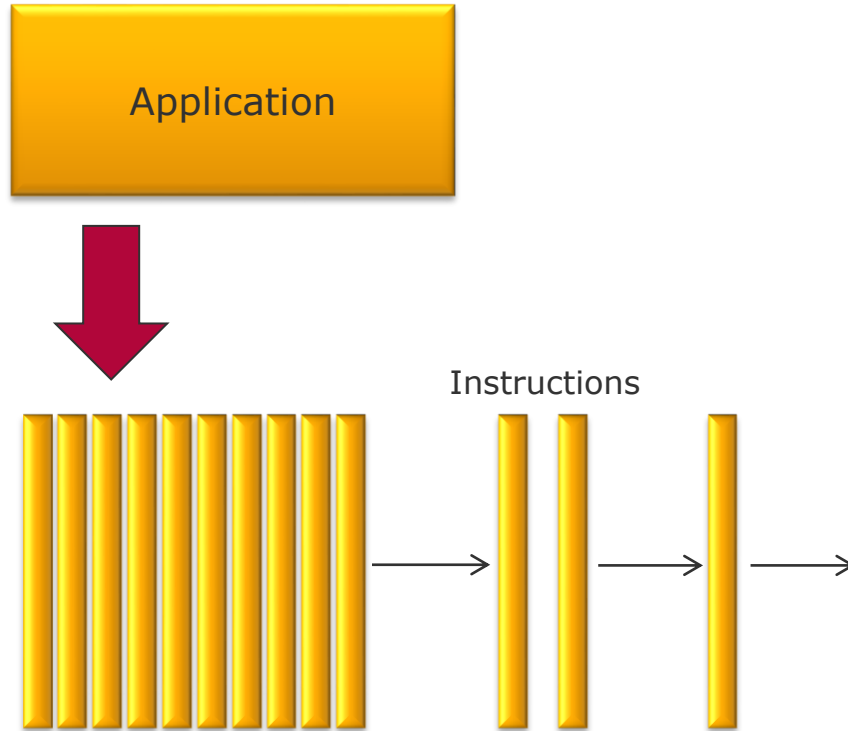- Meeting links for Thursday will also appear on the website

https://osm.hpi.de/parProg/2020/

**ParProg 2020**
**Introduction**

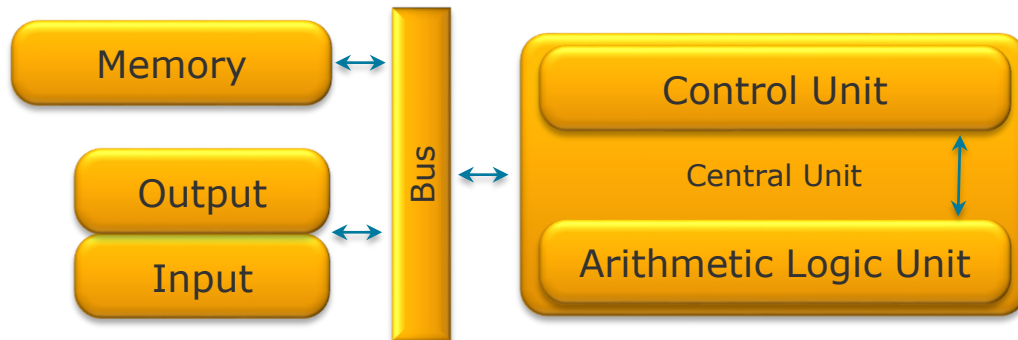Chart **5**

# Course Design
# Assignments

- Implementation of parallel algorithms to demonstrate the use of programming models discussed in the lecture

- Presented every ~3 Weeks

- Solved in teams of 2 persons

- Submission to https://osm.hpi.de/submit/

- At least 50% correct solutions enable participation in final exam

**ParProg 2020**
**Introduction**

Chart **6**

# Running Applications
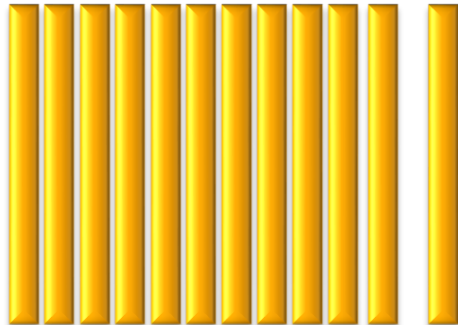


Application

Instructions

# Machine Model

- First computers had fixed programs (electronic calculator)
- **von Neumann architecture** (1945, for EDVAC project)
  - Instruction set for control flows stored in memory
  - Program is treated as data, which allows the exchange of code during runtime and self-modification
  - Introduced the **von Neumann bottleneck**
- CPUs are built from logic gates, which are built from transistors
- Multiple CPUs (SMP) were always possible, but exotic

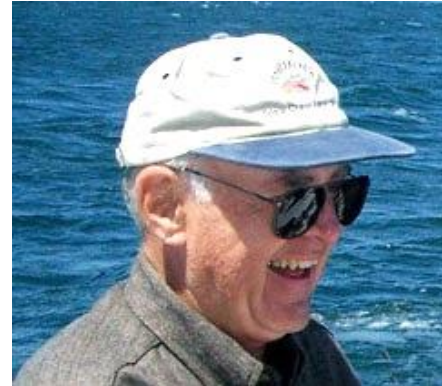# Three Ways of Doing Anything Faster [Pfister]

Application

Instructions

- Work Harder
  (clock speed)

- Work Smarter
  (optimization, caching)
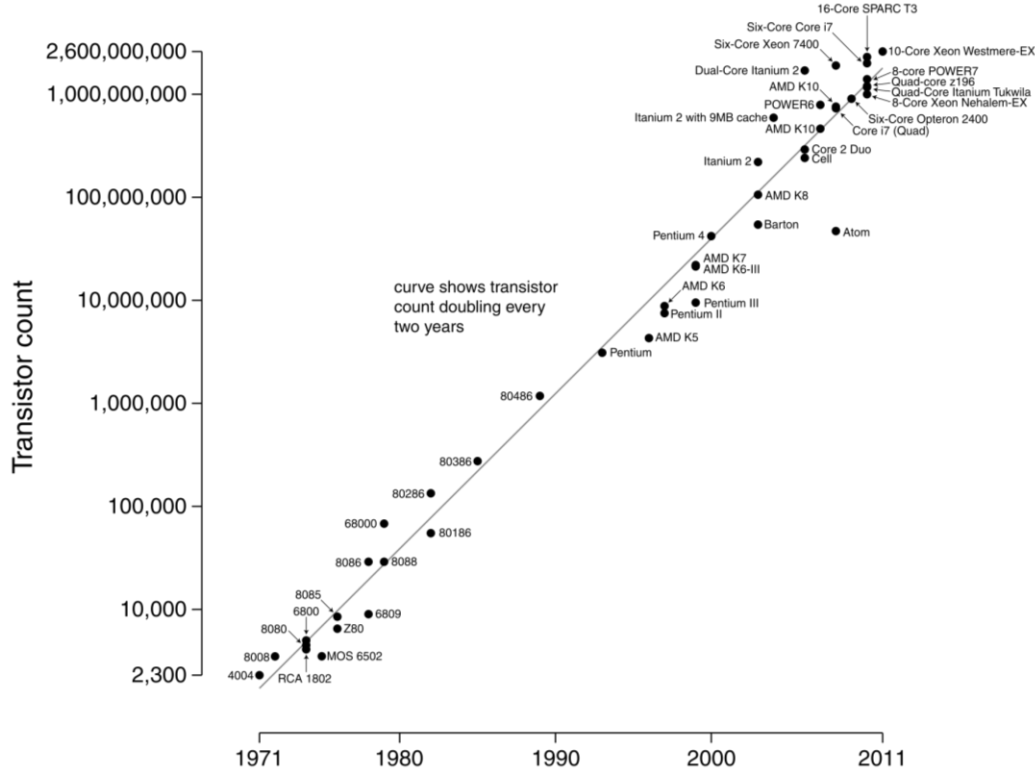
- Get Help
  (parallelization)

# Moore's Law

- *„...the number of transistors that can be inexpensively placed on an integrated circuit is increasing exponentially, doubling approximately every two years. ...“ (Gordon Moore, 1965)*

  - Rule of exponential growth
  - Applied to many IT hardware developments
  - Sometimes misinterpreted as performance indication
  - Has become a self-fulfilling prophecy
  - Comes to an end within the next 5-10 years



**ParProg 2020**
**Introduction**

Chart **10**
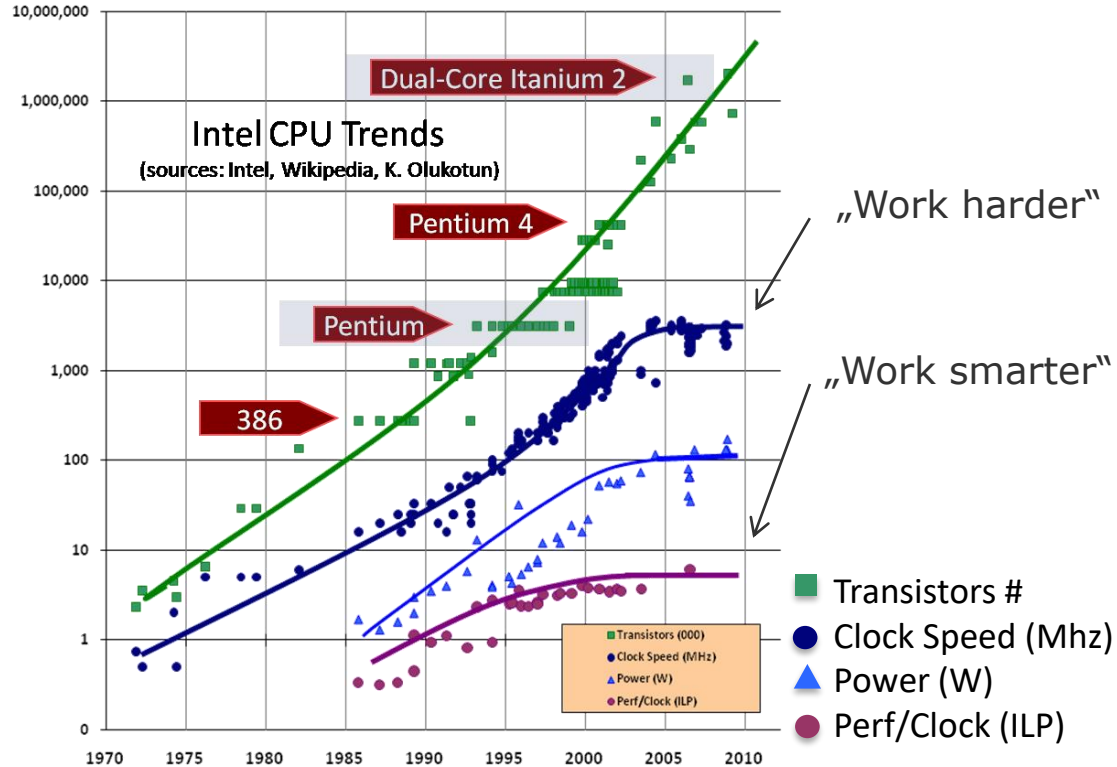
# Moore's Law

**ParProg 2020
Introduction**

Chart **11**

# Moore's Law vs. Software

- Gate's law: *"The speed of software halves every 18 months."*

- Wirth's law: *"Software is getting slower more rapidly than hardware becomes faster."*

- May's law: *"Software efficiency halves every 18 months, compensating Moore's Law."*

- Jevons paradox:
  *"Technological progress that increases the efficiency with which a resource is used tends to increase (rather than decrease) the rate of consumption of that resource."*

- Zawinski's Law of Software Envelopment:
  *"Every program attempts to expand until it can read mail.
  Those programs which cannot so expand are replaced by ones which can."*

# Processor Speed Development

[Herb Sutter, 2009]



Intel CPU Trends
(sources: Intel, Wikipedia, K. Olukotun)

Dual-Core Itanium 2

Pentium 4

Pentium

386

„Work harder"

„Work smarter"

- Transistors #
- Clock Speed (Mhz)
- Power (W)
- Perf/Clock (ILP)

Transistors (000)
Clock Speed (MHz)
Power (W)
Perf/Clock (ILP)
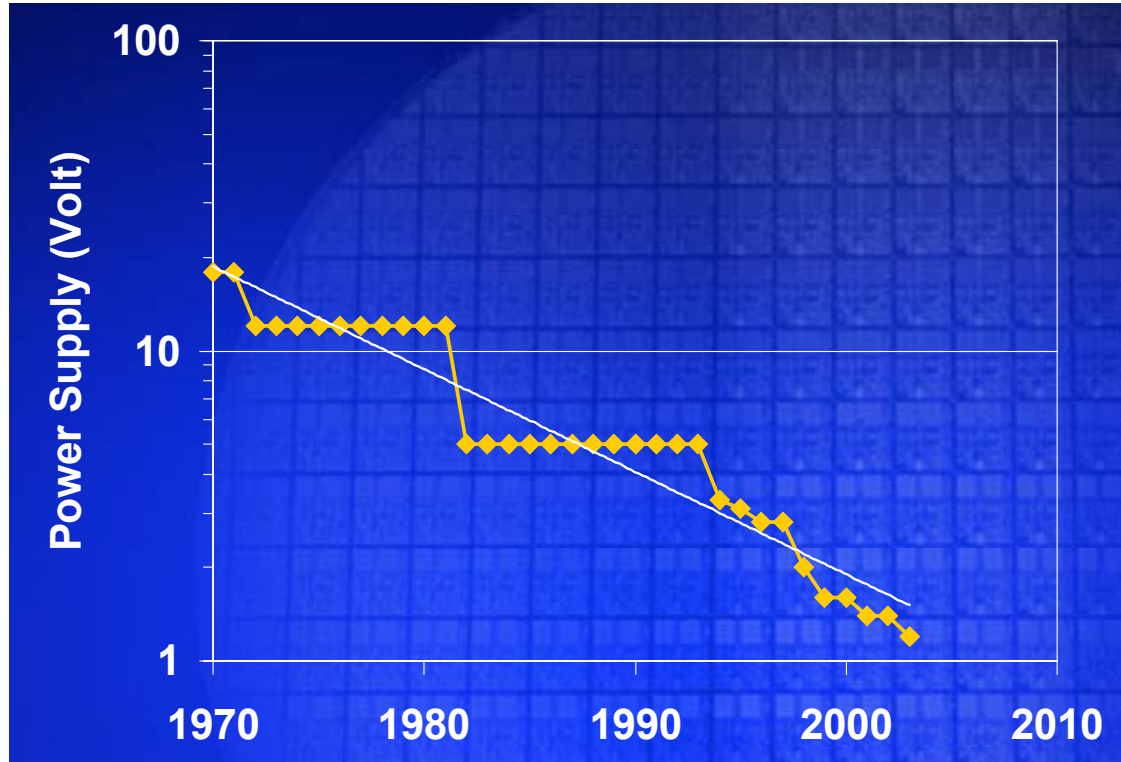
**ParProg 2020
Introduction**

Chart **13**

# A Physics Problem

- Power: Energy needed per time unit
  - Power density: Watt/mm$^2$ → Cooling
- **Static power**: Leakage of transistors while being inactive
- **Dynamic power**: Energy needed to switch a gate

**Dynamic Power ~
Number of Transistors (N) x Capacitance (C) x
Voltage$^2$ (V$^2$) x Frequency (F)**

- Moore's law: N goes up exponentially, C goes down with the size
- The trick
  - Bringing down V reduces energy consumption, quadratically
  - Don't use all the N for gates (e.g. caches)
  - Keeps the dynamic power increase moderate
  - We can happily increase F with N for faster computation
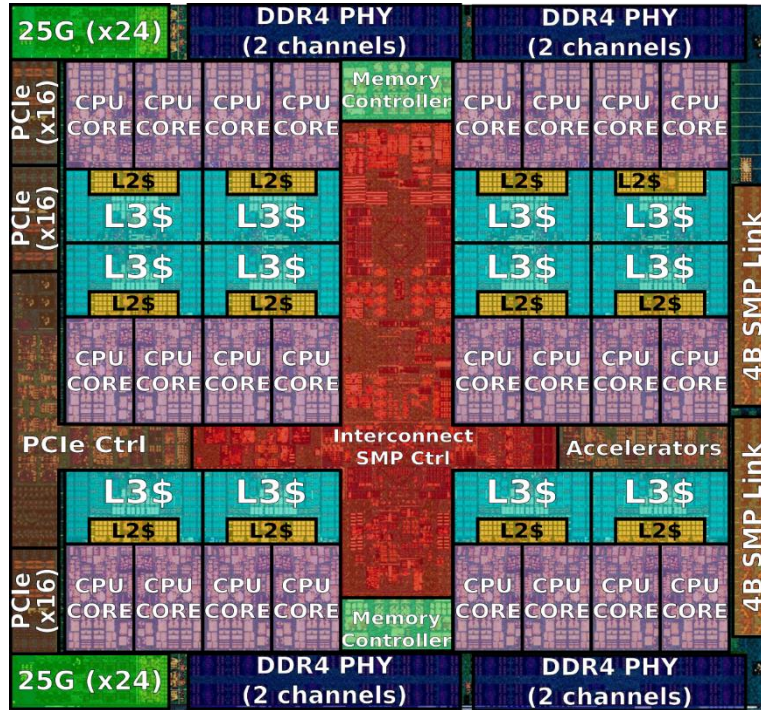
# Processor Supply Voltage

[Moore, ISSCC]



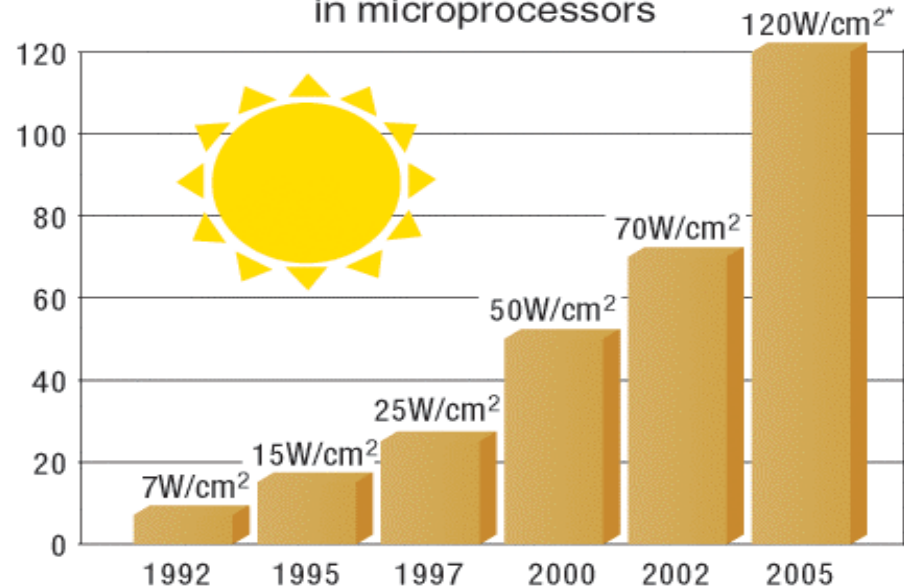**ParProg 2020
Introduction**

Chart **15**

# Transistor Usage



[https://en.wikichip.org/wiki/ibm/microarchitectures/power9]

**ParProg 2020
Introduction**

Chart **16**

# Power Density

**SIZZLING SEMICONDUCTORS**
growth of watts per square centimeter
in microprocessors

7W/cm² — 1992
15W/cm² — 1995
25W/cm² — 1997
50W/cm² — 2000
70W/cm² — 2002
120W/cm²* — 2005

*Could be higher, depends on level of integration.
SOURCE: HEWLETT-PACKARD LABS

# Power Density = Temperature

1st CPU     2nd CPU

[source: Devgan'05]



Power 4 server chip

cache

thermal profile during runtime

- Higher temperature leads to
  - Increased transistor leakage
  - Decreased transistor speed
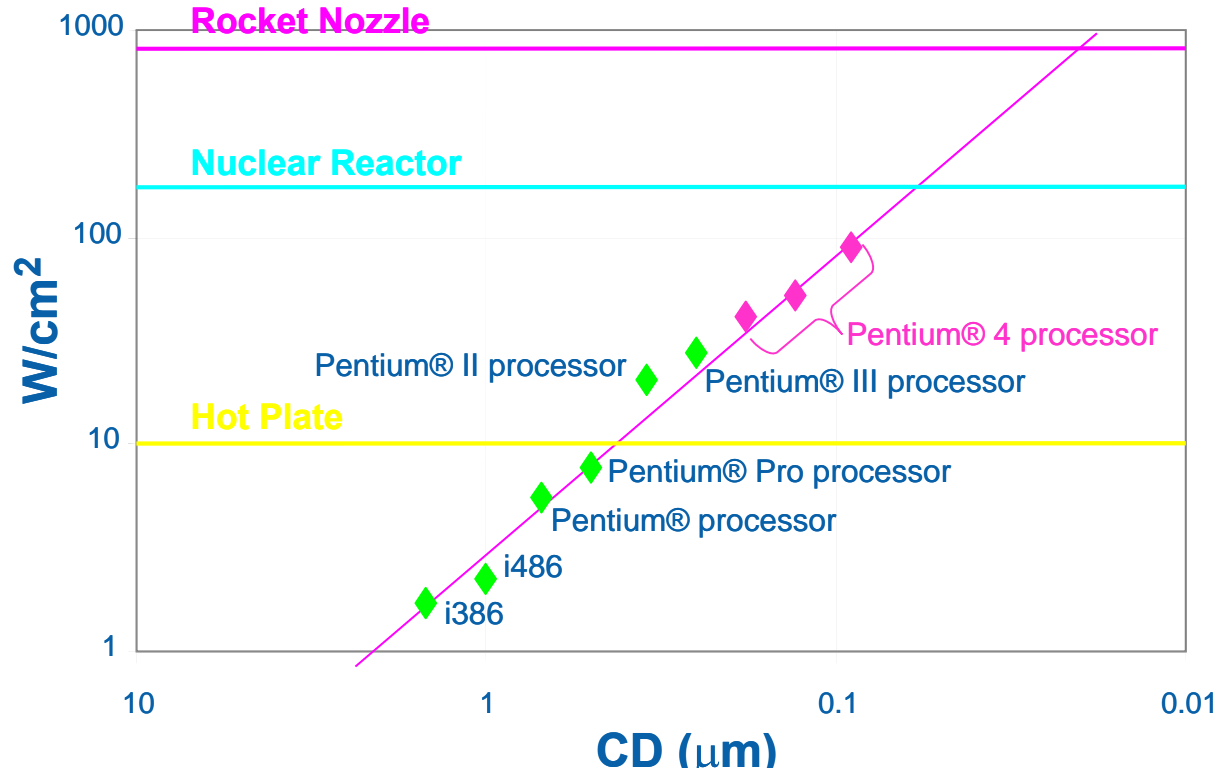  - Higher failure probability

# Power Density



[Kevin Skadron, 2007]

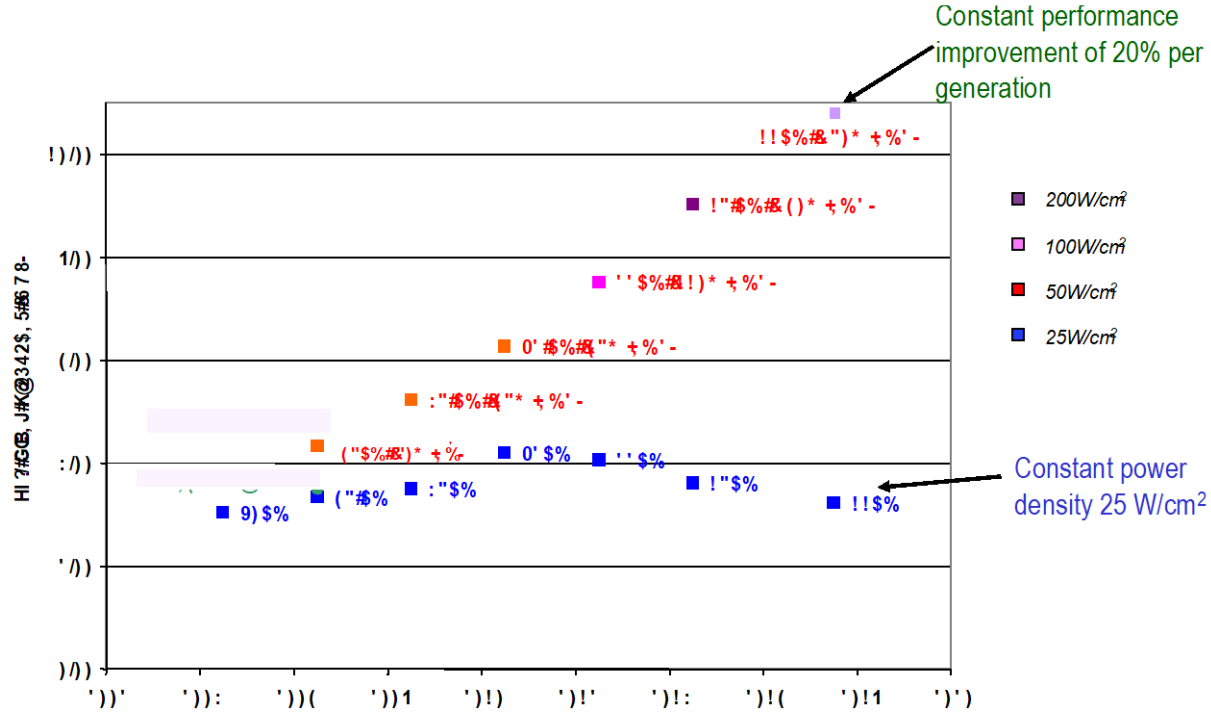**ParProg 2020
Introduction**

Chart **19**

# Power Density

[Taylor, 2009]



**Rocket Nozzle**

**Nuclear Reactor**

Pentium® II processor

Pentium® III processor

Pentium® 4 processor

**Hot Plate**

Pentium® Pro processor

Pentium® processor

i486

i386

$W/cm^2$

CD ($\mu$m)

1000
100
10
1

10   1   0.1   0.01

**ParProg 2020
Introduction**

Chart **20**

# Power Density



Constant performance improvement of 20% per generation

Constant power density 25 W/cm$^2$

Legend:
- 200W/cm$^2$
- 100W/cm$^2$
- 50W/cm$^2$
- 25W/cm$^2$

Source: *D. Frank, C. Tyberg,* IBM Research

# Second Problem: Leakage Increase

# A Physics Problem
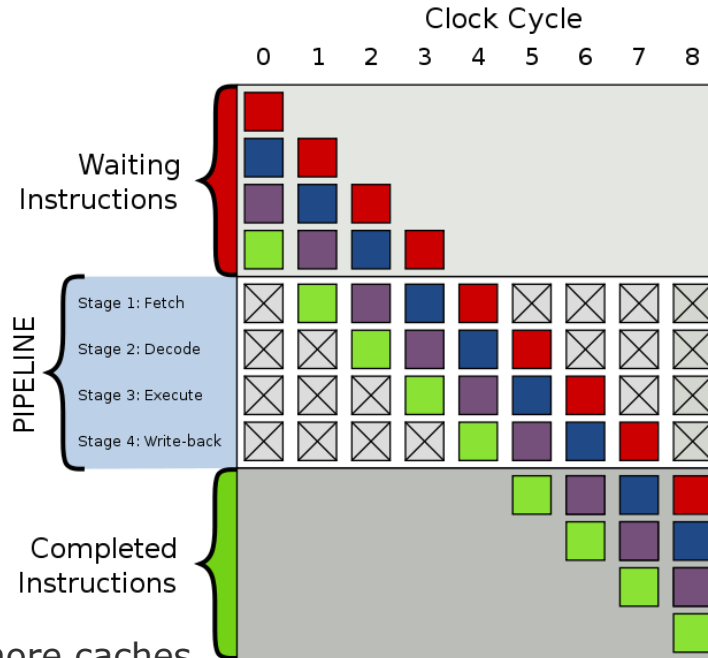
**Dynamic Power = N x C x V$^2$ x F**

- Even if we would keep F constant
  - N continues to increase exponentially → dynamic power
  - Increasing N sums up to more leakage → static power
- Cooling performance is constant (100-125 Celsius)
  - Static and dynamic power consumption has a limit
- Further reducing V for compensating an additionally increased F
  - Also makes the transistors slower
  - We can't do that endlessly, 0V is the limit
  - Strange physical effects
- Increasing the frequency is no longer possible
  → **"Power Wall"**
- Ok, so let's use the additional N for smarter processors

# Instruction Level Parallelism

- Increasing transistor count was also used for more gate logic in **instruction level parallelism (ILP)**
  - **Instruction pipelining**
    - Overlapped execution of serial instructions
  - **Superscalar execution**
    - Multiple execution units are used in parallel
  - **Out-of-order execution**
    - Reorder instructions that have no data dependency
  - **Speculative execution**
    - Control flow speculation, memory dependence prediction, branch prediction
- Today's processors are packed with ILP logic
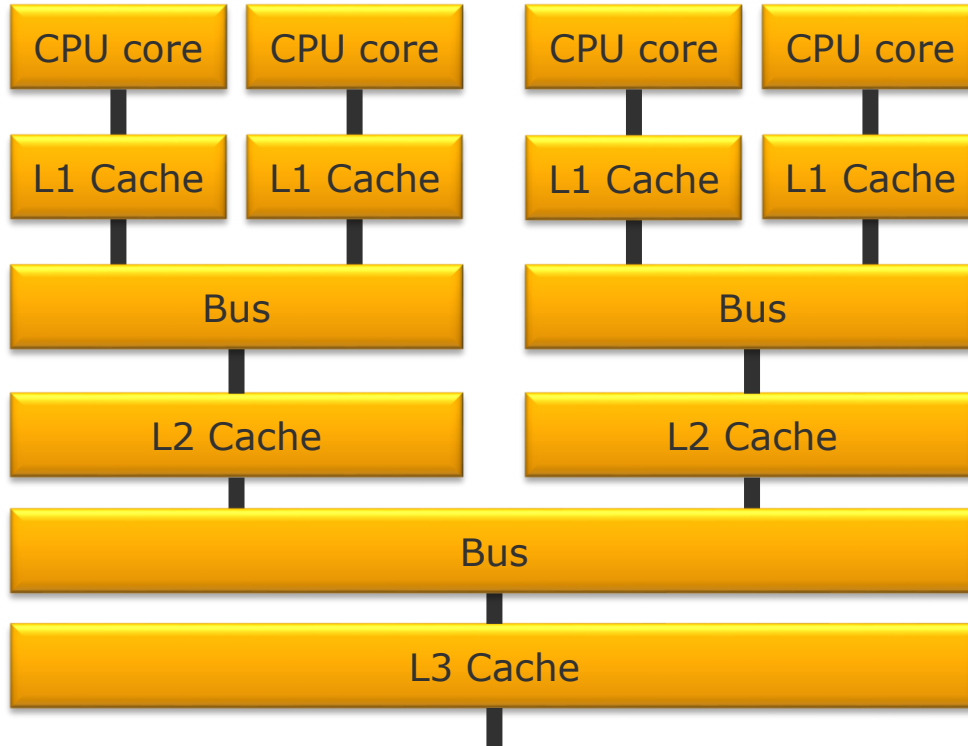
# The ILP Wall

- No longer cost-effective to dedicate new transistors to ILP mechanisms

- Deeper pipelines make the power problem worse

- High ILP complexity effectively reduces the processing speed for a given frequency (e.g. mispredictions)

- More aggressive ILP technologies too risky for products due to unknown real-world workloads

- → **"ILP wall"**

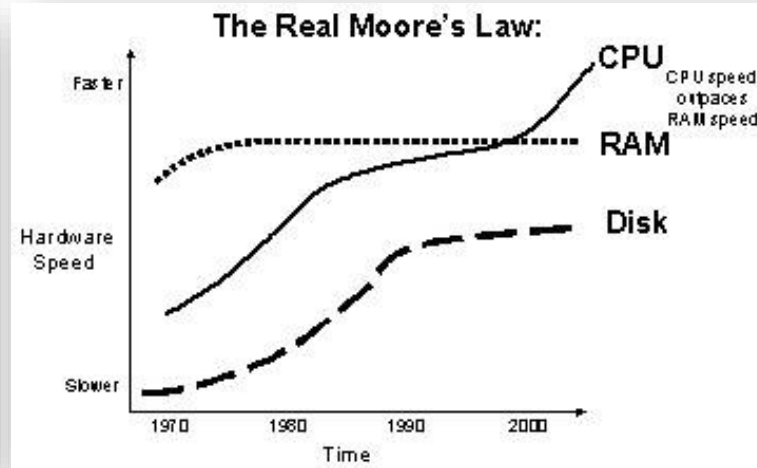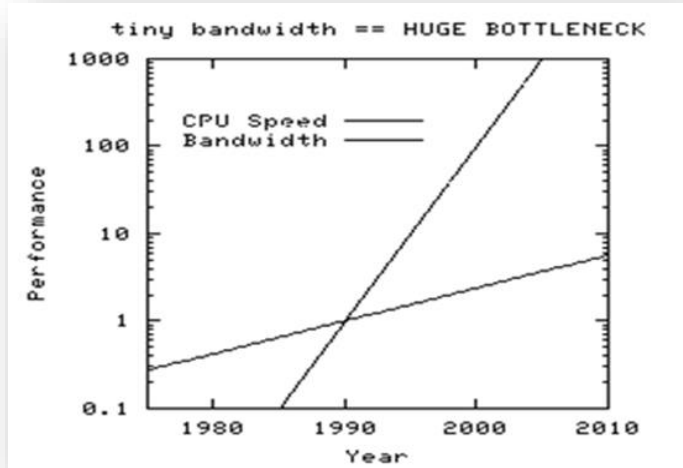- Ok, so let's use the additional N for more caches

[Wikipedia]



**ParProg 2020
Introduction**

Chart **25**

# Memory Hierarchy

# Memory Wall
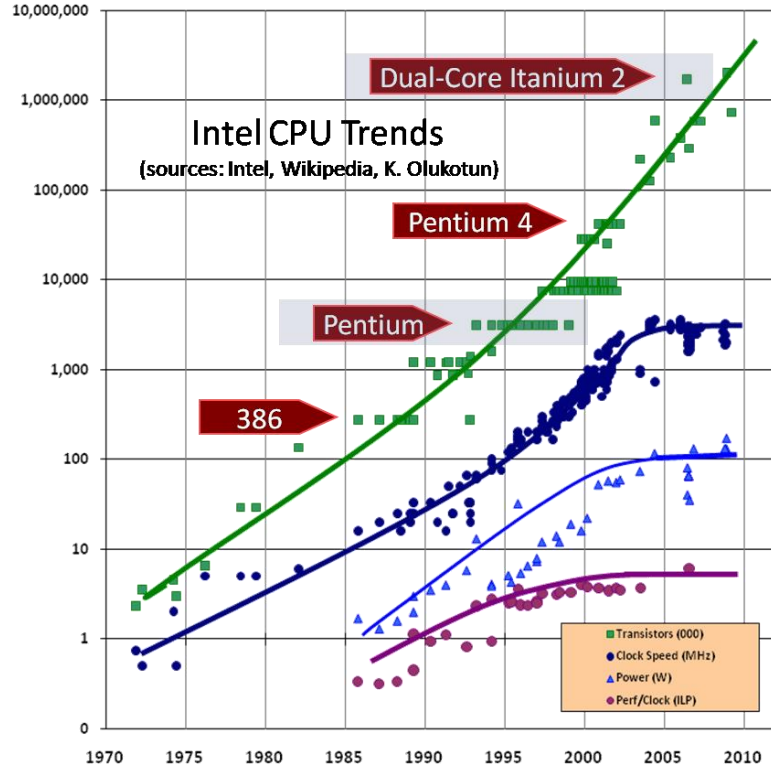
# Memory Wall

- Sandia National Labs investigated the speedup achievable by increasing parallelism (ILP, multiple processors) in 2009
- Example: Number of clerks behind a supermarket counter
  - Two clerks can serve more customers than one
  - 4 ? 8 ? 16 ? 32 ? 64 ? ... 1000 ?
- The problem: Shared memory is ‚shared‘
  - Memory bandwidth
    - Memory transfer speed is limited by the power wall
    - Memory transfer size is limited by the power wall
    - Putting memory into the processor is too costly
  - Bus contention
- Another problem: Memory need kept the pace of CPU speedup
- → **"Memory wall"**

# Processor Speed Development

- Clock speed curve flattened in 2003
  - Heat
  - Power consumption
  - Leakage
- 3-4 GHz since 2001 (!)
- Speeding up the serial instruction execution through clock speed improvements no longer works
- We stumbled into the **Many-Core Era**

[Herb Sutter, 2009]



**ParProg 2020 Introduction**

Chart **29**

# Conventional Wisdoms Replaced

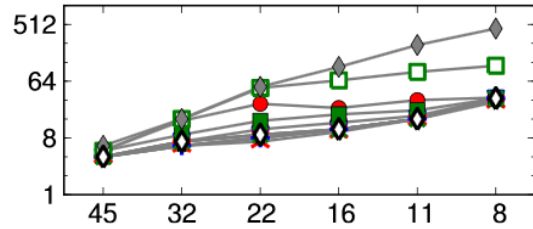| Old Wisdom | New Wisdom |
|---|---|
| Power is free, transistors are expensive | „Power wall" |
| Only dynamic power counts | Static leakage makes 40% of power |
| Multiply is slow, load-and-store is fast | „Memory wall" |
| Instruction-level parallelism gets constantly better via compilers and architectures | „ILP wall" |
| Parallelization is not worth the effort, wait for the faster uniprocessor | Performance doubling might now take 5 years due to physical limits |
| Processor performance improvement by increased clock frequency | Processor performance improvement by increased parallelism |

**ParProg 2020**
**Introduction**

Chart **30**

# Power Wall 2.0

- Power consumption increases with Moore's law, even under constant frequencies

- Cooling is a constant factor

  - Maximum temperature of 100-125 C

  - Hot spots make it worse

- Next-generation processors need to use less power

  - Lower the frequencies

  - Dynamic frequencies scaling (see latest Intel products)

  - Minimize ‚power per bit of I/O' [Skadron 2007]

  - Better cache locality, stop moving stuff around

  - Start to use specialized co-processors and accelerators

**ParProg 2020
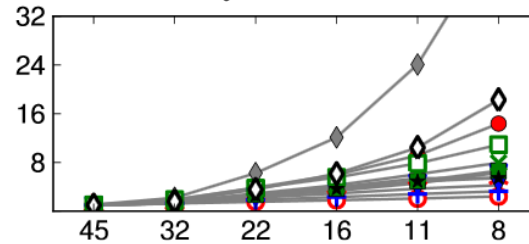Introduction**

Chart **31**

# Power Wall 2.0 = Dark Silicon

"**Dark Silicon and the End of Multicore Scaling**"

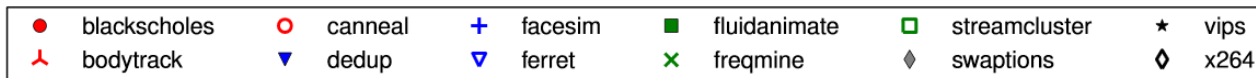by Hadi Esmaeilzadeh, Emily Blem, Renée St. Amant, Karthikeyan Sankaralingam, Doug Burger

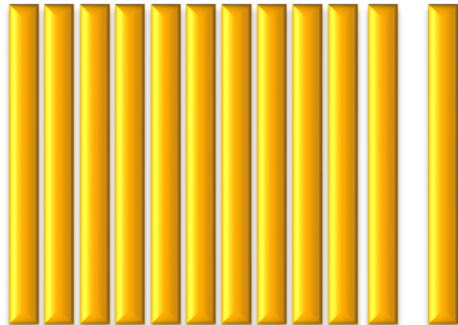**ParProg 2020 Introduction**

Chart **32**

# The Situation

- Hardware people
  - Number of transistors N is still increasing
  - Building larger caches no longer helps (memory wall)
  - ILP is out of options (ILP wall)
  - Voltage / power consumption is at the limit (power wall)
    - Some help with dynamic scaling approaches
  - Frequency is stalled (power wall)
  - Only possible offer is to use increasing N for more cores
- For faster software in the future ...
  - Speedup must come from the utilization of an increasing core count, since F is now fixed
  - Software must participate in the power wall handling, to keep F fixed
  - Software must tackle the memory wall

**ParProg 2020**
**Introduction**

Chart **33**

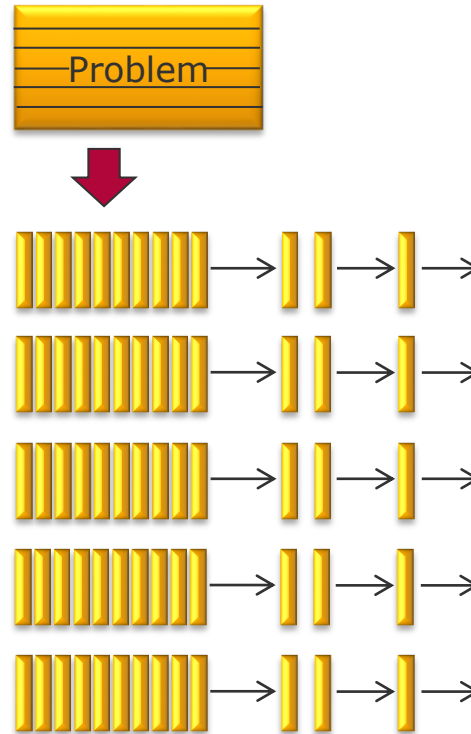# Three Ways of Doing Anything Faster [Pfister]

Application

Instructions

- Work Harder
  (clock speed)

- Work Smarter
  (optimization, caching)

- **Get Help
  (parallelization)**

# Getting Help

- Parallelization not only in computer science
  - Building construction, car manufacturing, large companies
- The basic idea is easy to understand
- Meanwhile tons of options for parallel processing
  - Languages, execution environments, patterns
- Parallelism is a hardware property that must be exploited by software
  - *„A parallel computer is a set of processors that are able to work cooperatively to solve a computational problem.*" (Foster 1995)

Problem

**ParProg 2020 Introduction**

Chart **35**

Thank you
for your attention!