

Digital Engineering • Universität Potsdan

Parallel Programming and Heterogeneous Computing

Introduction

Sven Köhler, Lukas Wenzel, Max Plauth, and Andreas Polze Operating Systems and Middleware Group



Running Applications



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]





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





Moore's Law





- 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

Hasso Plattner Institut

A Physics Problem

- Power: Energy needed per time unit
 - Power density: Watt/mm² → 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 Voltage² (V²) 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







[Moore, ISSCC]

ParProg 2020 Introduction

Chart 15



Transistor Usage





Chart 16

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

Power Density







SOURCE: HEWLETT-PACKARD LABS

Power Density = Temperature





thermal profile during runtime

- Higher temperature leads to
 - Increased transistor leakage

cache

- Decreased transistor speed
- Higher failure probability

Power Density





[Kevin Skadron, 2007]

ParProg 2020 Introduction

Chart 19

Power Density





[Taylor, 2009]







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

ASC





[www.ieeeghn.org]

ParProg 2020 Introduction

Chart 22



A Physics Problem

Dynamic Power = $N \times C \times V^2 \times F$

- Even if we would keep F constant
 - $_{\Box}$ $\,$ N continues to increase exponentially \rightarrow dynamic power
 - $\hfill\square$ Increasing N sums up to more leakage \rightarrow 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, OV 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

ParProg 2020 Introduction

• Today's processors are packed with ILP logic



Chart 25

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
- \rightarrow "ILP wall"
- Ok, so let's use the additional N for more caches

PIPELINE





[Wikipedia]



Memory Hierarchy





Memory Wall



HPI Hasso Plattner Institut

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





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

Chart 31

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



HPI Hasso Plattner Institut

Power Wall 2.0 = Dark Silicon



×

fregmine

۲

swaptions

٥

x264

ParProg 2020 Introduction

Chart **32**

- bodytrack
- ck V dedup

V

ferret

Chart 33

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



Three Ways of Doing Anything Faster [Pfister]



HPI Hasso Plattner Institut

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)







Digital Engineering • Universität Potsdam



Thank you for your attention!