



Friedrich-Alexander-Universität Erlangen-Nürnberg

Programming Techniques for Supercomputers: Modern processors: Single Core

Introduction

Basic technology trend / Moore's law Basic concept of single core architecture

Key single core features

Pipelining Superscalarity SingleInstructionMultipleData

Maximum In-Core Performance

<u>Prof. Dr. G. Wellein^(a,b)</u>, Dr. G. Hager^(a) ^(a)Erlangen National Center for High Performance Computing ^(b)Department für Informatik Friedrich-Alexander-Universität Erlangen-Nürnberg, Sommersemester 2025







Friedrich-Alexander-Universität Erlangen-Nürnberg

Programming Techniques for Supercomputers Modern processors: Single Core

Introduction

Basic technology trend / Moore's law Basic concept of core architecture Key single core features: Pipelining Superscalarity SingleInstructionMultipleData Maximum In-Core Performance



Multi-core today: Intel Xeon Sapphire Rapids (2023)

- Xeon "Sapphire Rapids" (Platinum/Gold/Silver/Bronze): Up to 60 cores running at 1.7+ GHz (+ "Turbo Mode" 4.8 GHz),
- Simultaneous Multithreading
 → reports as 120-way chip
- "Intel 7" process / up to 350 W
- Multi-die package (4 chips)
- Clock frequency: flexible ⁽²⁾



→ One memory domain per die





https://www.techpowerup.com/292204/intel-sapphire-rapids-xeon-4-tile-mcm-annotated





P

Friedrich-Alexander-Universität Erlangen-Nürnberg

Programming Techniques for Supercomputers Modern processors: Single Core

Introduction

Basic technology trend / Moore's law Basic concept of core architecture Key single core features: Pipelining Superscalarity SingleInstructionMultipleData Maximum In-Core Performance

Basic "stored program computer" concept – still in use

 Stored Program Computer" concept (Turing 1936)

 Similar designs on all modern systems



Flexibility! (Still) multiple potential bottlenecks

Stored Program Computer



From high level code to actual execution



General-purpose (cache based) microprocessor core



Introduction: From application to microprocessor core



 High Level Programming Language (e.g. C / C++ / Fortran): Aplication – portable

 Compiler translates program to Instruction set (architecture) (IA32, Intel 64, AMD64 a.k.a. x86, x86_64)

 Instruction Set Architecture (ISA): Hardware specific

Introduction: Instruction Set Paradigms

- In the beginning (60's): Complex Instruction Set Computers (CISC) :
 - Powerful & complex instructions
 - Instruction set is close to high-level programming language
 - Variable length of instructions Save storage!

MULT r0 * [a2] \rightarrow [a1]

Multiply content of address a2 with register content **r**0 and write back to address a1

- Mid 80's: Reduced Instruction Set Computer (RISC) evolved:
 - Fixed instruction length; enables pipelining and high clock frequencies
 - Uses simple instructions, e.g., above instruction is split into at least 3 instructions:

LOAD $[a2] \rightarrow r1;$ MULT $r0 + r1 \rightarrow r2;$ STORE $r2 \rightarrow [a1]$

Nowadays: RISC processor cores

Almost...

x86 CISC/RISC hybrid

 Current x86_64 processors (Intel, AMD): Compiler still generates CISC instructions; but processor core is RISC-like

Example:

- addsd xmm1, [rsi+rax*8] xmm1: register holding floating point data rsi, rax: register holding integer data
- \rightarrow combined address calculation, LD, and ADD instruction
 - 1. Calculate address rsi+rax*8
 - 2. Load double value from that address
 - 3. Add double value into register xmm1 (accumulate)

From high level code to machine execution (CISC-style)



May 6, 2025

From high level code to macro-/microcode execution



Key single-core features: Pipelining



Key single-core features: Superscalarity



Key single-core features: SIMD







Friedrich-Alexander-Universität Erlangen-Nürnberg

Programming Techniques for Supercomputers Modern processors: Single Core

Introduction

Basic technology trend / Moore's law Basic concept of core architecture Key single core features: Pipelining Superscalarity SingleInstructionMultipleData Maximum In-Core Performance



Pipelining of arithmetic/functional units

- Concept:
 - Split complex instruction into several simple / fast steps (stages)
 - Each step takes the same amount of time, e.g. a single cycle
 - Execute different steps on different instructions at the same time (in parallel)
- Benefit:
 - Pipeline can work on multiple instructions simultaneously (in parallel)
 - If pipeline is full one instruction completes every cycle \rightarrow Throughput: 1 inst./cy.
 - Enables faster clock speeds (simple steps/stages)
- Drawback:
 - Pipeline must be filled ("wind-up") → start-up "latency" = number of stages
 - Independent instructions required → complex instruction scheduling by hardware ("out-of-order") or compiler ("software-pipelining")
- Pipelining is widely used in modern computer architectures
- Pipelining addresses Instruction Level Parallelism

Interlude: Possible stages for Floating Point Multiply

Real numbers can be represented as mantissa and exponent in a "normalized" representation, e.g.: s*0.m * 10^e with

Sign s={-1,1} Mantissa m which does not contain 0 in leading digit Exponent e some positive or negative integer

```
    Multiply two real numbers r1*r2 = r3
r1=s1*0.m1 * 10<sup>e1</sup> , r2=s2*0.m2 * 10<sup>e2</sup> :
s1*0.m1 * 10<sup>e1</sup> * s2*0.m2 * 10<sup>e2</sup>
    → (s1*s2)* (0.m1*0.m2) * 10<sup>(e1+e2)</sup>
    → Normalize result: s3* 0.m3 * 10<sup>e3</sup>
```



A(2)

A(1)

First result is available after 5 cycles (=latency of pipeline)!

After that one instruction is completed in each cycle (N-1 cycles)!

Empty pipeline stages in Wind-up/-down phase!

Insert Sign

A(N-3)

A(N-4)

. . .

B(N

Pipelining: Latency, Throughput and Speed-Up

- Assume m-stage pipeline (pipeline latency: m cycles), fixed clock speed and N independent instructions to be executed
- Speed-up of pipelined (T_{pipe}) vs. non-pipelined (T_{seq}) execution time

$$\frac{T_{seq}}{T_{pipe}} = \frac{m \cdot N}{m + N - 1}$$

• **Pipeline throughput**, i.e. average instructions completed per cycle [inst./cy]:

$$\frac{N}{T_{pipe}} = \frac{N}{N + m - 1}$$

$$\Rightarrow \text{Speed-Up: } \frac{T_{seq}}{T_{pipe}} \approx m \qquad \text{for } N \gg m$$

$$\Rightarrow \text{Throughput: } \frac{N}{T_{pipe}} \approx 1\frac{\text{inst.}}{\text{cy}} \qquad \text{for } N \gg m$$

Throughput as function of pipeline stages



m = #pipeline stages

Efficient use of Pipelining

(Potential) dependencies within loop body may prevent efficient software pipelining or OOO execution, e.g.:



Dependency: do i=2,N a(i) = a(i-1) + send do



 Single core on Intel Xeon E5-2695 v3 ("Haswell") with clock speed fixed to 2.3 GHz (Compiler: -O3 -no-vec) – HW limit: 1 MULT instr./cy



Pipelining: Data dependencies – performance model



Pipeline utilization / performance improvement by unrelated workload:



 Single core on Intel Xeon E5-2695 v3 ("Haswell") with clock speed fixed to 2.3 GHz (Compiler: -O3 -no-vec) – HW limit: 1 MULT instr./cy



 Single core on Intel(R) Xeon(R) Platinum 8360Y CPU with clock speed fixed to 2.0 GHz (Compiler: -03 -no-vec) – HW limit: 2 MULT instr./cy



Single core on AMD EPYC 7543 CPU with clock speed fixed to 2.1 GHz (Compiler: -03 -no-vec) – HW limit: 2 MULT instr./cy



* s

Sometime the data dependencies are not that obvious..



- Data (register) dependency on sum (xmm1) → 1 F / m cy for above code! (assuming an ADD latency of m cycles, m=3 for Intel)
- How to enable pipelining here?

Increase pipeline utilization by "loop unrolling"

"2-way Modulo Variable Expansion" (we assume that N is even)





```
sum1=0.d0
sum2=0.d0
do i=1, N, 2
    sum1=sum1+A(i)
    sum2=sum2+A(i+1)
enddo
sum = sum1 + sum2
```



sum2 +=A(4)	
sum1 +=A(3)	

- m-way Modulo Variable Expansion (MVE) to get best performance!
- Sum is split up in m independent partial sums
- Optimal for Intel ADD: 3-way MVE

```
N_r = 3 * (N/3)
sum1=0.d0
sum2=0.d0
sum3=0.d0
do i=1, N_r, 3
    sum1=sum1+A(i)
    sum2=sum2+A(i+1)
    sum3=sum3+A(i+2)
enddo
do i=Nr+1, N
    sum1=sum1+A(i)
enddo

sum=sum1+sum2+sum3
```

- Compiler can do that, if it is allowed to do so...
 - High optimization levels
 - Compiler prefers powers of 2 for unrolling
- Reason: Computer's floating point arithmetic is not associative!

$$\left(\left((a+b)+c\right)+d\right)+e)+f\right) \neq (a+b)+(c+d)+(e+f)$$

- If you require binary exact results (-fp-model strict for Intel) the compiler is not allowed to do this transformation
- Beware additional latency due to reduction at the end
 - Final sum cannot be pipelined
 - High unrolling factor leads to high overhead
 - High unrolling may lead to register shortage

Pipelining: Available resources in modern CPUs

- Typical number of pipeline stages on modern cores:
 - 2-5 for most (important) hardware pipelines: LoaD; STore; MULT; ADD; FMA
 - >>10 for other floating point pipelines: DIVide/SQuareRooT
 - Many other other piplined ALUs, e.g. integer arithmetic, logical, shift, branch, address generation
- Most "older" x86 cores (AMD, Intel):
 - 1 MULT & 1 ADD floating point unit per processor core
 → Max. 1 MULT & 1 ADD instruction per cycle
- Latest Intel (Haswell, Broadwell, Skylake) & AMD (Zen+) cores:
 1 (AMD) or 2 Floating Point Fused MultiplyAdd (FMA) floating point units
 - FMA3 instruction: $s=s+a*b \rightarrow 1$ Input register (s) is overwritten
 - FMA4 instruction: $s=r+a*b \rightarrow No$ input register is modified
 - Typically 2 (1) FMA instruction per cycle for Intel (AMD) processors
 - On Intel: Per cycle up to 2 MULT or ADD instructions

Costs of arithmetic instructions: Intel Skylake processors



- Consequence: Avoid expensive instructions in hot spots!
- Other expensive math (transcendental, log,...) is done in libraries

Pipelining: The Instruction pipeline

 Besides arithmetic & functional units, instruction execution itself is pipelined also, e.g.: one instruction performs at least 3 steps:



Hardware Pipelining on processor (all units can run concurrently):

1	Fetch Instruction 1 from L1I			
2	Fetch Instruction 2	Decode		
	from L1I	Instruction 1		
	Fetch Instruction 3	Decode	Execute	
3	from L1I	Instruction 2	Instruction 1	
4	Fetch Instruction 4	Decode	Execute	
	from L1I	Instruction 3	Instruction 2	

. . .

- Non-predictable branches can stall this pipeline!
 - Hardware can predict conditional branches w/ high accuracy
- Each unit is pipelined itself (cf. Execute=Multiply Pipeline)

Pipelining: The Instruction pipeline

Problem: Unpredictable branches to other instructions



Pipelining summary

Pipelining tries to achieve

- Maximum instruction throughput (1 instr/cy in many cases)
- Hiding of instruction latency
- Prerequisites
 - Independent instructions
 - A lot of independent instructions for maximum efficiency ($N \gg m$)
 - Highest benefit if code & data are close to the core (L1 instr./data cache)
 - Conditional branches must be correctly predicted by hardware
- Drawbacks
 - Pipeline must be filled \rightarrow inefficient for N $\leq m$
 - Dependencies between pipelines may increase effective depth (see tutorial)
 - Unresolvable data dependencies are hazardous





Friedrich-Alexander-Universität Erlangen-Nürnberg

Programming Techniques for Supercomputers Modern processors: Single Core

Introduction

Basic technology trend / Moore's law Basic concept of core architecture Key single core features: Pipelining Superscalarity SingleInstructionMultipleData Maximum In-Core Performance



- Superscalar processors provide additional hardware (i.e. transistors) to execute multiple instructions per cycle!
- → Exploit Instrucion Level Parallelism (ILP)
- Parallel hardware components / pipelines are available to
 - fetch / decode / issues multiple instructions per cycle (typically 3 – 8 per cycle)
 - perform multiple integer / address calculations per cycle
 - perform multiple load (store) multiple instructions per cycle (e.g. one LD and one ST per cycle)
 - perform multiple floating point (FP) instructions per cycle (e.g., 2 floating point instructions/cycle, e.g. 1 MULT + 1 ADD)
- Parallelization of instruction stream" required
- Performance metrics quantifying superscalarity:

Instructions Per Cycle: IPC Cycles Per Instruction: CPI

Superscalar Processors – Instruction Level Parallelism



Multiple pipelines at work: Interleaving instructions

Example:

```
Fortran Code:
                     do i=1, N
                         a(i) = a(i) * c
                     end do
                                                                                      Assumed
                                                                                       Latencies
                                  Load operand to register (4 cycles)
load r1, a[i]
                                  Multiply a(i) with c (2 cycles); a[i],c in registers
mult r1 = c, r1
                                  Store result from register to mem./cache (2 cycles)
store a[i], r1
                                  Increase loop counter as long as i less or equal N (0 cycles)
branch.loop
                 Simple Pseudo Code:
                                                     Dependencies on r1
                 loop: load r1, a[i]
                                                       within one iteration
                         mult r1 = c, r1
                         store a[i], r1
                                                       across iterations
                         branch.loop
```

Superscalar & Pipelined Execution

a[i]=a[i]*c;N=12

Naive instruction issue Cycle 1 load a[1] Cycle 2 Cvcle 3 Cycle 4 Cycle 5 mult **a[1]**=c,**a[1]** Cycle 6 Cycle 7 store a[1] Cycle 8 Cycle 9 load a[2] Cycle 10 Cycle 11 Cycle 12 Cycle 13 mult a[2]=c,a[2] Cycle 14 Cycle 15 store a[2] Cycle 16 Cycle 17 load a[3] Cycle 18 Cycle 19

```
Simple Pseudo Code:
loop: load a[i]
    mult a[i] = c, a[i]
    store a[i]
    branch.loop
```

Instruction executed "in-order"

Total execution time: T= 12 * (4+2+2) cy = 96 cyIPC = 3/8 instr./cy

No pipelining and superscalarity!

CPI = 8/3 cy/instr.

time

Superscalar & Pipelined Execution

a[i]=a[i]*c;N=12

Optimized instruction issue



```
Simple Pseudo Code:
loop: load a[i]
    mult a[i] = c, a[i]
    store a[i]
    branch.loop
```

Assumptions:

- LD/MULT/ST can be executed in parallel!
- Instructions are perfectly reordered but dependecies (within loop iteration) are maintained!
- Register renaming required

Kernel:

Full pipelining and high superscalarity!

Reordering the instruction stream: Two options

- Software pipelining
 - Done by the compiler
 - Compiler reorders instructions
 - Requires deep insight into application (data dependencies) and processor (latencies of functional units)

```
<... prologue ...>
kernel: load a[i+6]
mult a[i+2] = c, a[i+2]
store a[i]
branch → kernel
<... epilogue ...>
```

- Required on "in-order" architectures
- Rarely used today (see right)

- Dynamic reordering of instructions at runtime
 - Done by the hardware
 - Out-of-order (OOO) execution
 - Instructions are executed when operands are available



All modern general-purpose CPUs do this

Register renaming

- Prerequisite for good OoO execution: "Bogus" register dependencies can be resolved
 - Hardware has "shadow registers" it can use to store intermediate values that are already "officially" overwritten



- Solution: Hardware assigns a new register with the same name as soon as the old value gets overwritten
- "Shadow copy" lives as long as necessary
 - Until no instructions in flight reference the register any more

Superscalar processors

Intel processors - qualitative view ("Intel Sandy Bridge")

