Parallel Computing



Master Informatics Eng.

2019/20 *A.J.Proença*

Multithreading & Data Parallelism

(some slides are borrowed, mod's in green)

AJProença, Advanced Architectures, MiEI, UMinho, 2019/20

Processor arch: beyond Instruction-Level Parallelism

- \sim
- When exploiting ILP, goal is to minimize CPI

Pipeline CPI (efficient to exploit loop-level parallelism) =>

- Ideal pipeline CPI +
- Structural stalls +
- Data hazard stalls +
- Control stalls +
- Memory stalls ... cache techniques ...
- > Multiple issue =>
 - find enough parallelism to keep pipeline(s) occupied
- Multithreading =>

➢ find additional ways to keep pipeline(s) occupied

Multithreading

- Performing multiple threads of execution in parallel
 - Replicate registers, PC/IP, etc.
 - Fast switching between threads
- Fine-grain multithreading / time-multiplexed
 MT
 - Switch threads after each cycle
 - Interleave instruction execution
 - If one thread stalls, others are executed
- Coarse-grain multithreading
 - Only switch on long stall (e.g., L2-cache miss)
 - Simplifies hardware, but doesn't hide short stalls (eg, data hazards)

Chapter 7 — Multicores, Multiprocessors, and Clusters — 3

Simultaneous Multithreading

- In multiple-issue dynamically scheduled processor
 - Schedule instructions from multiple threads
 - Instructions from independent threads execute when function units are available
 - Within threads, dependencies handled by scheduling and register renaming
 - Example: Intel Pentium-4 HT
 - Two threads: duplicated registers, shared function units and caches

HT: Hyper-Threading, Intel trade mark for their SMT implementation MT in Xeon Phi KNC: 4-way SMT with time-mux MT, **not HT**!



Multithreading Example





Chapter 7 — Multicores, Multiprocessors, and Clusters — 5

Multithreading Example



R

Chapter 7 — Multicores, Multiprocessors, and Clusters — 6



Article Talk

Hyper-threading

From Wikipedia, the free encyclopedia

*X*X

As seen before...

Internal architecture of Intel P6 processors

2

<u>Note:</u> "Intel P6" is the common µarch name for PentiumPro, Pentium II & Pentium III, which inspired Core, Nehalem and Phi



The pipelined functional units might have better use if shared among more threads =>

Note: white boxes are bubbles...



Processor arch: beyond Instruction-Level Parallelism

- \sim
- When exploiting ILP, goal is to minimize CPI

Pipeline CPI (efficient to exploit loop-level parallelism) =>

- Ideal pipeline CPI +
- Structural stalls +
- Data hazard stalls +
- Control stalls +
- Memory stalls ... cache techniques ...
- > Multiple issue =>
 - find enough parallelism to keep pipeline(s) occupied
- Multithreading =>

➢ find additional ways to keep pipeline(s) occupied

Insert data parallelism features

AJProença, Advanced Architectures, MiEI, UMinho, 2019/20

Instruction and Data Streams

		Data Streams	
		Single	Multiple
Instruction Streams	Single	SISD: Intel Pentium 4	SIMD : SSE instructions of x86
	Multiple	MISD: No examples today	MIMD: Intel Xeon e5345

SPMD: Single Program Multiple Data

- A parallel program on a MIMD computer
- Conditional code for different processors

Introduction

- SIMD architectures can exploit significant datalevel parallelism for:
 - matrix-oriented <u>scientific computing</u>
 - media-oriented <u>image</u> and <u>sound</u> processing
- SIMD is more energy efficient than MIMD
 - only needs to fetch one instruction per data operation
 - makes SIMD attractive for personal mobile devices
- SIMD allows programmers to continue to think sequentially



SIMD Parallelism

- Vector architectures (slides 5 to 18)
- SIMD & extensions (slides 19 to 24)
- Graphics Processor Units (GPUs) (next set)
- For x86 processors:
 - Expected grow:
 2 more cores/chip/year
 - SIMD width: 2x every 4 years
 - Potential speedup: SIMD 2x that from MIMD!





Introduction

Vector Architectures

Basic idea:

- Read sets of data elements (<u>gather</u> from memory) into "vector registers"
- Operate on those registers
- Store/<u>scatter</u> the results back into memory
- Registers are controlled by the compiler
 - Used to hide memory latency
 - Leverage memory bandwidth





AJProença, Advanced Architectures, MiEI, UMinho, 2016/17

Challenges

Start up time

- Latency of vector functional unit
- Assume the same as Cray-1
 - Floating-point add => 6 clock cycles
 - Floating-point multiply => 7 clock cycles
 - Floating-point divide => 20 clock cycles
 - Vector load => 12 clock cycles

Improvements:

- > 1 element per clock cycle (1)
- Non-64 wide vectors (2)
- IF statements in vector code (3)
- Memory system optimizations to support vector processors (4)
- Multiple dimensional matrices (5)
- Sparse matrices (6)
- Programming a vector computer (7)



Vector Programming (7)

- Compilers are a key element to give hints on whether a code section will vectorize or not
- Check if loop iterations have data dependencies, otherwise vectorization is compromised
- Vector Architectures have a too high cost, but simpler variants are currently available on off-the-shelf devices; however:
 - most do not support non-unit stride => care must be taken in the design of data structures
 - same applies for gather-scatter...



SIMD Extensions

- Media applications operate on data types narrower than the native word size
 - Intel SIMD Ext started with 64-bit wide vectors and grew to wider vectors and more facilit
 - Current AVX generation is 512-bit wide



- Limitations, compared to vector architectures:
 - Number of data operands encoded into op code
 - No sophisticated addressing modes (strided, scatter-gather, but...)
 - No mask registers



SIMD Implementations

- Implementations:
 - Intel MMX (1996)
 - Eight 8-bit integer ops or four 16-bit integer ops
 - Streaming SIMD Extensions (SSE) (1999)
 - Eight 16-bit integer ops
 - Four 32-bit integer/fp ops or two 64-bit integer/fp ops
 - Advanced Vector eXtensions (AVX) (2010)
 - Eight 32-bit fp ops or Four 64-bit fp ops (integer in AVX-2)
 - 512-bits wide in AVX-512 (and also in Larrabee & Phi-KC)
 - Operands <u>must be in consecutive and aligned</u> memory locations



Extensões de processamento vetorial no Intel 64



AJProença, Advanced Architectures, MiEI, UMinho, 2019/20

18

Additional features in Intel x86

公



from Marat Dukhan, 2014

Beyond Vector/SIMD architectures

\sim

- Vector/SIMD-extended architectures are hybrid approaches
 - mix (super)scalar + vector op capabilities on a single device
 - highly pipelined approach to reduce memory access penalty
 - tightly-closed access to shared memory: lower latency
- Evolution of Vector/SIMD-extended architectures
 - CPU cores with wider vectors and/or SIMD cores:
 - <u>DSP</u> VLIW cores with vector capabilities: **Texas Instruments** (...?)
 - <u>PPC</u> cores coupled with SIMD cores: **Cell** (past...) , **IBM Power BQC...**
 - <u>ARM64</u> cores coupled with SIMD cores: from Tegra to Parker (NVidia) (...?)
 - <u>x86</u> many-core: **Intel** MIC / Xeon KNL, **AMD** FirePro...
 - other many-core: **ShenWay** 260, **Adapteva** Epiphany-V...
 - coprocessors (require a host scalar processor): accelerator devices
 - on disjoint physical memories (e.g., Xeon KNC with PCI-Expr, PEZY-SC)
 - focus on SIMT/SIMD to hide memory latency: GPU-type approach
 - ISA-free architectures, code compiled to silica: FPGA