Lecture 20: Data Level Parallelism -- Introduction and Vector Architecture

CSE 564 Computer Architecture Summer 2017

Department of Computer Science and Engineering Yonghong Yan yan@oakland.edu

www.secs.oakland.edu/~yan

Very Important Terms

- Dynamic Scheduling → Out-of-order Execution
- Speculation → In-order Commit
- Superscalar → Multiple Issue

Techniques	Goals	Implementation	Addressing	Approaches
Dynamic Scheduling	Out-of- order execution	Reservation Stations, Load/Store Buffer and CDB	Data hazards (RAW, WAW, WAR)	Register renaming
Speculation	In-order commit	Branch Prediction (BHT/BTB) and Reorder Buffer	Control hazards (branch, func, exception)	Prediction and misprediction recovery
Superscalar/ VLIW	Multiple issue	Software and Hardware	To Increase CPI	By compiler or hardware

Last Lecture: Multithreading



CSE 564 Class Contents

- Introduction to Computer Architecture (CA)
- Quantitative Analysis, Trend and Performance of CA
 - Chapter 1
- Instruction Set Principles and Examples
 - Appendix A
- Pipelining and Implementation, RISC-V ISA and Implementation
 - Appendix C, RISC-V (riscv.org) and UCB RISC-V impl
- Memory System (Technology, Cache Organization and Optimization, Virtual Memory)
 - Appendix B and Chapter 2
 - Midterm covered till Memory Tech and Cache Organization
- Instruction Level Parallelism (Dynamic Scheduling, Branch Prediction, Hardware Speculation, Superscalar, VLIW and SMT)
 - Chapter 3
- Data Level Parallelism (Vector, SIMD, and GPU)
 - Chapter 4
- Thread Level Parallelism
 - Chapter 5

Topics for Data Level Parallelism (DLP)

- Parallelism (centered around ...)
 - Instruction Level Parallelism
 - Data Level Parallelism
 - Thread Level Parallelism
- DLP Introduction and Vector Architecture – 4.1, 4.2
- SIMD Instruction Set Extensions for Multimedia - 4.3
- Graphical Processing Units (GPU)
 -4.4
- GPU and Loop-Level Parallelism and Others - 4.4, 4.5, 4.6, 4.7

Finish in two/three sessions

Acknowledge and Copyright

- Slides adapted from
 - UC Berkeley course "Computer Science 252: Graduate Computer Architecture" of David E. Culler Copyright(C) 2005 UCB
 - UC Berkeley course Computer Science 252, Graduate Computer Architecture Spring 2012 of John Kubiatowicz Copyright(C) 2012 UCB
 - Computer Science 152: Computer Architecture and Engineering, Spring 2016 by Dr. George Michelogiannakis from UC Berkeley
 - Arvind (MIT), Krste Asanovic (MIT/UCB), Joel Emer (Intel/MIT), James Hoe (CMU), John Kubiatowicz (UCB), and David Patterson (UCB)

https://passlab.github.io/CSE564/copyrightack.html

Flynn's Classification (1966)

Broad classification of parallel computing systems

- based upon the number of concurrent Instruction

(or control) streams and **Data** streams



Michael J. Flynn:

- SISD: Single Instruction, Single Data
 - conventional uniprocessor
- SIMD: Single Instruction, Multiple Data
 - one instruction stream, multiple data paths
 - distributed memory SIMD (MPP, DAP, CM-1&2, Maspar)
 - shared memory SIMD (STARAN, vector computers)
- **MIMD:** Multiple Instruction, Multiple Data
 - message passing machines (Transputers, nCube, CM-5)
 - non-cache-coherent shared memory machines (BBN Butterfly, T3D)
 - cache-coherent shared memory machines (Sequent, Sun Starfire, SGI Origin)
- **MISD: Multiple Instruction, Single Data**
 - Not a practical configuration

7

Flynn's Taxonomy



https://en.wikipedia.org/wiki/Flynn%27s_taxonomy

More Categories

- Single program, multiple data (SPMD)
 - Multiple autonomous processors execute the program at independent points
 - Difference with SIMD: SIMD imposes a lockstep
 - Programs at SPMD can be at independent points
 - SPMD can run on general purpose processors
 - Most common method for parallel computing
- Multiple program, multiple data (MPMD)
 ≈ MIMD
 - Multiple autonomous processors simultaneously operating at least 2 independent programs



SIMD: Single Instruction, Multiple Data (Data Level Paralleism)

- SIMD architectures can exploit significant data-level parallelism for:
 - matrix-oriented scientific computing
 - media-oriented image and sound processors
- SIMD is more energy efficient than MIMD
 - Only needs to fetch one instruction per data operation processing multiple data elements
 - Makes SIMD attractive for personal mobile devices
- SIMD allows programmer to continue to think sequentially



SIMD Parallelism

- Three variations
 - Vector architectures
 - SIMD extensions
 - Graphics Processor Units (GPUs)
- For x86 processors:
 - Expect two additional cores per chip per year (MIMD)
 - SIMD width to double every four years
 - Potential speedup from SIMD to be twice that from MIMD!

Vector Architecture

VLIW vs Vector

 VLIW takes advantage of instruction level parallelism (ILP) by specifying multiple instructions to execute in parallel

1					
Int Op 1	Int Op 2	Mem Op 1	Mem Op 2	FP Op 1	FP Op 2

- Vector architectures perform the same operation on multiple data elements – single instruction
 - Data-level parallelism



Vector Programming Model



Control Information

- VLR limits the highest vector element to be processed by a vector instruction
 - VLR is loaded prior to executing the vector instruction with a special instruction
- Stride for load/stores:
 - Vectors may not be adjacent in memory addresses
 - E.g., different dimensions of a matrix
 - Stride can be specified as part of the load/store



Basic Structure of Vector Architecuture

- VMIPS
- eight 64-element vector registers
- all the functional units are vector functional units.
- The vector and scalar registers have a significant number of read and write ports to allow multiple simultaneous vector operations.
- A set of crossbar switches (thick gray lines) connects these ports to the inputs and outputs of the vector functional units.



VMIPS Vector Instructions

 Suffix 	Instruction	Operands	Function
– VV suffix	ADDVV.D ADDVS.D	V1,V2,V3 V1,V2,F0	Add elements of V2 and V3, then put each result in V1. Add F0 to each element of V2, then put each result in V1.
– VS suffix	SUBVV.D SUBVS.D SUBSV.D	V1,V2,V3 V1,V2,F0 V1,F0,V2	Subtract elements of V3 from V2, then put each result in V1. Subtract F0 from elements of V2, then put each result in V1. Subtract elements of V2 from F0, then put each result in V1.
Load/Store	MULVV.D MULVS.D	V1,V2,V3 V1,V2,F0	Multiply elements of V2 and V3, then put each result in V1. Multiply each element of V2 by F0, then put each result in V1.
– LV/SV	DIVVV.D DIVVS.D	V1,V2,V3 V1,V2,F0	Divide elements of V2 by V3, then put each result in V1. Divide elements of V2 by F0, then put each result in V1.
– LVWS/SVW	S IVSV.D	V1.F0.V2	Divide F0 by elements of V2, then but each result in V1.
Dogistors		R1.V1	Store vector register V1 into memory starting at address R1.
- Registers	LVWS	V1,(R1,R2)	Load V1 from address at R1 with stride in R2 (i.e., $R1 + i \times R2$).
– VLR (vector	SVWS	(R1,R2),V1	Store V1 to address at R1 with stride in R2 (i.e., $R1 + i \times R2$).
length	LVI	V1,(R1+V2)	Load V1 with vector whose elements are at R1 + V2(i) (i.e., V2 is an index).
register	SVI	(R1+V2),V1	Store V1 to vector whose elements are at R1 + V2(i) (i.e., V2 is an index).
register)		V1.R1	Create an index vector by storing the values 0, $1 \times R1$, $2 \times R1$,, $63 \times R1$ into V1.
 VM (vector mask) 	SVV.D SVS.D	V1,V2 V1,F0	Compare the elements (EQ, NE, GT, LT, GE, LE) in V1 and V2. If condition is true, put a 1 in the corresponding bit vector; otherwise put 0. Put resulting bit vector in vector- mask register (VM). The instruction SVS.D performs the same compare but using a scalar value as one operand.
	РОР	R1,VM	Count the 1s in vector-mask register VM and store count in R1.
	CVM		Set the vector-mask register to all 1s.
	MTC1 MFC1	VLR,R1 R1,VLR	Move contents of R1 to vector-length register VL. Move the contents of vector-length register VL to R1.
	MVTM MVFM	VM,FO FO,VM	Move contents of F0 to vector-mask register VM. Move contents of vector-mask register VM to F0.

Figure 4.3 The VMIPS vector instructions, showing only the double-precision floating-point operations. I addition to the vector registers, there are two special registers, VLR and VM, discussed below. These special register

Highlight of VMIPS Vector Instructions

- Vector operations have the letters "VV/VS" appended.
 - E.g. ADDVV.D is an addition of two double-precision vectors.
- Vector instructions input:
 - 1) a pair of vector registers (ADDVV.D) or
 - 2) a vector register and a scalar register (ADDVS.D).
 - » all operations use the same value in the scalar register as one input.
- LV/LVWS and SV/SVWS: vector load and vector store which load or store an entire vector of double-precision data.
 - One operand is the vector register to be loaded or stored;
 - The other operand, a GPR, is the starting address of the vector in memory.
 - LVWS/SVWS: For stride load/store
 - LVI/SVI: indexed load/store
- Two additional special-purpose registers:
 - Vector-length register: when the natural vector length is NOT 64
 - Vector-mask register: when loops involve IF statements.
- In-order scalar processor for vector architecture
 - Not out- of-order superscalar processors.
- Vectors naturally accommodate varying data sizes.
 - one view of a vector register size is 64 64-bit data elements, but 128 32-bit elements, 256 16-bit elements, and even 512 8-bit elements are equally valid views.

AXPY (64	ele	ments) and	(Y = a *) VMIPS	X + Y) in MIPS		
for (i=0; i<64; i++) Y[i] = a* X[i] + Y[i];			The starting ; are in Rx a	The starting addresses of X and Y are in Rx and Ry, respectively		
# instrs:	Loop:	L.D DADDIU L.D	F0,a R4,Rx,#512 F2,0(Rx) F2 F2 F0	;load scalar a ;last address to load ;load X[i]		
 6 vs ~600 Pipeline stalls 64x higher by MIPS 		MUL.D L.D ADD.D S.D DADDIU	F2,F2,F0 F4,0(Ry) F4,F4,F2 F4,9(Ry) Rx,Rx,#8	;a × X[i] ;load Y[i] ;a × X[i] + Y[i] ;store into Y[i] ;increment index to X		
Vector <i>chaining</i> (forwarding)		DADDIU DSUBU BNEZ	Ry,Ry,#8 R20,R4,Rx R20,Loop	; compute bound ; compute if done		
– V1, V2, V3 and	V4	L.D LV MULVS.D LV ADDVV.D SV	F0,a V1,Rx V2,V1,F0 V3,Ry V4,V2,V3 V4,Ry	<pre>;load scalar a ;load vector X ;vector-scalar multiply ;load vector Y ;add ;store the result</pre>		

Vector Memory-Memory versus Vector Register Machines

- Vector memory-memory instructions hold all vector operands in main memory
- The first vector machines, CDC Star-100 ('73) and TI ASC ('71), were memory-memory machines
- Cray-1 ('76) was first vector register machine



Vector Memory-Memory vs. Vector Register Machines

 Vector memory-memory architectures (VMMA) require greater main memory bandwidth, why?

- All operands must be read in and out of memory

- VMMAs make if difficult to overlap execution of multiple vector operations, why?
 - Must check dependencies on memory addresses
- VMMAs incur greater startup latency
 - Scalar code was faster on CDC Star-100 (VMM) for vectors < 100 elements</p>
- Apart from CDC follow-ons (Cyber-205, ETA-10) all major vector machines since Cray-1 have had vector register architectures
- (we ignore vector memory-memory from now on)

Vector Instruction Set Advantages

- Compact
 - one short instruction encodes N operations
- Expressive and predictable, tells hardware that these N operations:
 - are independent
 - use the same functional unit
 - access disjoint registers
 - access registers in same pattern as previous instructions
 - access a contiguous block of memory (unit-stride load/store)
 - access memory in a known pattern (strided load/store)

Scalable

– can run same code on more parallel pipelines (lanes)

History: Supercomputers

- Definition of a supercomputer:
 - Fastest machine in world at given task
 - A device to turn a compute-bound problem into an I/O bound problem
 - Any machine costing \$30M+
 - Any machine designed by Seymour Cray (originally)
- CDC6600 (Cray, 1964) regarded as first supercomputer
 - A vector machine
- <u>www.cray.com</u>: The Supercomputer Company

The Father of Supercomputing



Seymour Cray

Electrical engineer

Seymour Roger Cray was an American electrical engineer and supercomputer architect who designed a series of computers that were the fastest in the world for decades, and founded Cray Research which built many of these machines. Wikipedia

Born: September 28, 1925, Chippewa Falls, WI

Died: October 5, 1996, Colorado Springs, CO

Awards: Eckert–Mauchly Award

Parents: Seymour R. Cray, Lillian Cray

Education: University of Minnesota, Chippewa Falls High School

Fields: Applied mathematics, Computer Science, Electrical engineering

https://en.wikipedia.org/wiki/Seymour_Cray

http://www.cray.com/company/history/seymour-cray

Supercomputer Applications

- Typical application areas
 - Military research (nuclear weapons, cryptography)
 - Scientific research
 - Weather forecasting
 - Oil exploration
 - Industrial design (car crash simulation)
 - **Bioinformatics**
 - Cryptography
- All involve huge computations on large data sets
- In 70s-80s, Supercomputer = Vector Machine

Vector Supercomputers

- Epitomy: Cray-1, 1976
- **Scalar Unit**
 - Load/Store Architecture
- Vector Extension
 - Vector Registers
 - Vector Instructions

Implementation

- Hardwired Control
- Highly Pipelined Functional Units
- Interleaved Memory System
- No Data Caches
- No Virtual Memory



Cray-1 (1976)



memory bank cycle 50 ns *processor cycle* 12.5 ns (80MHz)

Today's Supercomputers (www.top500.org)

TOP 10 Sites for November 2016

For more information about the sites and systems in the list, click on the links or view the complete list.

Rmax Rpeak Power **Rank Site** System Cores (TFlop/s) (TFlop/s) (kW) National Supercomputing Sunway TaihuLight - Sunway MPP, Sunway 93,014.6 125,435.9 15,371 10,649,600 MIMD and Hybrid SW26010 260C 1.45GHz, Sunway Center in Wuxi China NRCPC National Super Computer Tianhe-2 (MilkyWay-2) - TH-IVB-FEP Cluster, 3,120,000 33,862.7 54,902.4 17,808 2 **MIMD**/vector or Center in Guangzhou Intel Xeon E5-2692 12C 2.200GHz. TH China Express-2. Intel Xeon Phi 31S1P **MIMD/GPU** NUDT DOE/SC/Oak Ridge Titan - Cray XK7, Opteron 6274 16C 560,640 17,590.0 27,112.5 8,209 National Laboratory 2.200GHz, Cray Gemini interconnect, NVIDIA United States K20x Cray Inc. Accelerator/Co-Processor System Share DOE/NNSA/LLNL Sequoia - BlueGene/Q, Power BQC 16C 1.60 1,572,864 17,173.2 20,132.7 7,890 United States GHz, Custom NVIDIA Tesla K40 IBM NVIDIA Tesla K80 Intel Xeon Phi 5110P 22.7% DOE/SC/LBNL/NERSC Cori - Cray XC40, Intel Xeon Phi 7250 68C 14,014.7 27,880.7 3,939 622,336 29.5% NVIDIA Tesla K20x United States 1.4GHz. Aries interconnect Intel Xeon Phi 5120D Crav Inc. Intel Xeon Phi 7120P Joint Center for Advanced Oakforest-PACS - PRIMERGY CX1640 M1. 556,104 13.554.6 24.913.5 2.719 NVIDIA 2050 **High Performance** Intel Xeon Phi 7250 68C 1.4GHz, Intel Omni-NVIDIA Tesla K20 10.2% Computing Path NVIDIA Tesla K20m Fujitsu Japan NVIDIA 2090 Others **RIKEN Advanced Institute** K computer, SPARC64 VIIIfx 2.0GHz, Tofu 705.024 10.510.0 11.280.4 12.660 for Computational Science interconnect (AICS) Fujitsu Japan Piz Daint - Cray XC50, Xeon E5-2690v3 12C 8 Swiss National 206,720 9,779.0 15,988.0 1,312 Supercomputing Centre 2.6GHz, Aries interconnect, NVIDIA Tesla [CSCS] P100 Switzerland Crav Inc. DOE/SC/Argonne National Mira - BlueGene/Q, Power BQC 16C 1.60GHz, 786,432 8,586.6 10,066.3 3,945 9

#1 of TOP500 as of Nov 2016



Vector Arithmetic Execution

- Use deep pipeline (=> fast clock) to execute element operations
- Simplifies control of deep pipeline because elements in vector are independent (=> no hazards!)



Six stage multiply pipeline

V3 <- v1 * v2

Vector Execution: Element Group



29

Vector Instruction Execution with Pipelined Functional Units



Vector Unit Structure (4 Lanes)



Vector Instruction Parallelism

Can overlap execution of multiple vector instructions



Complete 24 operations/cycle while issuing 1 short instruction/cycle

Vector Chaining

Vector version of register bypassing

```
- introduced with Cray-1
```



Vector Chaining Advantage

• Without chaining, must wait for last element of result to be written before starting dependent instruction



• With chaining, can start dependent instruction as soon as first result appears



Automatic Code Vectorization



Vector Length Register

- Vector length not known at compile time?
- Use Vector Length Register (VLR)
- Use strip mining for vectors over the maximum length (serialized version before vectorization by compiler)

```
low = 0;
VL = (n % MVL); /*find odd-size piece using modulo op % */
for (j = 0; j <= (n/MVL); j=j+1) { /*outer loop*/
for (i = low; i < (low+VL); i=i+1) /*runs for length VL*/
Y[i] = a * X[i] + Y[i] ; /*main operation*/
low = low + VL; /*start of next vector*/
VL = MVL; /*reset the length to maximum vector length*/
}
```



Vector Stripmining Problem: Vector registers have finite length

Solution: Break loops into pieces that fit in registers, "Stripmining"

	ANDI KI, N, $03 $ # N IIIOQ 04
	MTC1 VLR, R1
<pre>Eor (i=0; i<n; i++)<="" pre=""></n;></pre>	loop:
C[i] = A[i] + B[i];	LV V1, RA
A B C	DSLL R2, R1, 3 # Multiply by 8
	DADDU RA, RA, R2 # Bump pointer
	LV V2, RB
	DADDU RB, RB, R2
64 elements	ADDV.D V3, V1, V2
	SV V3, RC
	DADDU RC, RC, R2
	DSUBU N, N, R1 # Subtract elements
	LI R1, 64
	MTC1 VLR, R1
	BGTZ N, loop # Any more to do?

Vector Mask Registers

```
for (i = 0; i < 64; i=i+1)
if (X[i] != 0)
X[i] = X[i] - Y[i];
```

Use vector mask register to "disable" elements (1 bit per element):

V1,Rx	;load vector X into V1
V2,Ry	;load vector Y
F0,#0	;load FP zero into F0
V1,F0	;sets VM(i) to 1 if V1(i)!=F0
V1,V1,V2	;subtract under vector mask
Rx,V1	;store the result in X
	V1,Rx V2,Ry F0,#0 V1,F0 V1,V1,V2 Rx,V1

- GFLOPS rate decreases!
 - Vector operation becomes bubble ("NOP") at elements where mask bit is clear

Masked Vector Instructions

Simple Implementation

execute all N operations, turn off result writeback according to mask

M[7]=1 X[7]	Y[7]
M[6]=0 X[6]	Y[6]
M[5]=1 X[5]	Y[5]
M[4]=1 X[4]	Y[4]
M[3]=0 X[3]	Y[3]
M[2]=0 > M[1]=1 >	<[2] <[1]
M[0]=0	x[0]
Write Enable	Write data port

Density-Time Implementation

 scan mask vector and only execute elements with non-zero masks



Compress/Expand Operations

- Compress packs non-masked elements from one vector register contiguously at start of destination vector register
 - population count of mask vector gives packed vector length
- Expand performs inverse operation



Compress Expand

Used for density-time conditionals and also for general selection operations

- Must vectorize multiplication of rows of B with columns of D
 - Row-major: B: 1 double (8 bytes), and D: 100 doubles (800 bytes)
- Use non-unit stride
 - LDWS
- Bank conflict (stall) occurs when the same bank is hit faster than bank busy time:

- #banks / LCM(stride, #banks) < bank busy time</p>

Scatter-Gather

Sparse matrix:

– Non-zero values are compacted to a smaller value array (A[])

– indirect array indexing, i.e. use an array to store the index to value array (K[])

for (i = 0; i < n; i=i+1) A[K[i]] = A[K[i]] + C[M[i]];

• Use index vector:

LV	Vk, Rk	;load K
LVI	Va, (Ra+Vk)	;load A[K[]]
LV	Vm, Rm	;load M
LVI	Vc, (Rc+Vm)	;load C[M[]]
ADDVV.D	Va, Va, Vc	;add them
SVI	(Ra+Vk), Va	;store A[K[]]

(1.0)	0	5.0	0	0	0	0	0
0	3.0	0	0	0	0	11.0	0
0	0	0	0	9.0	0	0	0
0	0	6.0	0	0	0	0	0
0	0	0	7.0	0	0	0	0
2.0	0	0	0	0	10.0	0	0
0	0	0	8.0	0	0	0	0
0	4.0	0	0	0	0	0	12.0/

Memory operations

- Load/store operations move groups of data between registers and memory
 - Increased mem/instr ratio (intensity)
- Three types of addressing
 - <u>Unit stride</u>
 - » Contiguous block of information in memory
 - » Fastest: always possible to optimize this
 - <u>Non-unit</u> (constant) <u>stride</u>
 - » Harder to optimize memory system for all possible strides
 - » Prime number of data banks makes it easier to support different strides at full bandwidth
 - Indexed (gather-scatter)
 - » Vector equivalent of register indirect
 - » Good for sparse arrays of data
 - » Increases number of programs that vectorize

Interleaved Memory Layout

- Great for unit stride:
 - Contiguous elements in different DRAMs
 - Startup time for vector operation is latency of single read
- What about non-unit stride?
 - Above good for strides that are relatively prime to 8
 - Bad for: 2, 4



Avoiding Bank Conflicts

- Lots of banks int x[256][512]; for (j = 0; j < 512; j = j+1) for (i = 0; i < 256; i = i+1) x[i][j] = 2 * x[i][j];
- Even with 128 banks, since 512 is multiple of 128, conflict on word accesses
- SW: loop interchange or declaring array not power of 2 ("array padding")
- HW: Prime number of banks
 - bank number = address mod number of banks
 - address within bank = address / number of words in bank
 - modulo & divide per memory access with prime no. banks?
 - address within bank = address mod number words in bank
 - bank number? easy if 2^N words per bank

Finding Bank Number and Address within a bank

- <u>Problem</u>: Determine the number of banks, N_b and the number of words in each bank, N_w, such that:
 - given address x, it is easy to find the bank where x will be found, B(x), and the address of x within the bank, A(x).
 - for any address x, B(x) and A(x) are unique
 - the number of bank conflicts is minimized
- <u>Solution</u>: Use the Chinese remainder theorem to determine B(x) and A(x):

 $B(x) = x MOD N_b$

 $A(x) = x MOD N_w^{o}$ where N_b and N_w are co-prime (no factors)

- Chinese Remainder Theorem shows that B(x) and A(x) unique.
- Condition allows N_w to be power of two (typical) if N_b is prime of form 2^m-1.
- Simple (fast) circuit to compute (x mod N_b) when $N_b = 2^m 1$:
 - Since $2^k = 2^{k-m} (2^m-1) + 2^{k-m} \Rightarrow 2^k \text{ MOD } N_b = 2^{k-m} \text{ MOD } N_b = ... = 2^j \text{ with } j < m$
 - And, remember that: (A+B) MOD C = [(A MOD C)+(B MOD C)] MOD C
 - for every power of 2, compute single bit MOD (in advance)
 - B(x) = sum of these values MOD N_b (low complexity circuit, adder with ~ m bits)

Conclusion

- Vector is alternative model for exploiting ILP
 - If code is vectorizable, then simpler hardware, more energy efficient, and better real-time model than Out-of-order machines
 - Design issues include number of lanes, number of functional units, number of vector registers, length of vector registers, exception handling, conditional operations