# Vector Processors Architecutre

#### ENCS5331: Advanced Computer Architecture Fall 2024/2025 Instructor: Dr. Ayman Hroub

Special thanks to Prof. Onur Mutlu (ETHZ) for part of the slides

STUDENTS-HUB.com

#### Flynn's Taxonomy of Computers

- Mike Flynn, "Very High-Speed Computing Systems," Proc. of IEEE, 1966
- SISD: Single instruction operates on single data element
- SIMD: Single instruction operates on multiple data elements
  - ♦ Array processor

♦ Vector processor

MISD: Multiple instructions operate on single data element

♦ Closest form: systolic array processor, streaming processor

MIMD: Multiple instructions operate on multiple data elements (multiple instruction streams)

STUDENTS-HUB Son Multiprocessor and Multithreaded processor

# Data Parallelism

- Concurrency arises from performing the same operation on different pieces of data
  - ♦ Single instruction multiple data (SIMD)
  - $\diamond$  E.g., dot product of two vectors
- Contrast with data flow
  - ♦ Concurrency arises from executing different operations in parallel (in a data driven manner)
- Contrast with thread ("control") parallelism
  - ♦ Concurrency arises from executing different threads of control in parallel
- SIMD exploits operation-level parallelism on different data
  - ♦ Same operation concurrently applied to different pieces of data
  - ♦ A form of ILP where instruction happens to be the same across data

#### STUDENTS-HUB.com

# Loading/Storing Vectors from/to Memory

- Requires loading/storing multiple elements
- Elements separated from each other by a constant distance (stride)
  - $\diamond$  Assume stride = 1 for now
- Elements can be loaded in consecutive cycles if we can start the load of one element per cycle
  - ♦ Can sustain a throughput of one element per cycle
- Question: How do we achieve this with a memory that takes more than 1 cycle to access?
- Answer: Bank the memory; interleave the elements across banks

STUDENTS-HUB.com

4

#### Three Variations of SIMD

Vector Processors

- Graphics Processing Units (GPUs)
- Multimedia SIMD instruction set extensions

## Vector Processors (I)

- ✤ A vector is a one-dimensional array of numbers
- Many scientific/commercial programs use vectors

for (i = 0; i<=49; i++)

C[i] = (A[i] + B[i]) / 2

- A vector processor is one whose instructions operate on vectors rather than scalar (single data) values
- A single instruction operates on multiple data elements in consecutive time steps using the same space (PE)
- ✤ Basic requirements
  - $\diamond$  Need to load/store vectors  $\rightarrow$  vector registers (contain vectors)
  - $\diamond$  Need to operate on vectors of different lengths  $\rightarrow$  vector length register (VLEN)

• Stride: distance in memory between two elements of a vector STUDENTS-HUB.com

# Vector Stride Example: Matrix Multiply

✤ A and B matrices, both stored in memory in row-major order



$$\mathsf{A}_{4\mathsf{x}6} \: \mathsf{B}_{6\mathsf{x}10} \to \mathsf{C}_{4\mathsf{x}10}$$



Dot product of each row vector of A with each column vector of B

- Load A's row 0 (A0 through A5) into vector register V<sub>1</sub>
  - ♦ Each time, increment address by 1 to access the next column
  - ♦ Accesses have a stride of 1
- ✤ Load B's column 0 (B0 through B50) into vector register V<sub>2</sub>
  - $\diamond$  Each time, increment address by 10 to access the next row
  - ♦ Accesses have a stride of 10



STUDENTS-HUB.com

# Vector Processors (II)

- A vector instruction performs an operation on each element in consecutive cycles
  - ♦ Vector functional units are pipelined
  - ♦ Each pipeline stage operates on a different data element
- Vector instructions allow deeper pipelines
  - ♦ No intra-vector dependencies → no hardware interlocking needed within a vector
  - $\diamond~$  No control flow within a vector

#### Vector Architecture Basic Strucutre



STUDENTS-HUB.com

#### Storing Multiple Data Elements: Vector Registers

Each vector data register holds N M-bit values

 $\diamond$  Each register stores a vector





10

## Vector Processor Advantages

#### + No dependencies within a vector

- □ Pipelining & parallelization work really well
- □ Can have very deep pipelines (without the penalty of deep pipelines)
- + Each instruction generates a lot of work (i.e., operations)
  - Reduces instruction fetch bandwidth requirements
  - Amortizes instruction fetch and control overhead over many data
    - --> Leads to high energy efficiency per operation
- + No need to explicitly code loops
  - □ Fewer branches in the instruction sequence
- + Highly regular memory access pattern

STUDENTS-HUB.com

### Vector Processor Disadvantages

- -- Works (only) if parallelism is regular (data/SIMD parallelism)
  - ++ Vector operations
  - -- Very inefficient if parallelism is irregular
    - -- How about searching for a key in a linked list?

- -- Memory (bandwidth) can easily become a bottleneck, especially if
  - 1. compute/memory operation balance is not maintained
  - 2. data is not mapped appropriately to memory banks

# Vector Registers

- Each vector data register holds N M-bit values
- Vector control registers: VLEN, VSTR, VMASK
- Maximum VLEN can be N
  - $\diamond$  Maximum number of elements stored in a vector register
- Vector Mask Register (VMASK)

 $\diamond$  Indicates which elements of vector to operate on

 $\diamond$  Set by vector test instructions



STUDENTS-HUB.com

Uploaded By: Jibreel Bornat

13

# Vector Functional Units

- Use a deep pipeline to execute element operations
  - $\rightarrow$  fast clock cycle
- Control of deep pipeline is simple because elements in vector are independent

Six stage multiply pipeline



STUDENTS-Huderedic Kiste Asanovic

#### 15

# Memory Banking

- Memory is divided into banks that can be accessed independently; banks share address and data buses (to reduce memory chip pins)
- Can start and complete one bank access per cycle
- Can sustain N concurrent accesses if all N go to different banks



STUDENTS-Hick Berek Chiou

### Vector Memory System

- Next address = Previous address + Stride
- If (stride == 1) && (consecutive elements interleaved across banks) && (number of banks >= bank latency), then
  - $\diamond$  we can sustain 1 element/cycle throughput



# Scalar Code Example: Element-Wise Avg.

- rightarrow for i = 0 to 49
  - $\diamond$  C[i] = (A[i] + B[i]) / 2
- Scalar code (instruction and its latency in clock cycles)

MOVI R0 = 50	1
MOVA R1 = A	1
MOVA R2 = B	1 304 dynamic instructions
MOVA R3 = C	1
X: LD R4 = MEM[R1++]	11 ;autoincrement addressing
LD R5 = MEM[R2++]	11
ADD R6 = R4 + R5	4
SHFR R7 = R6 >> 1	1
ST MEM[R3++] = R7	11
DECBNZ R0, X	2 ;decrement and branch if NZ
STUDENTS-HUB.com	Uploaded E

# Scalar Code Execution Time (In Order)

- Scalar execution time on an in-order processor with 1 bank
  - ♦ First two loads in the loop cannot be pipelined: 2\*11 cycles
  - 4 + 50\*40 = 2004 cycles

- Scalar execution time on an in-order processor with 1 bank with 2 memory ports (two different memory accesses can be serviced concurrently) OR 2 banks (where arrays A and B are stored in different banks)
  - $\diamond$  First two loads in the loop can be pipelined: 1 + 11 cycles
  - 4 + 50\*30 = 1504 cycles

STUDENTS-HUB.com

#### Vectorizable Loops

- ✤ A loop is vectorizable if each iteration is independent of any other
- ✤ For i = 0 to 49
  - ♦ C[i] = (A[i] + B[i]) / 2
- Vectorized loop (each instruction and its latency):

MOVI VLEN = 50	1
MOVI VSTR = 1	1 7 dynamic instructions
VLD V0 = A	11 + VLEN – 1
VLD V1 = B	11 + VLEN – 1
VADD V2 = V0 + V1	4 + VLEN – 1
VSHFR V3 = V2 >> 1	1 + VLEN – 1
VSTC = V3	11 + VLEN – 1

STUDENTS-HUB.com

# Basic Vector Code Performance (I)

- Assume no chaining (no vector data forwarding)
  - i.e., output of a vector functional unit cannot be used as the direct input of another
  - The entire vector register needs to be ready before any element of it can be used as part of another operation
- 1 memory port (one address generator) per bank
- If a memory banks (word-interleaved: consecutive elements of an array are stored in consecutive banks)



STUDENTS-HUB.com

### Basic Vector Code Performance (II)

- Why 16 banks?
  - ♦ 11-cycle memory access latency
  - ♦ Having 16 (>11) banks ensures there are enough banks to overlap enough memory operations to cover memory latency
- The above assumes a unit stride (i.e., stride = 1)
  - ♦ Correct for our example program
- ✤ What if stride > 1?
  - ♦ How do you ensure we can access 1 element per cycle when memory latency is 11 cycles?

STUDENTS-HUB.com

21

### Vector Chaining

Vector chaining: Data forwarding from one vector functional unit to another



STUDENTS-HeideBredichiste Asanovic

#### Vector Code Performance - Chaining

Vector chaining: Data forwarding from one vector functional unit to another



STUDENTS-HUB.com

#### Vector Code Performance - Multiple Memory Ports

Chaining and 2 load ports, 1 store port in each bank



STUDENTS-HUB.com

Uploaded By: Jibreel Bornat

# Questions (I)

- What if # data elements > # elements in a vector register?
  - ♦ Idea: Break loops so that each iteration operates on # elements in a vector register
    - E.g., 527 data elements, 64-element VREGs
    - 8 iterations where VLEN = 64
    - I iteration where VLEN = 15 (need to change value of VLEN)
  - ♦ Called vector stripmining

# Questions (II)

- What if vector data is not stored in a strided fashion in memory? (irregular memory access to a vector)
  - ♦ Idea: Use indirection to combine/pack elements into vector registers
  - ♦ Called scatter/gather operations
  - Doing so also helps with avoiding useless computation on sparse vectors (i.e., vectors where many elements are 0)

#### Gather/Scatter Operations

Want to vectorize loops with indirect accesses:

```
for (i=0; i<N; i++)
    A[i] = B[i] + C[D[i]]</pre>
```

Indexed load instruction (Gather)

LV vD, rD	#	Load indices in D vector
LVI vC, rC, vD	#	Load indirect from rC base
LV vB, rB	#	Load B vector
ADDV.D vA,vB,vC	#	Do add
SV vA, rA	#	Store result

STUDENTS-HUB.com

27

### Gather/Scatter Operations

- Gather/scatter operations often implemented in hardware to handle sparse vectors (matrices) or indirect indexing
- Vector loads and stores use an index vector which is added to the base register to generate the addresses

#### ✤ Scatter example

Index Vector	Data Vector (to Store)	Stored Ve	ctor (in Memory)
0	3.14	Base+0	3.14
2	6.5	Base+1	Х
6	71.2	Base+2	6.5
7	2.71	Base+3	Х
		Base+4	Х
		Base+5	Х
		Base+6	71.2
JB.com		Base+7	2.71 Uploaded By: Jibre

STUDENTS-HU

eel Bornat

# Conditional Operations in a Loop

What if some operations should not be executed on a vector (based on a dynamically-determined condition)?

```
for (i=0; i<N; i++)</pre>
```

```
if (a[i] != 0) then b[i]=a[i]*b[i]
```

- Idea: Masked operations
  - VMASK register is a bit mask determining which data element should not be acted upon

```
VLD V0 = A
VLD V1 = B
VMASK = (V0 != 0)
VMUL V1 = V0 * V1
VST B = V1
```

This is predicated execution. Execution is predicated on mask bit.
 STUDENTS-HUB.com
 Uploaded By: Jibreel Bornat

#### Another Example with Masking

for (i	= 0; i	. < 64; ++i)					
if (a	a[i] >=	= b[i])	Steps to execute the loop in SIMD code				
c[i] = a[i] else c[i] = b[i]		.] = a[i] .] = b[i]	1. Compare A, B to get VMASK				
			2. Masked store of A into C				
A 1	B 2	VMASK 0	3. Complement VMASK				
2 3	2 2 10	1 1	4. Masked store of B into C				
4 -5 0	10 -4 -3	0 0 1					
6 -7	5 -8	1 1					

STUDENTS-HUB.com

#### Masked Vector Instructions

#### Simple Implementation

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

	M[7]=1	A[7]	B[7]	
	M[6]=0	A[6]	B[6]	
	M[5]=1	A[5]	B[5]	
	M[4]=1	A[4]	B[4]	
	M[3]=0	A[3]	B[3]	
	M[2]=0			
	M[0]=0	7	C[0]	
STI	Write En	able HUB.c	Write da	ata port

#### **Density-Time Implementation**

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



#### Stride and bank count

- ♦ As long as stride and bank count are *relatively prime* to each other and there are enough banks to cover bank access latency, we can sustain 1 element/cycle throughput
- Storage format of a matrix
  - Row major: Consecutive elements in a row are laid out consecutively in memory
  - Column major: Consecutive elements in a column are laid out consecutively in memory
  - $\diamond\,$  You need to change the stride when accessing a row versus column

# Bank Conflicts in Matrix Multiplication

✤ A and B matrices, both stored in memory in row-major order



$$\mathsf{A}_{4x6} \: \mathsf{B}_{6x10} \to \mathsf{C}_{4x10}$$

Dot product of each row vector of A with each column vector of B

#### Load A's row 0 into vector register V<sub>1</sub>

- $\diamond$  Each time, increment address by 1 to access the next column
- ♦ Accesses have a stride of 1
- Load B's column 0 into vector register V<sub>2</sub>
  - ♦ Each time, increment address by 10
  - ♦ Accesses have a stride of 10

Different strides can lead to bank conflicts

How do we minimize them?

STUDENTS-HUB.com

Bo	0	1	2	3	4	5	6	7	8	9
	10	11	12	13	14	15	16	17	18	19
	20									
	30									
	40									
	50									

# Minimizing Bank Conflicts

More banks

- More ports in each bank
- Better data layout to match the access pattern
  - $\diamond$  Is this always possible?
- Better mapping of address to bank
  - ♦ E.g., randomized mapping

STUDENTS-HUB.com

#### **Vector Instruction Execution**



#### Vector Unit Structure



#### Vector Instruction Level Parallelism

Can overlap execution of multiple vector instructions

- ♦ Example machine has 32 elements per vector register and 8 lanes
- ♦ Example with 24 operations/cycle (steady state) while issuing 1 vector instruction/cycle



#### Automatic Code Vectorization



#### Vector/SIMD Processing Summary

- Vector/SIMD machines are good at exploiting regular datalevel parallelism
  - ♦ Same operation performed on many data elements
  - ♦ Improve performance, simplify design (no intra-vector dependencies)
- Performance improvement limited by vectorizability of code
  - ♦ Scalar operations limit vector machine performance
  - ♦ Remember Amdahl's Law
- Many existing ISAs include (vector-like) SIMD operations
  - ♦ Intel MMX/SSEn/AVX, PowerPC AltiVec, ARM Advanced SIMD, RISC V V (RVV)

STUDENTS-HUB.com