

## HY425 Lecture 10: Vector processors

Dimitrios S. Nikolopoulos

University of Crete and FORTH-ICS

November 5, 2010

### Vector processors

#### Vector instructions

- ▶ a loop in an instruction
- ▶ explicit parallelism, programmer guarantees independence
- ▶ hazard check between blocks of operations
- ▶ exploit memory parallelism
- ▶ less loop overhead, control hazards

#### Benefits

- ▶ Scientific, engineering, multimedia applications
  - ▶ Computation on arrays dominates, data parallelism

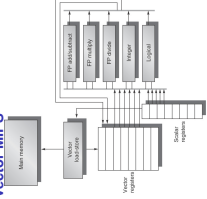
### Limitations of ILP

#### ILP walls

- ▶ Hard to exploit higher degrees of ILP
- ▶ Deeper pipelines, wider instruction issue
  - ▶ Increased hardware complexity with small performance gain
  - ▶ Heavy burden on software on statically scheduled processors

### Vector architecture

#### Vector MIPS



#### Features compared to scalar

- ▶ Vector registers (wide)
- ▶ Vector FUs, load/store units
- ▶ Parallel pipelines or lanes
- ▶ Vector ISA

## Example: SAXPY, DAXPY

### MIPS code

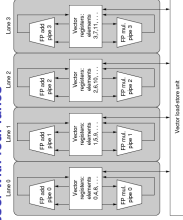
```
LD F0, A
ADDI R4, Rk, #512 ;last address to load
Loop: LD F2, 0(Rk) ;load X(1)
      MADD F2, F2, F0 ;a * X(1)
      ADDI R4, R4, 4 ;increment to X(2)
      SD 0(Rk), F2 ;store into Y(1)
      SD 0(Rk), F4 ;store into Y(1)
      ADDIU Rk, Rk, #8 ;increment address to X
      ADDIU Rk, Rk, #8 ;increment address to Y
      SDBU R20, R4, Rk ;compute bound
      SNEZ R20, Loop ;check if done
```

### VMIPS code

```
LD F0, A ;load scalar a
LV V1, Rk ;load vector X
MULVSD V2, V1, F0 ;vector-scalar multiply
ADDVSD V3, V2, V3 ;add vector y
SV Rk, V4 ;store the result
```

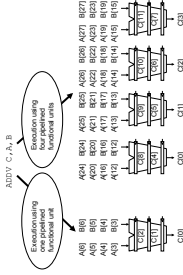
## Vector instruction-level parallel execution

### Vector processor with four lanes



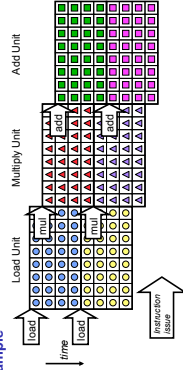
## Vector instruction execution

### One versus multiple pipelined units



## Multi-lane vector unit

### Example



## Vector processor performance factors

- ▶ Length of vector registers
- ▶ Vector operation initiation rate ( $\geq 1$  per cycle)
- ▶ Number of lanes (parallel pipelines)
- ▶ Vector instructions executing in same cycle – **convoy**
- ▶ Time to execute a single convoy – **chime**

```

1: LV V1, R0      ; load vector X
2: MULVSD V2, V1, R0 ; multiply
3: LV V3, R1, V1, R0 ; load vector Y
4: ADDV V4, V2, V3 ; add
5: SV R2, V4     ; store the result
    
```

## Vector processor performance measurement

### Vector instruction latency

- ▶ Pipelined vector FUs
- ▶ Convoys do not overlap due to dependences
- ▶ Vector startup time:
  - ▶ time until first result out from pipelined FU
- ▶ Rest of results come out one per cycle
- ▶  $S + VL - 1$ ,  $S$ : startup time,  $VL$ : vector length

## Example: DAXPY

### 4 convoys

Functional unit	Latency		
load/store unit	12		
multiply unit	7		
add unit	6		
Convoy	Start time	First result	Last result
1: LD	0	12	11+n
2: MULVSD, LV	12+n	24+n	23+2n
3: ADDV	24+2n	30+2n	29+3n
4: SV	30+3n	42+3n	41+4n

- ▶ Latency of convoy depends on slowest instruction in convoy
- ▶ Latency of second convoy is latency of slower LV instruction
- ▶ Shorter vector length implies more pipeline restarts in FUs

### Requirements from memory system

- ▶ Load/store units need word/cycle bandwidth from memory
- ▶ Hard to meet demand, even with advanced memory systems
  - ▶ Multiple memory banks, interleaving
  - ▶ High memory bank cycle time
  - ▶ Multiple loads/stores per clock need to be supported
- ▶ How many banks are needed to sustain throughput to load/store units?

## Vector load/store unit

## Example

- 6-cycle bank latency, starting address=136
- ▶ double-word (8-byte) bank interleaving

cycle	0	1	2	3	4	5	6	7
	Bank							
0	136	144						
1	busy	busy	152					
2	busy	busy	busy	160				
3	busy	busy	busy	busy	168			
4	busy	busy	busy	busy	busy	176		
5	busy	busy	busy	busy	busy	busy	184	
6	192							
7	busy	200						
8	busy	busy	208					
9	busy	busy	busy	216				
10	busy	busy	busy	busy	224			
11	busy	busy	busy	busy	busy	232		
12	busy	busy	busy	busy	busy	busy	240	
13	busy	busy	busy	busy	busy	busy	busy	248
14	256							
15	busy	264						
16	busy	busy	272					

## Stripmining

### Vector length

- ▶ Typically shorter than real vector sizes (2–64 elements)
- ▶ Loop transformation to vectorize code
- ▶ Partition loop into vector instructions

### Example

```

/* Original code */
for (i=1; i<n; i++)
    Y[i] = a * X[i] + Y[i]

/* Vectorized code */
int i;
int n;
void *a;
void *x;
void *y;
for (i=0; i < n/MVL; i++)
    Y[i*low; i < low+VL-1; i++)
        Y[i] = a * X[i] + Y[i]; /*main operation*/
low = low + VL; /*start of next vector*/
VL = MVL; /*reset the length to mask*/
    
```

## Vector processor performance metrics

- ▶ **FLOPS**: Floating point operations per second
- ▶  $R_n$ : FLOPS with  $MVL = n$
- ▶  $T_n$ : Time with  $MVL = n$
- ▶  $R_\infty$ : FLOPS with  $MVL = \infty$
- ▶  $M_{1/2}$ : Vector length to achieve  $\frac{R_\infty}{2}$  FLOPS
- ▶  $N_v$ : Vector length to achieve more FLOPS than scalar

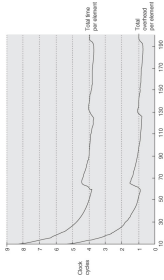
## Stripmining performance

- ▶  $T_{start}$ : vector pipeline startup cost for loop body
- ▶  $T_{loop}$ : scalar code per outer loop iteration
- ▶ **chimes**: number of chimes needed to execute conveys in the loop
- ▶ Pipelined vector units, no overlap of conveys

$$T_n = \left\lceil \frac{n}{MVL} \right\rceil \times (T_{loop} + T_{start}) + n \times \text{chimes} \quad (1)$$

## Stripmining performance (cont.)

- ▶ Overhead more than 50% of total time for short vectors
- ▶ Jumps indicate iterations of outer loop after stripmining



## Vector strides

### Matrix-matrix multiplication

```

for (i=0; i<100; i++)
  for (j=0; j<100; j++) {
    A[i][j] = 0.0;
    for (k=0; k<100; k++)
      A[i][j] = A[i][j] + B[i][k]*C[k][j]
  }
    
```

### Vectorization considerations

- ▶ Innermost loop over row of B and column of C
- ▶ Row-major array allocation in C
- ▶ Vector with non-adjacent elements for k-th column of C
- ▶ Vector registers pack data with non-unit strides in memory
  - ▶ Compare to caches where blocks store only unit-stride data

## Vector strides

### Instructions for strided data access

- ▶ **LWVS**: load vector with stride
- ▶ **SVWS**: store vector with stride
- ▶ Also known as **gather/scatter** operations

## Vector strides (cont.)

### Implications for memory system

- ▶ Non-unit strides may increase conflicts in memory banks
- ▶ Example:  $100 \times 100$  matrix-matrix multiplication, double-word (8-byte) interleaving, array type double, starting address of C = 136

cycle	0	1	2	3	4	5	6	7
	136							
0	busy	busy	936					
1	busy	busy	busy	busy				
2	busy	busy	busy	busy	busy			
3	busy	busy	busy	busy	busy	busy		
4	busy	busy	busy	busy	busy	busy	busy	
5	1736							
6	busy	busy	2536					
7	busy	busy	busy	busy	2536			
8	busy	busy	busy	busy	busy	busy		
9	busy	busy	busy	busy	busy	busy	busy	
10	busy	busy	busy	busy	busy	busy	busy	
11	busy	busy	busy	busy	busy	busy	busy	
12	3036							
13	busy	busy	3036					
14	busy	busy	busy	busy	3036			
15	busy	busy	busy	busy	busy	busy	busy	
16	busy	busy	busy	busy	busy	busy	busy	



## Scatter-gather operations

### Indirection arrays

- ▶ Sparse matrix codes access arrays indirectly (e.g.  $A[B[i]]$ )
- ▶ Indirection arrays implemented with index vector registers
- ▶ **Gather**: Collect elements pointed to by index vector
- ▶ **Scatter**: Distribute elements to locations pointed to by index vector

```

for (i=0; i<64; i++)
  A[K(i)] = A[I(i)] - C[M(i)]
    
```

### VMIPS gather (LV), scatter (SV) instructions

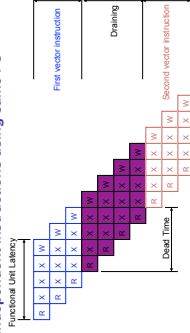
Assume  $R_{16}$ ,  $R_{18}$ ,  $R_{16}$ ,  $R_{16}$  have starting addresses of  $A$ ,  $K$ ,  $C$ ,  $M$ .

```

LV  V0, R16 ; load K
LV  V0, (R0+V0) ; load A(K(i))
LV  V0, R18 ; load M
LV  V0, (R0+V0) ; load M(i)
ADDV D, V0, V0, V0 ; add 3x
SVI  (R0+V0), V0 ; store A(K(i))
    
```

## Pipeline startup and draining latency

### Two independent instructions using same FU



## Vector processor performance metrics

- ▶ **FLOPS**: Floating point operations per second
- ▶  $R_n$ : FLOPS with  $MVL = n$
- ▶  $T_n$ : Time with  $MVL = n$
- ▶  $R_{\infty}$ : FLOPS with  $MVL = \infty$
- ▶  $M_{1/2}$ : Vector length to achieve  $\frac{R_{\infty}}{2}$  FLOPS
- ▶  $N_v$ : Vector length to achieve more FLOPS than scalar

$$T_n = \left\lceil \frac{n}{MVL} \right\rceil \times (T_{loop} + T_{start}) + n \times T_{chime} \quad (3)$$

## DAXPY example

### Execution with 500 MHz dual-issue vector processor

LV V1,Rx MULVS V2,V1,F0 Convoy 1: Chained load and multiply  
LV V3,Ry ADDVD V4,V2,V3 Convoy 2: Chained load and add  
SV V4,Ry Convoy 3: Store result

### Assumptions

- ▶  $MVL=64$ ,  $LV/SV$  FU latency=12,  $ADDVD$  FU latency=6,  $MULVS$  FU latency=7,  $T_{loop} = 15$
- ▶  $T_{start} = 12 + 7 + 12 + 6 + 12 = 49$
- ▶  $T_{loop} = 15$
- ▶  $T_{chime} = 3$
- ▶  $T_n = \left\lceil \frac{n}{64} \right\rceil (49 + 15) + 3n$
- ▶  $T_n \leq 4n + 64$





## Intel SSE extensions

### Pentium II onwards

- ▶ 8–16 128-bit registers (XMM registers), plus 8–16 64-bit MMX registers
- ▶ Registers enable packed data (2 double-words, 4 words, 8 half-words, 16 chars) operations and/or scalar operations
- ▶ Single-cycle SSE ALU operations
- ▶ SIMD memory-register load/store operations
  - ▶ prefetching
  - ▶ streaming stores

## Intel SSE extensions

### Example: inner product

```
typedef float v8f; // attribute... (mode(V8F)); // floating point vector type
float x[k]; float y[k]; // operand vectors of length k
float inner_product = 0.0; temp[4]; // zero the accumulator
v8f acc; // zero the accumulator
acc = _builtin_is32_xorps(acc, acc); // zero the accumulator
for (int i = 0; i < (k - 3); i += 4) {
    X = _builtin_is32_loadps(x+i); // Load groups of four floats
    Y = _builtin_is32_loadps(y+i);
    acc = _builtin_is32_accps(acc, _builtin_is32_mulps(X, Y));
}
_builtin_is32_storeps(temp, acc); // add the accumulated value
inner_product = temp[0] + temp[1] + temp[2] + temp[3];
for (; i < k; i++) // add up the remaining floats
    inner_product += x[i] * y[i];
```

xorps: bit-wise xor of single-precision floating point values  
loadps: move four unaligned packed single-precision floats from memory  
mulps, accps: multiply, add single-precision floats