



CS 152 Computer Architecture and Engineering

Lecture 17: Vectors Part II

Krste Asanovic

Electrical Engineering and Computer Sciences
University of California, Berkeley

<http://www.eecs.berkeley.edu/~krste>

<http://inst.cs.berkeley.edu/~cs152>

April 1, 2010

CS152, Spring 2010



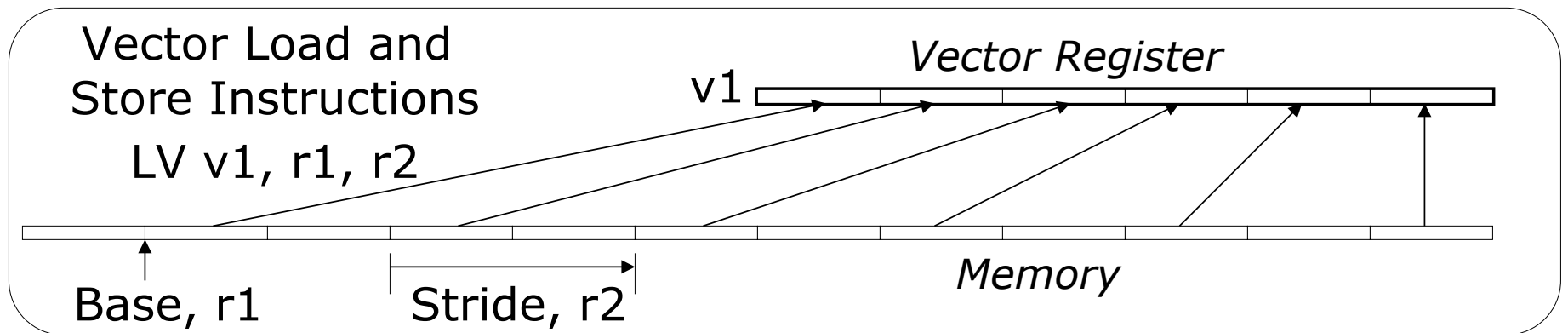
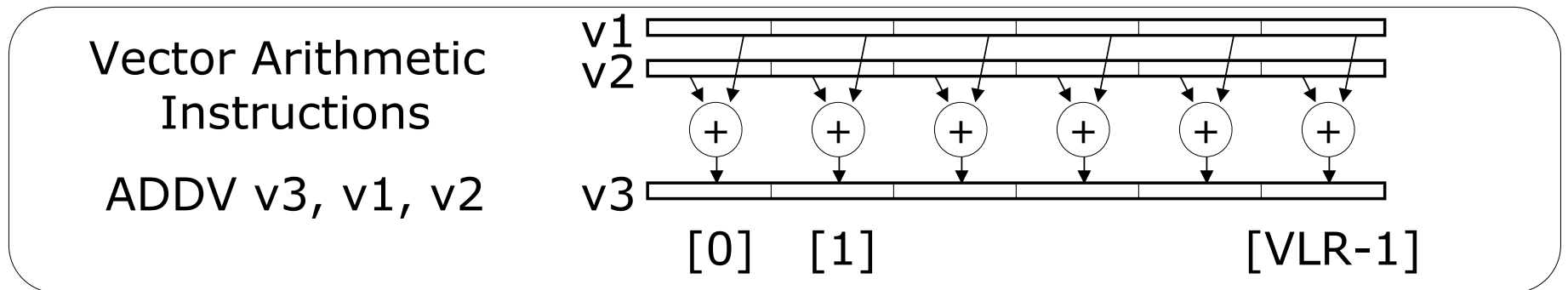
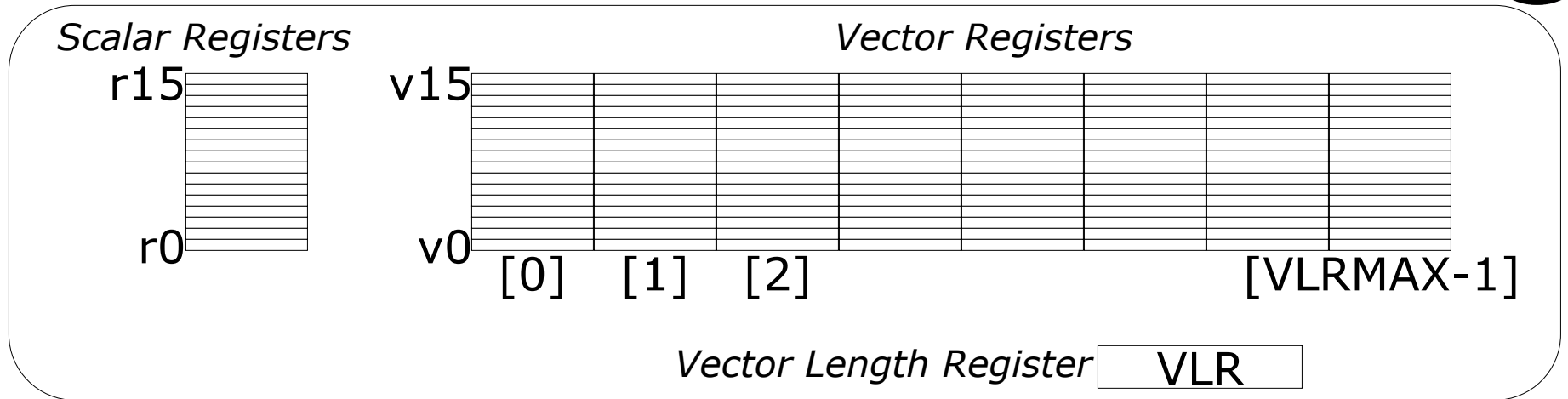
Last Time: Vector Supercomputers

Epitomized by 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



Vector Programming Model





Vector Code Example

# C code	# Scalar Code	# Vector Code
<pre>for (i=0; i<64; i++) C[i] = A[i] + B[i];</pre>	<pre>LI R4, 64 loop: L.D F0, 0(R1) L.D F2, 0(R2) ADD.D F4, F2, F0 S.D F4, 0(R3) DADDIU R1, 8 DADDIU R2, 8 DADDIU R3, 8 DSUBIU R4, 1 BNEZ R4, loop</pre>	<pre>LI VLR, 64 LV V1, R1 LV V2, R2 ADDV.D V3, V1, V2 SV V3, R3</pre>



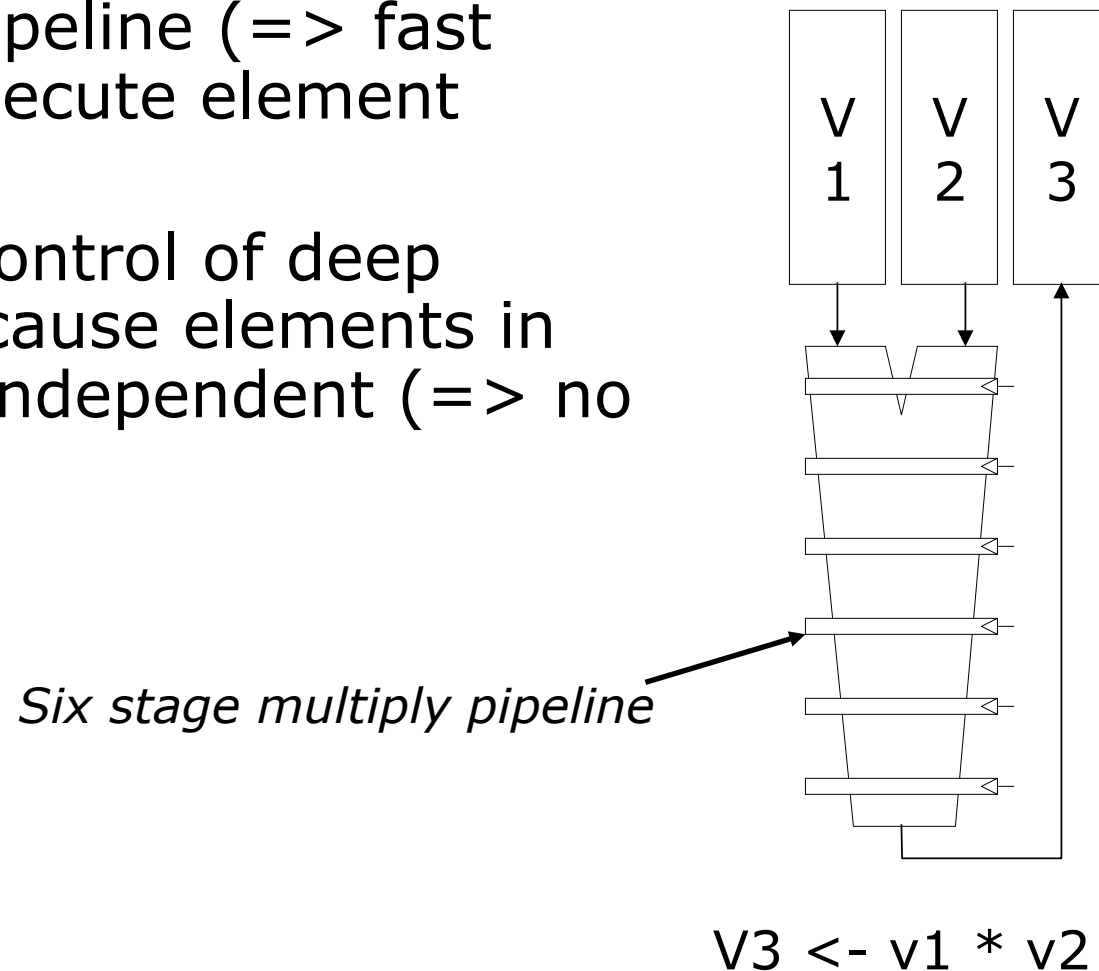
Vector Instruction Set Advantages

- Compact
 - one short instruction encodes N operations
- Expressive, 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*)



Vector Arithmetic Execution

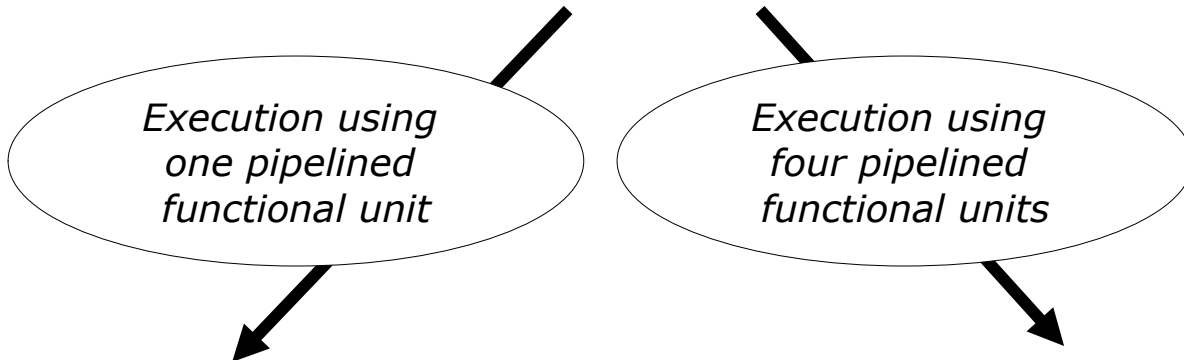
- Use deep pipeline (\Rightarrow fast clock) to execute element operations
- Simplifies control of deep pipeline because elements in vector are independent (\Rightarrow no hazards!)



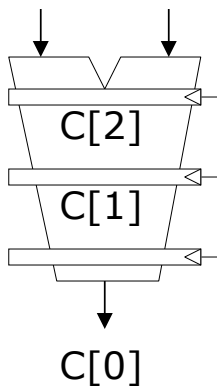


Vector Instruction Execution

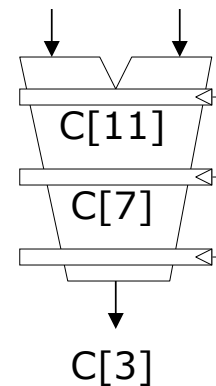
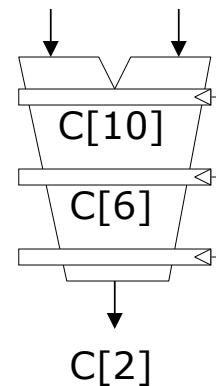
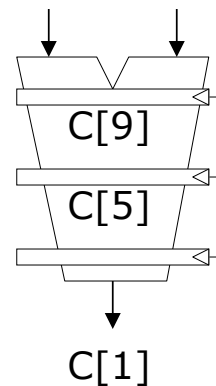
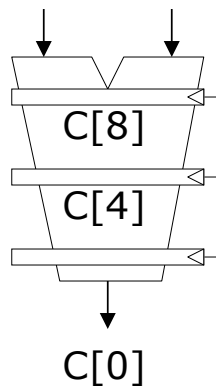
ADDV C,A,B



A[6] B[6]
A[5] B[5]
A[4] B[4]
A[3] B[3]



A[24] B[24] A[25] B[25] A[26] B[26] A[27] B[27]
A[20] B[20] A[21] B[21] A[22] B[22] A[23] B[23]
A[16] B[16] A[17] B[17] A[18] B[18] A[19] B[19]
A[12] B[12] A[13] B[13] A[14] B[14] A[15] B[15]

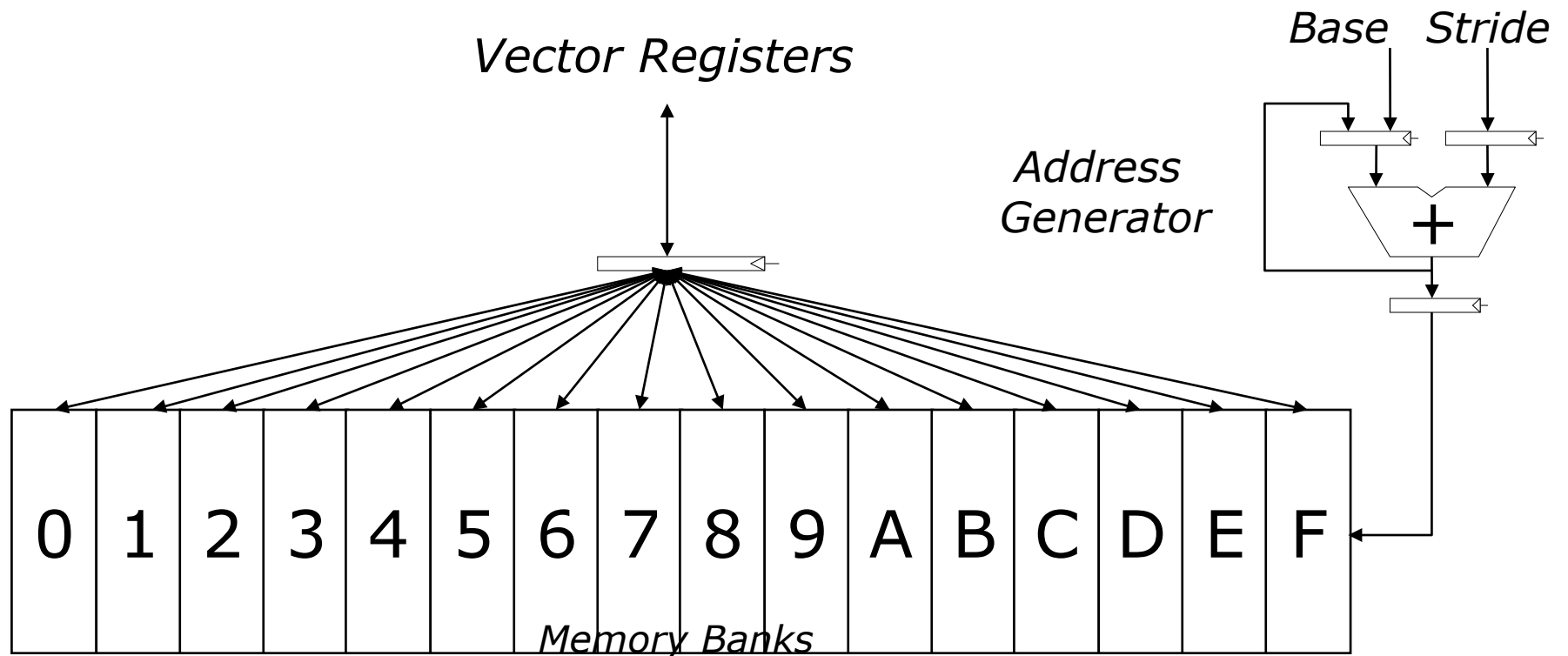




Vector Memory System

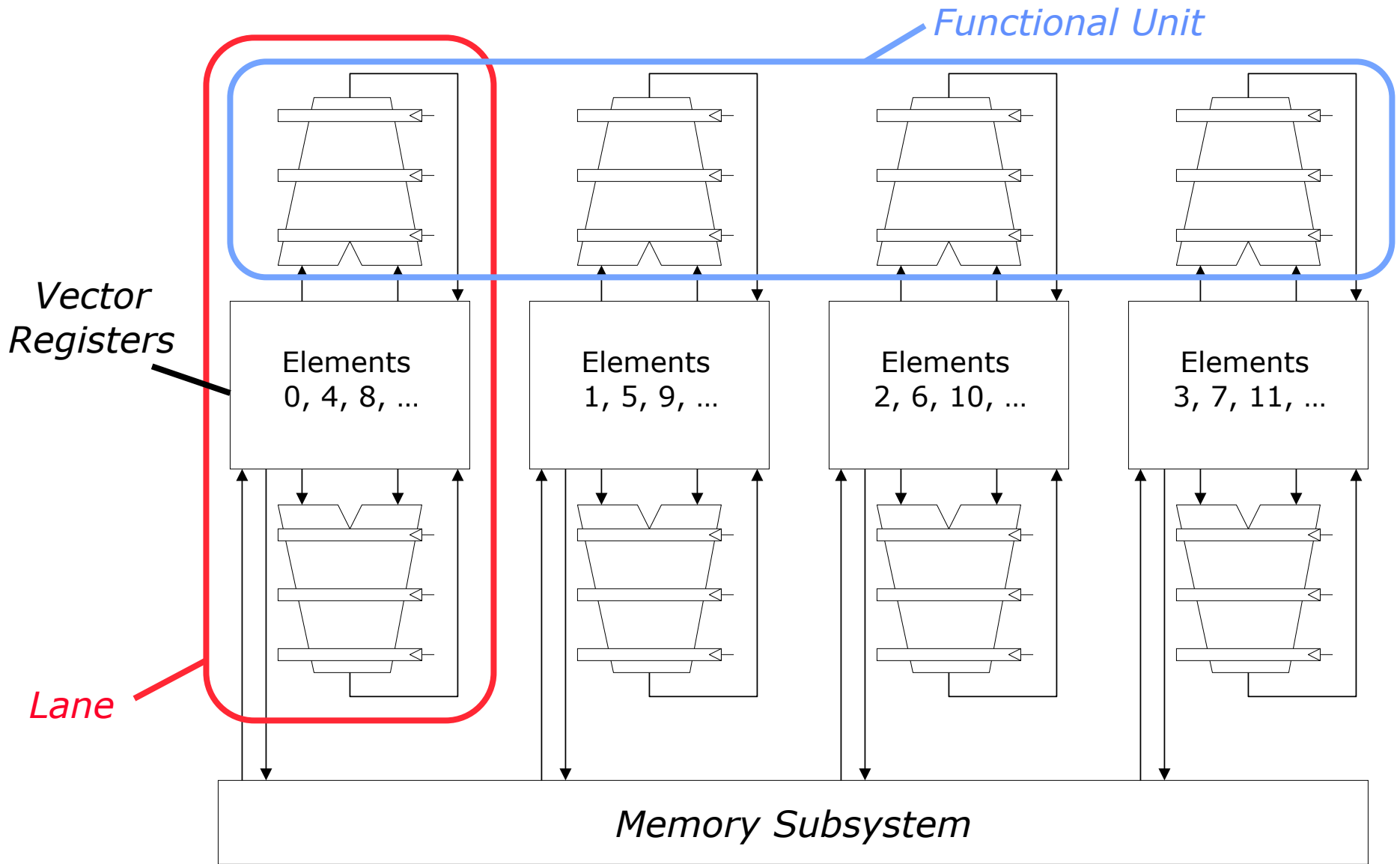
Cray-1, 16 banks, 4 cycle bank busy time, 12 cycle latency

- *Bank busy time*: Time before bank ready to accept next request



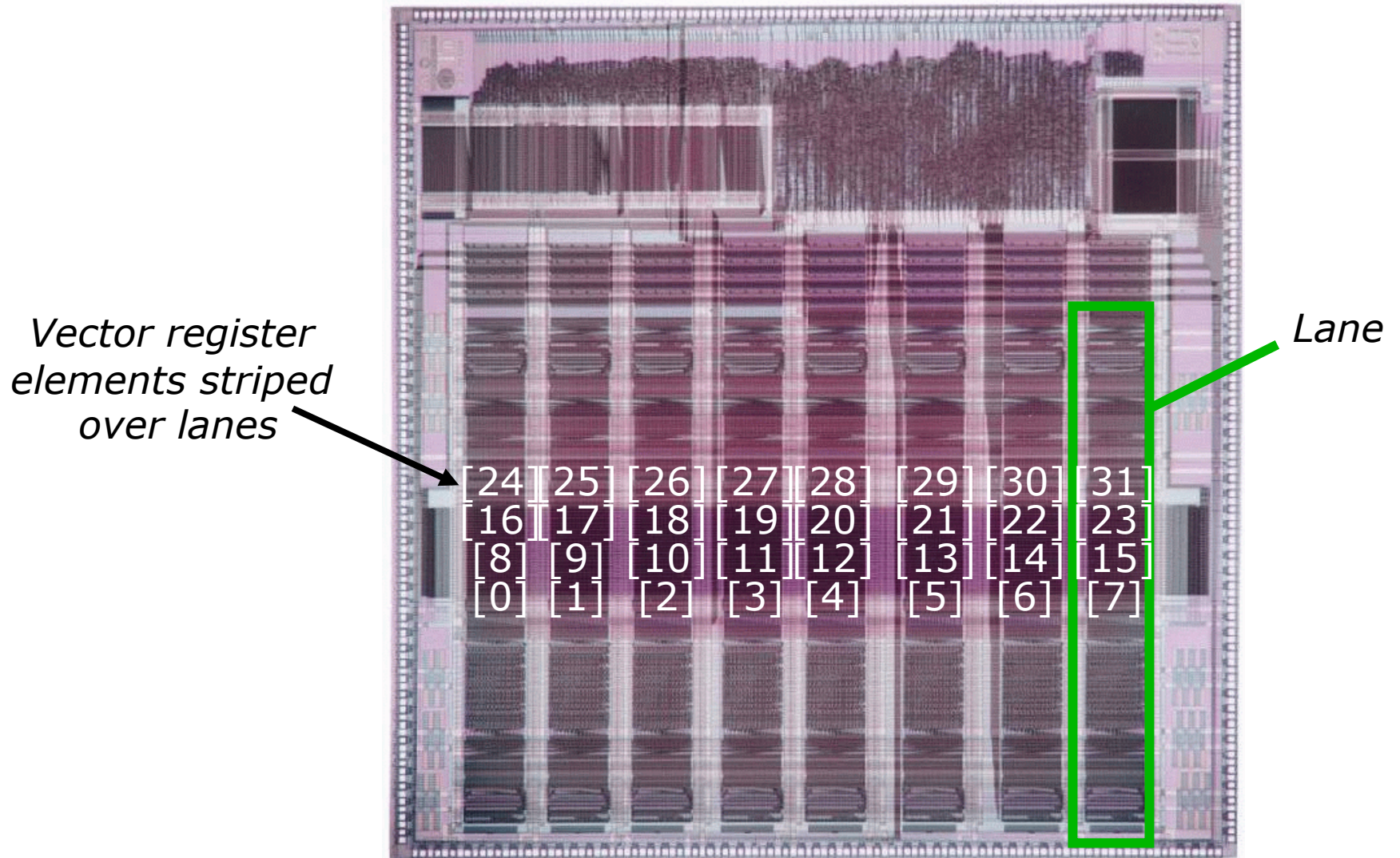


Vector Unit Structure





T0 Vector Microprocessor (UCB/ICSI, 1995)

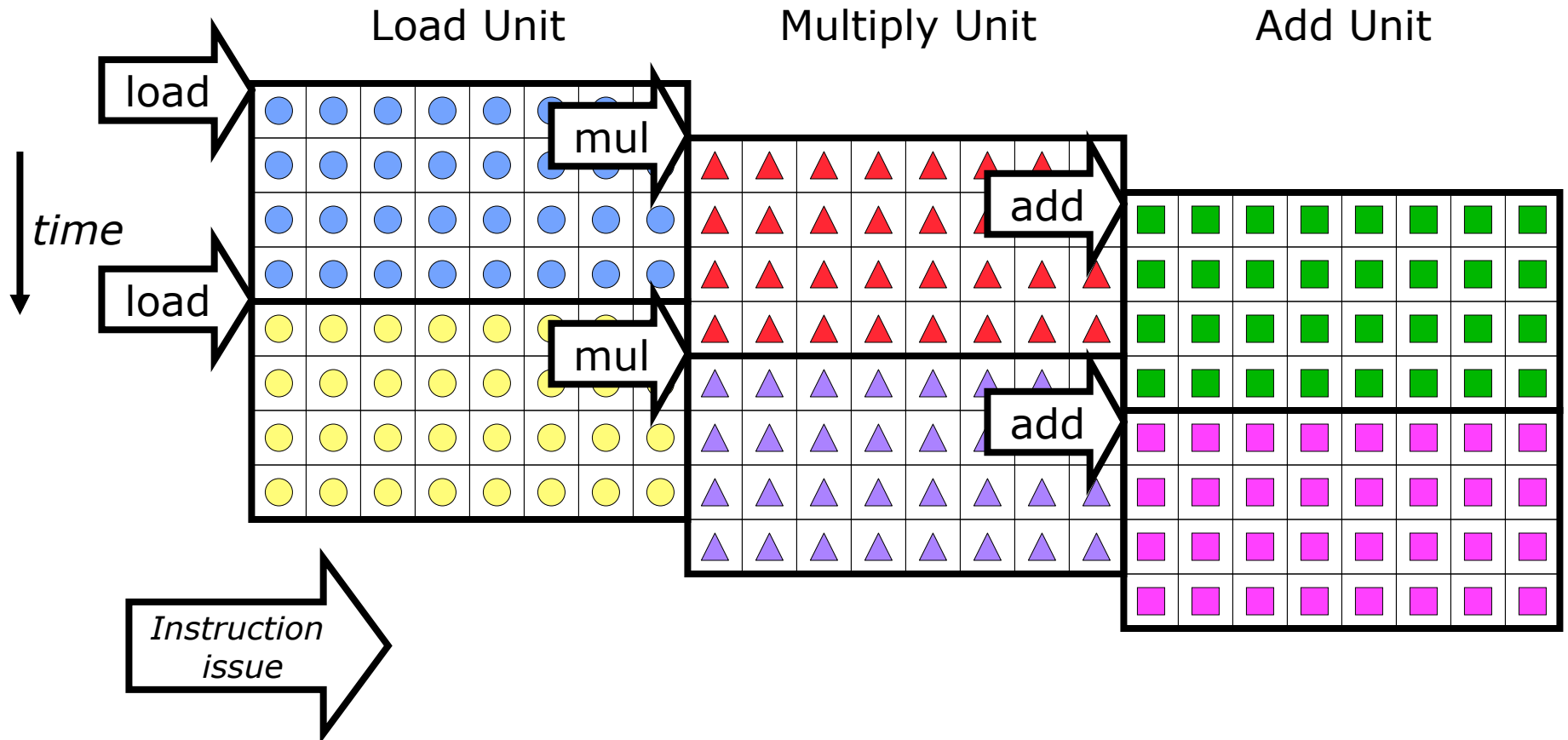




Vector Instruction Parallelism

Can overlap execution of multiple vector instructions

- example machine has 32 elements per vector register and 8 lanes



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

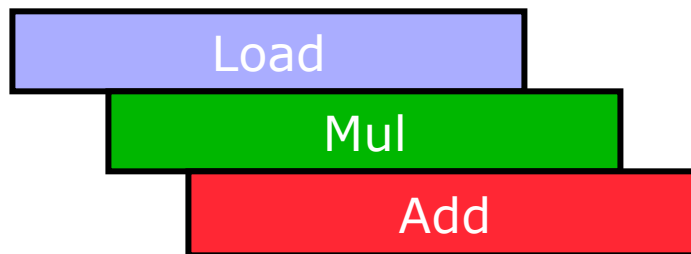


Vector Chaining Advantage

- 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



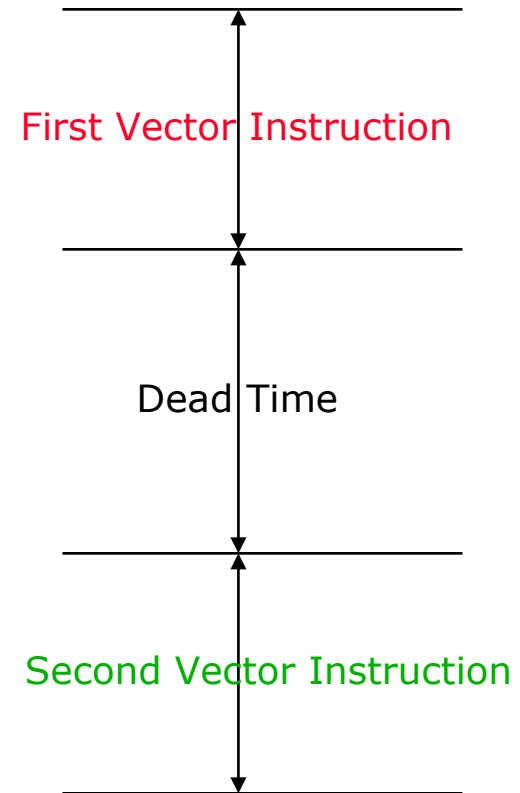
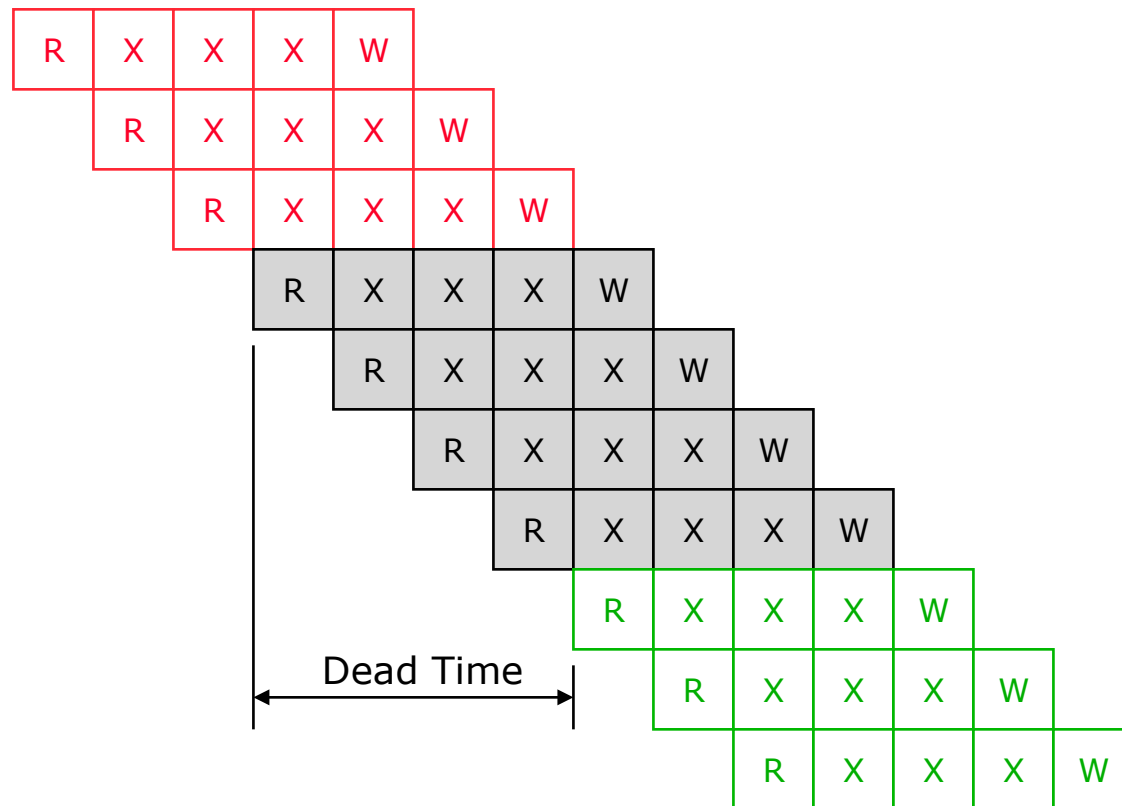
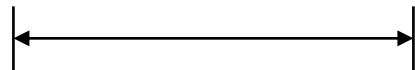


Vector Startup

Two components of vector startup penalty

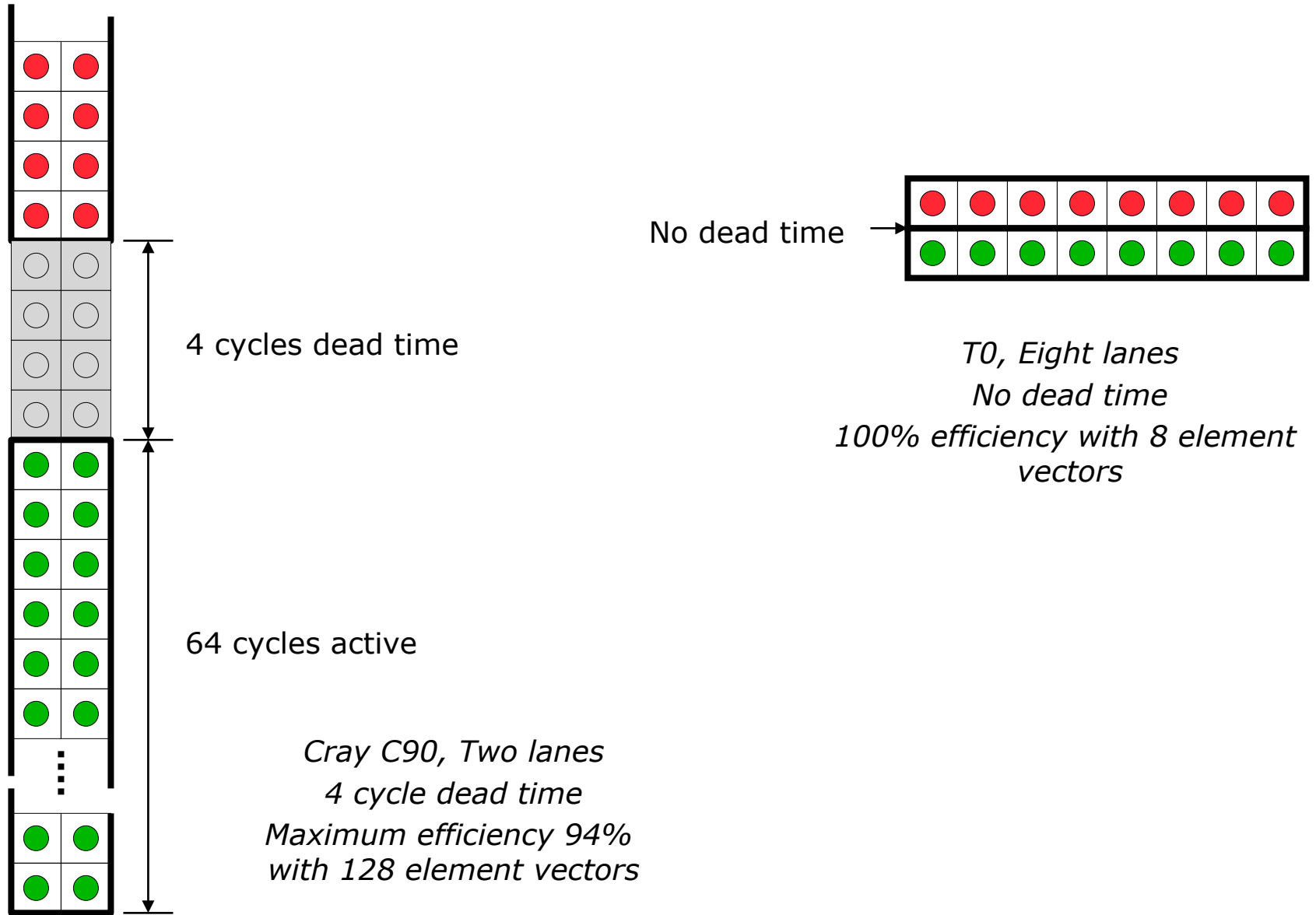
- functional unit latency (time through pipeline)
- dead time or recovery time (time before another vector instruction can start down pipeline)

Functional Unit Latency





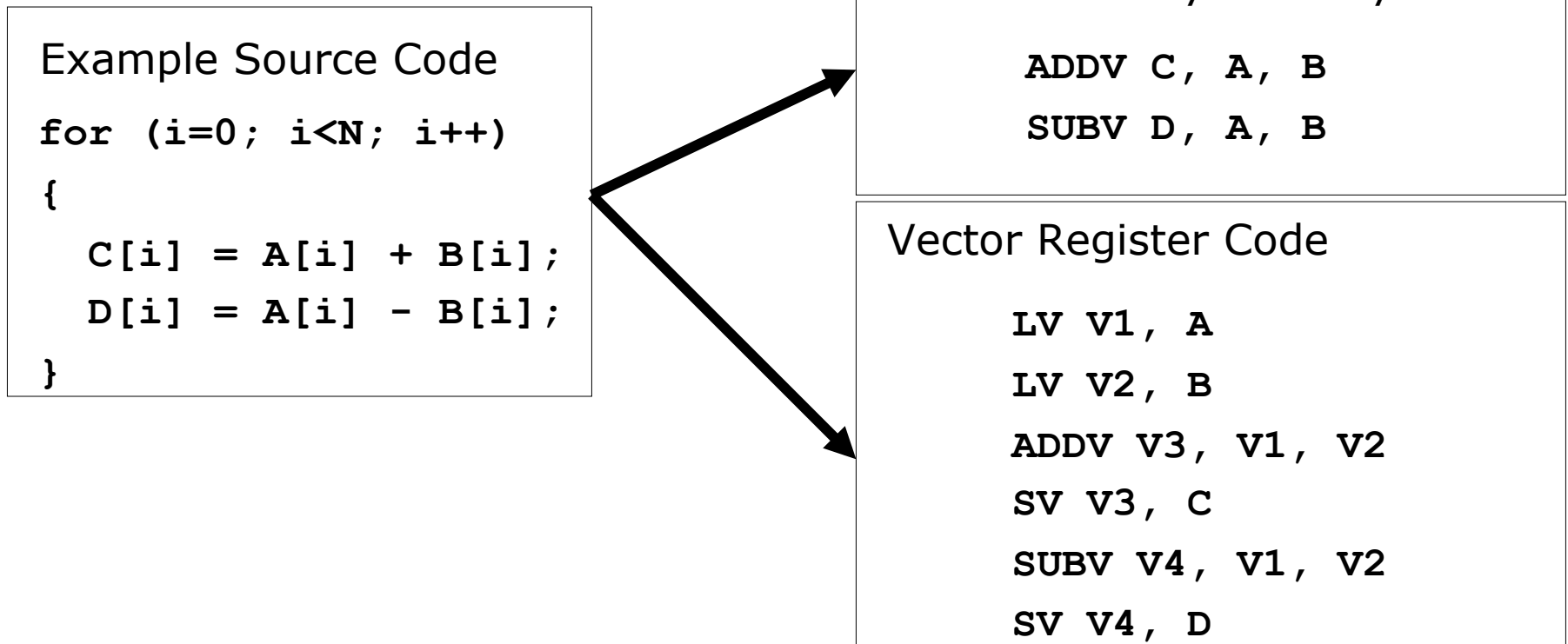
Dead Time and Short Vectors





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
- VMMA makes it difficult to overlap execution of multiple vector operations, why?
 - Must check dependencies on memory addresses
- VMMA incurs greater startup latency
 - Scalar code was faster on CDC Star-100 for vectors < 100 elements
 - For Cray-1, vector/scalar breakeven point was around 2 elements

⇒ *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)



CS152 Administrivia

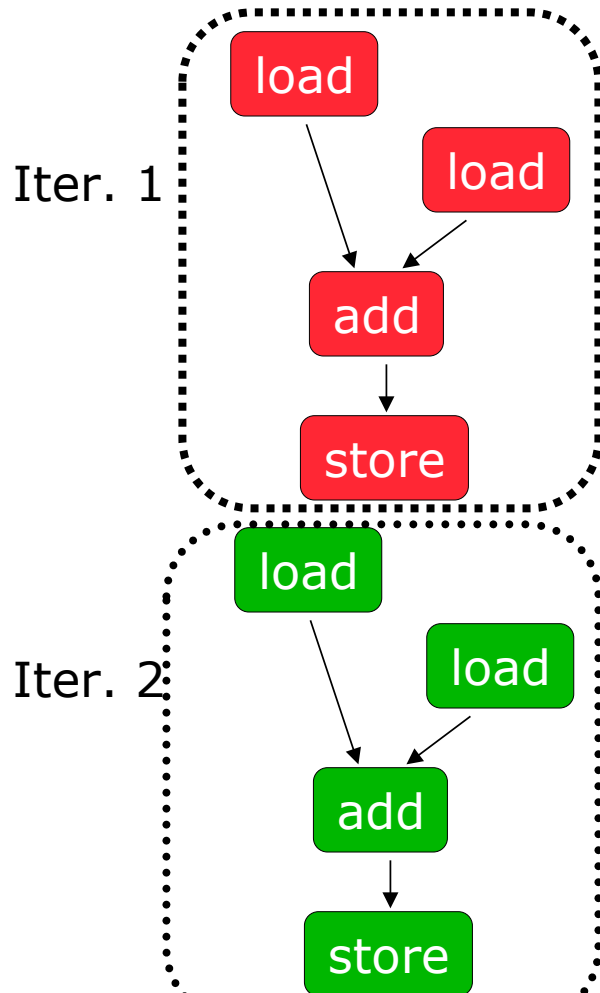
- Quiz 4, Tue Apr 13



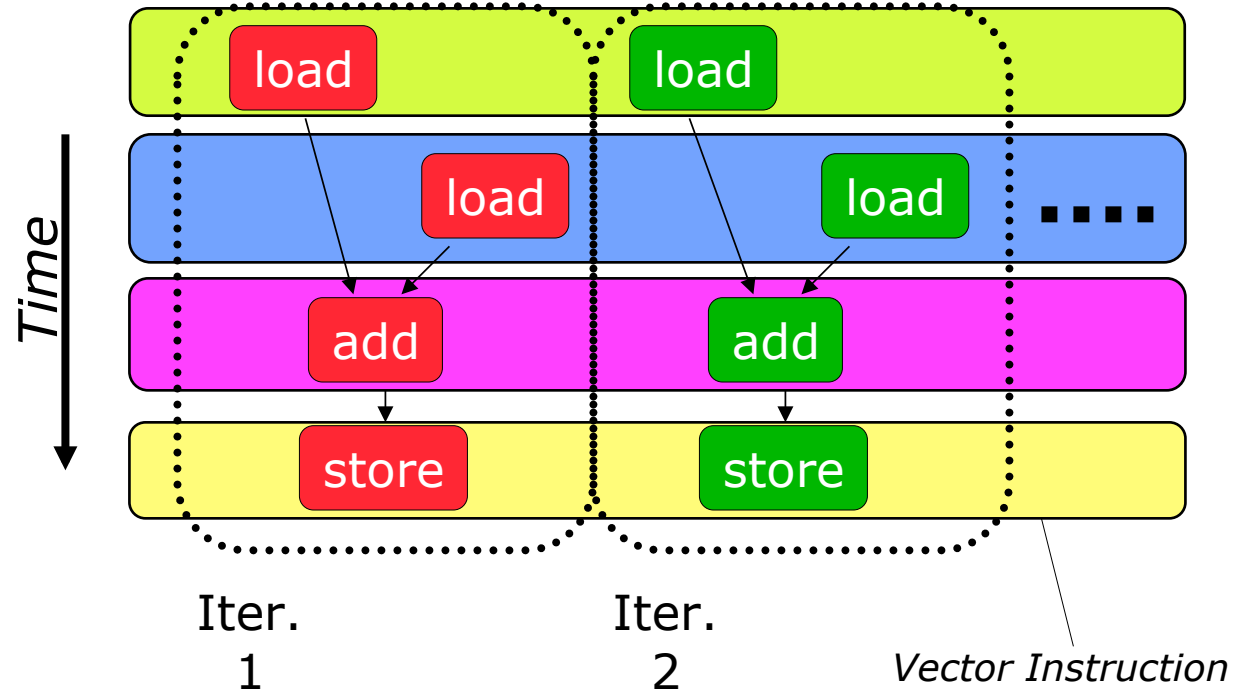
Automatic Code Vectorization

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

Scalar Sequential Code



Vectorized Code



Vectorization is a massive compile-time reordering of operation sequencing requires extensive loop dependence analysis

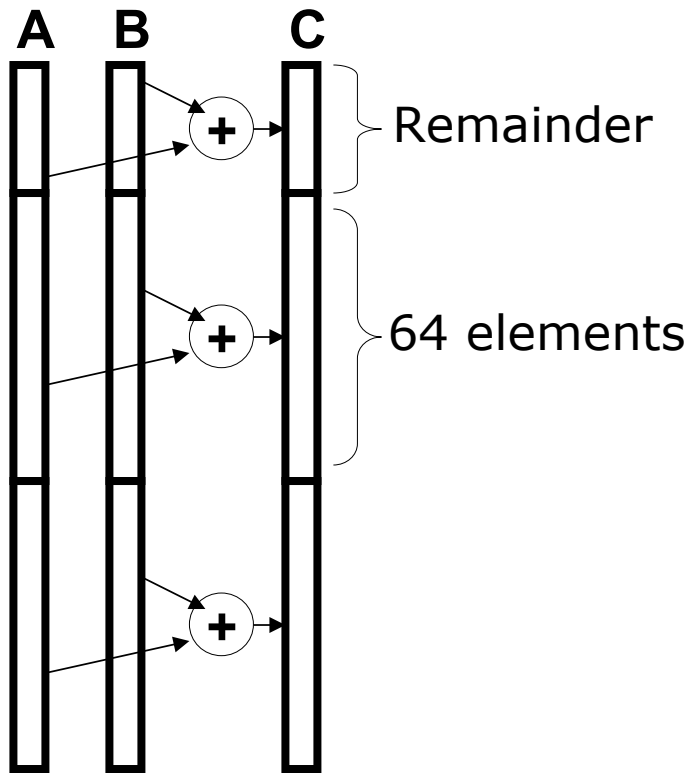


Vector Stripmining

Problem: Vector registers have finite length

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

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



```

ANDI R1, N, 63      # N mod 64
MTC1 VLR, R1       # Do remainder
loop:
LV V1, RA
DSLL R2, R1, 3     # Multiply by 8
DADDU RA, RA, R2   # Bump pointer
LV V2, RB
DADDU RB, RB, R2
ADDV.D V3, V1, V2
SV V3, RC
DADDU RC, RC, R2
DSUBU N, N, R1    # Subtract elements
LI R1, 64
MTC1 VLR, R1     # Reset full length
BGTZ N, loop     # Any more to do?

```



Vector Conditional Execution

Problem: Want to vectorize loops with conditional code:

```
for (i=0; i<N; i++)
    if (A[i]>0) then
        A[i] = B[i];
```

Solution: Add vector *mask* (or *flag*) registers

- vector version of predicate registers, 1 bit per element

...and *maskable* vector instructions

- vector operation becomes NOP at elements where mask bit is clear

Code example:

```
CVM                # Turn on all elements
LV vA, rA          # Load entire A vector
SGTVS.D vA, F0    # Set bits in mask register where A>0
LV vA, rB          # Load B vector into A under mask
SV vA, rA          # Store A back to memory under mask
```

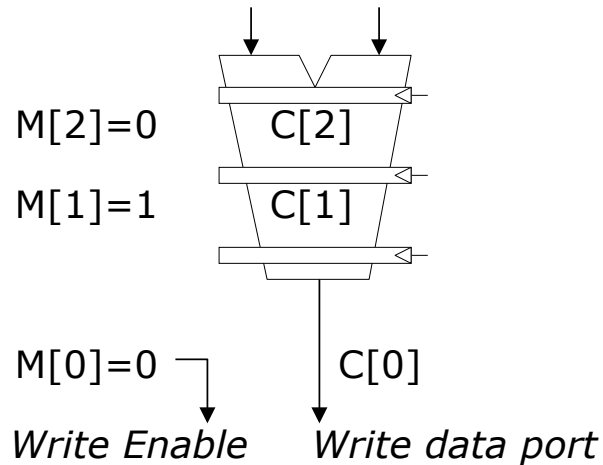


Masked Vector Instructions

Simple Implementation

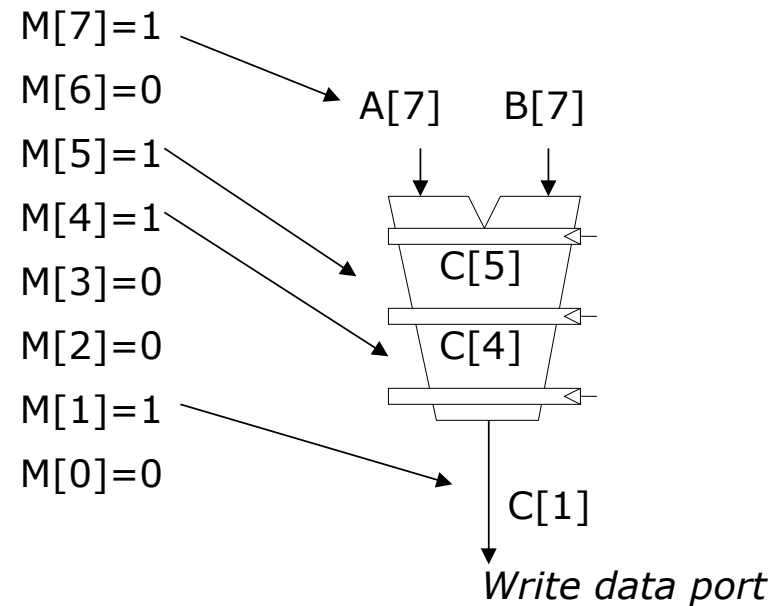
- execute all N operations, turn off result writeback according to mask

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]



Density-Time Implementation

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





Vector Reductions

Problem: Loop-carried dependence on reduction variables

```
sum = 0;
for (i=0; i<N; i++)
    sum += A[i]; # Loop-carried dependence on sum
```

Solution: Re-associate operations if possible, use binary tree to perform reduction

```
# Rearrange as:
sum[0:VL-1] = 0 # Vector of VL partial sums
for(i=0; i<N; i+=VL) # Stripmine VL-sized chunks
    sum[0:VL-1] += A[i:i+VL-1]; # Vector sum
# Now have VL partial sums in one vector register
do {
    VL = VL/2; # Halve vector length
    sum[0:VL-1] += sum[VL:2*VL-1] # Halve no. of partials
} while (VL>1)
```



Vector Scatter/Gather

Want to vectorize loops with indirect accesses:

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

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
```




Vector Scatter/Gather

Scatter example:

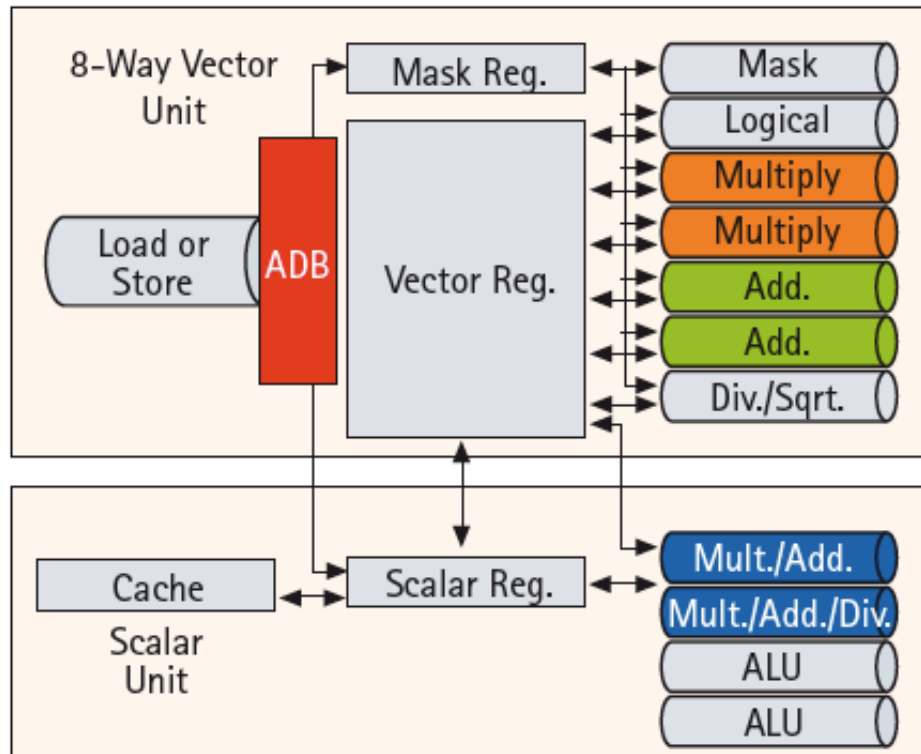
```
for (i=0; i<N; i++)  
    A[B[i]]++;
```

Is following a correct translation?

```
LV vB, rB          # Load indices in B vector  
LVI vA, rA, vB     # Gather initial A values  
ADDV vA, vA, 1     # Increment  
SVI vA, rA, vB     # Scatter incremented values
```



A Modern Vector Super: NEC SX-9 (2008)



- 65nm CMOS technology
- Vector unit (3.2 GHz)
 - 8 foreground VRegs + 64 background VRegs (256x64-bit elements/VReg)
 - 64-bit functional units: 2 multiply, 2 add, 1 divide/sqrt, 1 logical, 1 mask unit
 - 8 lanes (32+ FLOPS/cycle, 100+ GFLOPS peak per CPU)
 - 1 load or store unit (8 x 8-byte accesses/cycle)
- Scalar unit (1.6 GHz)
 - 4-way superscalar with out-of-order and speculative execution
 - 64KB I-cache and 64KB data cache

- Memory system provides 256GB/s DRAM bandwidth per CPU
- Up to 16 CPUs and up to 1TB DRAM form shared-memory *node*
 - total of 4TB/s bandwidth to shared DRAM memory
- Up to 512 nodes connected via 128GB/s network links (message passing between nodes)

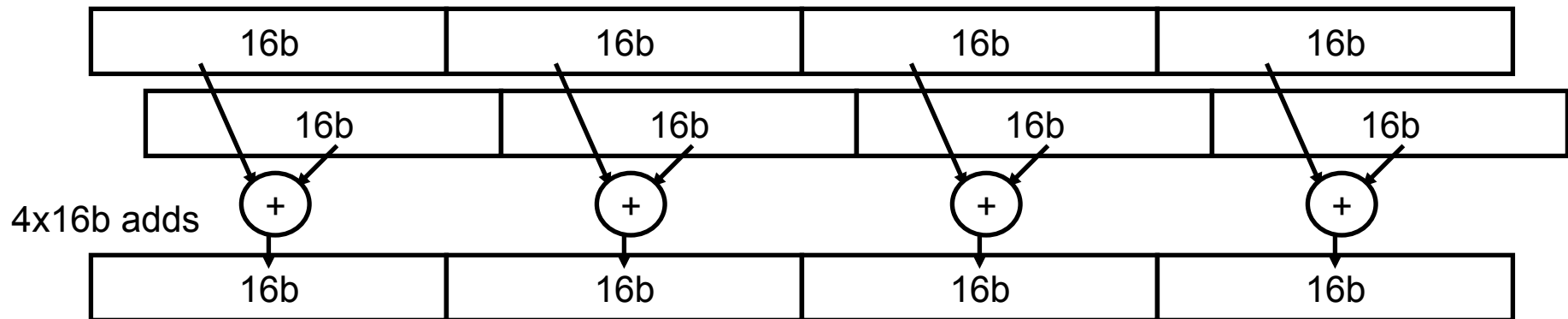
(See also Cray X1E in Appendix F)



Multimedia Extensions (aka SIMD extensions)



- Very short vectors added to existing ISAs for microprocessors
- Use existing 64-bit registers split into 2x32b or 4x16b or 8x8b
 - This concept first used on Lincoln Labs TX-2 computer in 1957, with 36b datapath split into 2x18b or 4x9b
 - Newer designs have 128-bit registers (PowerPC AltiVec, Intel SSE2/3/4)
- Single instruction operates on all elements within register





Multimedia Extensions versus Vectors

- Limited instruction set:
 - no vector length control
 - no strided load/store or scatter/gather
 - unit-stride loads must be aligned to 64/128-bit boundary
- Limited vector register length:
 - requires superscalar dispatch to keep multiply/add/load units busy
 - loop unrolling to hide latencies increases register pressure
- Trend towards fuller vector support in microprocessors
 - Better support for misaligned memory accesses
 - Support of double-precision (64-bit floating-point)
 - New Intel AVX spec (announced April 2008), 256b vector registers (expandable up to 1024b)



Graphics Processing Units (GPUs)

- Original GPUs were dedicated fixed-function devices for generating 3D graphics
- More recently, GPUs have been made more programmable, so called “General-Purpose” GPUs or GP-GPUs.
- Base building block of modern GP-GPU is very similar to a vector machine
 - e.g., NVIDIA G80 series core (NVIDIA term is Streaming Multiprocessor, SM) has 8 “lanes” (NVIDIA term is Streaming Processor, SP). Vector length is 32 elements (NVIDIA calls this a “warp”).
- Currently machines are built with separate chips for CPU and GP-GPU, but future designs will merge onto one chip
 - Already happening for smartphones and tablet designs



Acknowledgements

- These slides contain material developed and copyright by:
 - Arvind (MIT)
 - Krste Asanovic (MIT/UCB)
 - Joel Emer (Intel/MIT)
 - James Hoe (CMU)
 - John Kubiatowicz (UCB)
 - David Patterson (UCB)
- MIT material derived from course 6.823
- UCB material derived from course CS252