

# CS184c: Computer Architecture [Parallel and Multithreaded]

Day 11: May10, 2001  
Data Parallel  
(SIMD, SPMD, Vector)



CALTECH cs184c Spring2001 -- DeHon

## Today

- Data Parallel
  - Model
  - Application
  - Resources
  - Architectures
    - Abacus
    - T0

CALTECH cs184c Spring2001 -- DeHon

## Data Parallel Model

- Perform same computation on multiple, distinct data items
- SIMD
  - recall simplification of general array model
  - every PE get same instruction
    - feed large number of PEs with small instruction bandwidth

CALTECH cs184c Spring2001 -- DeHon

CS184a

## Architecture Instruction Taxonomy

Control Threads (PCs)					
<i>pinsts</i> per Control Thread					
Instruction Depth					
Granularity					
Architecture/Examples					
0	0	0	n/a	Hardwired Functional Unit (e.g. ECC/EDC Unit, FP MPY)	
	n	1	w	FPGA	
		$n_v \cdot 1$	$w$	Reconfigurable ALUs	
1	1	c	$w$	Bitwise SIMD	
		$n_v \cdot w$	$w$	Traditional Processors	
		$n_v \cdot w$	$w$	Vector Processors	
	n	c	1	DPGA	
m	n	8	16	PADDI	
		c	w	VLIW	
	1	n	1	1	HSRA/SCORE
		c	$n_v \cdot w$	$w$	MSIMD
		c	1	1	VEGA
m	1	8	16	PADDI-2	
		c	w	MIMD (traditional)	

CALTECH cs184c Spring2001 -- DeHon

## Example

- Operations on vectors
  - vector sum
  - dot, cross product
  - matrix operations
- Simulations / finite element...
  - same update computation on every site
- Image/pixel processing
  - compute same thing on each pixel

CALTECH cs184c Spring2001 -- DeHon

## Model

- Zero, one, infinity
  - good model has unbounded number of processors
  - user allocates virtual processors
  - folded (as needed) to share physical processor se

CALTECH cs184c Spring2001 -- DeHon

## How do an if?

- Have large set of data
- How do we conditionally deal with data?

CALTECH cs184c Spring2001 -- DeHon

## ABS example

- Saw hoops had to jump through to compute absolute value w/out conditional branching

CALTECH cs184c Spring2001 -- DeHon

## Key: Local State

- Set state during computation
- Use state to modify transmitted instruction
  - Could simply be  $PE.op(inputs, state)$
  - Often mask
    - select subset of processors to operate
    - like predicated operations in conventional processor

CALTECH cs184c Spring2001 -- DeHon

## Local State Op

- Consider 4-LUT with two states
  - w/ local state bit, can implement a 3-LUT function with one state bit
  - state bit is 4th input to LUT can decide which operation to perform

CALTECH cs184c Spring2001 -- DeHon

## ABS with Mask

- $\text{Tmp} = \text{val} < 0$
- $\text{rval} = \text{val}$
- $\text{mask tmp} == \text{true}$   
     $\text{rval} = -(\text{val})$
- $\text{unmask}$

CALTECH cs184c Spring2001 -- DeHon

## Model

- Model remains
  - all PEs get same operation
  - compute on local state with operation

CALTECH cs184c Spring2001 -- DeHon

## Synchronization

- Strong SIMD model
  - all operations move forward in lock-step
  - don't get asynchronous advance
  - don't have to do explicit synchronization

CALTECH cs184c Spring2001 -- DeHon

## Communications

- Question about how general
- Common, low-level
  - nearest-neighbor
  - cheap, fast
  - depends on layout...
  - effect on virtual processors and placement?

CALTECH cs184c Spring2001 -- DeHon

## Communications

- General network
  - allow model with more powerful shuffling
  - how rich? (expensive)
  - wait for longest operation to complete?
- Use Memory System?

CALTECH cs184c Spring2001 -- DeHon

## Memory Model?

- PEs have local memory
- Allow PEs global pointers?
- Allow PEs to dereference arbitrary addresses?
  - General communications
  - Including conflicts on PE/bank
    - potentially bigger performance impact in lock-step operation
- Data placement important

CALTECH cs184c Spring2001 -- DeHon



## Vector Model

- Primary data structure
- Memory access very predictable
  - easy to get high performance on
    - e.g. burst memory fetch, banking
  - one address and get stream of data

CALTECH cs184c Spring2001 -- DeHon

## How effect control flow?

- Predicated operations take care of local flow control variations
- Sometimes need to effect entire control stream
- E.g. relaxation convergence
  - compute updates to refine some computation
  - until achieve tolerance

CALTECH cs184c Spring2001 -- DeHon

## Flow Control

- Ultimately need one bit (some digested value) back at central controller to branch upon
- How get?
  - Pick some value calculated in memory?
  - Produce single, aggregate result

CALTECH cs184c Spring2001 -- DeHon

## Reduction Value

- Example: summing-or
  - Or together some bit from all Pes
    - build reduction tree....log depth
  - typical usage
    - processor asserts bit when find solution
    - processor deassert bit when solution quality good enough
      - detect when all processors done

CALTECH cs184c Spring2001 -- DeHon

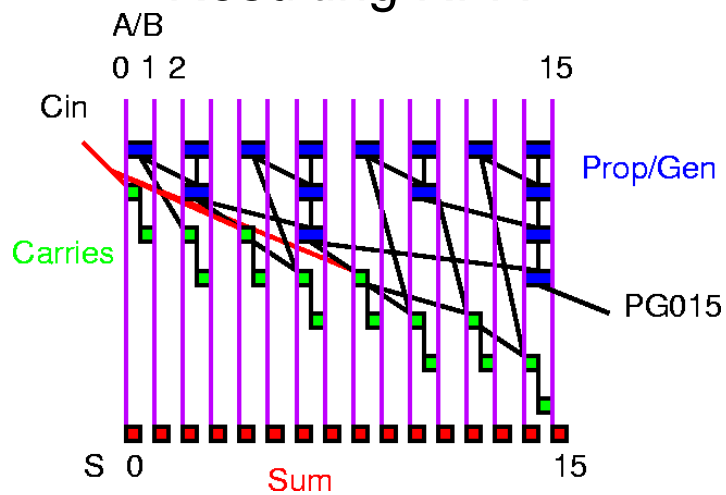
## Key Algorithm: Parallel Prefix

- Often will want to calculate some final value on aggregate
  - dot product: sum of all pairwise products
  - Karl showed us: saturating sums
    - for example in ADPCM compression
  - Already saw in producing log-depth carries

CALTECH cs184c Spring2001 -- DeHon

CS184a

### Resulting RPA



CALTECH cs184c Spring2001 -- DeHon

## Parallel Prefix

- Calculate **all** intermediate results in log depth
  - e.g. all intermediate carries
  - e.g. all sums to given point in vector
- More general than tree reduction
  - tree reduction (sum, or, and) uses commutativity
  - parallel prefix only requires associativity

CALTECH cs184c Spring2001 -- DeHon

## Parallel Prefix...

- Count instances with some property
- Parsing
- List operations
  - pointer jumping, find length, matching

CALTECH cs184c Spring2001 -- DeHon

# Resources

CALTECH cs184c Spring2001 -- DeHon

## Contrast VLIW/SS

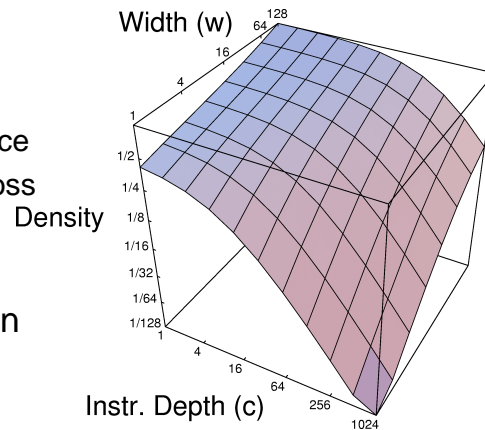
- Single instruction shared across several ALUs
  - (across more bits)
- Significantly lower control
- Simple/predictable control flow
- Parallelism (of data) in model

CALTECH cs184c Spring2001 -- DeHon

CS184a

## Peak Densities from Model

- Only 2 of 4 parameters
  - small slice of space
  - 100× density across
- Large difference in peak densities
  - large design space!



CALTECH cs184c Spring2001 -- DeHon

CS184a

## Calibrate Model

FPGA	model $w = 1, d = c = 1, k = 4$	<b>880K<math>\lambda^2</math></b>
	Xilinx 4K	<b>630K<math>\lambda^2</math></b>
	Altera 8K	<b>930K<math>\lambda^2</math></b>
SIMD	model $w = 1000, c = 0, d = 64, k = 3$	<b>170K<math>\lambda^2</math></b>
	Abacus	<b>190K<math>\lambda^2</math></b>
Processor model	model $w = 32, d = 32, c = 1024, k = 2$	<b>2.6M<math>\lambda^2</math></b>
	MIPS-X	<b>2.1M<math>\lambda^2</math></b>

CALTECH cs184c Spring2001 -- DeHon

# Examples

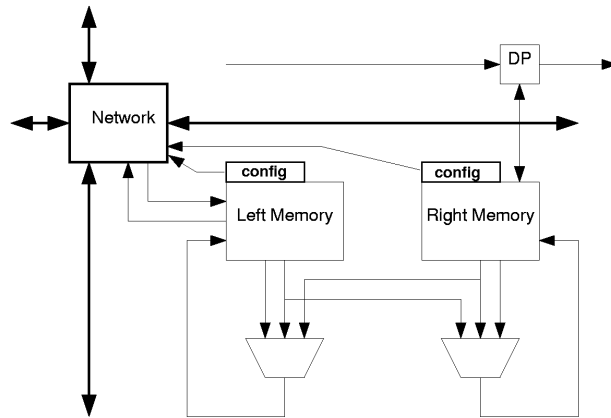
CALTECH cs184c Spring2001 -- DeHon

## Abacus: bit-wise SIMD

- Collection of simple, bit-processing units
- PE:
  - 2x3-LUT (think adder bit)
  - 64 memory bits, 8 control config
  - active (mask) register
- Network: nearest neighbor with bypass
- Configurable word-size
- [Bolotski et. al. ARVLSI'95]

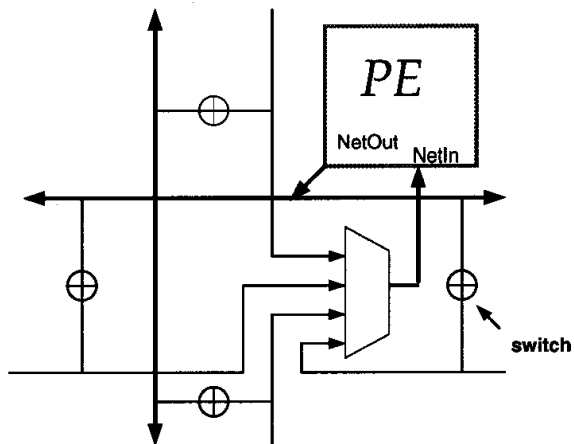
CALTECH cs184c Spring2001 -- DeHon

## Abacus: PE



CALTECH cs184c Spring2001 -- DeHon

## Abacus: Network



CALTECH cs184c Spring2001 -- DeHon





## Abacus: bit-wise SIMD

- High raw density:
  - 660 ALU Bit Ops/ $\lambda^2$ -s
- Do have to synthesize many things out of several operations
- Nearest neighbor only

CALTECH cs184c Spring2001 -- DeHon

## Abacus: Cycles

Operation	8-bit		16-bit		32-bit	
	Cycles	GOPS	Cycles	GOPS	Cycles	GOPS
Add	4	4.0	4	2.0	5	0.7
Shift	2	8.0	2	4.0	2	2.0
Accumulate	3	5.2	3	2.6	3	1.3
Move	3	5.2	4	2.0	6	0.6
Compare	6	2.6	11	0.6	12	0.2
Multiply ( $16 \times 16$ )					180	0.03

Algorithm	Cycles	Time ( $\mu$ sec)
Edge Detection $\sigma = 1.6$	380	3
Optical Flow, $\Delta = 2$ , $5 \times 5$ region	3000	24
Surface Reconstruction (1 iteration)	370	3

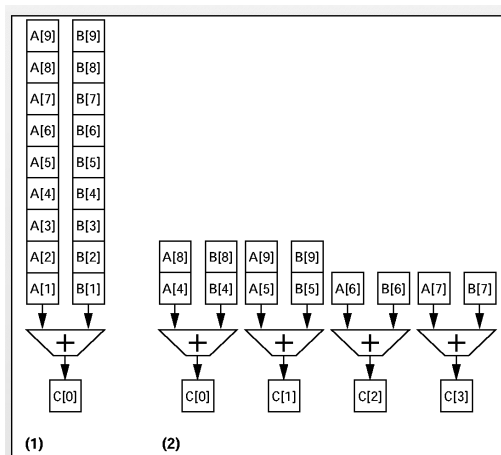
CALTECH cs184c Spring2001 -- DeHon

## T0: Vector Microprocessor

- Word-oriented vector pipeline
- Scalable vector abstraction
  - vector ISA
  - size of physical vector hardware abstracted
- Communication mostly through memory
- [Asanovic et. al., IEEE Computer 1996]
- [Asanovic et. al., Hot Chips 1996]

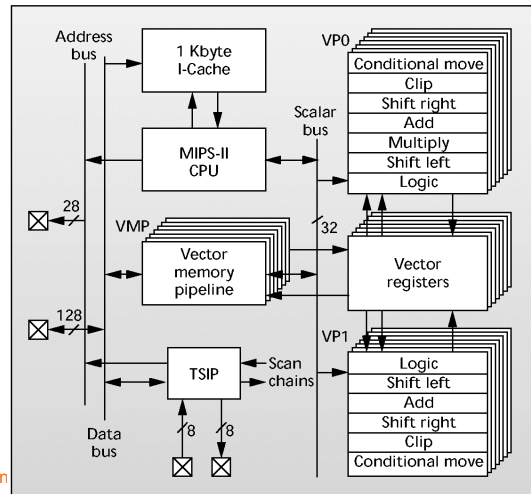
CALTECH cs184c Spring2001 -- DeHon

## Vector Scaling



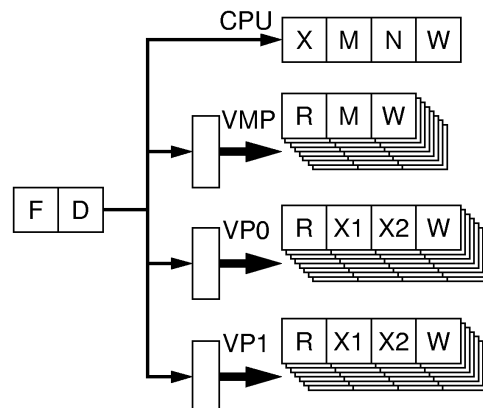
CALTECH cs184c Spring2001 -- DeHon

# T0 Microarchitecture



CALTECH cs184c Spring

# T0 Pipeline



CALTECH cs184c Spring2001 -- DeHon

## T0 ASM example

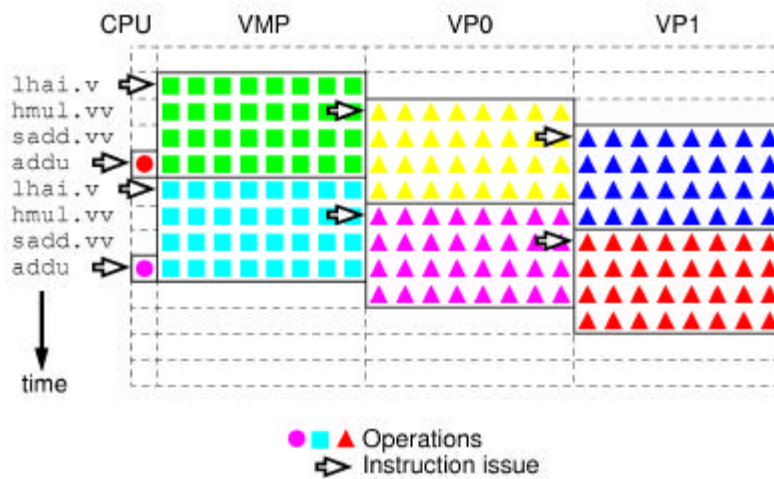
```

lhai.v vv1, t0, t1      # Vector load.
hmul.vv vv4, vv2, vv3   # Vector mul.
sadd.vv vv7, vv5, vv7   # Vector add.
addu t2, -1             # Scalar add.
lhai.v vv2, t0, t1      # Vector load.
hmul.vv vv5, vv1, vv3   # Vector mul.
sadd.vv vv8, vv4, vv8   # Vector add.
addu t7, t4             # Scalar add.

```

CALTECH cs184c Spring2001 -- DeHon

## T0 Execution Example



CALTECH

## T0: Vector Microprocessor

- Higher raw density than (super)scalar microprocessors
  - 22 ALU Bit Ops/ $\lambda^2$ -s (vs. <10)
- Clean ISA, scaling
  - contrast VIS, MMX
- Easy integration with existing  $\mu$ P/tools
  - assembly library for vector/matrix ops
  - leverage work in vectorizing compilers

CALTECH cs184c Spring2001 -- DeHon

## Big Ideas

- Model for computation
  - enables programmer think about machine capabilities a high level
  - abstract out implementation details
  - allow scaling/different implementations
- Exploit structure in computation
  - use to reduce hardware costs

CALTECH cs184c Spring2001 -- DeHon