

# CS184b: Computer Architecture (Abstractions and Optimizations)

Day 24: May 25, 2005  
Heterogeneous Computation  
Interfacing



Caltech CS184 Spring2005 -- DeHon

## Previously

- Homogenous model of computational array
  - single word granularity, depth, interconnect
  - all post-fabrication programmable
- Understand tradeoffs of each

Caltech CS184 Spring2005 -- DeHon

2

## Today

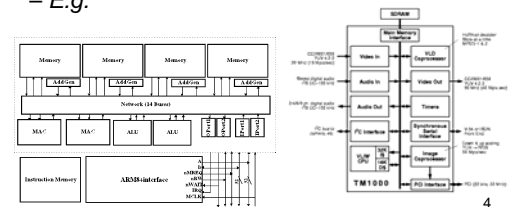
- Heterogeneous architectures
  - Why?
- Focus in on Processor + Array hybrids
  - Motivation
  - Single Threaded Compute Models
  - Architecture
  - Examples

Caltech CS184 Spring2005 -- DeHon

3

## Why?

- Why would we be interested in heterogeneous architecture?
  - *E.g.*

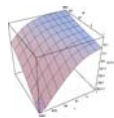


Caltech CS184 Spring2005 -- DeHon

4

## Why?

- Applications have a mix of characteristics
- Already accepted
  - seldom can afford to build most general (unstructured) array
    - bit-level, deep context,  $p=1$
  - => are picking some structure to exploit
- May be beneficial to have portions of computations optimized for different structure conditions.

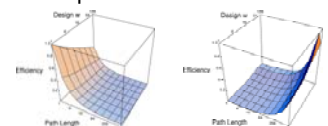


Caltech CS184 Spring2005 -- DeHon

5

## Examples

- Processor+FPGA
- Processors or FPGA add
  - multiplier or MAC unit
  - FPU
  - Motion Estimation coprocessor



Caltech CS184 Spring2005 -- DeHon

## Optimization Prospect

- Cost capacity for composite than either pure
  - $(A_1+A_2)T_{12} < A_1T_1$
  - $(A_1+A_2)T_{12} < A_2T_2$

## Optimization Prospect Example

- Floating Point
  - Task: I integer Ops + F FP-ADDs
  - $A_{\text{proc}}=125M\lambda^2$
  - $A_{\text{FPU}}=40M\lambda^2$
  - I cycles / FP Ops = 60
  - $125(I+60F) \square 165(I+F)$ 
    - $(7500-165)/40 = I/F$
    - $183 \approx I/F$

## Motivational: Other Viewpoints

- Replace interface glue logic
- IO pre/post processing
- Handle real-time responsiveness
- Provide powerful, application-specific operations
  - possible because of previous observation

## Wide Interest

- PRISM (Brown)
- PRISC (Harvard)
- DPGA-coupled uP (MIT)
- GARP, Pleiades, ... (UCB)
- OneChip (Toronto)
- REMARC (Stanford)
- NAPA (NSC)
- E5 etc. (Triscend)
- Chameleon
- Quicksilver
- Excalibur (Altera)
- Virtex+PowerPC (Xilinx)
- Stretch

## Pragmatics

- Tight coupling important
  - numerous (anecdotal) results
    - we got 10x speedup...but were bus limited
      - would have gotten 100x if removed bus bottleneck
- **Speed Up =  $T_{\text{seq}}/(T_{\text{accel}} + T_{\text{data}})$** 
  - e.g.  $T_{\text{accel}} = 0.01 T_{\text{seq}}$
  - $T_{\text{data}} = 0.10 T_{\text{seq}}$

## Key Questions

- How do we co-architect these devices?
- What is the compute model for the hybrid device?

## Compute Models

- Unaffected by array logic (interfacing)
- Dedicated IO Processor
  - Specialized multithreaded
- Instruction Augmentation
  - Special Instructions / Coprocessor Ops
  - VLIW/microcoded extension to processor
  - Configurable Vector unit
- Memory ⇒ memory coprocessor

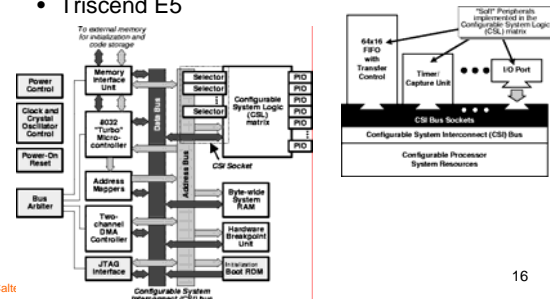
## Interfacing

## Model: Interfacing

- Logic used in place of
  - ASIC environment customization
  - external FPGA/PLD devices
- Example
  - bus protocols
  - peripherals
  - sensors, actuators
- Case for:
  - Always have some system adaptation to do
  - Modern chips have capacity to hold processor + glue logic
  - reduce part count
  - Glue logic vary
  - value added must now be accommodated on chip (formerly board level)

## Example: Interface/Peripherals

- Triscend E5



## IO Processor

## Model: IO Processor

- Array dedicated to servicing IO channel
  - sensor, lan, wan, peripheral
- Provides
  - protocol handling
  - stream computation
    - compression, encrypt
- Looks like IO peripheral to processor
- Maybe processor can map in
  - as needed
  - physical space permitting
- Case for:
  - many protocols, services
  - only need few at a time
  - dedicate attention, offload processor

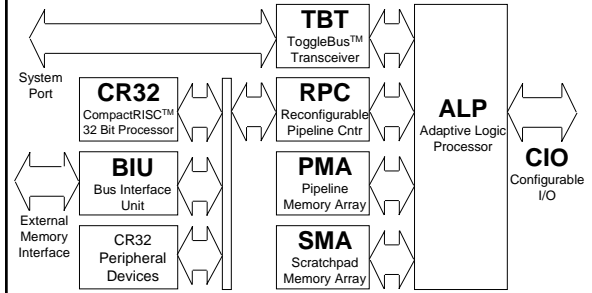
## IO Processing

- Single threaded processor
  - cannot continuously monitor multiple data pipes (src, sink)
  - need some minimal, local control to handle events
  - for performance or real-time guarantees , may need to service event rapidly
  - *E.g.* checksum (decode) and acknowledge packet

Caltech CS184 Spring2005 -- DeHon

19

## NAPA 1000 Block Diagram

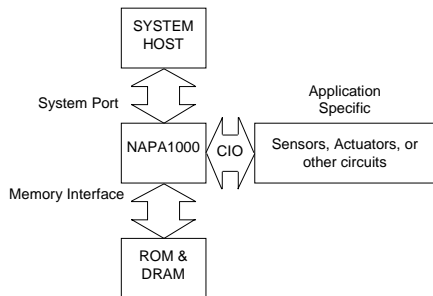


Source: National Semiconductor

Caltech CS184 Spring2005 -- DeHon

20

## NAPA 1000 as IO Processor



Source: National Semiconductor

Caltech CS184 Spring2005 -- DeHon

21

## Instruction Augmentation

Caltech CS184 Spring2005 -- DeHon

22

## Model: Instruction Augmentation

- Observation: Instruction Bandwidth
  - Processor can only describe a small number of basic computations in a cycle
    - 1 bits  $\rightarrow 2^l$  operations
  - This is a small fraction of the operations one could do even in terms of  $w \otimes w \rightarrow w$  Ops
    - $w2^{2(2^w)}$  operations

Caltech CS184 Spring2005 -- DeHon

23

## Model: Instruction Augmentation (cont.)

- Observation: Instruction Bandwidth
  - Processor could have to issue  $w2^{2(2^w - 1)}$  operations just to describe some computations
  - An *a priori* selected base set of functions could be very bad for some applications

Caltech CS184 Spring2005 -- DeHon

24

## Instruction Augmentation

- Idea:
  - provide a way to augment the processor's instruction set
  - with operations needed by a particular application
  - close semantic gap / avoid mismatch

## Instruction Augmentation

- What's required:
  - some way to fit augmented instructions into instruction stream
  - execution engine for augmented instructions
    - if programmable, has own instructions
  - interconnect to augmented instructions

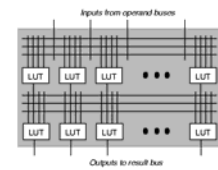
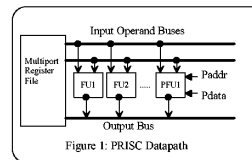
## PRISC

- How integrate into processor ISA?

[Razdan+Smith: Harvard]

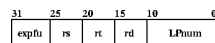
## PRISC

- Architecture:
  - couple into register file as "superscalar" functional unit
  - flow-through array (no state)



## PRISC

- ISA Integration
  - add expfu instruction
  - 11 bit address space for user defined expfu instructions
  - fault on pfu instruction mismatch
    - trap code to service instruction miss
  - all operations occur in clock cycle
  - easily works with processor context switch
    - no state + fault on mismatch pfu instr



## PRISC Results

- All compiled
- working from MIPS binary
- <200 4LUTs ?
  - 64x3
- 200MHz MIPS base

Optimization	CPS	EQN	EXP	GCC	LI	SC
PFU-expression	9	0	48	13	4	12
PFU-table-lookup	0	0	0	0	0	0
PFU-predication	0	1	0	13	0	0
PFU-jump	10	0	47	103	0	35
PFU-loop	0	3	0	4	0	0
TOTAL	19	4	95	133	4	47

Table 1: Static PFU optimization instances in SPECint92.

	CPS	EQN	EXP	GCC	LI	SC
Speedup	1.15	1.91	1.16	1.10	1.06	1.12

Table 2: Cycle count speedup for a PRISC-1 microarchitecture with a single PFU resource. The speedup for each application is an arithmetic average (as defined by SPEC) of all of the data sets for that application.

## Chimaera

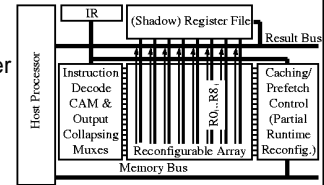
- Start from PRISC idea
  - integrate as functional unit
  - no state
  - RFUOPs (like expfu)
  - stall processor on instruction miss, reload
- Add
  - manage multiple instructions loaded
  - more than 2 inputs possible

[Hauck: Northwestern]  
Caltech CS184 Spring2005 -- DeHon

31

## Chimaera Architecture

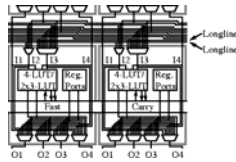
- “Live” copy of register file values feed into array
- Each row of array may compute from register values or intermediates (other rows)
- Tag on array to indicate RFUOP



Caltech CS184 Spring2005 -- DeHon

## Chimaera Architecture

- Array can compute on values as soon as placed in register file
- Logic is combinational
- When RFUOP matches
  - stall until result ready
    - critical path
      - only from late inputs
  - drive result from matching row

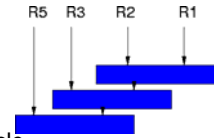


Caltech CS184 Spring2005 -- DeHon

33

## Chimaera Timing

- If presented
  - R1, R2
  - R3
  - R5
  - can complete in one cycle
- If R1 presented last
  - will take more than one cycle for operation



Caltech CS184 Spring2005 -- DeHon

34

## Chimaera Results

- Speedup
- Compress 1.11
  - Eqntott 1.8
  - Life 2.06 (160 hand parallelization)

[Hauck/FCCM97]  
Caltech CS184 Spring2005 -- DeHon

35

## Instruction Augmentation

- Small arrays with limited state
  - so far, for automatic compilation
    - reported speedups have been small
  - open
    - discover less-local recodings which extract greater benefit

Caltech CS184 Spring2005 -- DeHon

36

## GARP Motivation

- Single-cycle flow-through
  - not most promising usage style
- Moving data through RF to/from array
  - can present a limitation
    - bottleneck to achieving high computation rate

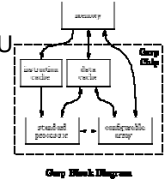
[Hauser+Wawrzynek: UCB]

Caltech CS184 Spring2005 -- DeHon

37

## GARP

- Integrate as coprocessor
  - similar bwtd to processor as FU
  - own access to memory
- Support multi-cycle operation
  - allow state
  - cycle counter to track operation
- Fast operation selection
  - cache for configurations
  - dense encodings, wide path to memory



Caltech CS184 Spring2005 -- DeHon

38

## GARP

- ISA -- coprocessor operations
  - issue `gaconfig` to make a particular configuration resident (may be active or cached)
  - explicitly move data to/from array
    - 2 writes, 1 read (like FU, but not 2W+1R)
  - processor suspend during coproc operation
    - cycle count tracks operation
  - array may directly access memory
    - processor and array share memory space
      - cache/mmu keeps consistent between
    - can exploit streaming data operations

Caltech CS184 Spring2005 -- DeHon

39

## GARP

- Processor Instructions

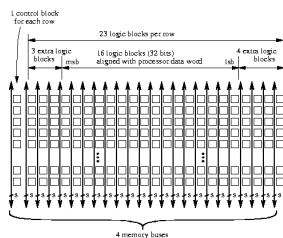
Instruction	Interlock?	Description
<code>gaconf reg</code>	yes	Load (or switch to) configuration at address given by <code>reg</code> .
<code>mtga reg, array-row-reg, count</code>	yes	Copy <code>reg</code> value to <code>array-row-reg</code> and set array clock counter to <code>count</code> .
<code>mfga reg, array-row-reg, count</code>	yes	Copy <code>array-row-reg</code> value to <code>reg</code> and set array clock counter to <code>count</code> .
<code>gabump reg</code>	no	Increase array clock counter by value in <code>reg</code> .
<code>gastop reg</code>	no	Copy array clock counter to <code>reg</code> and stop array by zeroing clock counter.
<code>gacinv reg</code>	no	Invalidate cache copy of configuration at address given by <code>reg</code> .
<code>cfga reg, array-control-reg</code>	no	Copy value of array control register <code>array-control-reg</code> to <code>reg</code> .
<code>gasave reg</code>	yes	Save all array data state to memory at address given by <code>reg</code> .
<code>garestore reg</code>	yes	Restore previously saved data state from memory at address given by <code>reg</code> .

Caltech CS184 Spring2005 -- DeHon

40

## GARP Array

- Row oriented logic
  - denser for datapath operations
- Dedicated path for
  - processor/memory data
- Processor not have to be involved in array↔memory path



Caltech CS184 Spring2005 -- DeHon

41

## GARP Hand Results

Table 2. Garp's speedups over Ultrasparc for hand-coded functions.

Function	Data size	Speedup	Limiting factor
Image median filter	640 × 480 pixels	43	Compute throughput
DES (ECB mode)	1 Mbyte	18.7	Compute throughput
Image dithering	640 × 480 pixels	17.0	Compute throughput
<code>strlen</code>	1,024 chars	14.2	Memory bandwidth
<code>strlen</code>	16 chars	1.84	Overhead
Sort	2 Mbytes	2.2	Scattered memory accesses

[Callahan, Hauser, Wawrzynek. IEEE Computer, April 2000]

Caltech CS184 Spring2005 -- DeHon

42

## GARP Compiler Results

Table 1. Kernels from a wavelet image compression program.

Kernel	Percentage of original software execution time	Iteration Interval	No. of queues used	ILP (average operations per nonstall cycle)	No. of executions	Average no. of compute cycles per execution (including stalls)	Average no. of overhead cycles per execution	Net speedup over MIPS only
forward_wavelet_696	18.2	2	2	10.0	448	1,176	114	2.1
forward_wavelet_647	13.8	2	2	10.0	448	310	91	3.1
int_image_354	12.8	1	2	8.0	1	65,852	564	12.7
forward_wavelet_711	10.1	2	2	7.0	448	241	59	4.9
entropy_encode_544	10.0	1	1	5.0	1	65,538	989	9.9
forward_wavelet_674	9.3	1	3	13.0	448	128	76	6.6
block_quantize_411	5.5	2	0	5.5	320	353	56	2.8
entropy_encode_557	3.9	6	0	2.8	3,262	31	24	1.4
RLE_encode_509	3.8	1	1	11.0	774	22	48	4.6

[Callahan, Hauser, Wawrzynek. IEEE Computer, April 2000]

Caltech CS184 Spring2005 -- DeHon

43

## Common Theme

- To get around instruction expression limits
  - define new instruction in array
    - many bits of config ... broad expressability
    - many parallel operators
  - give array configuration short “name” which processor can callout
    - ...effectively the address of the operation

Caltech CS184 Spring2005 -- DeHon

44

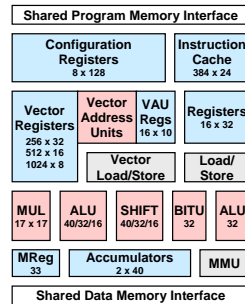
## VLIW/microcoded Model

- Similar to instruction augmentation
- Single tag (address, instruction)
  - controls a number of more basic operations
    - E.g. Silicon Spice SpiceEngine
- Some difference in expectation
  - can sequence a number of different tags/operations together

Caltech CS184 Spring2005 -- DeHon

45

## SpiceEngine™ Architecture



- **1 KB Vector Registers**
  - Fast structured data access
  - Reduce memory demand
  - Preload vector registers to hide memory latency
- **10 Parallel ops/cycle**
  - Fixed-point DSP arithmetic
  - Configurable 16/32/40 bits
  - ITU saturation arithmetic
- **Wide configuration issue**
  - Compiler generates, preloads
  - Dynamic reconfiguration

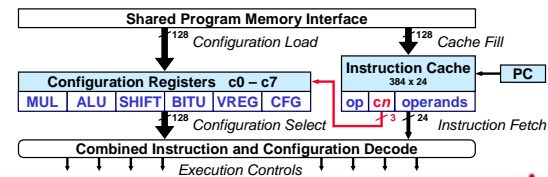
[Nicholls/MPF 2001]

46 MPF 2001

BROADCOM

## Instruction / Configuration Control

- Fetch 24-bit instructions from cached shared memory
  - Instruction field *cn* selects a Configuration Register
- Issue combined instruction and 128-bit configuration *cn*
  - Compiler/program preloads configuration registers from memory
- Dynamically reconfigure and control execution units



47 MPF 2001

[Nicholls/MPF 2001]

BROADCOM

## Vector and Shared Memory

Caltech CS184 Spring2005 -- DeHon

48



## Configurable Vector Unit Model

- Perform vector operation on datastreams
- Setup spatial datapath to implement operator in configurable hardware
- Potential benefit in ability to chain together operations in datapath
- May be way to use GARP/NAPA?
- OneChip (to come...)

Caltech CS184 Spring2005 -- DeHon

49

## Observation

- All single threaded
  - limited to parallelism
    - instruction level (VLIW, bit-level)
    - data level (vector/stream/SIMD)
  - no task/thread level parallelism
    - except for IO dedicated task parallel with processor task

Caltech CS184 Spring2005 -- DeHon

50

## Scaling

- Can scale
  - number of inactive contexts
  - number of PFUs in PRISC/Chimaera
    - but still limited by single threaded execution (ILP)
    - exacerbate pressure/complexity of RF/interconnect
- Cannot scale
  - number of active resources
  - and have automatically exploited

Caltech CS184 Spring2005 -- DeHon

51

## Processor/FPGA run in Parallel?

- What would it take to let the processor and FPGA run in parallel?
  - And still get reasonable program semantics?

Caltech CS184 Spring2005 -- DeHon

52

## Modern Processors

- Deal with
  - variable delays
  - dependencies
  - multiple (unknown to compiler) func. units
- Via
  - register scoreboarding
  - runtime dataflow (Tomasulo)

Caltech CS184 Spring2005 -- DeHon

53

## Dynamic Issue

- PRISC (Chimaera?)
  - register→register, work with scoreboard
- GARP
  - works with memory system, so register scoreboard not enough

Caltech CS184 Spring2005 -- DeHon

54

## OneChip Memory Interface [1998]

- Want array to have direct memory→memory operations
- Want to fit into programming model/ISA
  - w/out forcing exclusive processor/FPGA operation
  - allowing decoupled processor/array execution

[Jacob+Chow: Toronto]  
Caltech CS184 Spring2005 -- DeHon

55

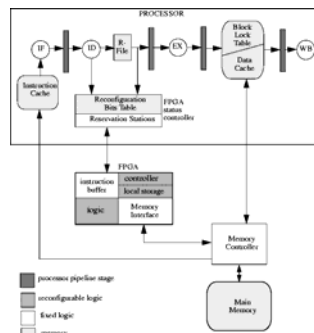
## OneChip

- Key Idea:
  - FPGA operates on memory→memory regions
  - make regions explicit to processor issue
  - scoreboard memory blocks
    - Synchronize exclusive access

Caltech CS184 Spring2005 -- DeHon

56

## OneChip Pipeline

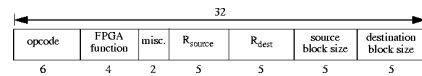


Caltech CS184 Spring2005

57

## OneChip Instructions

- Basic Operation is:
  - FPGA MEM[Rsource]→MEM[Rdst]
  - block sizes powers of 2



- Supports 14 "loaded" functions
  - DPFGA/contexts so 4 can be cached

Caltech CS184 Spring2005 -- DeHon

58

## OneChip

- Basic op is: FPGA MEM→MEM
- no state between these ops
- coherence is that ops appear sequential
- could have multiple/parallel FPGA Compute units
  - scoreboard with processor and each other
- single source operations?
- can't chain FPGA operations?

Caltech CS184 Spring2005 -- DeHon

59

## Model Roundup

- Interfacing
- IO Processor (Asynchronous)
- Instruction Augmentation
  - PFU (like FU, no state)
  - Synchronous Coproc
  - VLIW
  - Configurable Vector
- Memory⇒memory coprocessor

Caltech CS184 Spring2005 -- DeHon

60

## Big Ideas

- Exploit structure
  - area benefit to
  - tasks are heterogeneous
  - mixed device to exploit
- Instruction description
  - potential bottleneck
  - custom “instructions” to exploit

Caltech CS184 Spring2005 -- DeHon

61

## Big Ideas

- Spatial
  - denser raw computation
  - supports definition of powerful instructions
    - assign short name --> descriptive benefit
    - build with spatial --> dense collection of active operators to support
  - efficient way to support
    - repetitive operations
    - bit-level operations

Caltech CS184 Spring2005 -- DeHon

62

## Big Ideas

- Model
  - for heterogeneous composition
  - preserving semantics
  - limits of sequential control flow
  - decoupled execution
  - avoid sequentialization / expose parallelism w/in model
    - extend scoreboard/locking to memory
    - important that memory regions appear in model
  - tolerate variations in implementations
  - support scaling

Caltech CS184 Spring2005 -- DeHon

63