## Octavo: An FPGA-Centric Processor Architecture

Charles Eric LaForest J. Gregory Steffan ECE, University of Toronto

FPGA 2012, February 24

## Easier FPGA Programming

- We focus on overlay architectures
  - Nios, MicroBlaze, Vector Processors
    - These inherited their architectures from ASICs
  - Easy to use with existing software tools
  - Performance penalty
  - ASIC architectures poor fit to FPGA hardware!
- ASIC ≠ FPGA
  - ASIC: transistors, poly, vias, metal layers
  - FPGA: LUTs, BRAMs, DSP Blocks, routing
    - Fixed widths, depths, other discretizations

## FPGA-centric processor design?

## How do FPGAs Want to Compute?

Hardware (Stratix IV)	Width (bits)	Fmax (MHz)
DSP Blocks	36	480
Block RAMs	36	550
ALUTs	1	800
Nios II/f	32	230

# What processor architecture best fits the underlying FPGA?

### **Research Goals**

- 1. Assume threaded data parallelism
- 2. Run at maximum FPGA frequency
- 3. Have high performance
- 4. Never stall
- 5. Aim for simple, minimal ISA
- 6. Match architecture to underlying FPGA

#### Result: Octavo

- 10 stages, 8 threads, 550 MHz
- Family of designs
  - Word width (8 to 72 bits)
  - Memory depth (2 to 32k words)
  - Pipeline depth (8 to 16 stages)

## Snapshot of work-in-progress

## **Designing Octavo**

#### High-Level View of Octavo



Unified registers and RAM

#### Octavo vs. Classic RISC



- All memories unified (no loads/stores)
- How to pipeline Octavo?

## Design For Speed: Self-Loop Characterization

#### Self-Loop Characterization

- Connect module outputs to inputs

   Accounts for the FPGA interconnect
- Pipeline loop paths to absorb delays
- Pointed to other limits than raw delay
  - Minimum clock pulse widths
    - DSP Blocks: 480 MHz
    - BRAMs: 550 MHz

#### We measured some surprising delays...

#### **BRAM Self-Loop Characterization** Μ Μ Μ Μ Μ Μ Μ Μ 398 MHz 531 MHz 656 MHz 710 MHz (routing!) Μ Μ Ň Μ Μ Μ Μ Ň Μ

Must connect BRAMs using registers 11

## **Building Octavo: Memory**

## Building Octavo: Memory





Replicated "scratchpad" memories with I/O while still exceeding 550 MHz limit.

## **Building Octavo: ALU**



Fully pipelined (4 stages)
 – Never stalls



- Multiplication
  - Uses DSP Blocks
  - Must overcome their 480 MHz limit...

## Building Octavo: Multiplier

One multiplier is wide enough but too slow



• Two multipliers working at half-speed

- Send data to both multipliers in alternation



## Octavo: Putting It All Together



**5**¦

6¦

**7**¦

8<sup>¦</sup>

g



0¦

- 10 stages

3¦

2¦

4¦

- Actually 8 stages with one exception (more later)
- No result forwarding or pipeline interlocks
- Scalar, Single-Issue, In-Order, Multi-Threaded



**5**¦

**6**¦

**7**¦

8<sup>¦</sup>

g



3¦

2¦

4

- Indexed by current thread PC
- Provides a 3-operand instruction
- On-chip BRAMs only

#### Octavo 0 2¦ 3 **5**: A/B A/B

- A and B Memories
  - Receive operand addresses from instruction

6

**7**¦

8¦

g

- Provide data operands to ALU and Controller
  - Some addresses map to I/O ports
- On-chip BRAMs only

6

**7**¦

8

g



- Pipeline Registers
  - Avoid an odd number of stages
  - Separate BRAMs for best speed
    - Predicted by BRAM self-loop characterization
    - Unusual but essential design constraint



#### Controller

- Receives opcode, source/destination operands
- Decides branches
- Provides current PC of next thread to I memory



• ALU

- Receives opcode and data
- Writes result to all memories



Longest mandatory loop: 8 stages

 Along A/B memories and ALU
 Fill with 8 threads to avoid stalls



- Special case longest loop: 10 stages
  - Along instruction memory and ALU
  - Does not affect most computations
    - Adds a delay slot to subroutine and loop code

## **Results: Speed and Area**

## Experimental Framework

- Quartus 10.1 targeting Stratix IV (fastest)
  - Optimize and place for speed
  - Average speed over 10 placement runs
- Varied processor parameters:
  - Word width
  - Memory depth
  - Pipeline depth
- Measure Frequency, Area, and Density

## Maximum Operating Frequency

#### Maximum Operating Frequency



Faster

Maximum Operating Frequency



#### Maximum Operating Frequency



## Area Density

## Area Density



"Sweet spot" 72 bits, 1024 words



## Designing Octavo: Lessons & Future Work

#### Lessons

- Soft-processors can hit BRAM Fmax
   Octavo: 8 threads, 10 stages, 550 MHz
- Self-loop characterization for modules

   Helps reason about their pipelining
   Shows true operating envelopes on FPGA
- Octavo spans a large design space
   Significant range of widths, depths, stages

#### **Consider FPGA-centric architecture!**

#### Future Work

