

WaveScalar: the Executive Summary

Dataflow machine

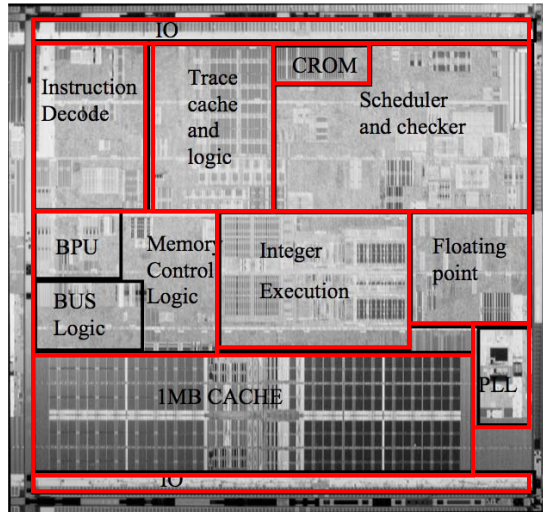
- good at exploiting ILP
- dataflow parallelism + traditional coarser-grain parallelism
 - cheap thread management
- low operand latency because of a hierarchical organization
- memory ordering enforced through **wave-ordered memory**
 - no special languages

WaveScalar

Additional motivation:

- increasing disparity between computation (fast transistors) & communication (long wires)
- increasing circuit complexity
- decreasing fabrication reliability

Monolithic von Neumann Processors



A phenomenal success today.
But in 2016?

⊗ Performance
Centralized processing & control
Long wires
e.g., operand broadcast networks

⊗ Complexity
40-75% of "design" time is design verification

⊗ Defect tolerance
1 flaw -> paperweight

Autumn 2006

CSE P548

3

WaveScalar's Microarchitecture

Good performance via distributed microarchitecture 😊

- hundreds of PEs
- dataflow execution – no centralized control
- short point-to-point (producer to consumer) operand communication
- organized hierarchically for fast communication between neighboring PEs
- scalable

Low design complexity through simple, identical PEs 😊

- design one & stamp out thousands

Defect tolerance 😊

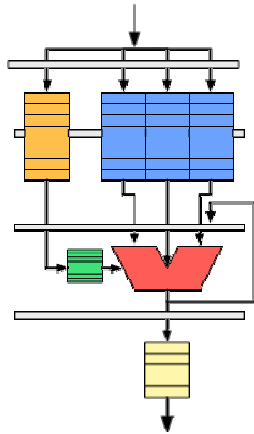
- route around a bad PE

Autumn 2006

CSE P548

4

Processing Element



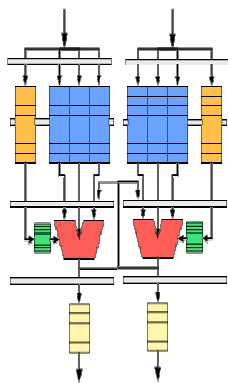
- Simple, small (.5M transistors)
- 5-stage pipeline (receive input operands, match tags, instruction schedule, execute, send output)
- Holds 64 (decoded) instructions
- 128-entry token store
- 4-entry output buffer

Autumn 2006

CSE P548

5

PEs in a Pod



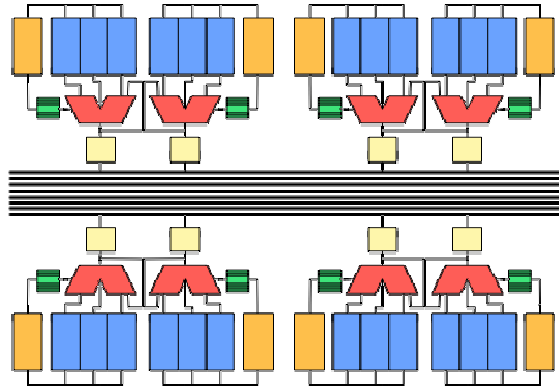
- Share operand bypass network
- Back-to-back producer-consumer execution across PEs
- Relieve congestion on intra-domain bus

Autumn 2006

CSE P548

6

Domain

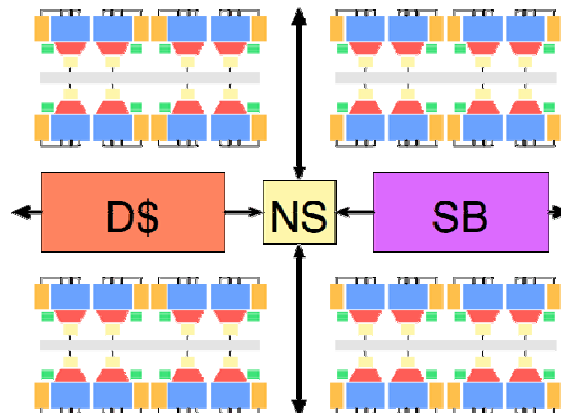


Autumn 2006

CSE P548

7

Cluster



Autumn 2006

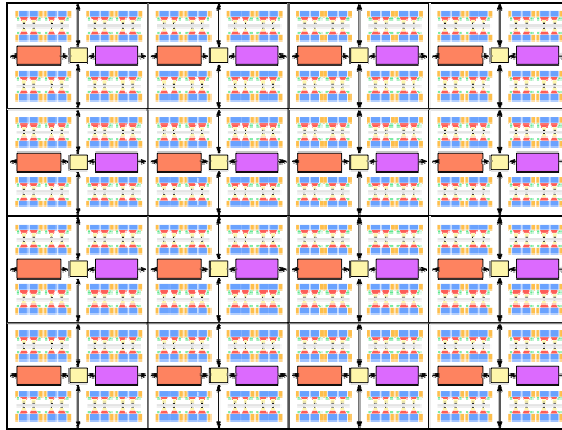
CSE P548

8

WaveScalar Processor

Long distance
communication

- dynamic routing
- grid-based network
- 2-cycle hop/cluster



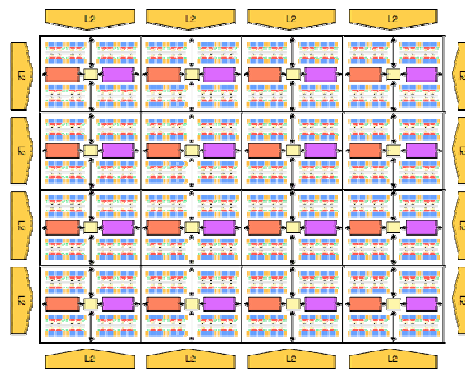
Autumn 2006

CSE P548

9

Whole Chip

- Can hold 32K instructions
- Normal memory hierarchy
- Traditional directory-based cache coherence
- ~400 mm² in 90 nm technology
- 1GHz.
- ~85 watts



Autumn 2006

CSE P548

10

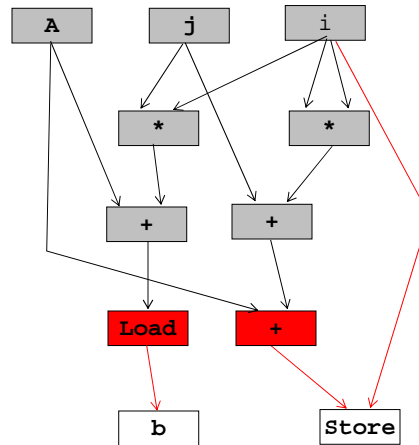
WaveScalar Instruction Placement

Place instructions in PEs to maximize data locality & instruction-level parallelism.

- Instruction placement algorithm based on a performance model that captures the important performance factors
- Carve the dataflow graph into segments
 - a particular depth to make chains of dependent instructions that will be placed in the same pod
 - a particular width to make multiple independent chains that will be placed in different, but near-by pods
- Snakes segments across PES in the chip on demand
- K-loop bounding to prevent instruction “explosion”

Example to Illustrate the Memory Ordering Problem

`A[j + i*i] = i;`
`b = A[i*j];`



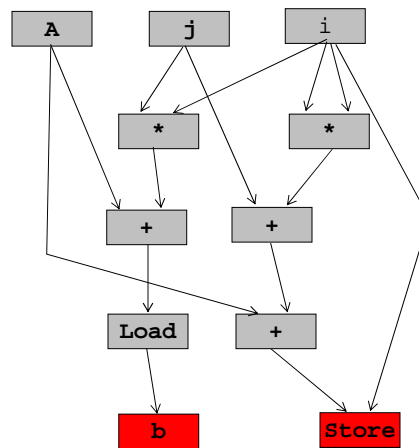
Autumn 2006

CSE P548

13

Example to Illustrate the Memory Ordering Problem

`A[j + i*i] = i;`
`b = A[i*j];`

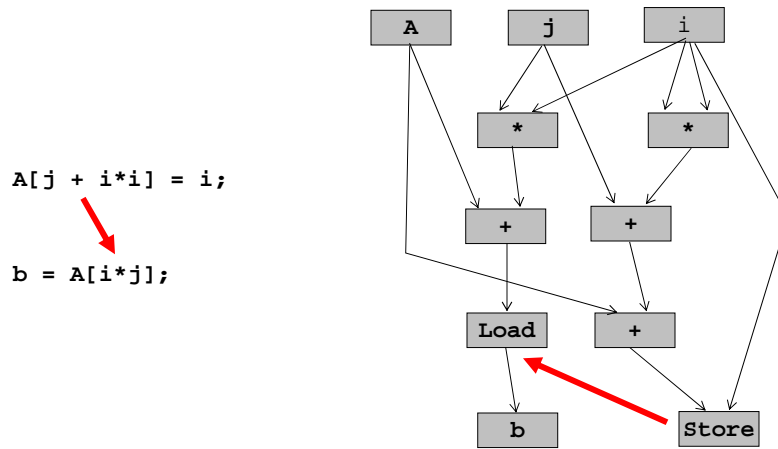


Autumn 2006

CSE P548

14

Example to Illustrate the Memory Ordering Problem



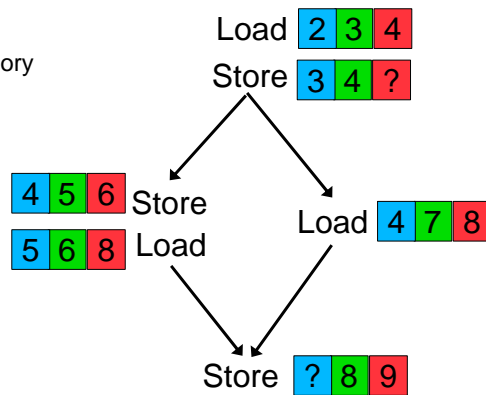
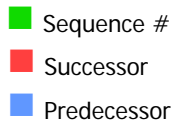
Autumn 2006

CSE P548

15

Wave-ordered Memory

- Compiler annotates memory operations

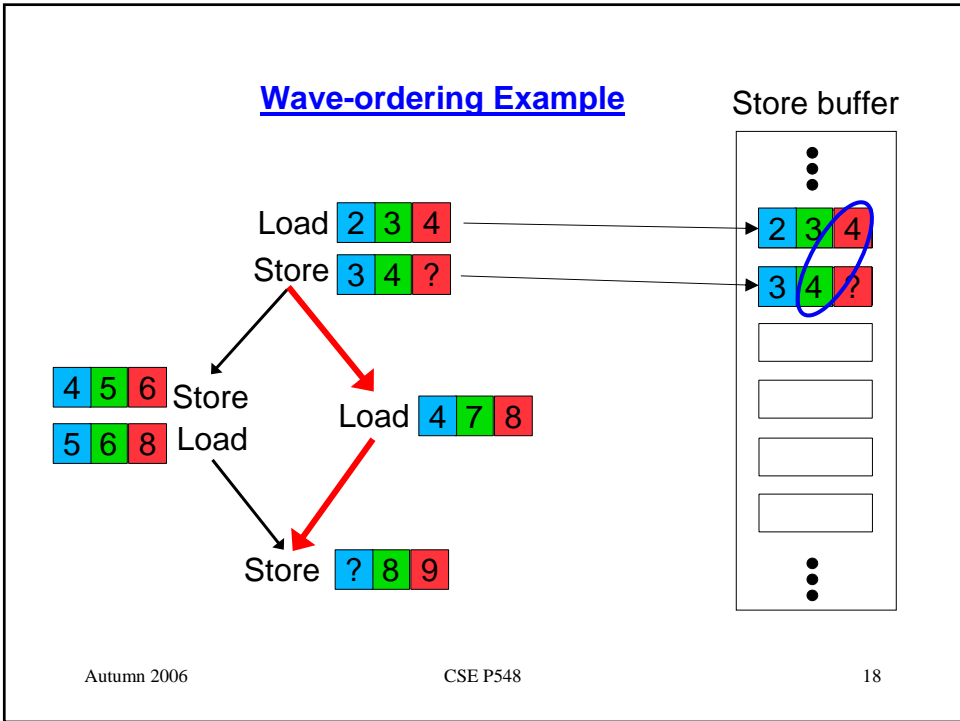
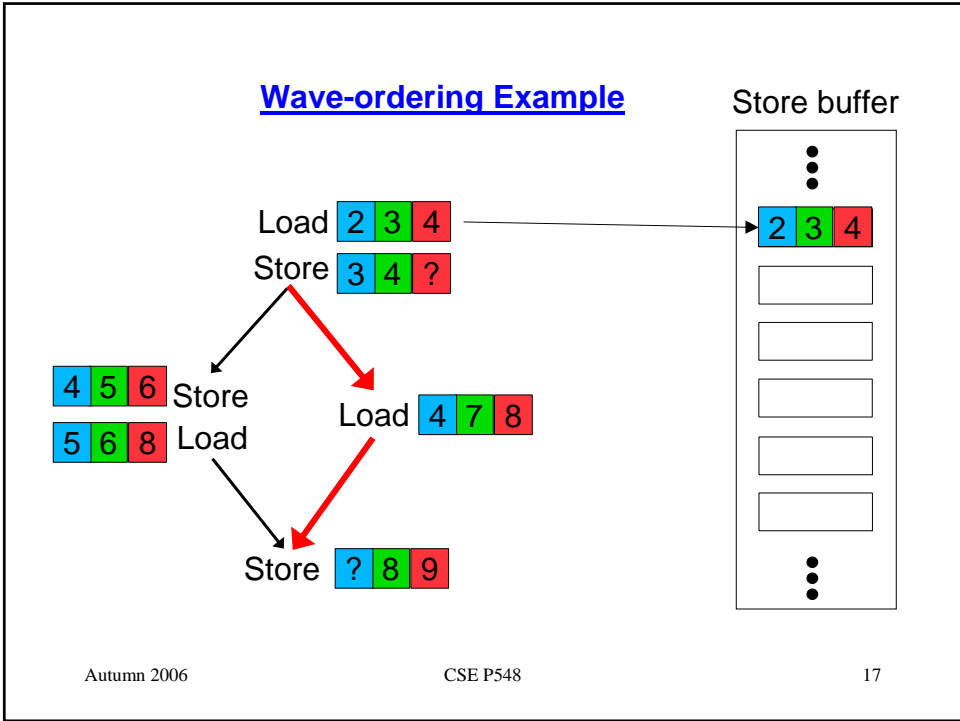


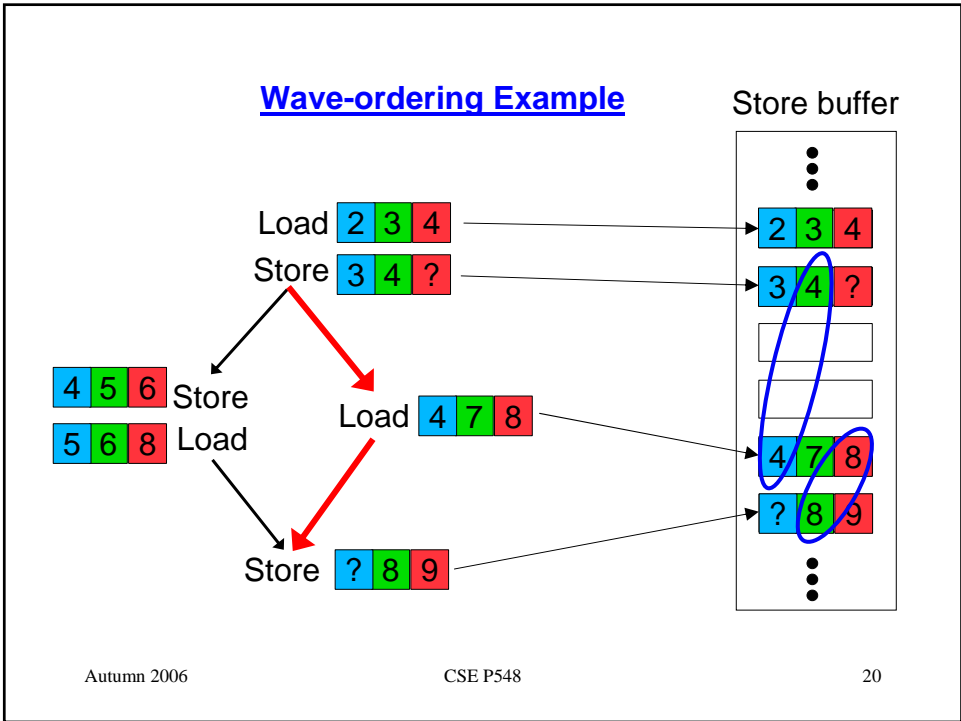
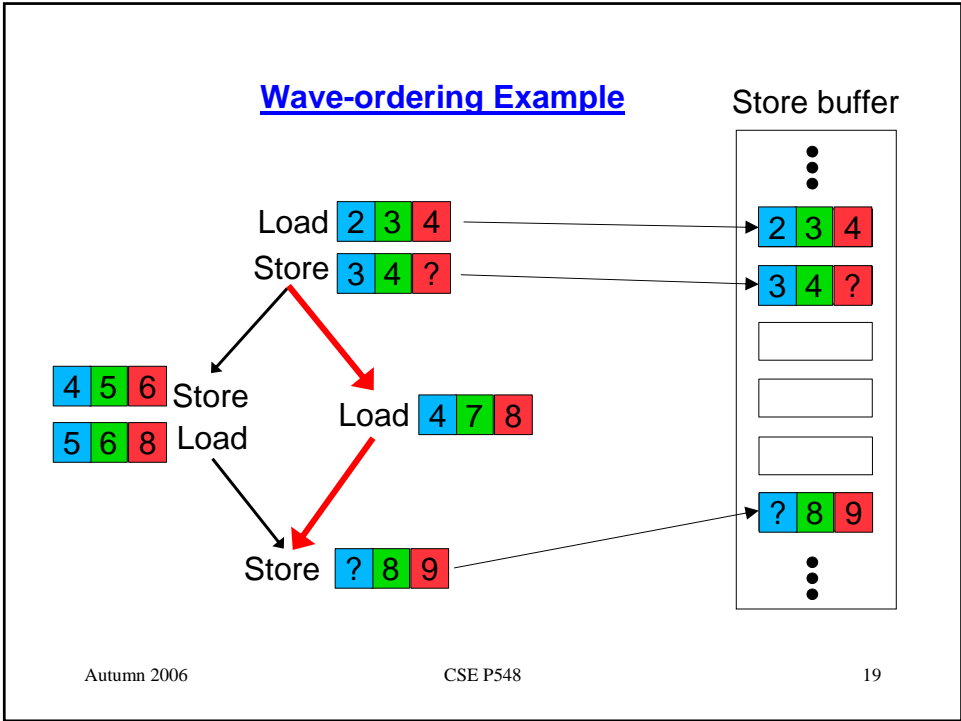
- Send memory requests in any order
- Hardware reconstructs the correct order

Autumn 2006

CSE P548

16





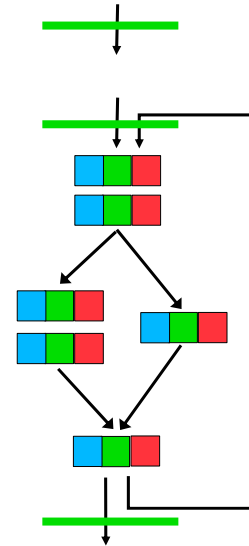
Wave-ordered Memory

Waves are loop-free sections of the dataflow graph

Each dynamic wave has a **wave number**
Wave number is incremented between waves

Ordering memory:

- wave-numbers
- sequence number within a wave



Autumn 2006

CSE P548

21

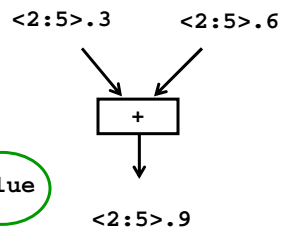
WaveScalar Tag-matching

WaveScalar tag

- thread identifier
- wave number

Token: **tag** & **value**

<ThreadID:Wave#:InstructionID> value

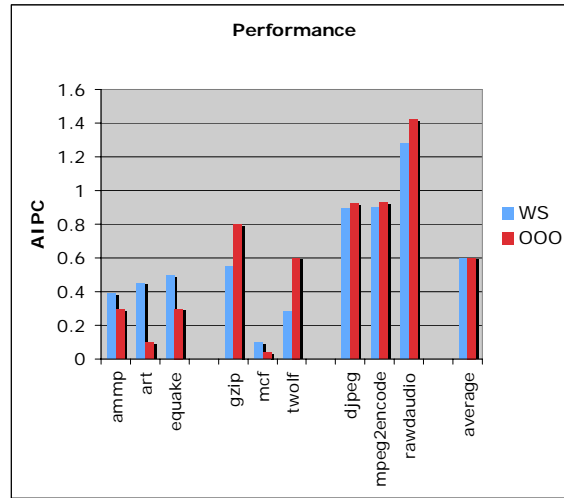


Autumn 2006

CSE P548

22

Single-thread Performance

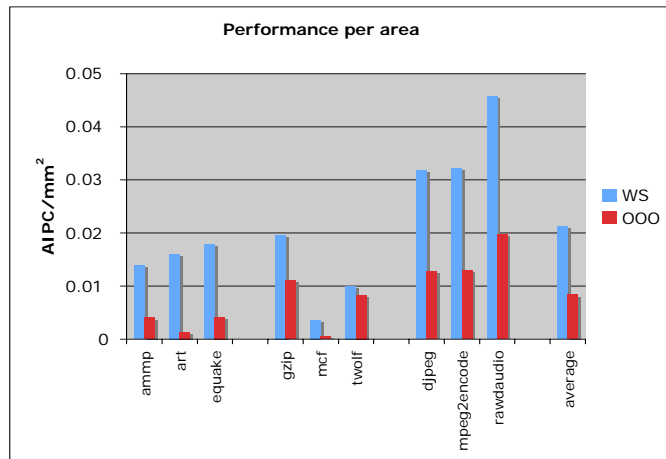


Autumn 2006

CSE P548

23

Single-thread Performance per Area



Autumn 2006

CSE P548

24

Multithreading the WaveCache

Architectural-support for WaveScalar threads

- instructions to start & stop memory orderings, i.e., threads
- memory-free synchronization to allow exclusive access to data (thread communicate instruction)
- fence instruction to force all previous memory operations to fully execute (to allow other threads to see the results of this one's memory ops)

Combine to build threads with multiple granularities

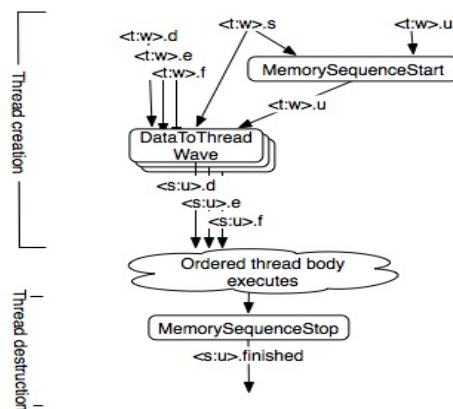
- coarse-grain threads: 25-168X over a single thread; 2-16X over CMP, 5-11X over SMT
- fine-grain, dataflow-style threads: 18-242X over single thread
- combine the two in the same application: 1.6X or 7.9X -> 9X

Autumn 2006

CSE P548

25

Creating & Terminating a Thread

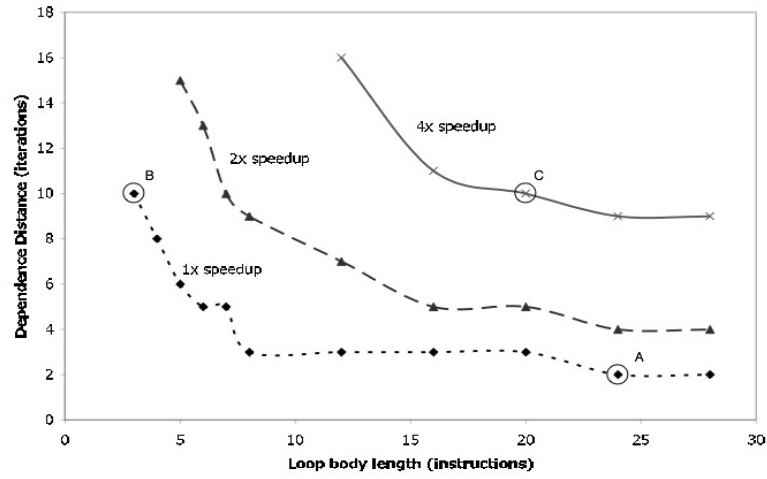


Autumn 2006

CSE P548

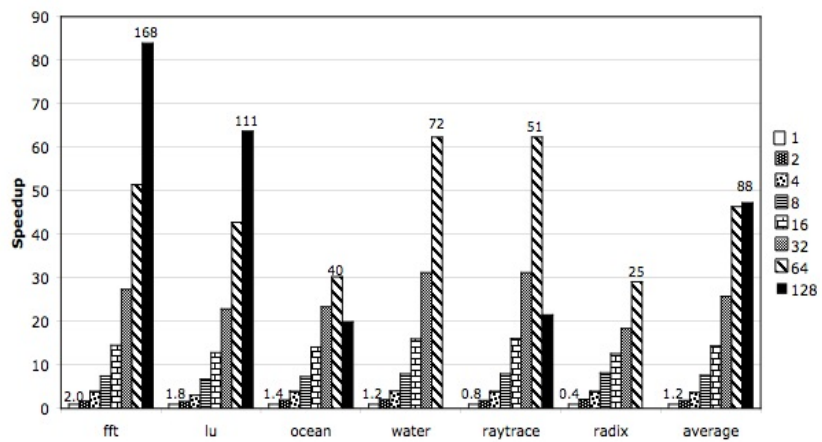
26

Thread Creation Overhead



27

Performance of Coarse-grain Parallelism

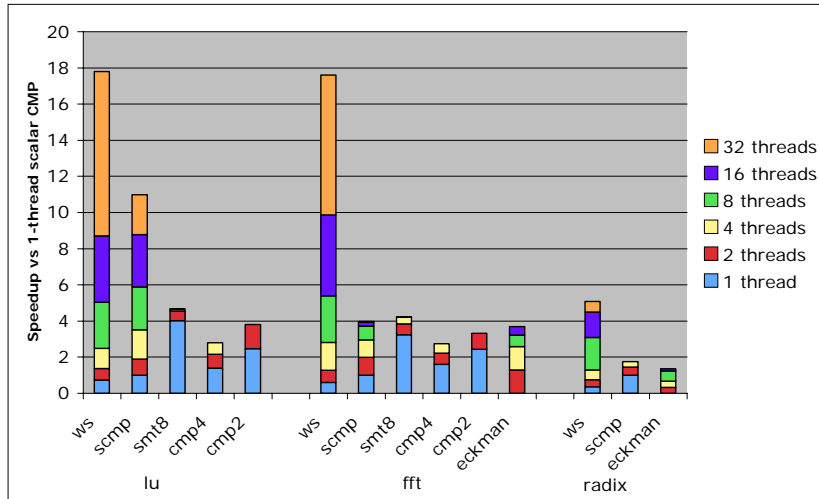


Autumn 2006

CSE P548

28

CMP Comparison



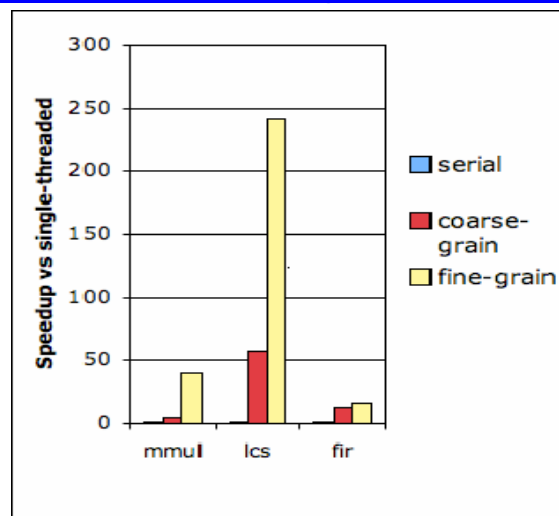
Autumn 2006

CSE P548

29

Relies on:
Cheap synchronization
Load once, pass data (not load/compute/store)

Performance of Fine-grain Parallelism



Autumn 2006

CSE P548

30

Building the WaveCache

RTL-level implementation

- some didn't believe it could be built in a normal-sized chip
- some didn't believe it could achieve a decent cycle time and load-use latencies
- Verilog & Synopsys CAD tools

Different WaveCache's for different applications

- 1 cluster: low-cost, low power, single-thread or embedded
 - 42 mm² in 90 nm process technology, 2.2 AIPC on Splash2
- 16 clusters: multiple threads, higher performance: 378 mm² , 15.8 AIPC

Board-level FPGA implementation

- OS & real application simulations