Using Fine Grain Multithreading for Energy Efficient Computing

Alex Gontmakher, Technion Avi Mendelson, Intel Labs Assaf Schuster, Technion

Instruction Level Parallelism

X

Program

Dataflow Execution





Practical Execution: constrained by processor resources



How Out-of-Order Works





- Sliding window over the program
 - Independent instructions within the window executed
- Much of the parallelism is beyond the horizon!
 - Good for extremely fine granularity!
- All instructions within the window must be checked for dependencies
 - High complexity
 - Increased energy consumption

Who's The/A Culprit?					
	MIPS		Alpha		
$2\mathbf{v}$	R5K	R10K	21164	21264	
Pipeline	InO	000	InO	000	
SPECInt95	5.5	10.1	15.4	27.7	
SPECFP95	5.5	9.71	21.1	58.7	
Power (W)	10	30	28	91	
Clock Rate	200	250	500	500	
Technology (μ)	0.35	0.35	0.35	0.35	
Transistors	$3.6\mathrm{M}$	$6.7\mathrm{M}$	9.3M	$15.2\mathrm{M}$	
Transistors (Logic)	$)$ $0.8 \mathrm{M}$	2.3M	2.6M	6M	
Physical Regs		64		80	



- S. Hily, A. Seznec: Out-of-order execution may not be cost-effective on SMT processors
 - **Single thread**: inorder 46% slower than OOO
 - 4 threads: inorder only 15% slower than OOO

 SUN's Niagara processor: Massive SMT support with inorder pipelines, especially for low-power computing (throughputoriented)

Thread Based Parallelism

Coarse Granularity



Fine Granularity



Parallelism expressed during compilation

- Communication through shared variables
- OS-based runtime support

Relatively high threading overhead

 Good only for coarse granularity!

The Main Idea





Inthreads: An extremely lightweight threading mechanism

- Hardware-based runtime support
- Communication through shared registers

Medium-level granularity

- Beyond the horizon for OOO
- Too fine-grain for regular threads

Inthreads: Programming Model

Lightweight architecture

- Fixed number of threads
- Shared registers
- Synchronization instructions



Code regions explicitly belong to specific threads

- Opportunity for better compiler optimizations

Fits within a function call frame. Function calls?

- inlining
- suspend/resume

Inthreads: Programming Model 2



- Threads share the registers cooperatively
 - <u>Thread-private variables</u> use different registers in each thread
 - <u>Shared variables</u> must be allocated to same register in all threads
 - Accesses to shared variables must be protected by synchronization
- Memory (+register) consistency model: Data-Race-Free-1 (DRF1)
 - Software: no data races
 - Processor: obeys instruction ordering



Inthreads vs. SMT: Fast Comm



Inthreads ISA by Example





Condition Registers



More Instructions:

inth.clr - clear a condition variable
inth.kill - kills a given thread

How lightweight it really is?

- Two independent tasks executed in parallel
 - Expected top speedup: 2x
 - Speedup 1x => Overhead is the same as task size





000 Execution Out Of Rename Commit Decode Fetch ssue Order InO Execution Decode Fetch Issue Inorder Execution Inorder/ Decode Fetch Issue Wait NO MT

Fetch Stage



Jump Address



Instruction Issue



Dynamic Instruction Mix





Inorder+Inthreads: less unnecessary work!

Execution Time

- Inorder is much worse than OOO...
- But Inorder+Inthreads isn't!
 - Although by different methods, Inorder+Inthreads achieves the same latency-tolerance as OOO



Energy, Energy-Delay Results

Processor mostly idle, waiting for memory. Little progress occurs, but power is consumed.



How IPC Affects Energy

Energy consumption per period

- X
- Divide execution into periods 1000 cycles long
- For each period, plot the energy consumed as a function of the IPC



Conclusions



OOO: **fast single task** ↔ **high-power**

MT: energy efficient ↔ throughput-oriented

Inthreads:

Shared Register Threads + Data-Race-Free-1



Best of both

Questions

