Fine-Grained Multithreading Filling Pipes with Instructions from Alternate Sources

Superscalar Issue Example

- Multiple pipelines with different functions
- Issue instructions as available from prefetch
- 8 pipes of depth 10
- To keep example simple, assume no stalls

Br	FP	FP	L/S	Br	Int	L/S	Int	Int	Int
O SO SO S	000000	00000	0000000	00000	200000	000000	000000		00000

FP	
FP	
Int	
Int	
Int	
L/S	
L/S	
Br	



L/S	FP	FP	FP	Br	FP	FP	L/S	Br	Int

FP	
FP	
Int	
Int	
Int	
L/S	
L/S	
Br	



L/S L	J/S L/	S FP	FP	L/S	FP	FP	FP	Br	





Br	Int	Br	L/S	L/S	L/S	FP	FP	L/S	FP	
000000000000000000000000000000000000000	000000		000000		00000	00000	0000	00000	0000	50000



Int Int	Int	Br Int	Br	L/S	L/S	L/S	FP	











FP	FP	L/S	L/S	L/S	Int	Int	Int	Int	Br	





FP	FP	BR	FP	FP	FP	L/S	L/S	L/S	Int	





FP	FP	BR	FP	FP	BR	FP	FP	FP	L/S

FP		
FP		
Int		
Int		
INT		
1 /C		
L/ J		
1/9		
Br		
		1212



FP	FP	BR	FP	FP	BR	FP	FP	BR	FP	





Int	Int	Int	FP	FP	BR	FP	FP	BR	FP
					00000				



34 slots filled/80 available = 42.5% utilization

Simple Multithreading

- Prefetch two threads at once
- Simplest source is user and supervisor
- Issue from thread with most ready ops
- Do not issue from stalled thread



FP	FP	Br	Int	FP	L/S	Int	Br	FP	FP
Int	Int	FP	L/S	Br	Int	L/S	Int	Int	Int

FP				
FP				
Int				
Int				
Int				
L/S				
1/0				
L/				
₽r				
	00000000000	2202020202020	KOKOKOKOK	40X0X0X0X0-D



FP	FP	Br	Int	FP	L/S	Int	Br	FP	FP
FP	FP	FP	L/S	Br	Int	L/S	Int	Int	Int



L/S	L/S	Br	L/S	LS	FP	FP	Br	Int	FP
L/S	FP	FP	L/S	FP	FP	FP	L/S	Br	Int





L/S	L/S	Br	L/S	LS	FP	FP	Br	Int	FP	
Br	Int	Br	L/S	L/S	L/S	FP	FP	L/S	FP	

FP	
FP	
Int	
Int	
Int	
L/S	
L/S	
Br	



L/S	Int	Int	Int	L/S	L/S	Br	L/S	LS	FP
Br	Int	Br	L/S	L/S	L/S	FP	FP	L/S	FP

FP	
FP	
Int	
Int	
Int	
L/S	
L/S	
Br	



L/S	Int	Int	Int	L/S	L/S	Br	L/S	LS	FP	
Br	Int	Br	L/S	L/S	L/S	FP	FP	L/S	FP	



Assume the red L/S stalls 4 cycles



L/S	Int	Int	Int	L/S	Int	Int	Int	L/S	L/S
Br	Int	Br	L/S	L/S	L/S	FP	FP	L/S	FP

FP	
FP	
Int	
Int	
Int	
L/S	
1 /C	
Rr	



Br	L/S	Int	Int	Int	L/S	Int	Int	Int	L/S
Br	Int	Br	L/S	L/S	L/S	FP	FP	L/S	FP





L/SIntFPBrL/SIntIntBrIntBrL/SL/SL/SFPFPL/S	L/S FP
Br Int Br L/S L/S FP FP L/S	FP
Br Int Br L/S L/S FP FP L/S	FP
	000000

FP	
FP	
Int	
Int	
Int	
L/S	
L/S	
Br	



Int	L/S	L/S	L/S	L/S	Int	FP	FP	Br	L/S
Br	Int	Br	L/S	L/S	L/S	FP	FP	L/S	FP

FP		
FP		
Int		
Int		
INT		
1/0		
L/S		
1 /오		
Br		



FD	FD	Br	т./s	т./s	Tnt	Tnt	т./с	т./s	т./s
Br	Int	Br	L/S	L/S	L/S	FP	FP	L/S	FP

FP		
FP		
Int		
Int		
Int		
L/S		
L/S		
Br		



<u> </u>									100					100		- N.		- N.	- A.	100	- C	100	100	100		10 A							- No.				÷																
	9		20							9		2	2	2		27	7	2	1	1	1	27	1	20	10				5	-		-	27	1	20	10					2	20	10				27	1	20			20	
									27		-		27	3	-11	1	47	1	1	1	47	1	47	Y-	20	1		5		\smile				27	1	10	1		\sim		27	1		\sim		9			1		$ \sim$	-10-	
								9		-		20		27	2	27	2	-17	7	47	1	47	1	20	1	10	\sim		\sim				1	1	47	1-		\sim		9		-17	1	\sim	\sim		-11		10	\sim		27	\sim
							9		-1	2	-17		-17		-11		-17	15	47	1	47	Y:	-17	1	-	1		\sim		\smile	\sim	9		-17			\sim		\sim		-	1		\sim		\sim		<u> </u>	1		\leq		60
				2		\leq		-		-1		-1		-11	20	-1	32	-1	2	-		1	n-	-	$\gamma \sim$		\sim		\sim		\simeq	\sim	1		27			\sim		9		40	\geq		\sim		-			\sim		27	
					2		31		1		1		-		-1		-1		1	Y.	1		1		1	\mathbb{N}		\simeq		\simeq	\sim	1	25	1		10	\geq		\sim		1		10	\mathbb{M}	\sim	M		H.C	\geq		M		
				\simeq		1		1		1				-1		1	\mathbf{x}	-1	22	1		1		1		1	\mathbb{M}		\simeq		\simeq	\mathbf{O}			-		10	\simeq		1		1		$\mathbf{\cap}$	\simeq		M		10	\simeq	\mathbf{O}	H.	
		- C.	÷.		-		Ψ.		÷.		4		Υ.		-		-		-		1	11-	11	10-	£1.			\sim		\sim			115	-	- 10-	1	1-		\sim	- 10	<u> </u>	10-	÷.,	\sim			- 15	-		<u> </u>		- 12%	
														=1						-	100																	0															

Int

Br

Int

L/S

L/S

L/S

L/S L/S

L/S

 \mathbf{FP}

L/S

Int

L/S

Br

FP

Int

FP

Int

Br

Int



40 slots filled/80 available = 50% utilization (even with the extra overhead of a stall)

Notes

- Could have issued from both threads
- Likely that dependences in stalled thread would limit forward progress
- Some of boost in utilization was due to the second thread having a better schedule
- System threads rarely have FP ops
- Assumed no other hazards
- Real situation likely to have lower utilization



Cost of Multithreading

- Extra prefetch queue and logic per thread, more complex scheduling
- Thread tags for instructions in pipelines (more complex commit)
- Causes functional units to run hotter
- OS must manage threads

Separate register bank(s) needed for each thread, including PC, PSW, etc. Switching logic to select register bank (and register update unit stations)

Simultaneous Multithreading

Issue the maximum mix of ready ops from the two threads

Requires analyzing a larger set of pending ops

FP	FP	Br	Int	FP	L/S	Int	Br	FP	FP	
FP	FP	FP	L/S	Br	Int	L/S	Int	Int	Int	

80000

							555555555555555555555555555555555555555
FP							
FP							
Int							
Int							
Int							
L/S							
L/S							10909090 10909090
	50505050	858585	158582	258585	158585	102020	5555-5



										0000000000
										000000000
Br	L/S	L/S	FP	FP	Br	Int	FP	L/S	Int	
										0000
L/S	FP	FP	L/S	FP	FP	FP	L/S	Br	Int	
										000
										000
										0000
										000
										000
										000
										000

FP	
FP	
Int.	
Int	
Int	
L/S	
L/S	
D	



Int	Int	L/S	L/S	Br	L/S	L/S	FP	FP	Br
								the second second	

Int	Br	L/S	L/S	L/S	FP	FP	L/S	FP	FP
858585									





Int	Int	Br	Int	Br	L/S	L/S	L/S	FP	FP

Int	Int	Int	L/S	L/S	Br	L/S	L/S	FP	FP
252525	000000	000000	20000	200000	000000	000000	00000	000000	000000

FP	
FP	
Int	
Int	
Int	
L/S	
L/S	
Br	



Int	Int	Int	L/S	Int	Int	Int	Int	L/S	L/S
Int	Int	Br	Int	Br	L/S	L/S	L/S	FP	FP

FP	
FP	
Int.	
Int	
Int	
L/S	
L/S	
Br	



Tnt	ED	ΓD	Bro	T/C	Tat	The	Tnt	т /с	Tnt
Inc	FP	ΓP	BL	Г\Р	LUC			<u>с, т</u>	
Int	Int	Int	Int	Br	Int	Br	L/S	L/S	L/S

000000000000000000000000000000000000000	
FP	
FP	
Int	
Int	
Int	
L/S	
L/S	
Br	



L/S L/S Int FP FP Br L/S Int Int L/S Int Int Int Int Br Int Br L/S L/S										
L/S Int Int Int Int Br Int Br L/S L/S	L/S	L/S	L/S	Int	FP	FP	Br	L/S	Int	Int
	L/S	Int	Int	Int	Int	Br	Int	Br	L/S	L/S

FP	
FP	
Int	
Int	
Int	
L/S	
L/S	
Br	
	, , , , , , , , , , , , , , , , , , ,



FP	Br	L/S	L/S	Int	Int	L/S	L/S	L/S	L/S
							and the second second	and the second s	

L/S	L/S	Int	Int	Int	Int	Br	Int	Br	L/S
						202020			

FP		
FP		
Int		
Int		
Int		
L/S		
L/S		
Br		
	-0	



FP	FP	Br	L/S	L/S	Int	Int	L/S	L/S	L/S	
FP	FP	L/S	L/S	L/S	Int	Int	Int	Int	Br	

FP			
FP			
Int			
Int			
Int			
L/S			
L/S			
Br			83333



	THE	505 ± 50	525532	OPTO	СПАР	с П, Р,	THE	THE	цір	ŏS
FP	FP	Br	FP	FP	FP	L/S	L/S	L/S	Int	
										000
										000
										000
										00
										000
										Š S S S S S S S S S S S S S S S S S S S
										000
KOKOKO										0

2202020	<u>zozozo</u>	202020	2020202	00000	020202	020202	0202020	2202020	202020
1545454	545454	525252	525252	262626	262626				
Tnt	Tnt	FP	FP	Br	T./S	T./S	Tnt	Tnt	T./S
THE	T11.C	2000	2000	OPTO		CHI C	THE	THE	ц,р
JUNUNU	HUNUNU	JUNUNU	JUAUAU	U AUAUA	UNUNUN	UAUAUA	UnUnUnu	JAUAUAU	hununu



Int Int L/S

Br



Br	FP	FP	FP	Br	FP	FP	FP	L/S	L/S
						00000			
	000000	000000			000000	000000			0000000

L/S	Int	L/S	Int	Int	FP	FP	Br	L/S	L/S
5252525	0000000	000000	00000	20000	200000	00000	00000	325252	000000



58 slots filled/80 available = 72.5% utilization

Notes

- Original thread is still not done
- Higher throughput, lower per-thread rate
- More pressure on all units
- Memory busy on most cycles
- Branch predictor complications
- Exception handling more complex



Tullsen ISCA 1995 Simultaneous Multithreading: Maximizing On-Chip Parallelism



SVI Introduced

- latency hiding
- Exploited in fine-grained way for computation in HEP
- Further developed in Tera
- thread on each cycle

Multithreading originated at least in early CDC machines (early 60s) for I/O

SMT adds ability to go beyond switching threads on each cycle, and blend

Models

- Single issue per thread per cycle
- Dual, quad, full issue per thread per cycle
- Compared with fine-grained multithreading

Limited range of connections between thread pipes and functional units

Results





Results



Compared to Multicore

Purpose of Test	Common Elements	Spe
Unlimited FUs: equal	Test A: $FUs = 32$	SM: 8
total issue bandwidth, equal number of register	Reg sets = 8	MP: 8
sets (processors or	Test B: $FUs = 16$	SM: 4
threads)	Issue bw = 4 Reg sets = 4	MP: 4
	Test C: $FUs = 16$	SM: 4
	Issue bw $= 8$ Reg sets $= 4$	MP: 4
Unlimited FUs: Test A,	Test D:	SM: 8
but limit SM to 10 FUs	Issue $bw = 8$ Reg sets = 8	MP: 8
Unequal Issue BW: MP	Test E:	SM: 8
has up to four times the	FUs = 32 Reg sets = 8	MP: 8
total issue bandwidth	Test F:	SM: 4
	FUs = 16 Reg sets = 4	MP: 4
FU Utilization: equal	Test G:	SM- 8
FUs, equal issue bw,	FUs = 8	MD. 2
unequal reg sets	Issue $BW = 8$	MP: 2

Figure 5: Results for the various multiprocessor vs. simultaneous multithreading comparisons. The multiprocessor always has one functional unit of each type per processor. In most cases the SM processor has the same total number of each FU type as the MP.

cific Configuration	Throughput (instructions/cycle)
3 thread, 8-issue	6.64
3 1-issue procs	5.13
thread, 4-issue	3.40
1-issue procs	2.77
thread, 8-issue	4.15
2-issue procs	3.44
thread, 8 issue, 10 FU	6.36
l-issue procs, 32 FU	5.13
8 thread, 8-issue	6.64
3 4-issue procs	6.35
thread, 8-issue	4.15
4-issue procs	3.72
3 thread, 8-issue	5.30
2 4-issue procs	1.94

Discussion

Mahlke ISCA 1995 ILP Processors



Predication

- Added operand that specifies predicate for execution of an instruction
- If satisfied, instruction completes, otherwise squashed
- Full predication adds capability to all instructions
- Partial supports some (e.g., conditional moves)
 - Depends on speculation

Implementation

- Set of predicate registers
- Predicate define instructions
- Instructions refer to defined predicates
- Similar to Itanium

Different from ARM-32, where each instruction can conditionally set a condition, and be predicated on existing conditions (ARM-64 more limited)

Compilation

- Uses hyperblocks (connected series of basic blocks) with single entry, multiple exit
- Control flow between blocks eliminated by if-conversion
- - Significant code expansion (roughly double)
- Then peephole optimize

Create fully predicated code, then transform to partially predicated code

Simple 8-issue Model



2-Branch Speculation



Simple 4-issue Model



Realistic Caches



Figure 11: Effectiveness of full and partial predicate support for an 8-issue, 1-branch processor with 64K instruction and data caches.

Effect on Instruction Count

Benchmark	Superblk			
008.espresso	489M			
022.li	31M			
023.eqntott	1030M			
026.compress	90M			
052.alvinn	3574M			
056.ear	11225M			
072.sc	91M			
ссср	3701K			
cmp	932K			
eqn	44M			
grep	1282K			
lex	36M			
qsort	44M			
wc	1493K			
yacc	43M			

Table 2: Dynamic instruction count comparison.

Cond. Move	Full Pred.			
812M (1.66)	626M(1.28)			
38M (1.23)	32M (1.04)			
1230M(1.19)	885M (0.86)			
128M(1.41)	108M(1.20)			
4003M (1.12)	3603M (1.01)			
13838M (1.23)	11073M (0.99)			
85M (0.93)	75M (0.83)			
5077K (1.37)	3855K(1.04)			
1422K(1.53)	922K (0.99)			
49M (1.11)	44M (0.99)			
2467K (1.92)	1647K(1.28)			
75M (2.10)	46M (1.29)			
70M (1.61)	49M (1.11)			
2999K (2.01)	1526K(1.02)			
66M(1.53)	50M (1.16)			

Effect on Rest of Branches

Benchmark	Superblock			Conditional Move			Full Predication		
	BR	MP	MPR	BR	MP	MPR	\mathbf{BR}	MP	MPR
008.espresso	75M	3402K	4.55%	38M	2066K	5.38%	33M	1039K	3.15%
022.li	7457K	774K	10.38%	6169K	694K	11.25%	6110K	702K	11.5%
023.eqntott	315M	42M	13.47%	53M	6732K	12.66%	51M	6931K	13.57%
026.compress	12M	1344K	10.9%	9269K	864K	9.32%	9240K	867K	9.38%
052.alvinn	463M	1091K	0.24%	74M	896K	1.23%	74M	1032K	1.38%
056.ear	1539M	66M	4.3%	443M	16M	3.52%	442M	15M	3.4%
072.sc	22M	1232K	5.49%	11M	1044K	9.19%	11M	934K	8.26%
ссср	921K	66K	7.19%	537K	65K	12.17%	534K	65K	12.15%
cmp	530K	4395	0.83%	26K	31	0.12%	26K	31	0.12%
eqn	7470K	1612K	8.2%	4506K	514K	11.4%	4495K	511K	11.37%
grep	663K	9660	1.46%	171K	20K	11.7%	171K	20K	11.73%
lex	14M	232K	1.65%	3070K	$201 \mathrm{K}$	6.55%	$3030 \mathrm{K}$	196K	6.46%
\mathbf{qsort}	8847K	1332K	15.06%	6092K	597K	9.79%	6066K	610K	10.06%
wc	478K	33K	6.85%	224K	57	.025%	$224 \mathrm{K}$	57	.025%
yacc	12M	517K	4.31%	5944K	445K	7.48%	5900K	431K	7.31%

BRanches, MisPredicts, MisPredict Rate

Discussion