ARCHITECTURE OF

COMPUTER SYSTEMS LECTURE 13 - VLIW MACHINES AND STATICALLY SCHEDULED ILP

LASTIIMEIN

Unified physical register file machines remove data values from ROB

- All values only read and written during execution
- Only register tags held in ROB
- Allocate resources (ROB slot, destination physical register, memory reorder queue location) during decode
- Issue window can be separated from ROB and made smaller than ROB (allocate in decode, free after instruction completes)
- Free resources on commit
- Speculative store buffer holds store values before commit to allow load-store forwarding
- Can execute later loads past earlier stores when addresses known, or predicted no dependence

SUPERSCALAR CONTROL LOGIC SGALLING

Previously
Issued
Instructions

Lifetime L

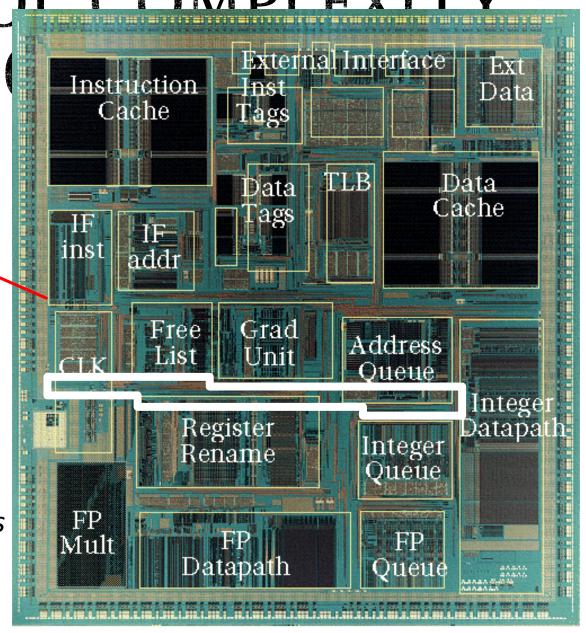
- Each issued instruction must somehow check against W*L instructions, i.e., growth in hardware ∞ W*(W*L)
- For in-order machines, L is related to pipeline latencies and check is done during issue (interlocks or scoreboard)
- For out-of-order machines, L also includes time spent in instruction buffers (instruction window or ROB), and check is done by broadcasting tags to waiting instructions at write back (completion)
- As W increases, larger instruction window is needed to find enough parallelism to keep machine busy => greater L

=> Out-of-order control logic grows faster than W^2 ($\sim W^3$)

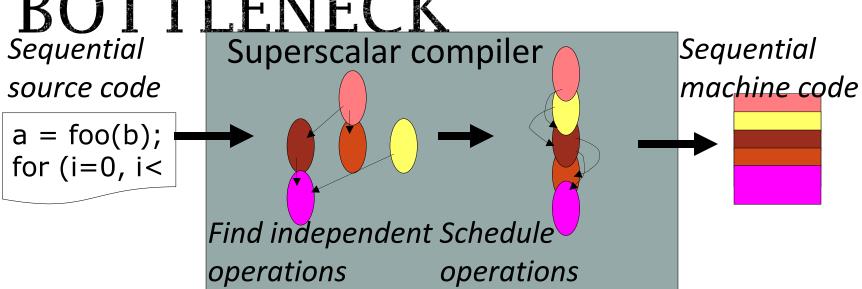
CONTRO MIPS R1

> Control Logic

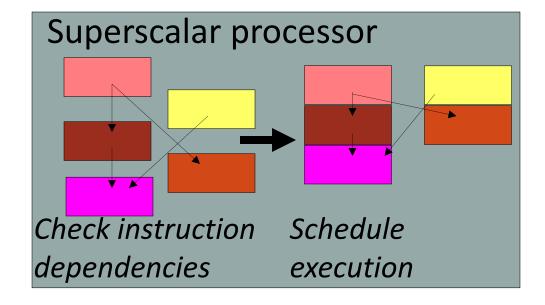
[SGI/MIPS Technologies Inc., 1995]

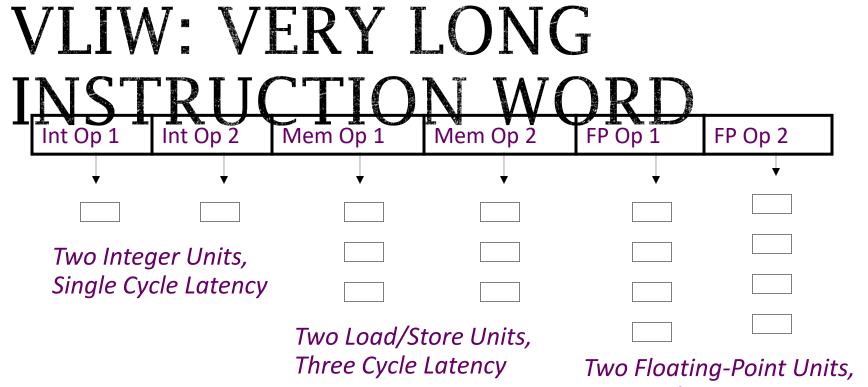


SEQUENTIAL ISA DOTTI ENIECV









- Multiple operations packed in to one instruction
- Each operation slot is for a fixed function
- Constant operation latencies are specified
- Architecture requires guarantee of:
 - Parallelism within an instruction => no cross-operation RAW check
 - No data use before data ready => no data interlocks

EARLY VLIW FPSAP120B (1976)

- scientific attached array processor
- first commercial wide instruction machine
- hand-coded vector math libraries using software pipelining and loop unrolling

• Multiflow Trace (1987)

- commercialization of ideas from Fisher's Yale group including "trace scheduling"
- available in configurations with 7, 14, or 28 operations/instruction
- 28 operations packed into a 1024-bit instruction word

Cydrome Cydra-5 (1987)

- 7 operations encoded in 256-bit instruction word
- rotating register file

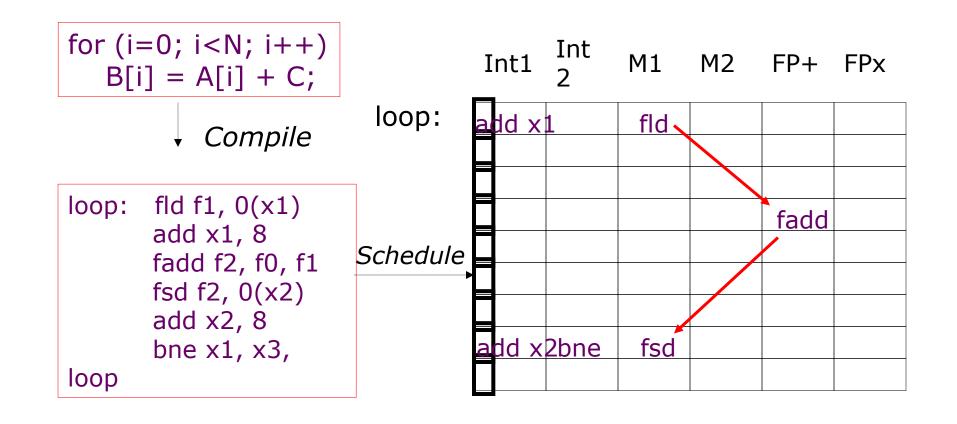
VLIW COMPILER RESPONSIBILITIES

Schedule operations to maximize parallel execution

Guarantees intra-instruction parallelism

- Schedule to avoid data hazards (no interlocks)
 - Typically separates operations with explicit NOPs

LOOP EXECUTION



How many FP ops/cycle?

1 fadd / 8 cycles = 0.125

LOOP UNROLLING

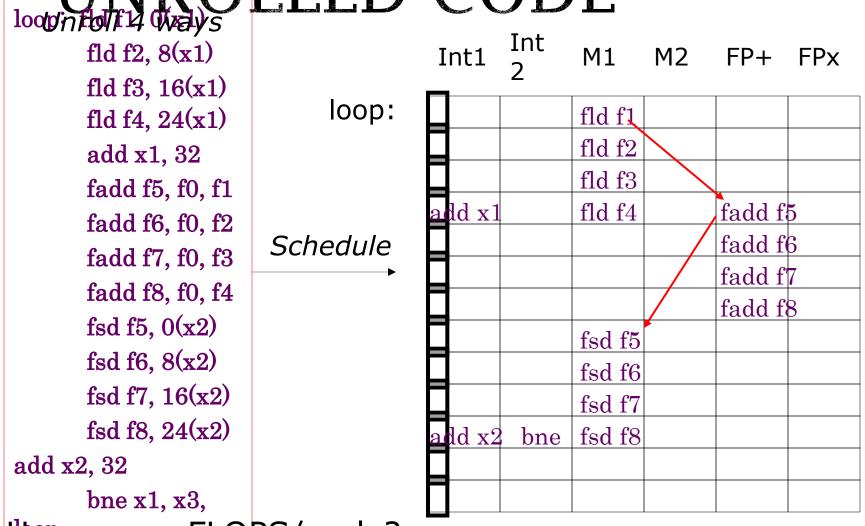
```
for (i=0; i<N; i++)
B[i] = A[i] + C;
```

Unroll inner loop to perform 4 iterations at once

```
for (i=0; i<N; i+=4)
{
    B[i] = A[i] + C;
    B[i+1] = A[i+1] + C;
    B[i+2] = A[i+2] + C;
    B[i+3] = A[i+3] + C;
}
```

Need to handle values of N that are not multiples of unrolling factor with final cleanup loop

UNROLLED CODE

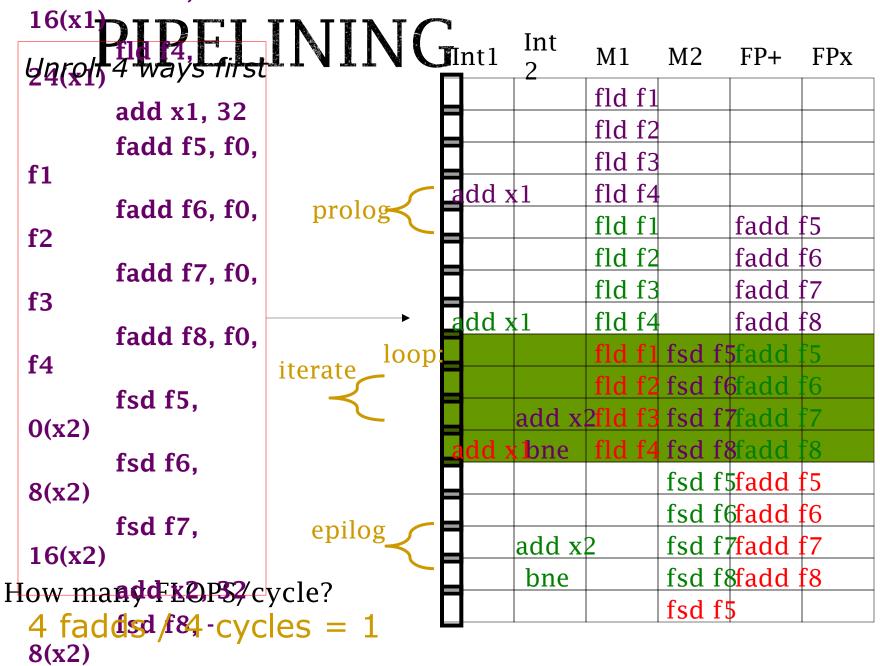


How many FLOPS/cycle?

4 fadds / 11 cycles = 0.36

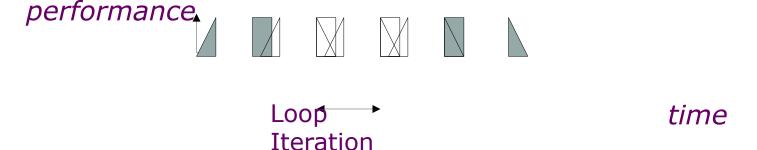
SMATIWAKE

la--a --1 --7



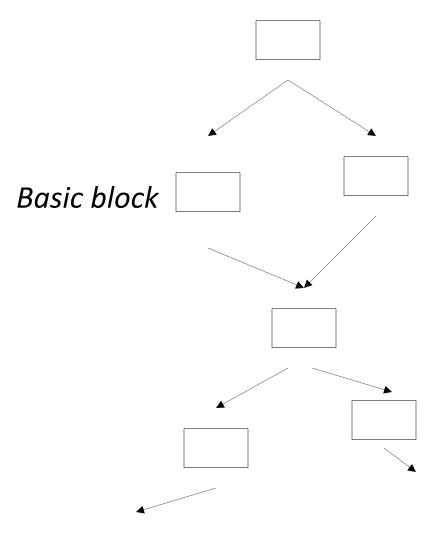
PIPELINING VS. LOOP UNROPPUNG Wind-down overhead performance Startup overhead Loop Iteration time

Software Pipelined



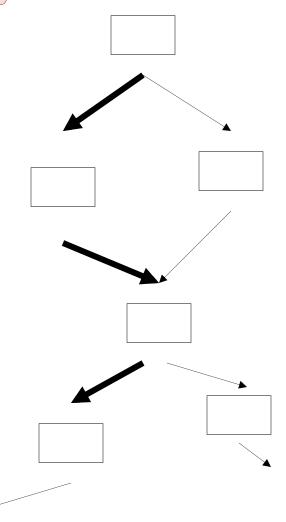
Software pipelining pays startup/wind-down costs only once per loop, not once per iteration

WHAT IF THERE ARE NO LOOPS?



- Branches limit basic block size in control-flow intensive irregular code
- Difficult to find ILP in individual basic blocks

SCHEDULING [FISHER, ELLIS]

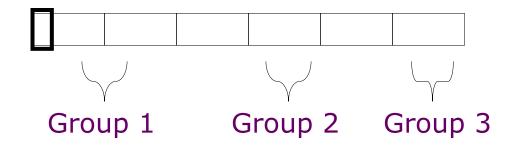


- Pick string of basic blocks, a *trace*, that represents most frequent branch path
- Use <u>profiling feedback</u> or compiler heuristics to find common branch paths
- Schedule whole "trace" at once
- Add fixup code to cope with branches jumping out of trace

"CLASSIC" VLIW

- Object-code compatibility
 - have to recompile all code for every machine, even for two machines in same generation
- Object code size
 - instruction padding wastes instruction memory/cache
 - loop unrolling/software pipelining replicates code
- Scheduling variable latency memory operations
 - caches and/or memory bank conflicts impose statically unpredictable variability
- Knowing branch probabilities
 - Profiling requires an significant extra step in build process
- Scheduling for statically unpredictable branches
 - optimal schedule varies with branch path

VLIW INSTRUCTION ENCODING



- Schemes to reduce effect of unused fields
 - Compressed format in memory, expand on I-cache refill
 - used in Multiflow Trace
 - introduces instruction addressing challenge
 - Mark parallel groups
 - used in TMS320C6x DSPs, Intel IA-64
 - Provide a single-op VLIW instruction
 - Cydra-5 UniOp instructions

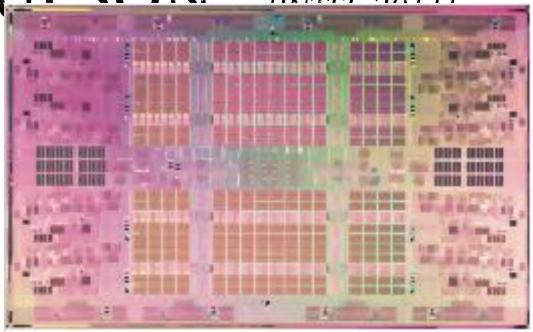
INTEL ITANIUM, EPIC IA-64

- EPIC is the style of architecture (cf. CISC, RISC)
 - Explicitly Parallel Instruction Computing (really just VLIW)
- IA-64 is Intel's chosen ISA (cf. x86, MIPS)
 - IA-64 = Intel Architecture 64-bit
 - An object-code-compatible VLIW
- Merced was first Itanium implementation (cf. 8086)
 - First customer shipment expected 1997 (actually 2001)
 - McKinley, second implementation shipped in 2002
 - Recent version, Poulson, eight cores, 32nm, announced 2011



EIGHT CORE ITANIUM

P() [CAN" [[[[]]]]]



- 8 cores
- 1-cycle 16KB L1 I&D caches
- 9-cycle 512KB L2 I-cache
- 8-cycle 256KB L2 D-cache
- 32 MB shared L3 cache
- 544mm² in 32nm CMOS
- Over 3 billion transistors

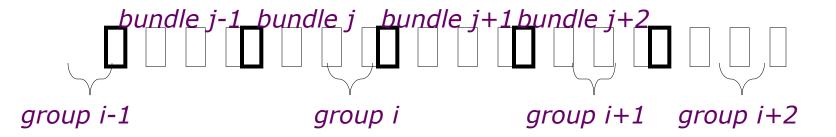
- Cores are 2-way multithreaded
- 6 instruction/cycle fetch
 - Two 128-bit bundles
- Up to 12 insts/cycle execute

IA-64 INSTRUCTION FORMAT

Instruction 2 Instruction 1 Instruction 0 Template

128-bit instruction bundle

- Template bits describe grouping of these instructions with others in adjacent bundles
- Each group contains instructions that can execute in parallel





IA-64 REGISTERS

- 128 General Purpose 64-bit Integer Registers
- 128 General Purpose 64/80-bit Floating Point Registers
- 64 1-bit Predicate Registers
- GPRs "rotate" to reduce code size for software pipelined loops
 - Rotation is a simple form of register renaming allowing one instruction to address different physical registers on each iteration

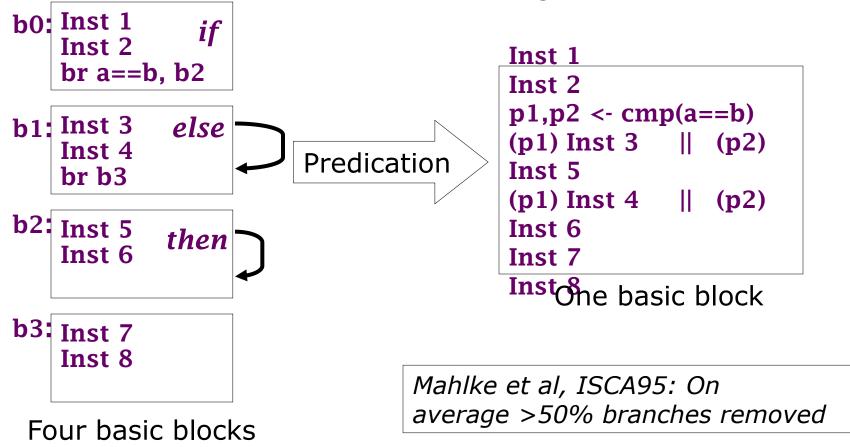


IA-04 PREDICATED

Problem: Mispredicted branches limit ILP

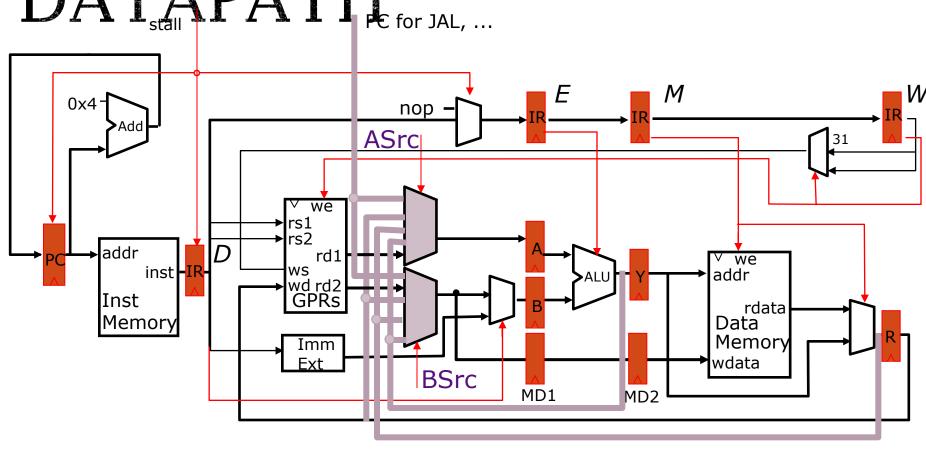
Solution: Eliminate hard to predict branches with predicated execution

- Almost all IA-64 instructions can be executed conditionally under predicate
- Instruction becomes NOP if predicate register false





FULLY BYPASSED DATAPATH for JAL, ...



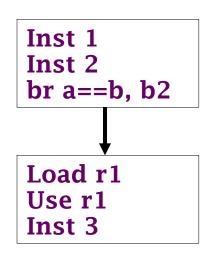
Where does predication fit in?



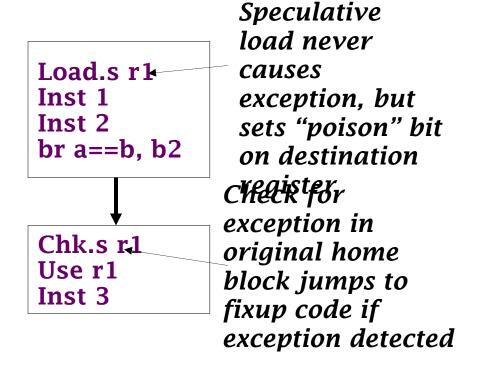
EXECUTION

Problem: Branches restrict compiler code motion

Solution: Speculative operations that don't cause exceptions



Can't move load above branch because might cause spurious exception



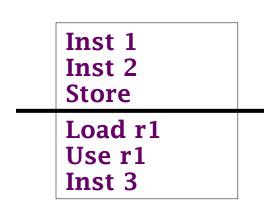
Particularly useful for scheduling long latency loads early



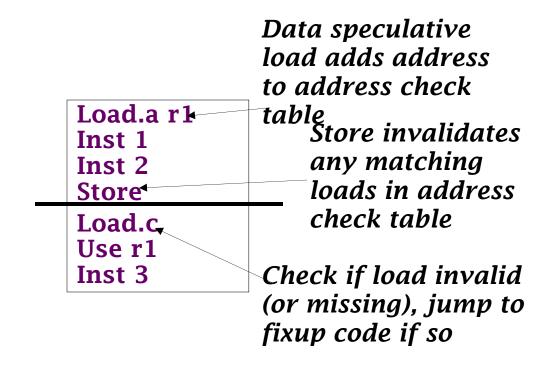
IA-04 DATA SPECULATION

Problem: Possible memory hazards limit code scheduling

Solution: Hardware to check pointer hazards



Can't move load above store because store might be to same address



Requires associative hardware in address check table



LIMITS OF STATIC

- Unprédictable branches
- Variable memory latency (unpredictable cache misses)
- Code size explosion
- Compiler complexity
- Despite several attempts, VLIW has failed in general-purpose computing arena (so far).
 - More complex VLIW architectures close to in-order superscalar in complexity, no real advantage on large complex apps.

Successful in embedded DSP market

 Simpler VLIWs with more constrained environment, friendlier code.



ACKNOWLEDGEMENTS

- These slides contain material developed and copyright by:
 - Arvind (MIT)
 - Krste Asanovic (MIT/UCB)
 - Joel Emer (Intel/MIT)
 - James Hoe (CMU)
 - John Kubiatowicz (UCB)
 - David Patterson (UCB)
- MIT material derived from course 6.823
- UCB material derived from course CS252