







| Loop-Level Parallelism | |
|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---|
| Exploit loop-level parallelism by "unrolling loop" (which increases BB size) either by Static methods via loop unrolling by compiler Dynamic methods via branch prediction | ١ |
| • Determining instruction dependence is critical to Loop Level Parallelism | |
| If 2 instructions are <u>Parallel</u>, they can execute simultaneously in a pipeline of arbitrary depth without causing any stalls (assuming no structural hazards) <u>Dependent</u>, they are not parallel and must be executed in order, although they may often be partially overlapped | |
| 5 | |



<section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item>













| | | | | - . | |
|---------|-------------|-----------------------------------|--------------|------------------|--|
| 1 Loop: | LD | F0,0(R1) | ;F0=vector e | lement | |
| 2 | stall | | | | |
| 3 | ADDD | F4 , F0 , F2 | ;add scalar | in F2 | |
| 4 | stall | | | | |
| 5 | stall | | | | |
| 6 | SD | 0(R1), F4 | ;store resul | t | |
| 7 | SUBI | R1,R1,8 | ;decrement p | ointer 8B (DW) | |
| 8 | BNEZ | R1,Loop | ;branch R1!= | zero | |
| 9 | stall | | ;delayed bra | nch slot | |
| Instr | uction | Instru | ction | Stall latency in | |
| ED AL | cing result | using i | result | clock cycles | |
| FP AL | Uon | Store | double | 2 | |
| Load | double | FP AL | J op | 1 | |
| Load | double | Store | double | 0 | |
| Integ | er op | Intege | r op | 0 | |



Unroll Loop Four Times (straightforward way)

| 2 100p. | ADDD | F4, F0, F2 | 2 cycles stall | |
|---------|----------|----------------------------|--------------------------|------------|
| 3 | SD | 0(R1),F4 | drop SUBI & BNEZ; | |
| 4 | LD | F6, <mark>-8</mark> (R1) | | |
| 5 | ADDD | F8,F6,F2 | | |
| 6 | SD | -8(R1),F8 | drop SUBI & BNEZ; | |
| 7 | LD | F10, <mark>-16</mark> (R1) | | |
| 8 | ADDD | F12,F10,F2 | | |
| 9 | SD | -16(R1),F12 | drop SUBI & BNEZ; | Rewrite |
| 10 | LD | F14,- <mark>24</mark> (R1) | | |
| 11 | ADDD | F16,F14,F2 | | ισορ το |
| 12 | SD | -24(R1),F16 | | minimize |
| 13 | SUBI | R1,R1, <mark>#32</mark> | ;alter to 4*8 | at all a 2 |
| 14 | BNEZ | R1,LOOP | | stalls ? |
| 15 | NOP | | | |
| 15 + | 4 x (1+2 |) = 27 clock cycle | es, or 6.8 per iteration | |
| | | Assumes R1 is mu | ultiple of 4 | 16 |







Where are the name dependencies? F0,0(R1) 1 Loop: LD 2 ADDD F4,F0,F2 3 SD 0(R1),F4 ;drop SUBI & BNEZ 4 LD F0,-8(R1) 5 ADDD F4,F0,F2 6 SD -8(R1),F4 ;drop SUBI & BNEZ 7 LD F0, -16(R1)8 ADDD F4,F0,F2 9 SD -16(R1), F4;drop SUBI & BNEZ 10 LD F0,-24(R1) 11 ADDD F4,F0,F2 12 SD -24 (R1), F4 R1,R1,<mark>#32</mark> 13 SUBI ;alter to 4*8 14 BNEZ R1,LOOP 15 NOP How can we remove them?

Where are the name dependencies? 1 Loop: LD F0,0(R1) 2 ADDD F4,F0,F2 3 SD 0(R1),F4 ;drop SUBI & BNEZ 4 LD F6,-8(R1) 5 ADDD F8, F6, F2 6 SD -8(R1),F8 ;drop SUBI & BNEZ 7 LD F10, -16(R1)8 ADDD F12,F10,F2 9 SD -16(R1),F12 ;drop SUBI & BNEZ F14, -24(R1)10 LD 11 ADDD F16,F14,F2 12 SD -24 (R1), F16 13 SUBI R1,R1,#32 ;alter to 4*8 14 BNEZ R1,LOOP 15 NOP Called "register renaming" 21





Five (5) Loop Unrolling Decisions

• Requires understanding how one instruction depends on another and how the instructions can be changed or reordered given the dependences.

• Main steps:

- 1. Determine loop unrolling useful by finding that loop iterations were independent
- 2. Use different registers to avoid unnecessary constraints forced by using same registers for different computations
- 3. Eliminate the extra test and branch instructions and adjust the loop termination and iteration code
- 4. Determine that loads and stores in unrolled loop can be interchanged by observing that loads and stores from different iterations are independent
 - Transformation requires analyzing memory addresses and finding that they do not refer to the same address
- 5. Schedule the code, preserving any dependences needed to yield the same result as the original code 24









<section-header><text><text><section-header><text><section-header><text>

















Advanced Techniques: 1. Correlated Branch Prediction (2-level Predictors)

- Hypothesis: recent branches are correlated; that is, behavior of recently executed branches affects prediction of current branch
 Keep track of the behavior of previous branches, and use that to
- Keep track of the behavior of previous branches, and use that predict the behavior of the current branch
- Idea: record *m* most recently executed branches (represents a path through the program) as taken or not taken, and use that pattern to select the proper *n*-bit branch history table
- In general, (*m*,*n*) predictor means record last *m* branches to select between 2^m history tables, each with *n*-bit counters ° Thus, old 2-bit BHT is a (0,2) predictor
- Global Branch History: *m*-bit shift register keeping T/NT status of last *m* branches





















Dynamic Branch Prediction Summary

- Prediction becoming important part of execution
- Branch History Table: 2 bits for loop accuracy
- Correlation: recently executed branches correlated with next branch
 - ° Either different branches,
 - $^{\circ}$ Or different executions of same branches
- Tournament predictors take insight to next level, by using multiple predictors
 - Usually one based on global information and one based on local information, and combining them with a selector
 - $^\circ\,$ Tournament predictors using \approx 30K bits were in processors like the Power5 and Pentium 4 (circa 2006)





- Dynamic scheduling Hardware rearranges the instruction execution to reduce stalls while maintaining data flow and exception behavior
- It handles cases when dependences unknown at compile time
- It allows the processor to tolerate unpredictable delays such as cache misses, by executing other code while waiting for the miss to resolve
- It allows code that is compiled for one pipeline to run efficiently on a different pipeline
- It simplifies the compiler
- Hardware speculation Technique with significant performance advantages, builds on dynamic scheduling 50

| HW Schemes: Instruction Parallelism |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Key idea: Allow instructions behind stall to proceed |
| DIVD F0,F2,F4 ADDD F10,F0,F8 SUBD F12,F8,F14 |
| Enables out-of-order execution and allows out-of-order completion (e.g., SUBD) |
| In a dynamically scheduled pipeline, all instructions still pass through issue stage in order (in-order issue) |
| Will distinguish when an instruction begins execution and when it completes execution; between 2 times, the instruction is in execution |
| 51 |











<section-header> basic Functions of Some Elements Load buffers have 3 functions: hold components of effective address until it is computed track outstanding loads that are waiting on the memory hold the results of the completed loads that are waiting for the CDB Store buffers have 3 functions: hold components of effective address until it is computed hold destination memory addresses of outstanding stores that are waiting for the data value to store until memory unit is available Hold address and value to store until memory unit is available he reservation stations, and the store buffers















































Distribution of the hazard detection logic Distributed reservation stations and the CDB If multiple instructions waiting on single result, & each instruction has other operand, then instructions can be released simultaneously by broadcast on CDB If a centralized register file were used, the units would have to read their results from the registers when register buses are available **2. Elimination of stalls for WAW and WAR hazards**









