



BSC~ Microsoft Research
Centre

ILP: Instruction Level Parallelism

Tassadaq Hussain

Riphah International University

Barcelona Supercomputing Center
Universitat Politècnica de Catalunya

Consultancy for:
FYPs and Future Career Guidance.
Engineering Workshops, Master
and Ph.D. thesis.
Design and Develop Industrial
Digital Systems. www.ucerd.com



Introduction

- Pipelining become universal technique in 1985
 - Overlaps execution of instructions
 - Exploits “Instruction Level Parallelism”
- Beyond this, there are two main approaches:
 - Hardware-based dynamic approaches
 - Used in server and desktop processors
 - Not used as extensively in Parallel Microprogrammed Processors (PMP)
 - Compiler-based static approaches
 - Not as successful outside of scientific applications

Basic Block

The amount of parallelism available within a basic block—a straight-line code sequence with no branches in except to the entry and no branches out except at the exit—is quite small.

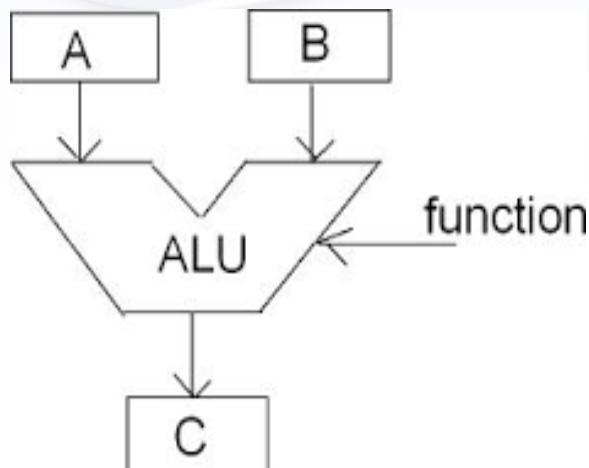
Basic Block

The amount of parallelism available within a basic block—a straight-line code sequence with no branches in except to the entry and no branches out except at the exit—is quite small.

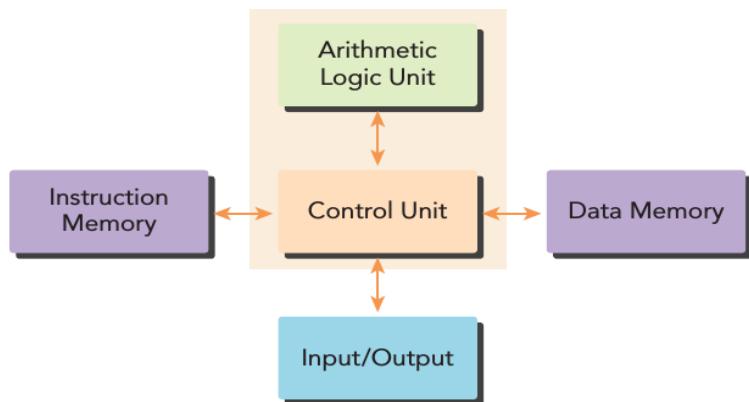
Instruction-Level Parallelism

- When exploiting instruction-level parallelism, goal is to maximize CPI
 - Pipeline CPI =
 - Ideal pipeline CPI +
 - Structural stalls +
 - Data hazard stalls +
 - Control stalls
- Parallelism with basic block is limited
 - Typical size of basic block = 3-6 instructions
 - Must optimize across branches

Basic introduction of Microprocessor



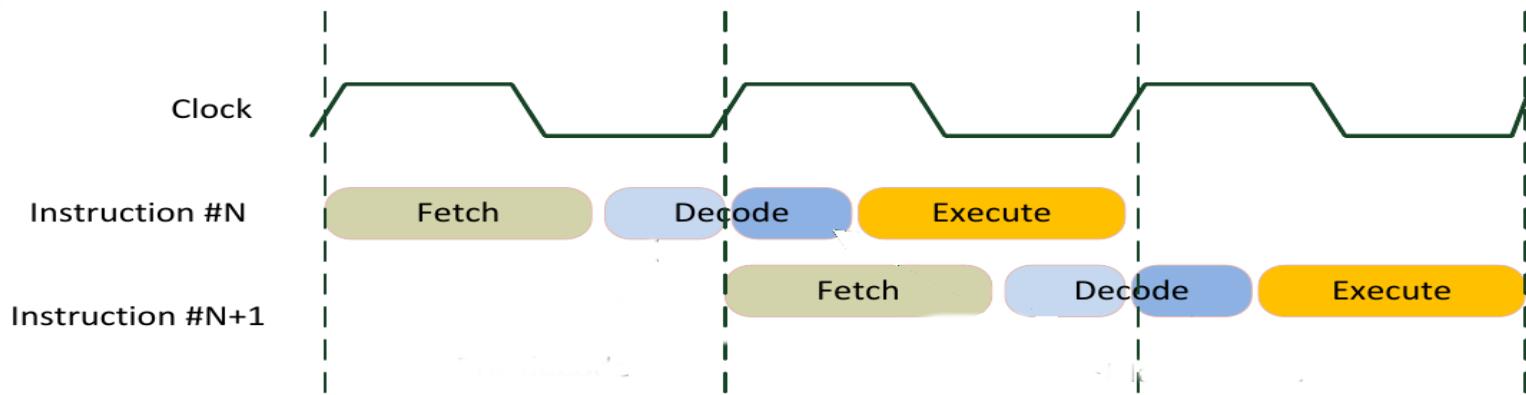
Function	Inputs A	Input B	Output
0	A	B	$A+B$
1	A	B	$A-B$



Processor: Performance Improvement

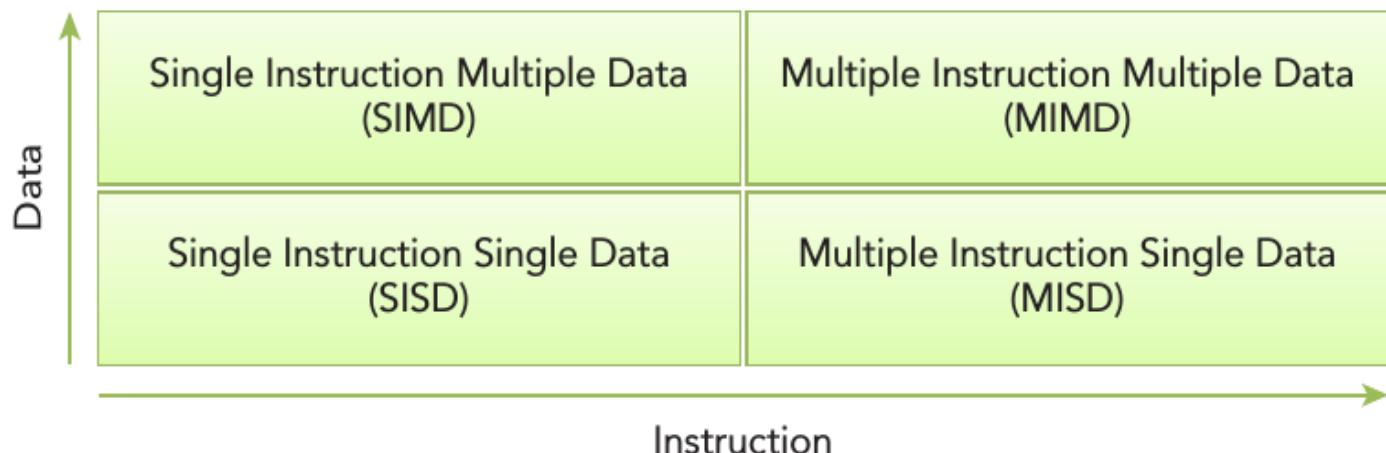
Instruction Level Parallelism
Task Level Parallelism
Data Level Parallelism

Speed and Performance



Processor Architectures

- Single Instruction Single Data (SISD)
- Single Instruction Multiple Data (SIMD)
- Multiple Instruction Single Data (MISD)
- Multiple Instruction Multiple Data (MIMD)



Data Dependence

- Loop-Level Parallelism
 - Unroll loop statically or dynamically
 - Use SIMD (vector processors and GPUs)
- Challenges:
 - Data dependency
 - Instruction j is data dependent on instruction i if
 - Instruction i produces a result that may be used by instruction j
 - Instruction j is data dependent on instruction k and instruction k is data dependent on instruction i
 - Dependent instructions cannot be executed simultaneously

Data Dependence

- Dependencies are a property of programs
- Pipeline organization determines if dependence is detected and if it causes a stall

- Data dependence conveys:
 - Possibility of a hazard
 - Order in which results must be calculated
 - Upper bound on exploitable instruction level parallelism

- Dependencies that flow through memory locations are difficult to detect

Name Dependence

A name dependence occurs when two instructions use the same register or memory location.

- Two instructions use the same name but no flow of information
 - Not a true data dependence, *but is a problem when reordering instructions*
 - **Antidependence:** instruction j writes a register or memory location that instruction i reads
 - Initial ordering (i before j) must be preserved
 - **Output dependence:** instruction i and instruction j write the same register or memory location
 - Ordering must be preserved
- To resolve, use renaming techniques

Other Factors

- Data Hazards
 - Read after write (RAW)
 - Write after write (WAW)
 - Write after read (WAR)
- Control Dependence
 - Ordering of instruction i with respect to a branch instruction
 - Instruction control dependent on a branch cannot be moved before the branch so that its execution is no longer controlled by the branch
 - An instruction not control dependent on a branch cannot be moved after the branch so that its execution is controlled by the branch

Examples

- Example 1:

DADDU R1,R2,R3
BEQZ R4,L
DSUBU R1,R1,R6

L: ...

OR R7,R1,R8

- Example 2:

DADDU R1,R2,R3
BEQZ R12,skip
DSUBU R4,R5,R6
DADDU R5,R4,R9

skip:

OR R7,R8,R9

- OR instruction dependent on DADDU and DSUBU

- Assume R4 isn't used after skip
 - Possible to move DSUBU before the branch

Compiler Techniques for Exposing ILP

- Pipeline scheduling
 - Separate dependent instruction from the source instruction by the pipeline latency of the source instruction

- Example:

```
for (i=999; i>=0; i=i-1)  
    x[i] = x[i] + s;
```

Instruction producing result	Instruction using result	Latency in clock cycles
FP ALU op	Another FP ALU op	3
FP ALU op	Store double	2
Load double	FP ALU op	1
Load double	Store double	0

Pipeline Stalls

Loop:

```
L.D F0,0(R1)
stall
ADD.D F4,F0,F2
stall
stall
S.D F4,0(R1)
DADDUI R1,R1,#-8
stall (assume integer load latency is 1)
BNE R1,R2,Loop
```

Instruction producing result	Instruction using result	Latency in clock cycles
FP ALU op	Another FP ALU op	3
FP ALU op	Store double	2
Load double	FP ALU op	1
Load double	Store double	0

Pipeline Scheduling

Scheduled code:

Loop: L.D F0,0(R1)
DADDUI R1,R1,#-8
ADD.D F4,F0,F2
stall
stall
S.D F4,8(R1)
BNE R1,R2,Loop

Instruction producing result	Instruction using result	Latency in clock cycles
FP ALU op	Another FP ALU op	3
FP ALU op	Store double	2
Load double	FP ALU op	1
Load double	Store double	0

Loop Unrolling

- Loop unrolling
 - Unroll by a factor of 4 (assume # elements is divisible by 4)
 - Eliminate unnecessary instructions

Loop:

```
L.D F0,0(R1)
ADD.D F4,F0,F2
S.D F4,0(R1) ;drop DADDUI & BNE
L.D F6,-8(R1)
ADD.D F8,F6,F2
S.D F8,-8(R1) ;drop DADDUI & BNE
L.D F10,-16(R1)
ADD.D F12,F10,F2
S.D F12,-16(R1) ;drop DADDUI & BNE
L.D F14,-24(R1)
ADD.D F16,F14,F2
S.D F16,-24(R1)
DADDUI R1,R1,#-32
BNE R1,R2,Loop
```

- note: number of live registers vs. original loop

Loop Unrolling/Pipeline Scheduling

- Pipeline schedule the unrolled loop:

Loop:

- L.D F0,0(R1)
- L.D F6,-8(R1)
- L.D F10,-16(R1)
- L.D F14,-24(R1)
- ADD.D F4,F0,F2
- ADD.D F8,F6,F2
- ADD.D F12,F10,F2
- ADD.D F16,F14,F2
- S.D F4,0(R1)
- S.D F8,-8(R1)
- DADDUI R1,R1,#-32
- S.D F12,16(R1)
- S.D F16,8(R1)
- BNE R1,R2,Loop

Strip Mining

- Unknown number of loop iterations?
 - Number of iterations = n
 - Goal: make k copies of the loop body
 - Generate pair of loops:
 - First executes $n \bmod k$ times
 - Second executes n / k times
 - “Strip mining”

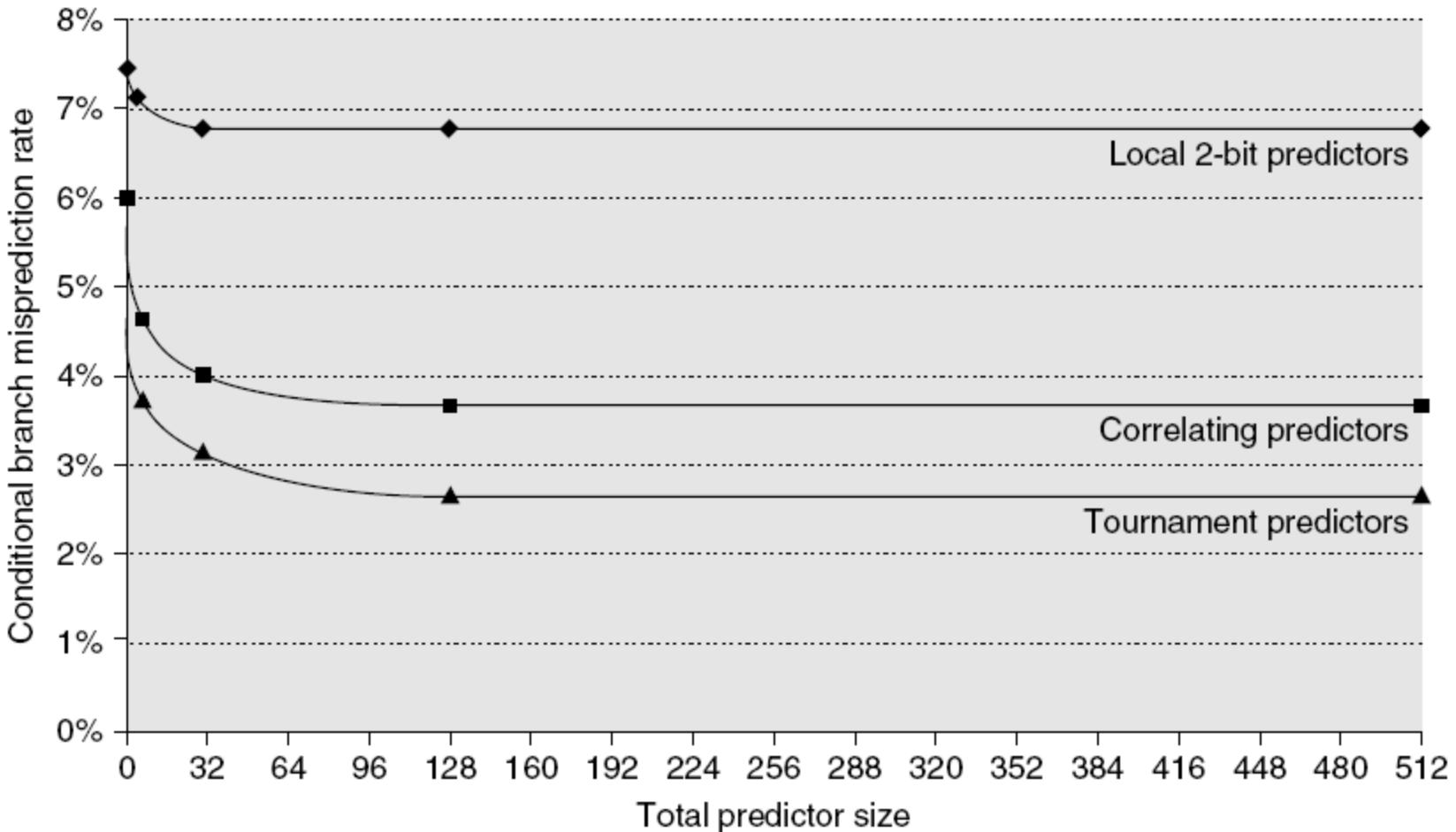
5 Loop Unrolling Decisions

- Requires understanding how one instruction depends on another and how the instructions can be changed or reordered given the dependences:
 1. Determine loop unrolling useful by finding that loop iterations were independent (except for maintenance code)
 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

Branch Prediction

- Basic 2-bit predictor:
 - For each branch:
 - Predict taken or not taken
 - If the prediction is wrong two consecutive times, change prediction
- Correlating predictor:
 - Multiple 2-bit predictors for each branch
 - One for each possible combination of outcomes of preceding n branches
- Local predictor:
 - Multiple 2-bit predictors for each branch
 - One for each possible combination of outcomes for the last n occurrences of this branch
- Tournament predictor:
 - Combine correlating predictor with local predictor

Branch Prediction Performance



Dynamic Scheduling

- Rearrange order of instructions to reduce stalls while maintaining data flow
- Advantages:
 - Compiler doesn't need to have knowledge of microarchitecture
 - Handles cases where dependencies are unknown at compile time
- Disadvantage:
 - Substantial increase in hardware complexity
 - Complicates exceptions

Dynamic Scheduling

- Dynamic scheduling implies:
 - Out-of-order execution
 - Out-of-order completion
- Creates the possibility for WAR and WAW hazards
- Tomasulo's Approach
 - Tracks when operands are available
 - Introduces register renaming in hardware
 - Minimizes WAW and WAR hazards

Register Renaming

- Example:

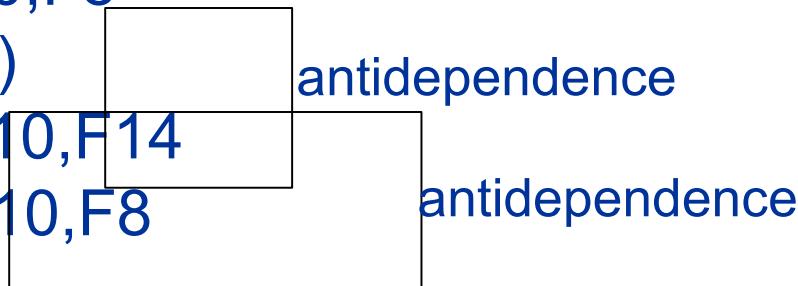
DIV.D F0,F2,F4

ADD.D F6,F0,F8

S.D F6,0(R1)

SUB.D F8,F10,F14

MUL.D F6,F10,F8



+ name dependence with F6

Register Renaming

- Example:

DIV.D F0,F2,F4

ADD.D S,F0,F8

S.D S,0(R1)

SUB.D T,F10,F14

MUL.D F6,F10,T

- Now only RAW hazards remain, which can be strictly ordered

Register Renaming

- Register renaming is provided by reservation stations (RS)
 - Contains:
 - The instruction
 - Buffered operand values (when available)
 - Reservation station number of instruction providing the operand values
 - RS fetches and buffers an operand as soon as it becomes available (not necessarily involving register file)
 - Pending instructions designate the RS to which they will send their output
 - Result values broadcast on a result bus, called the common data bus (CDB)
 - Only the last output updates the register file
 - As instructions are issued, the register specifiers are renamed with the reservation station
 - May be more reservation stations than registers

Hardware-Based Speculation

- Execute instructions along predicted execution paths but only commit the results if prediction was correct
- Instruction commit: allowing an instruction to update the register file when instruction is no longer speculative
- Need an additional piece of hardware to prevent any irrevocable action until an instruction commits
 - I.e. updating state or taking an execution

Reorder Buffer

- Reorder buffer – holds the result of instruction between completion and commit
- Four fields:
 - Instruction type: branch/store/register
 - Destination field: register number
 - Value field: output value
 - Ready field: completed execution?
- Modify reservation stations:
 - Operand source is now reorder buffer instead of functional unit

Reorder Buffer

- Register values and memory values are not written until an instruction commits
- On misprediction:
 - Speculated entries in ROB are cleared
- Exceptions:
 - Not recognized until it is ready to commit

Multiple Issue and Static Scheduling

- To achieve CPI < 1, need to complete multiple instructions per clock
- Solutions:
 - Statically scheduled superscalar processors
 - VLIW (very long instruction word) processors
 - dynamically scheduled superscalar processors

Multiple Issue

Common name	Issue structure	Hazard detection	Scheduling	Distinguishing characteristic	Examples
Superscalar (static)	Dynamic	Hardware	Static	In-order execution	Mostly in the embedded space: MIPS and ARM, including the ARM Coretex A8
Superscalar (dynamic)	Dynamic	Hardware	Dynamic	Some out-of-order execution, but no speculation	None at the present
Superscalar (speculative)	Dynamic	Hardware	Dynamic with speculation	Out-of-order execution with speculation	Intel Core i3, i5, i7; AMD Phenom; IBM Power 7
VLIW/LIW	Static	Primarily software	Static	All hazards determined and indicated by compiler (often implicitly)	Most examples are in signal processing, such as the TI C6x
EPIC	Primarily static	Primarily software	Mostly static	All hazards determined and indicated explicitly by the compiler	Itanium

VLIW Processors

- Package multiple operations into one instruction
- Example VLIW processor:
 - One integer instruction (or branch)
 - Two independent floating-point operations
 - Two independent memory references
- Must be enough parallelism in code to fill the available slots

VLIW Processors

- Disadvantages:
 - Statically finding parallelism
 - Code size
 - No hazard detection hardware
 - Binary code compatibility

Dynamic Scheduling, Multiple Issue, and Speculation

- Modern microarchitectures:
 - Dynamic scheduling + multiple issue + speculation
- Two approaches:
 - Assign reservation stations and update pipeline control table in half clock cycles
 - Only supports 2 instructions/clock
 - Design logic to handle any possible dependencies between the instructions
 - Hybrid approaches
- Issue logic can become bottleneck