# Instruction level parallelism

ILP:- overlapping inst during the execution.

There are 2 approaches to exploit ILP

1) An approach that relies on h/w to help discover & exploit parallelism dynamically

2) An approach that relies on software technology to find parallelism Statistically during compile time.

many Techniques are used to exploit parallelism among instructions in different blocks. The block may contain straight line codes, many ALU & I/O instructions. usually 15%-25%.

35-30% of code in a pgm has branch instructions.

Since these inst are likely to depend on each other. the amount of parallel overlap we can exploit within a basic block is likely to be less than the avg block size. To obtain Substantial performance enhancement, we must exploit ILP across multiple basic blocks.

The Simple way to exploit parallelism is iterations in a loop.

Eg    for (i=0; i<1000, i++)

      x[i] = x[i] + y[i]

This is also known as loop level parallelism.

Every iteration can overlap with any other iteration, although within each loop iteration there is little or no opportunity for overlap.

There are different techniques are used to convert loop level parallelism into instruction level parallelism.

# Data dependence and Hazards

Determining how one instruction depends on another is critical to determining how much parallelism exists in a program and how much parallelism can be exploited in a pgm.
In particular, to exploit instruction-level, we have to determine which instruction can be executed in parallel.

If two inst are in parallel, they can executed & simultaneously without causing stalls, then they are independent inst.

If two inst have to be executed, then it is not possible to executed completely overlapping 2 inst, then they are dependent inst

Dependents inst are executed in order by partially overlapping

There are diffrent types of dependence among the instructions

---

## 1) Data dependence:-

An inst j is data dependence on insti such that.

1) inst i produces the result that can be used by inst j

2) inst k is dependent on j & j is dependent on i.

known as chain of dependences. The dependence chain can be as long as the entire program.

For eg. consider the mips code

```
LOOP:        LD      R0, O(R1)
             ADD     R2, R0, R3
             SD      R2, O(R1)
             DADDUI  R1, R1, #-8
             BNE     R1, R2, Loop.
```

There is data dependence
b/w  LD    R0, O(R1)
     ADD   R2, R0, R3
     SD    RD, O(R1)

Data dependence

nce dependence must be preserved for the correct execution If 2 inst are dependent, they must be executed in order & can't execute simultaneously

Executing simultaneously causes pipeline interlock to detect hazard & stall, there by reducing the overlap.

Dependence are the properties of programs, whether a dependence <u>results in actual hazard being detected</u> or whether that hazard actually <u>causes a stall</u>, these are the properties of the pipeline organization.

A data dependence convey 3 things

1) The possibility of a hazard 2) The order in which results must be calculated. 3) The upper bound on how much parallelism can be exploited.

A dependence can limit the amount of ILP. A data dependence can over come 2 ways

1) maintaining the dependence but avoiding a <u>hazard</u>

2) <u>Eliminating</u> the dependence by transforming the code

### Name dependence :-

A name dependence occurs when two inst uses same register or memory location called a name, but there is no flow of data b/w those inst associated with that name.

There are 2 types of name dependence between i that preced inst j.

i) Antidependence :- b/w inst i and j occurs when j write a regr or memory that inst i reads, the original order must be preserved. eg    Sd R1, 0(R2)

DADDUI R2, R2 #8.

2) 'output dependence:- when both inst tries to write to a same regr/memory, the ordering must be preserved such that first value gets updated by i then j.

They are juust name dependencies not opposed to true data dependences, since there is no value being transmitted b/w the instructions. So there inst can be reordered & even be executed parallely.

## Data Hazards :-

A Exists whenever there is data dependence b/w two instructions, they are so close that changing the order ourlapping the inst may change the order of accessing the operands involved in the dependence.

program order must be maintained while preserving the dependence. Executing the inst in the order which it was executed sequentially.

The main goal is
    i) Preserve the pgming order, if it's output affects the result of the program.

2) maintaing the data dependence.
3) Detecting and avoiding the hazards & stalls.

The possible data hazards are.

1) RAW:- inst j tries to read the operand which was before.
inst i writes it.

       ADD R1, R2, R3
       ~~id~~ ~~R4,~~ ~~10~~ ~~R~~ SUB R4, R1, R5.
pgm order shld be preserved such that sub reads the updated data

WAW:- inst j waits tries to write an operand, where inst i is writing. The write end up being performed in the wrong order, leaving the value i instead of j in the destination. WAW occurs in pipelines that write in more than one pipe stage.

3) WAR:- (write after read):- inst j tries to write an operand before inst i reads it. So i incorrectly gets a new value. It occurs in pipeline where it performs write in pipestage early before other inst reads it.

RAR is not a hazard

## Control dependence:-

It defines ordering of inst i related to branch inst so that inst i should be executed in correct order and when it should be. Every inst, except the first basic block of the pgm is control dependent.

inst under the branch should should be executed in order and branch dependence should be preserved.

if P1 < S1;   S1 is control dependent on P1
         y

if P2 < S2;y   S2 is control dependent on P2.

There are 2 constraints imposed by control dependence

i) inst i inside the branch can not be put before the branch, so that it's pgming order is violated below

ii) Any inst above the branch should not be executed inside the branch, because those inst are not control dependent.

when processor executes these inst, compiler it shld make sure that program is executing in the order and also control dependence must be performed.

but if changing the order of the program doesn't the o/p of the program then we can compromise. for the control dependence.

The main critical properties to maintain control dependence are i) data flow   ii) Exceptions

Preserving the exceptions doesn't means that during the Execution inst shld not cause new exceptions

Eg
```
        DADDU  R1, R2, R3
        BEQZ   R1, L1
        LD     R4, 10(R1)

L1:
```
here even though, there is no data depence that exist directly since after ADD R1 value is not changed but putting LD before BEQZ cause memory protection violation exception. here kd shld exelē only iff

it is only the name dependence, branch is not taken.

because to preserve control dependence, we shld not put LD before

BEQZ.

here. BEQ is testing R4 when R4 isn't at all used by any inst.

2) Data flow:

```
        SUB  R1, R2, R3

              BEQZ R4, L1
        ADD  R1, R6, R5

L1:   .'_ _ _ _ _ _

LP:   _ '_ _ _ _ _ .

        OR.  R7, R1, R3
```

& but OR inst is data dependent on both SUB & ADD i.e if branch takes place, then OR shld hold the R1 of SUB else the result produced by ADD.

So here the control depend is note the main criteria but preserving data depender and pgming order is importan

ADD R1, R2, R3
Sub R4, R1, R2
BEQZR R, loop
Sub R5, R6, R7
loop: _
loop ADD R7, R8, R9

if Branch is not taken then
Sub can be exelē
before branch but Speculation this is required

Scanned with CamScanner

Basic compiler Techniques for exploiting ILP.

These Techniques can be used to Schedul instructions to exploit parallelism.

Basic Pipelining Scheduling and loop unrolling :-

To keep the pipelining full, parallelism among inst can be obtained by executing the sequence of unrelated inst that can be overlapped in the pipeline.

To avoid pipeline stall, dependent inst must be executed separated from the source inst by a distance equal to the pipeline latency of that source instruction.

The compiler ability to schedule the instructions depends on the amount of ILP available in the program and on the latency of the functional units in the pipeline.

Eg consider for (i=9999; i>=0; i--)

$$X[i] = X[i] + S;$$

Each iteration is independent but dependence exist inside the instruction itself. To see how this iteration execute first we will connvert this into assembly level language.

        Loop: MOV   F4, 0(R1)      highest addree
        ADD   F4, F4, R2      X[i]+S
        SD    F4 0(R1)
        BNE   R3
        DADDUI R1, R1, #-8.
        BNE   R1, R2, Loop      lowest element addree

The unscheduled loop with stalls is given below.
For this we consider the latency of FP operations.

| Inst producing result | Inst using result | latency in CC |
|---|---|---|
| FP ALU OP | FP ALU | 3 |
| FP ALU OP | store | 2 |
| load | FP ALU operation | 1 |
| Load | Store. | 0 |

```
1   LD    F4, 0(R1)
2        Stall
3   ADD.D  F4, F4, R3
4        Stall
5        Stall
6   S.D.   F4, 0(R1)
         Stall
7   DADDUI R F4, F4, #-8.
8   Stall
9   BNE   R1, R2.
```

This will take 7 cc for 1 iteration
for $4 = 18 \times 4 = 22$. Each element
is stored in 7 cc but actual work
LD, DADD, SD laty 3 cc.

So DADDUI & BNE are the staly
are overhead.

load 3
add 3  } 2
store (3)   (4) cc is extra

only 3 is enough.

---

**Scheduled loop**

```
LD   F4, 0(R1)      — 1
DADDUI R1, R1, #-8. 2
ADD.D  F4, F4, R3   —3
       Stall        —4
       stall        —5
SID  F4, 0(R1)      — 6
BNE  R1, R2         —7
```

Here we can complete
within 7 cc.

but 2 stalls &
DADD is a overhead
&
BNE

Another scheme to increase the no of instructions ⑤
relative to branch & overhead mgt is loop unrolling.
Unrolling simply replicates the loop body multiple times, by
adjusting the loop termination code correctly.

Loop unrolling is used to improve the scheduling, It
allows multiple iterations to be Scheduled together.
In this case, we can eliminate the data stall by creating
additional independent instructions. To create successfully
these addition mgt is not possible if we use the same
registers, because it may also affect the performance of the
Scheduling ∴ different registers can used for each iteration.

unrolled loop for the previous ex

        LD   F0, 0(R1)
        ADD.D  F0, F0, R2.
        S.D   F0, 0(R1).
    ——  LD   F1, @@R -8(R1)
        ADD.D  F1, F1, R3
        SD   F1, -8(R1)
        ————————————————
        LD   F2, 16(R1)
        ADDID  F2, F2, R4
        S.D   F2, -16(R1)
        ————————————————
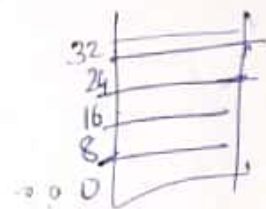        LD   F3, -24(R1)
        ADD  F9, F3, R5
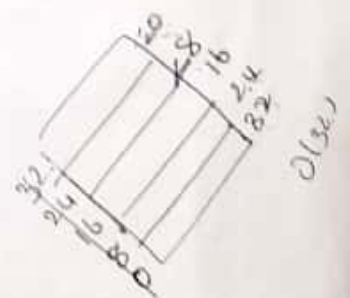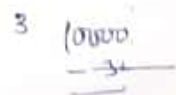        SD   F3, -24(R1)
        ————————————————
        DADDUI  R1, R1, #0-2
        BNE  R1, R2, loop.

here the unrolled loop without Schedule is

```
Loop   ┌LD    F0, 0(R1)
    1  ┤ADD   F1, F0, F2
    2  └SD    F1, 0(R1)

       LD    F3, -8(R1)
       ADD   F4, F3, F5
       SD    F4, -8(R1)

       LD    F6, -16(R1)
       ADD   F7, F6, F8
       SD    F7 -16(R1)

       LD    F9, -24(R1)
       ADD   F10, F9, F11
       SD    F10, -24(R1)

       DADDUI R1, R1, #32
    1  ─BNE   R1, R2, loop
```

LD — ADD — SD
  1       2

$3 \times 4 = 12$ stalls inside loop
+ 14 instruction
+ 1 ...... branch
─────────
27 CC is required

loop iteration i.e. Termination code
is adjusted with ld & sd &
different regrs are used to avoid
hazards

It optimizes the code by using
Substitution and Simplifications

If Actually we do not know the
upper bound on the array; i.e it
is difficult to known; it can
known during the compilation time

Suppose if the size is N & we should unroll k copies of each
iteration, then instead of Single large unrolling rod, Sim generate
a pair of consecutive loops first loop executes $n/k$ time
the second unrolled loop executes $n/k$ times. For large value of
H most of execution time is spent per loop unrolling

unrolling loops improves the performance of loop over pipeling
b eventhough code size is large.

```
loop   LD F0, 0(A)          SD F1, 0(R1)
       LD F3, -8(R1)        SD F4, -8(R1)
       LD F6, -16(R1)       SD F7, -16(R1)
       LD F9, -24(R1)       SD F10, -24(R1)
       ADD F1, F0, F2       DADDUI R1, R1, #32
       ADD F4, F3, F5       SD F10, -24(R1)
       ADD F7, F6, F8       BNE R1, R2, loop
       ADD F10, F9, F11
```

This loop will
work without
any stalls.

ere the unrolled loop without Schedule is

```
Loop
  1 ⌐LD    F0, 0(R1)           LD ─ ADD ─ SD      3×4 =12 Stalls inside loop
    ⌐ADD   F1, F0, F2              1     2            + 14 instruction
  2 ⌐SD    F1, 0(R1)                                  +  1 Sector Paadi[
                                                     ───────────────────
     LD    F3, -8(R1)                                  27 CC is required

     ADD   F4, F3, F5            loop iteration i.e. Termination code
                                 is adjusted with ld & sd &
     SD    F4, -8(R1)            different regrs are used to avoid
                                 hazards
     LD    F6, -16(R1)

     ADD   F7, F6, F8            It optimizes the code by using
                                 Substitution and Simplifications
     SD    F7 -16(R1)


     LD    F9, -24(R1)           If Actually we do not know the
                                 upper bound on the array i.e it
     ADD   F10, F9, F11          is difficult to known; we can
                                 known during the compilation time
     SD    F10, -24(R1)

     DADDUI R1, R1, #32

  1  ⌐ BNE R1, R2, loop
```

Suppose if the size is M & we should unroll k copies of each iteration, then instead of Single large unrolling nod, SIM generate a pair of consecutive loops first loop executes n/k times. the second unrolled loop executes n/k times. For large value of M most of execution time is Spent for loop unrolling unrolling loops improves the performance of loop over pipeling b eventhough code size is large.

```
loop   LD F0, 0(A)         SD F1, 0(R1)          This loop will
       LD F3, -8(R1)       SD F4, -8(R1)         work without
       LD F6, -16(R1)      SD F7, -16(R1)        any stalls.
       LD F9, -24(R1)      SD F10, -24(R1)
       ADD F1, F0, F2      DADDUI R1, R1, #32
       ADD F4, F3, F5      SD F10, -24(R1)
       ADD F7, F6, F8      BNE R1, R2, loop.
       ADD F10, F9, F11
```