# 1 Branch Prediction I

A processor implements an *in-order* pipeline with multiple stages. Each stage completes in a single cycle. The pipeline stalls upon fetching a conditional branch instruction and resumes execution once the condition of the branch is evaluated. There is no other case in which the pipeline stalls.

## 1.1 Part I: Microbenchmarking

You create a microbenchmark as follows to explore the pipeline characteristics:

```
LOOP1:
    SUB R1, R1, #1  // R1 = R1 - 1
    BGT R1, LOOP1   // Branch to LOOP1 if R1 > 0

LOOP2:
    B   LOOP2       // Branch to LOOP2
                    // Repeats until program is killed
```

The microbenchmark takes one input value `R1` and runs until it is killed (e.g., via an external interrupt).

You carefully run the microbenchmark using three different input values as summarized in Table 1. You terminate the microbenchmark using an external interrupt such that each run is guaranteed to execute the same number of *dynamic instructions*. Unfortunately, your testing infrastructure does *not* give you the actual number of instructions executed.

| Initial R1 Value | Number of Cycles Taken |
|---|---|
| 4 | 51 |
| 8 | 63 |
| 16 | 87 |

Table 1: Microbenchmark results.

Using this information, you need to determine the following three experiment characteristics. *Clearly show all work to receive full points!*

1. How many dynamic instructions are executed?

2. How many stages are in the pipeline?

3. For how many cycles does a conditional branch instruction cause a stall?

## 1.2 Part II: Performance Enhancement

To improve performance, the architects add a *mystery* branch prediction mechanism. They keep the rest of the design exactly the same as before. You re-run the microbenchmark for the same number of total dynamic instructions with the new design, and you find that with $R1 = 4$, the microbenchmark executes in 48 cycles.

Based on this given information, determine which of the following branch prediction mechanisms could be the *mystery* branch predictor implemented in the new version of the processor. For each branch prediction mechanism below, you should circle the configuration parameters that makes it match the performance of the mystery branch predictor.

a) **Static Branch Predictor**

Could this be the mystery branch predictor:     YES     NO

If YES, for which configuration below is the answer YES? Pick an option for each configuration parameter.

  I) Static Prediction Direction

    Always taken          Always not taken

Explain:

b) **Last Time Branch Predictor**

Could this be the mystery branch predictor?

YES                    NO

If YES, for which configuration is the answer *YES*? Pick an option for each configuration parameter.

   I) Initial Prediction Direction

      `Taken`               `Not taken`

  II) Local for each branch instruction (PC-based) or global (shared among all branches) history?

      `Local`               `Global`

Explain:

```



```

c) **Backward taken, Forward not taken (BTFN)**
Could this be the mystery branch predictor?

`YES`                  `NO`

Explain:

```



```

d) **Forward taken, Backwards not taken (FTBN)**
Could this be the mystery branch predictor?

`YES`                  `NO`

Explain:

```



```

e) **Two-bit Counter Based Prediction** (using saturating arithmetic)

Could this be the mystery branch predictor?

YES                              NO

If YES, for which configuration is the answer *YES*? Pick an option for each configuration parameter.

I) Initial Prediction Direction

00 (Strongly not taken)        01 (Weakly not taken)
10 (Weakly taken)              11 (Strongly taken)

II) Local for each branch instruction (i.e., PC-based, without any interference between different branches) or global (i.e., a single counter shared among all branches) history?

Local                Global

Explain:

## 2   Branch Prediction II

Assume a processor that implements an ISA with eight registers (R0-R7). In this ISA, the main memory is byte-addressable and each word contains 4 bytes. The processor employs a branch predictor to reduce the overhead of the branches. The ISA implements the instructions given in the following table:

| Instructions | Description |
|---|---|
| la $R_i$, Address | load the effective *Address* into $R_i$ |
| move $R_i$, $R_j$ | $R_i \leftarrow R_j$ |
| move $R_i$, $(R_j)$ | $R_i \leftarrow \text{Memory}[R_j]$ |
| move $(R_i)$, $R_j$ | $\text{Memory}[R_i] \leftarrow R_j$ |
| li $R_i$, Imm | $R_i \leftarrow \text{Imm}$ |
| add $R_i$, $R_j$, $R_k$ | $R_i \leftarrow R_j + R_k$ |
| addi $R_i$, $R_j$, Imm | $R_i \leftarrow R_j + \text{Imm}$ |
| cmp $R_i$, $R_j$ | Compare: Set sign flag, if $R_i < R_j$; set zero flag, if $R_i = R_j$ |
| cmp $R_i$, $(R_j)$ | Compare: Set sign flag, if $R_i < \text{Memory}[R_j]$; set zero flag, if $R_i = \text{Memory}[R_j]$ |
| cmpi $R_i$, Imm | Compare: Set sign flag, if $R_i < \text{Imm}$; set zero flag, if $R_i = \text{Imm}$. |
| jg label | Jump to the target address if **both** of sign and zero flags are zero. |
| jnz label | Jump to the target address if zero flag is zero. |
| halt | Stop executing instructions. |

The processor executes the following program. Answer the questions below related to the accuracy of the branch predictors that the processor can potentially implement.

```
1            la R0, Array
2            move R6, R0
3            li R1, 4
4            move R5, R1
5            move R7, R1
6            move R2, R0
7            addi R2, R2, 4
8    Loop:
9            move R3, (R2)
10           cmp R3, (R0)
11           jg Next_Iteration
12           move R4, (R0)
13           move (R0), R3
14           move (R2), R4
15   Next_Iteration:
16           addi R0, R0, 4
17           addi R2, R2, 4
18           addi R1, R1, -1
19           cmpi R1, 0
20           jnz Loop
21           move R1, R7
22           addi R5, R5, -1
23           move R0, R6
24           move R2, R0
25           addi R2, R2, 4
26           cmpi R5, 0
27           jnz Loop
28           halt
29   .data
30   Array: word 5, 20, 1, -5, 34
```

(a) What would be the prediction accuracy using a shared one-bit-history (last time) branch predictor for *all* the branches? The initial state of the predictor is "taken".

(b) What would be the prediction accuracy using a shared two-bit-history (two-bit counter) branch predictor for *all* the branches? Assume that the initial state of the two-bit-history branch predictor is "weakly taken". The "weakly taken" state transitions to "weakly not-taken" state on misprediction. Similarly, the "weakly not-taken" taken state transitions to "weakly taken" state on misprediction. A correct prediction in one of the "weak" states transitions the state to the corresponding "strong" state.

(c) What would be the prediction accuracy using two-bit-history (two-bit counter) branch predictor for *each* branch? The initial state is "weakly taken" and the state transitions are the same as in part (b).

# 3    Delayed Branching I

You are designing an ISA that uses delayed branch instructions. You are trying to decide how many instructions to place into the branch delay slot. How many branch delay slots would you need for the following different implementations? Explain your reasoning briefly.

(a) An in-order processor where conditional branches resolve during the 4th stage

(b) An out-of-order processor with 64 unified reservation station entries where conditional branches resolve during the 2nd cycle of branch execution. The processor has 15 pipeline stages until the start of the execution stages

# 4 Delayed Branching II

A machine has a five-stage pipeline consisting of fetch, decode, execute, mem and write-back stages. The machine uses delay slots to handle control dependences. Jump targets, branch targets and destinations are resolved in the execute stage.

a) What is the number of delay slots needed to ensure correct operation?



b) Which instruction(s) in the assembly sequences below would you place in the delay slot(s), assuming the number of delay slots you answered for part(a)? Clearly rewrite the code with the appropriate instruction(s) in the delay slot(s).

I)
```
ADD R5 <- R4, R3
OR  R3 <- R1, R2
SUB R7 <- R5, R6
J   X
```
Delay Slots
```
LW  R10 <- (R7)
ADD R6 <- R1, R2
X:
```

**Solution:**



II)
```
ADD R5 <- R4, R3
OR  R3 <- R1, R2
SUB R7 <- R5, R6
BEQ R5 <- R7, X
```
Delay Slots
```
LW  R10 <- (R7)
ADD R6 <- R1, R2
X:
```

**Solution:**



III) 
```
ADD R2 <- R4, R3
OR  R5 <- R1, R2
SUB R7 <- R5, R6
BEQ R5 <- R7, X
```
Delay Slots

```
LW  R10 <- (R7)
ADD R6 <- R1, R2
X:
```

**Solution:**



c) Can you modify the pipeline to reduce the number of delay slots (without introducing branch prediction)? Clearly state your solution and explain why.
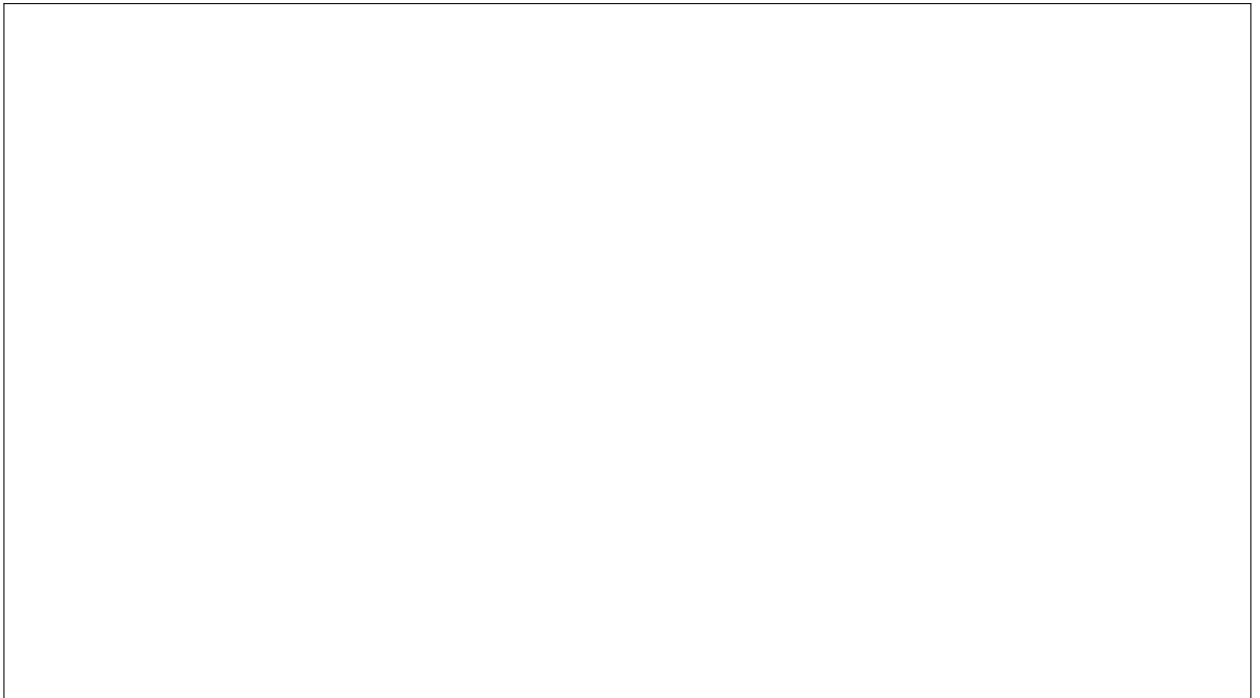
# 5   VLIW I

You are using a tool that transforms machine code that is written for the MIPS ISA to code in a VLIW ISA. The VLIW ISA is identical to MIPS except that multiple instructions can be grouped together into one VLIW instruction. Up to N MIPS instructions can be grouped together (N is the machine width, which depends on the particular machine). The transformation tool can reorder MIPS instructions to fill VLIW instructions, as long as loads and stores are not reordered relative to each other (however, independent loads and stores can be placed in the same VLIW instruction).

You give the tool the following MIPS program (we have numbered the instructions for reference below):

```
(01) lw   $t0 ← 0($a0)
(02) lw   $t2 ← 8($a0)
(03) lw   $t1 ← 4($a0)
(04) add  $t6 ← $t0, $t1
(05) lw   $t3 ← 12($a0)
(06) sub  $t7 ← $t1, $t2
(07) lw   $t4 ← 16($a0)
(08) lw   $t5 ← 20($a0)
(09) srlv $s2 ← $t6, $t7
(10) sub  $s1 ← $t4, $t5
(11) add  $s0 ← $t3, $t4
(12) sllv $s4 ← $t7, $s1
(13) srlv $s3 ← $t6, $s0
(14) sllv $s5 ← $s0, $s1
(15) add  $s6 ← $s3, $s4
(16) add  $s7 ← $s4, $s6
(17) srlv $t0 ← $s6, $s7
(18) srlv $t1 ← $t0, $s7
```

(a) Draw the dataflow graph of the program. Represent instructions as numbered nodes (01 through 18) and flow dependencies as directed edges (arrows).

(b) When you run the tool with its settings targeted for a particular VLIW machine, you find that the resulting VLIW code has 9 VLIW instructions. What minimum value of N must the target VLIW machine have?

(c) Write the MIPS instruction numbers (from the code above) corresponding to each VLIW instruction, for this value of N. When there is more than one MIPS instruction that could be placed into a VLIW instruction, choose the instruction that comes earliest in the original MIPS program.

| | MIPS Instr No | MIPS Instr No | MIPS Instr No | MIPS Instr No | MIPS Instr No | MIPS Instr No | MIPS Instr No | MIPS Instr No | MIPS Instr No | MIPS Instr No |
|---|---|---|---|---|---|---|---|---|---|---|
| VLIW Instr.1: | | | | | | | | | | |
| VLIW Instr.2: | | | | | | | | | | |
| VLIW Instr.3: | | | | | | | | | | |
| VLIW Instr.4: | | | | | | | | | | |
| VLIW Instr.5: | | | | | | | | | | |
| VLIW Instr.6: | | | | | | | | | | |
| VLIW Instr.7: | | | | | | | | | | |
| VLIW Instr.8: | | | | | | | | | | |
| VLIW Instr.9: | | | | | | | | | | |

(d) You find that the code is still not fast enough when it runs on the VLIW machine, so you contact the VLIW machine vendor to buy a machine with a larger machine-width "N". What minimum value of N would yield the maximum possible performance (i.e., the fewest VLIW instructions), assuming that all MIPS instructions (and thus VLIW instructions) complete with the same fixed latency and assuming no cache misses?

(e) Write the MIPS instruction numbers corresponding to each VLIW instruction, for this optimal value of N. Again, as in part (c) above, pack instructions such that when more than one instruction can be placed in a given VLIW instruction, the instruction that comes first in the original MIPS code is chosen.

|  | MIPS Instr No | MIPS Instr No | MIPS Instr No | MIPS Instr No | MIPS Instr No | MIPS Instr No | MIPS Instr No | MIPS Instr No | MIPS Instr No | MIPS Instr No |
|---|---|---|---|---|---|---|---|---|---|---|
| VLIW Instr.1: | | | | | | | | | | |
| VLIW Instr.2: | | | | | | | | | | |
| VLIW Instr.3: | | | | | | | | | | |
| VLIW Instr.4: | | | | | | | | | | |
| VLIW Instr.5: | | | | | | | | | | |
| VLIW Instr.6: | | | | | | | | | | |
| VLIW Instr.7: | | | | | | | | | | |
| VLIW Instr.8: | | | | | | | | | | |
| VLIW Instr.9: | | | | | | | | | | |

(f) A competing processor design company builds an in-order superscalar processor with the same machine-width N as the width you found above in part(b). The machine has the same clock frequency as the VLIW processor. When you run the original MIPS program on this machine, you find that it executes slower than the corresponding VLIW progra m on the VLIW machine in part (b). Why could this be the case?

(g) When you run some other program on this superscalar machine, you find it runs faster than the corresponding VLIW program on the VLIW machine. Why could this be the case?

# 6 VLIW II

Explain the motivation for VLIW in one sentence.

<br><br><br><br>

You are the human compiler for a VLIW machine whose specifications are as follows:

- There are 3 **fully** pipelined functional units (ALU, MU and FPU).

- Integer Arithmetic Logic Unit (ALU) has a 1-cycle latency.

- Memory Unit (MU) has a 2-cycle latency.

- Floating Point Unit (FPU) has a 3-cycle latency, and can perform either FADD or FMUL (floating point add / floating point multiply) on **floating point registers**.

- This machine has **only** 4 integer registers (r1 .. r4) and 4 floating point registers (f1 .. f4)

- The machine does not implement hardware interlocking or data forwarding.

a) For the given assembly code on the next page, fill **Table 1** (on the next page) with the appropriate VLIW instructions for only one iteration of the loop (The C code is also provided for your reference). Provide the VLIW instructions that lead to the **best** performance. Use the minimum number of VLIW instructions. Table 1 should **only** contain instructions provided in the assembly example. For all the instruction tables, show the NOP instructions you may need to insert. Note that BNE is executed in the **ALU**.

The base addresses for A, B, C are stored in r1, r2, r3 respectively. The address of the last element in the array C[N-1] is stored in r4, where N is an integer multiplier of 10! (read: 10 factorial).

**C Code**

```
float A[N];
float C[N];
int   B[N];
...  // code to initialize A and B
for (int i=0; i<N; i++)
    C[i] = A[i] * A[i] + B[i];
```

**Assembly Code**

```
loop: LD   f1, 0 (r1)
      LD   f2, 0 (r2)
      FMUL f1, f1, f1
      FADD f1, f1, f2
      ADDI r3, r3, 4
      ST   f1, -4, (r3)
      ADDI r1, r1, 4
      ADDI r2, r2, 4
      BNE  r3, r4, loop
```

| VLIW Instruction | ALU | MU | FPU |
|---|---|---|---|
| 1 | | | |
| 2 | | | |
| 3 | | | |
| 4 | | | |
| 5 | | | |
| 6 | | | |
| 7 | | | |
| 8 | | | |
| 9 | | | |
| 10 | | | |
| 11 | | | |
| 12 | | | |
| 13 | | | |
| 14 | | | |
| 15 | | | |
| 16 | | | |
| 17 | | | |
| 18 | | | |
| 19 | | | |
| 20 | | | |
| 21 | | | |
| 22 | | | |
| 23 | | | |
| 24 | | | |
| 25 | | | |

Table 1

What is the performance in Ops/VLIW instruction (Operations/VLIW instruction) for this design? An operation here refers to an instruction (in the Assembly Code), excluding NOPs.

b) Assume now we decide to unroll the loop once. Fill **Table 2** with the new VLIW instructions. You should optimize for latency first, then instruction count. **You can choose to use different offsets, immediates and registers, but you may not use any new instructions**.

| VLIW Instruction | ALU | MU | FPU |
|---|---|---|---|
| 1 | | | |
| 2 | | | |
| 3 | | | |
| 4 | | | |
| 5 | | | |
| 6 | | | |
| 7 | | | |
| 8 | | | |
| 9 | | | |
| 10 | | | |
| 11 | | | |
| 12 | | | |
| 13 | | | |
| 14 | | | |
| 15 | | | |
| 16 | | | |
| 17 | | | |
| 18 | | | |
| 19 | | | |
| 20 | | | |
| 21 | | | |
| 22 | | | |
| 23 | | | |
| 24 | | | |
| 25 | | | |

Table 2

What is the performance in Ops/VLIW instruction for this design?

c) Assume now we have **unlimited registers** and the loop is fully optimized (unrolled to the best performance possible). What is the performance in Ops/cycle for this design? Show your work and explain **clearly** how you arrived at your answer. You are not required to draw any tables, but you may choose to do so to aid your explanation. ( Hint: trace the dependent instructions )

# 7  Systolic Arrays

Figure 1 shows a systolic array processing element.

Each processing element takes in two inputs, M and N, and outputs P and Q. Each processing element also contains an "accumulator" R that can be read from and written to. The initial value of the "accumulator" is 0.

Figure 2 shows a systolic array composed of 9 processing elements. The smaller boxes are the inputs to the systolic array and the larger boxes are the processing elements. You will program this systolic array to perform the following calculation:

$$
\begin{bmatrix} c_{00} & c_{01} & c_{02} \\ c_{10} & c_{11} & c_{12} \\ c_{20} & c_{21} & c_{22} \end{bmatrix}
=
\begin{bmatrix} a_{00} & a_{01} & a_{02} \\ a_{10} & a_{11} & a_{12} \\ a_{20} & a_{21} & a_{22} \end{bmatrix}
\times
\begin{bmatrix} b_{00} & b_{01} & b_{02} \\ b_{10} & b_{11} & b_{12} \\ b_{20} & b_{21} & b_{22} \end{bmatrix}
$$

In each time cycle, each processing element will take in its two inputs, perform any necessary actions, and write on its outputs. The time cycle labels on the input boxes determine which time cycle the inputs will be fed into their corresponding processing elements. Any processing element input that is not driven will default to 0, and any processing element that has no output arrow will have its output ignored.

After all the calculations finish, each processing element's "accumulator" will hold one element of the final result matrix, arranged in the correct order.

(a) Please describe the operations that each individual processing element performs, using mathe- matical equations and the variables M, N, P, Q and R.
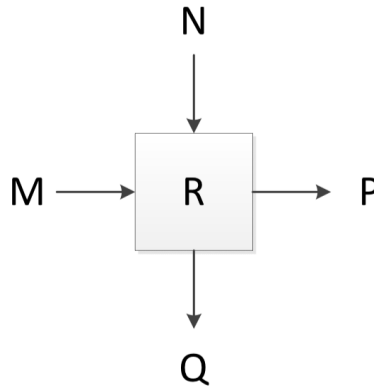


Figure 1: A systolic array processing element

P =

Q =

R =

(b) Please fill in all 30 input boxes in Figure 2 so that the systolic array computes the correct matrix multiplication result described on the previous page. (Hint: Use $a_{ij}$ and $b_{ij}$.)
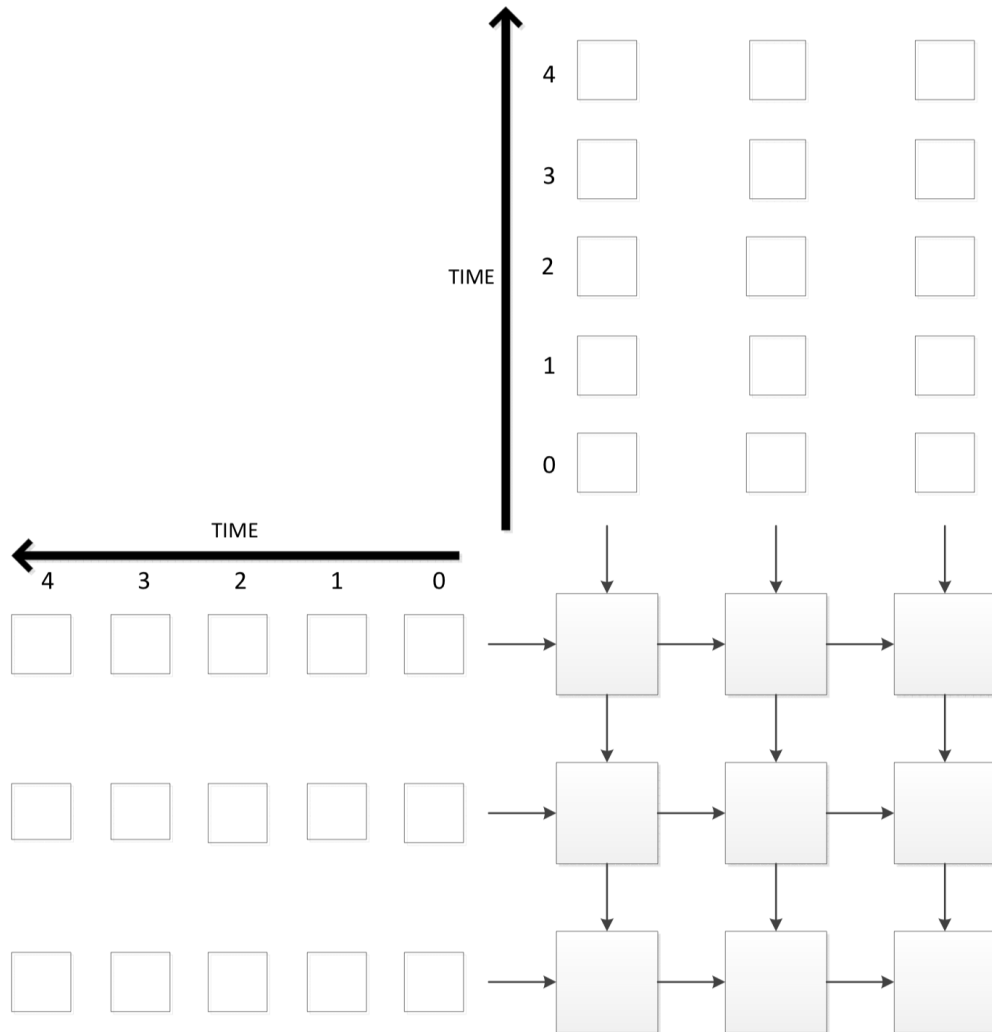


Figure 2: A systolic array

# 8    BONUS: Branch Prediction

Assume the following piece of code that iterates through a large array populated with **completely (i.e., truly) random** positive integers. The code has four branches (labeled B1, B2, B3, and B4). When we say that a branch is *taken*, we mean that the code *inside* the curly brackets is executed.

```
for (int i=0; i<N; i++) { /* B1 */
  val = array[i];        /* TAKEN PATH for B1 */
  if (val % 2 == 0) {    /* B2 */
    sum += val;          /* TAKEN PATH for B2 */
  }
  if (val % 3 == 0) {    /* B3 */
    sum += val;          /* TAKEN PATH for B3 */
  }
  if (val % 6 == 0) {    /* B4 */
    sum += val;          /* TAKEN PATH for B4 */
  }
}
```

(a) Of the four branches, list all those that exhibit *local correlation*, if any.

<br/>

(b) Which of the four branches are *globally correlated*, if any? Explain in less than 20 words.
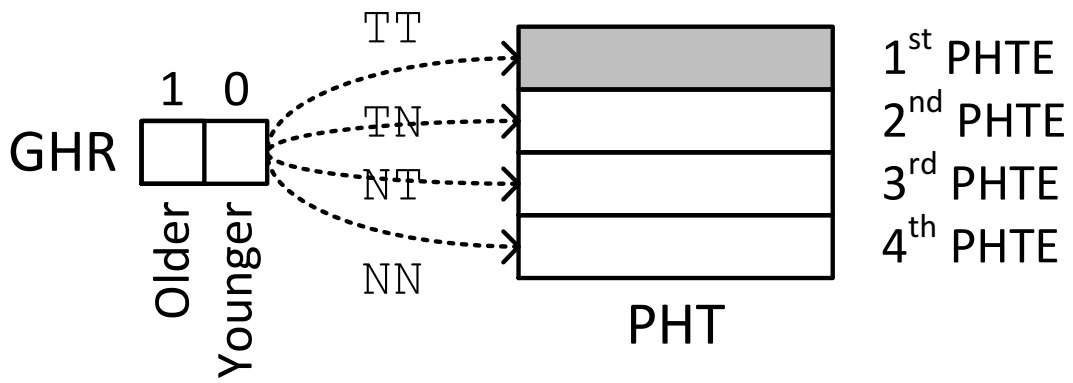
<br/>

Now assume that the above piece of code is running on a processor that has a global branch predictor. The global branch predictor has the following characteristics.

- Global history register (GHR): 2 bits.

- Pattern history table (PHT): 4 entries.

- Pattern history table entry (PHTE): 11-bit signed saturating counter (possible values: -1024–1023)

- Before the code is run, all PHTEs are initially set to 0.

- As the code is being run, a PHTE is incremented (by one) whenever a branch that corresponds to that PHTE is taken, whereas a PHTE is decremented (by one) whenever a branch that corresponds to that PHTE is not taken.

(c) After 120 iterations of the loop, calculate the **expected** value for only the first PHTE and fill it in the shaded box below. (Please write it as a base-10 value, rounded to the nearest one's digit.)

*Hint. For a given iteration of the loop, first consider, what is the probability that both B1 and B2 are taken? Given that they are, what is the probability that B3 will increment or decrement the PHTE? Then consider...*
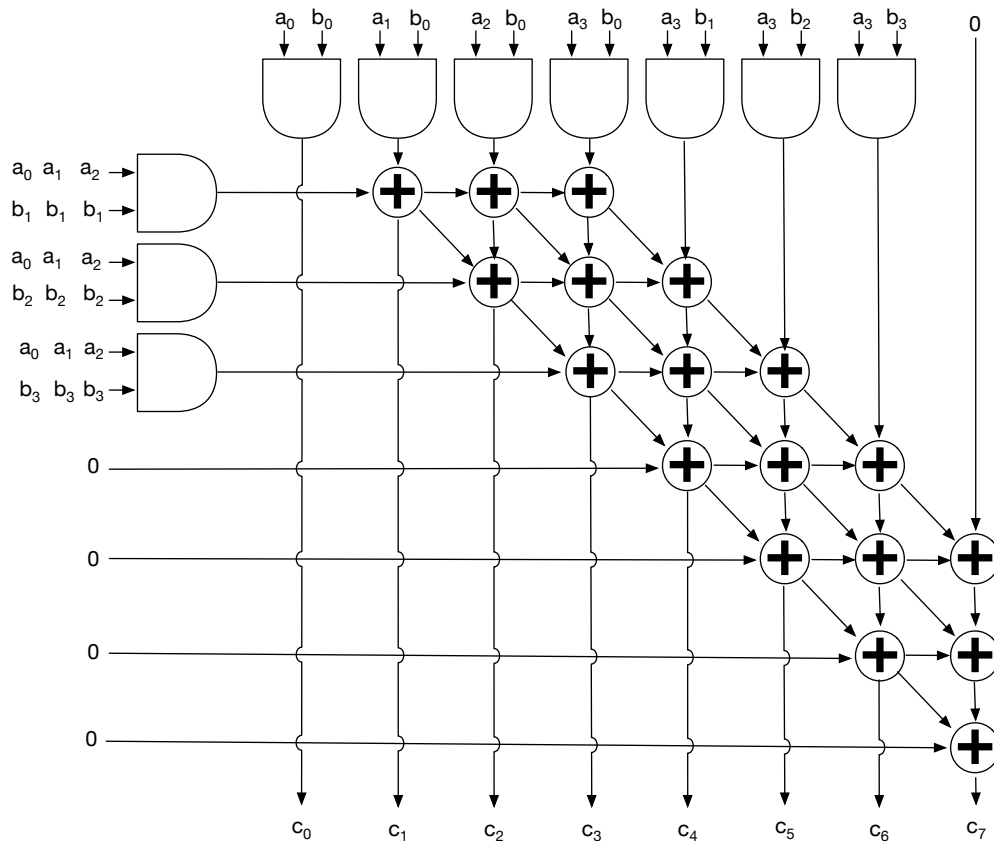
Show your work.

# 9 BONUS: Systolic Arrays

The following diagram is a systolic array that performs the multiplication of two 4-bit binary numbers ($\mathbf{a}$ and $\mathbf{b}$). For example, if $\mathbf{a}$=1110 and $\mathbf{b}$=1011, the result of the multiplication is $\mathbf{c}$=10011010:

$$
\begin{array}{r}
1011 \\
\times\ 1110 \\
\hline
0000 \\
1011 \\
1011 \\
+\ 1011 \\
\hline
10011010
\end{array}
\qquad (1)
$$

The input to the systolic arrays is through the AND gates. The figure shows which bits of the two numbers $\mathbf{a}$ and $\mathbf{b}$ are inserted into each AND gate. However, the figure does *not* indicate in which cycle each input is issued. Make the folowing assumptions:

- The latency of each adder is one cycle.

- Vertical arrows propagate the sum to the next adder.

- Diagonal arrows propagate the carry to the next adder.

- Horizontal arrows propagate the output of the AND gates in each row.

- An adder adds the value of its three inputs (vertical, diagonal and horizontal inputs)

- An adder can hold a value for only one cycle.

(a) How many cycles does it take to perform *one* multiplication of two 4-bit binary numbers in this systolic array? Indicate 1) in which cycle each bit is inputted in the systolic array and 2) in which cycle each bit of the result is produced.

(b) How many cycles does it take to perform $N$ consecutive multiplications of two 4-bit binary numbers in this systolic array?