



2024年度(令和6年)版

Ver. 2024-11-01a

Course number: CSC.T363

# コンピューターアーキテクチャ Computer Architecture

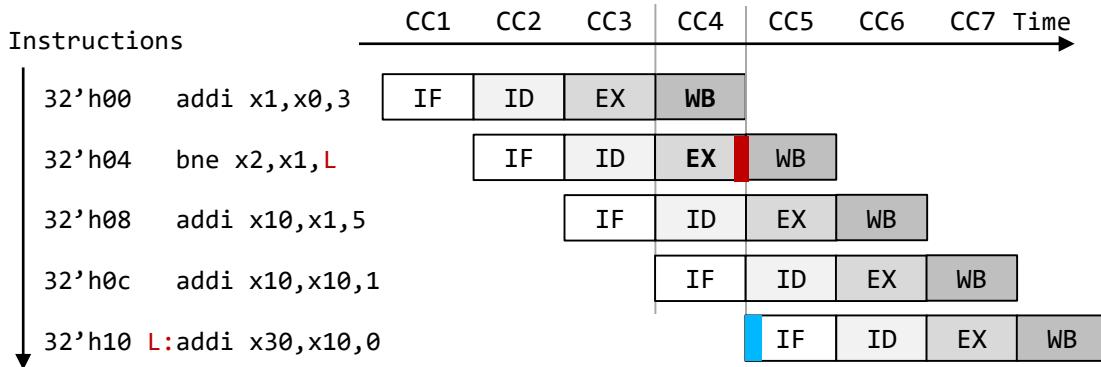
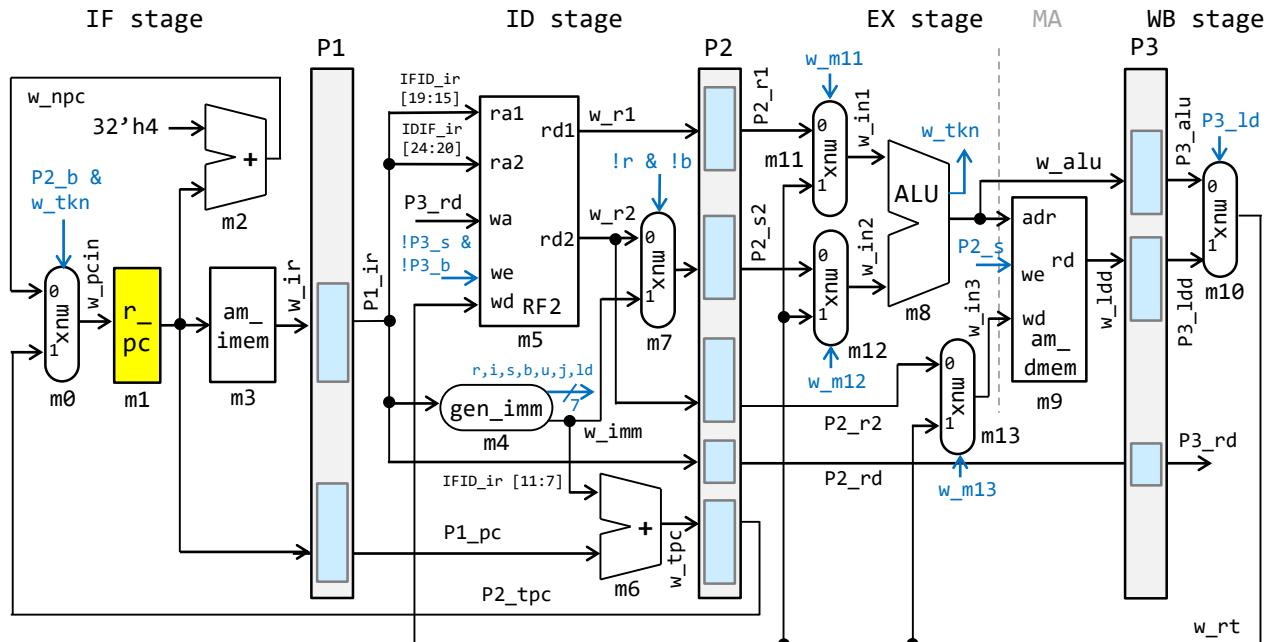
## 9. 分岐予測 (2) Branch Prediction (2)

[www.arch.cs.titech.ac.jp/lecture/CA/](http://www.arch.cs.titech.ac.jp/lecture/CA/)  
Tue 13:30-15:10, 15:25-17:05  
Fri 13:30-15:10

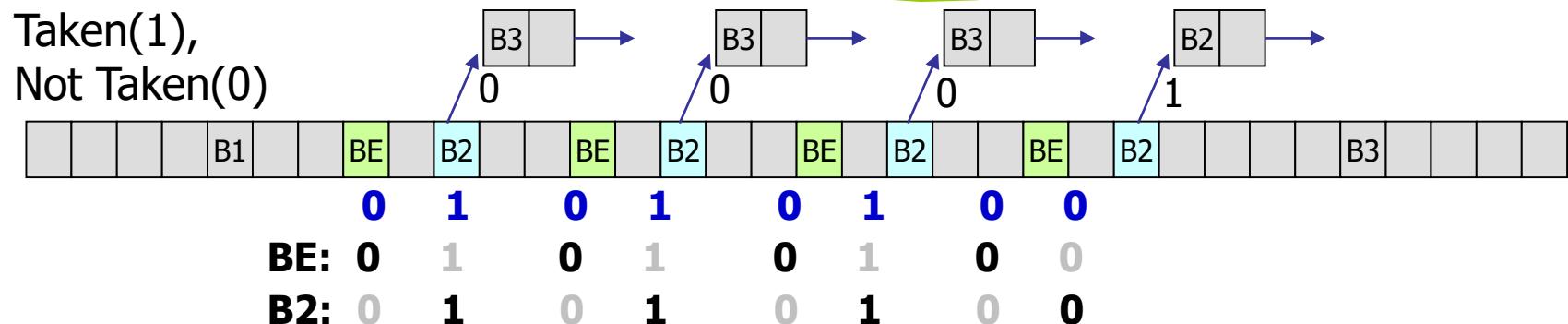


吉瀬 謙二 情報工学系  
Kenji Kise, Department of Computer Science  
kise\_at\_c.titech.ac.jp

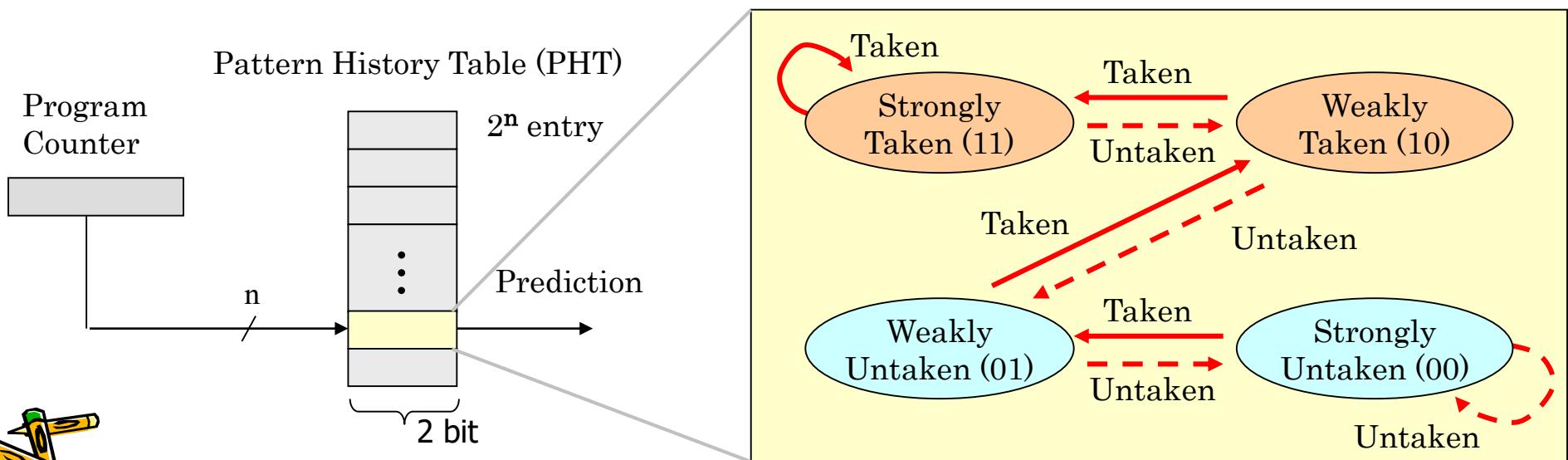
# proc8: 4-stage pipelining processor



# Bimodal branch predictor (ISCA 1981)

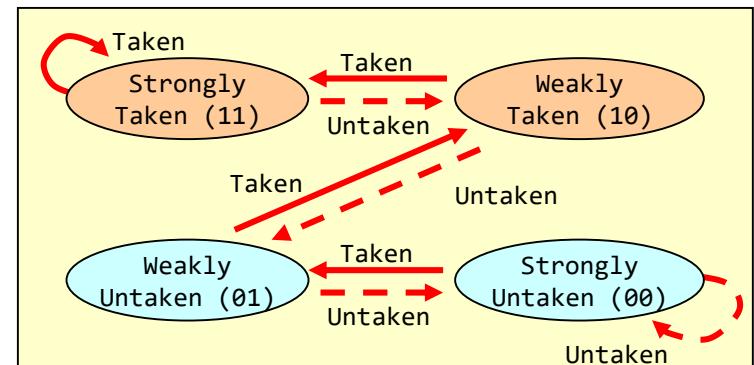


- 分岐アドレス(プログラムカウンタ)毎に履歴を切り替える
  - 分岐アドレスによりパターン履歴表(PHT)のインデックスを作成
- パターン履歴表は2ビットカウンタの配列.



# Bimodal branch predictor (ISCA 1981)

```
1 module m_bimodal (
2     input wire w_clk, w_we,
3     input wire [4:0] w_wadr, w_radr,
4     input wire w_tkn,
5     output wire w_pred
6 );
7     reg [1:0] mem [0:31];
8     wire [1:0] w_data = mem[w_radr];
9     assign w_pred = w_data[1];
10    wire [1:0] w_cnt = mem[w_wadr];
11    always @(posedge w_clk) if (w_we)
12        mem[w_wadr] <= (w_cnt < 3 & w_tkn) ? w_cnt + 1 :
13                                (w_cnt > 0 & !w_tkn) ? w_cnt - 1 : w_cnt;
14    integer i; initial for (i=0; i<32; i=i+1) mem[i] = 1;
15 endmodule
```

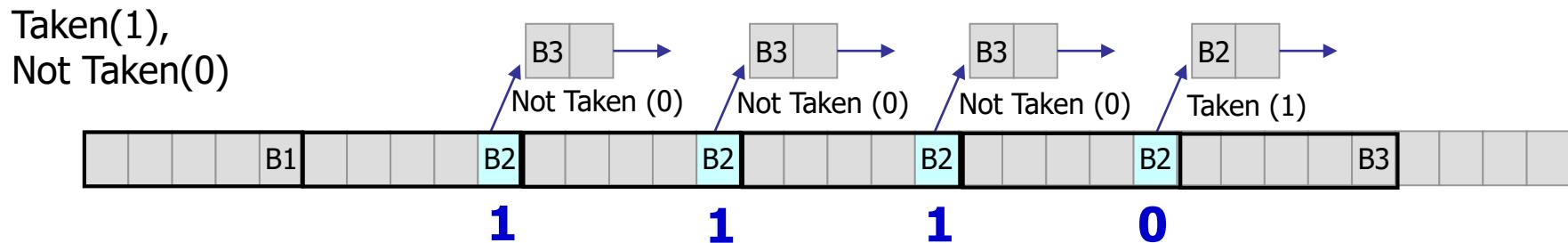


# サンプルプログラム

```
1 `MM[0]={12'd0,5'd0,3'h0,5'd10,7'h13};           // 00    addi x10,x0,0
2 `MM[1]={12'd0,5'd0,3'h0,5'd3,7'h13};           // 04    addi x3,x0,0
3 `MM[2]={12'd101,5'd0,3'h0,5'd1,7'h13};          // 08    addi x1,x0,101
4 `MM[3]={7'd0,5'd3,5'd10,3'h0,5'd10,7'h33};       // 0c    L:add x10,x10,x3
5 `MM[4]={12'd1,5'd3,3'h0,5'd3,7'h13};           // 10    addi x3,x3,1
6 `MM[5]={~12'd0,5'd1,5'd3,3'h1,5'b11001,7'h63}// 14    bne x3,x1,L
7 `MM[6]=32'h00050f13;                           // 18    HALT
```

```
1 `MM[0] ={12'd4, 5'd0, 3'h0,5'd1, 7'h13};        // 00    addi x1,x0,4
2 `MM[1] ={12'd101,5'd0, 3'h0,5'd2, 7'h13};        // 04    addi x2,x0,101
3 `MM[2] ={12'd0, 5'd0, 3'h0,5'd10,7'h13};         // 08    addi x10,x0,0
4 `MM[3] ={12'd0, 5'd0, 3'h0,5'd4 ,7'h13};         // 0c    addi x4,x0,0
5 `MM[4] ={12'd0, 5'd0, 3'h0,5'd3 ,7'h13};         // 10    M:addi x3,x0,0
6 `MM[5] ={12'd1, 5'd3 ,3'h0,5'd3, 7'h13};         // 14    L:addi x3,x3,1
7 `MM[6] ={12'd1, 5'd10,3'h0,5'd10,7'h13};        // 18    addi x10,x10,1
8 `MM[7] ={~12'd0,5'd1,5'd3,3'h1,5'b11001,7'h63}// 1c    bne x3,x1,L
9 `MM[8] ={12'd1, 5'd4 ,3'h0,5'd4, 7'h13};         // 20    addi x4,x4,1
10 `MM[9] ={12'd1, 5'd10,3'h0,5'd10,7'h13};        // 24    addi x10,x10,1
11 `MM[10]={~12'd0,5'd2,5'd4,3'h1,5'b01001,7'h63}// 28    bne x4,x2,M
12 `MM[11]=32'h00050f13;                           // 2c    HALT
```

# Branch history (分岐履歴)



## B2の分岐履歴

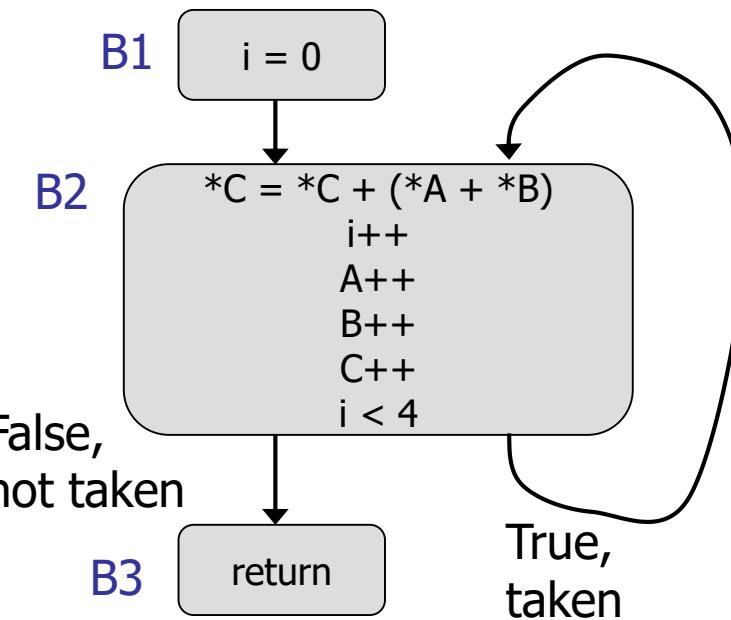
1110 ?

11101 ?

111011 ?

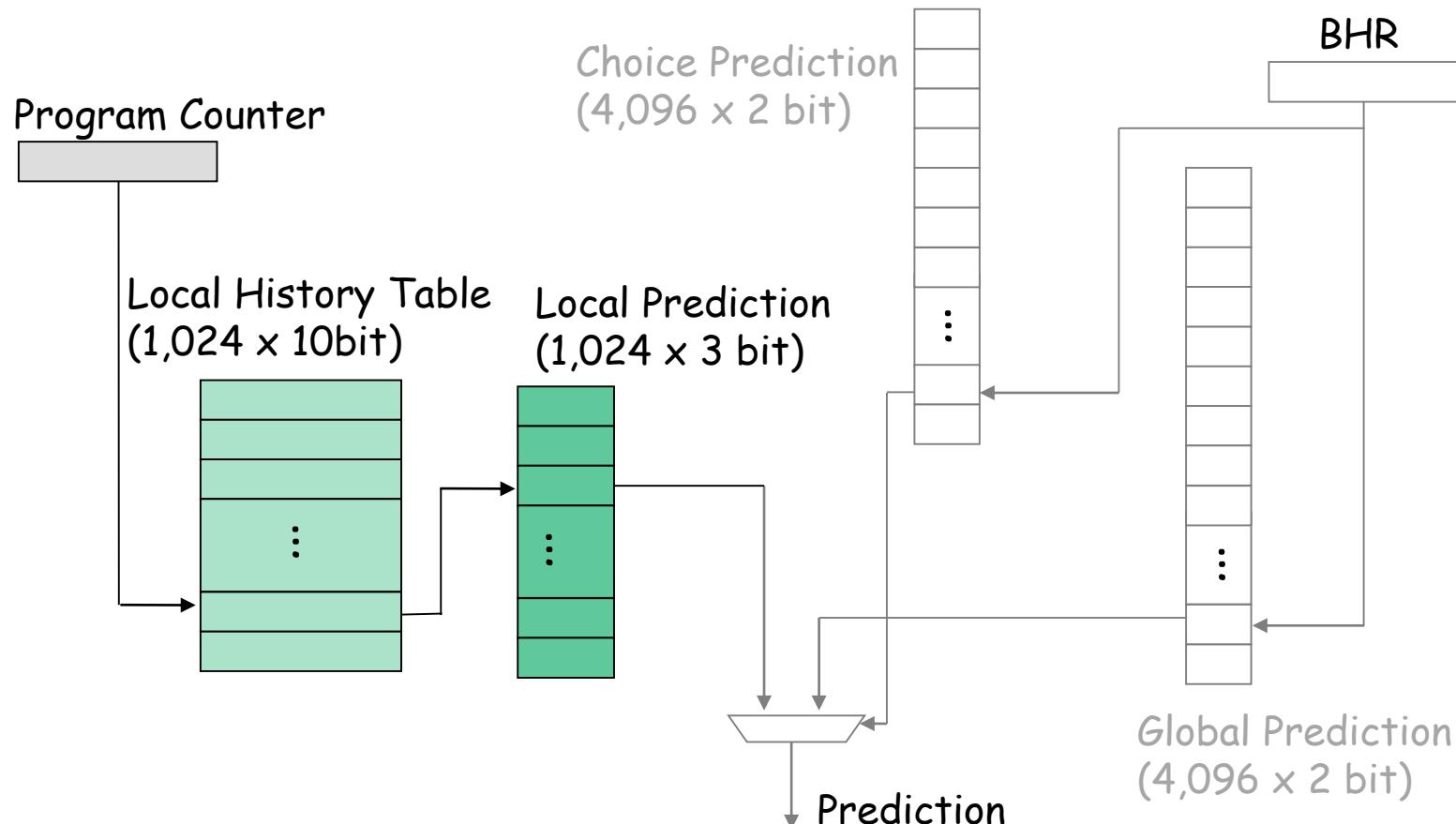
1110111 ?

11101110 ?

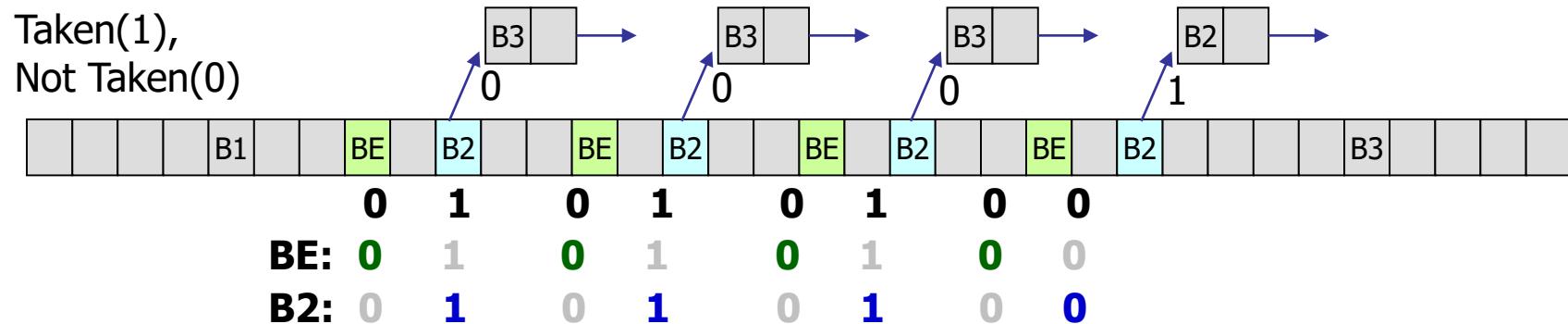


# Alpha 21264's hybrid branch predictor in 1996

- A hybrid of local prediction and global prediction implemented in DEC Alpha 21264 which was the state-of-the art commercial processor.
- A choice predictor is used as a meta-predictor



# ローカル分岐履歴とグローバル分岐履歴



BEの分岐履歴

0000 ?

00000 ?

000000 ?

0000000 ?

00000000 ?

ローカル分岐履歴

B2の分岐履歴

1110 ?

11101 ?

111011 ?

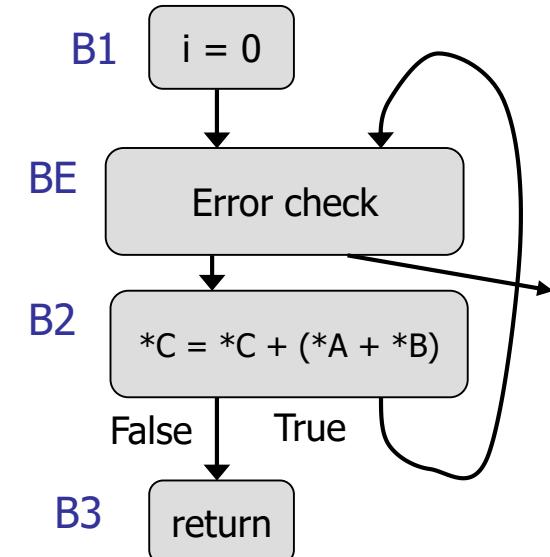
1110111 ?

11101110 ?

ローカル分岐履歴

B2とBEの分岐履歴

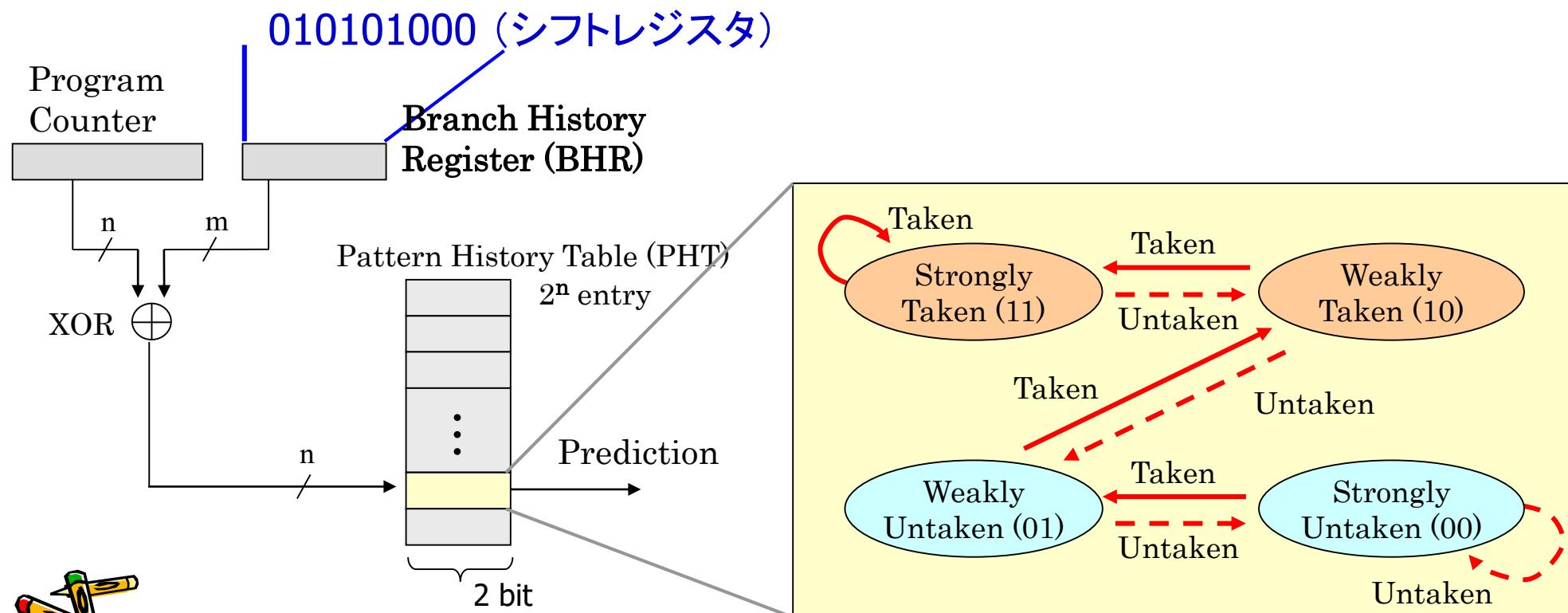
010101000 ?



グローバル分岐履歴

# Gshare (TR-DEC 1993)

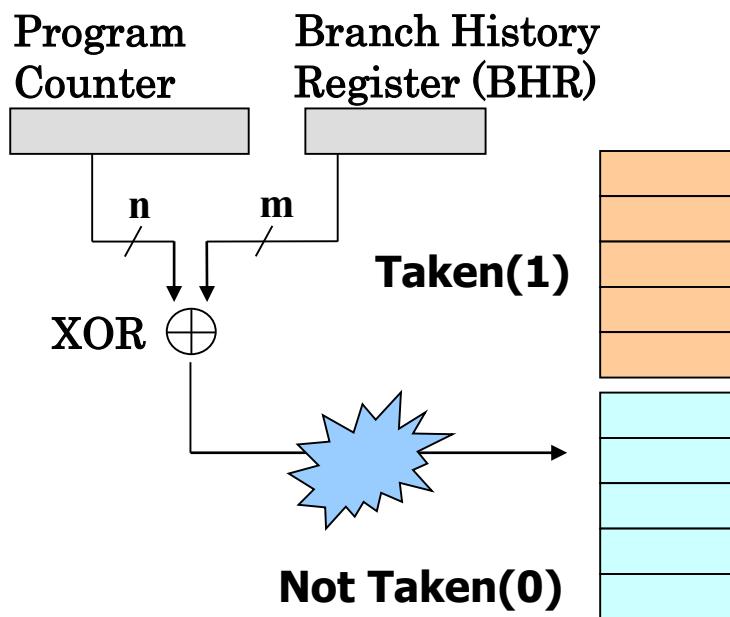
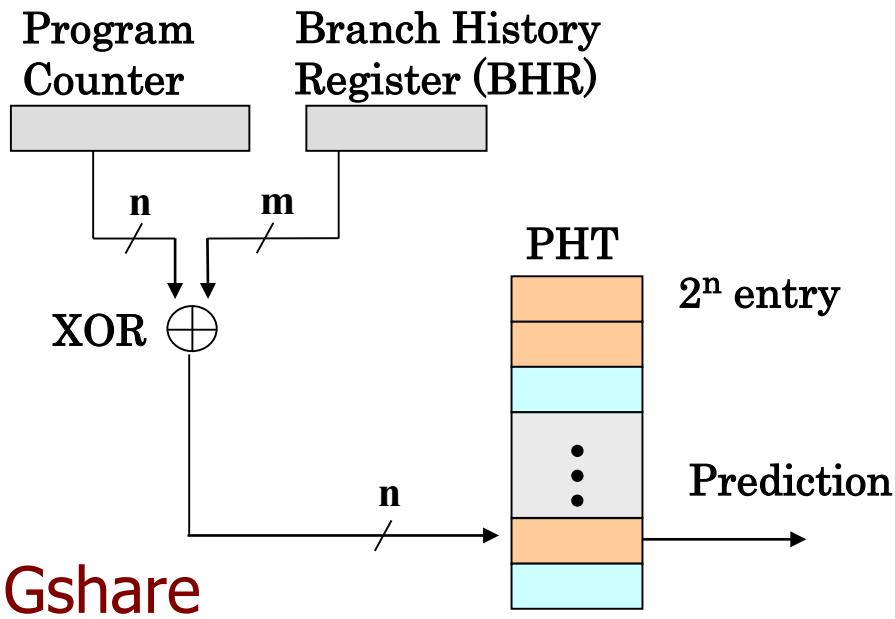
- グローバル分岐履歴と分岐アドレスとの排他的論理和によりパターン履歴表へのインデックスを作成
  - パターン履歴表は2ビット飽和型カウンタの配列で、選択された2ビットカウンタの値により分岐方向を予測(**bimodal**と同じ)
  - 分岐結果を用いて、予測に利用したカウンタを更新



# Gshare の改良

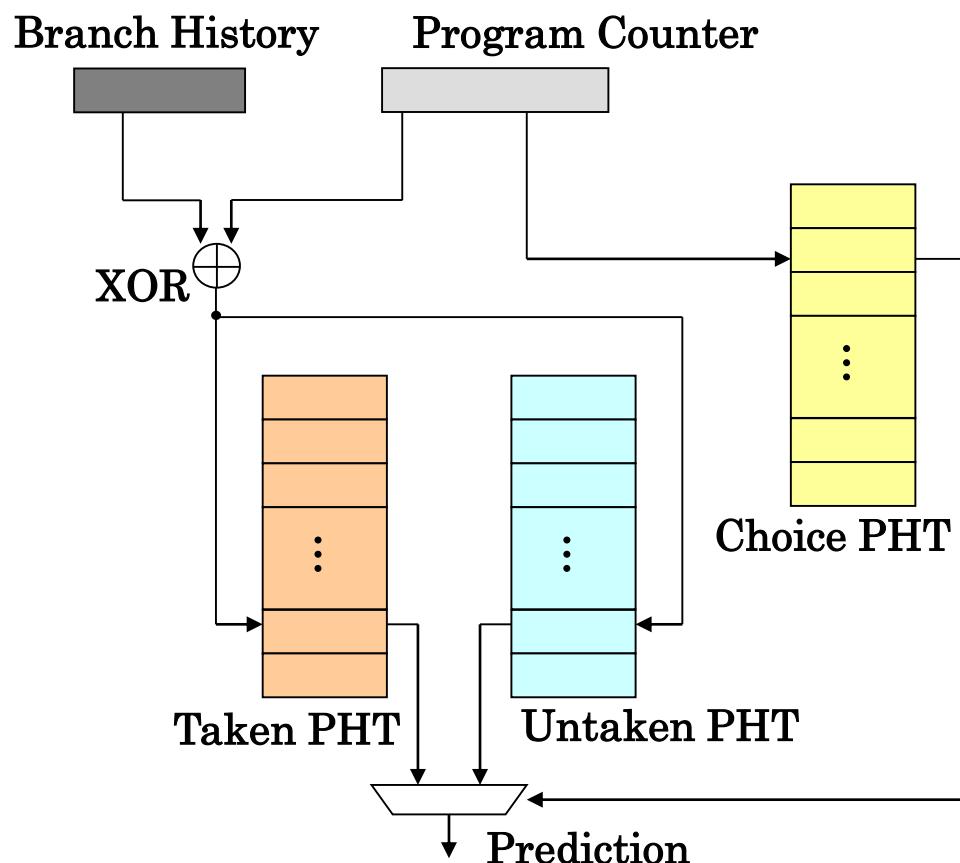


- PHTの競合が発生して性能が低下
- PCとBHRによって特定される予測(成立, 不成立)には偏りが存在するので,これらを別のテーブルに格納することで競合の悪影響を緩和
  - 分岐成立に偏っているもの
  - 分岐不成立に偏っているもの



# Bimode (MICRO 1997)

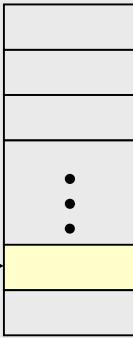
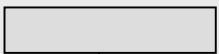
- 傾りを利用して競合の悪影響を緩和
  - 分岐成立に偏っているものをTaken PHTに格納
  - 分岐不成立に偏っているものをUntaken PHTに格納
- Choice PHT の内容で、どちらのテーブルを利用するか選択
- インデックスを工夫
  - Choice PHT は命令アドレス
  - Taken PHT, Untaken PHT は命令アドレスと分岐履歴



# 幾つかの分岐予測器

Pattern History Table

Program Counter



Prediction

(b) Bimodal

PC

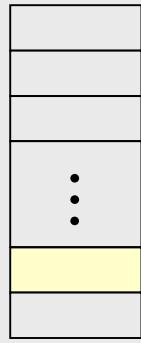


Branch History Register (BHR)

Pattern History Table

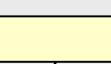
XOR

(c) Gshare



Prediction

2 bit



Prediction

(a) 2ビットカウンタ方式

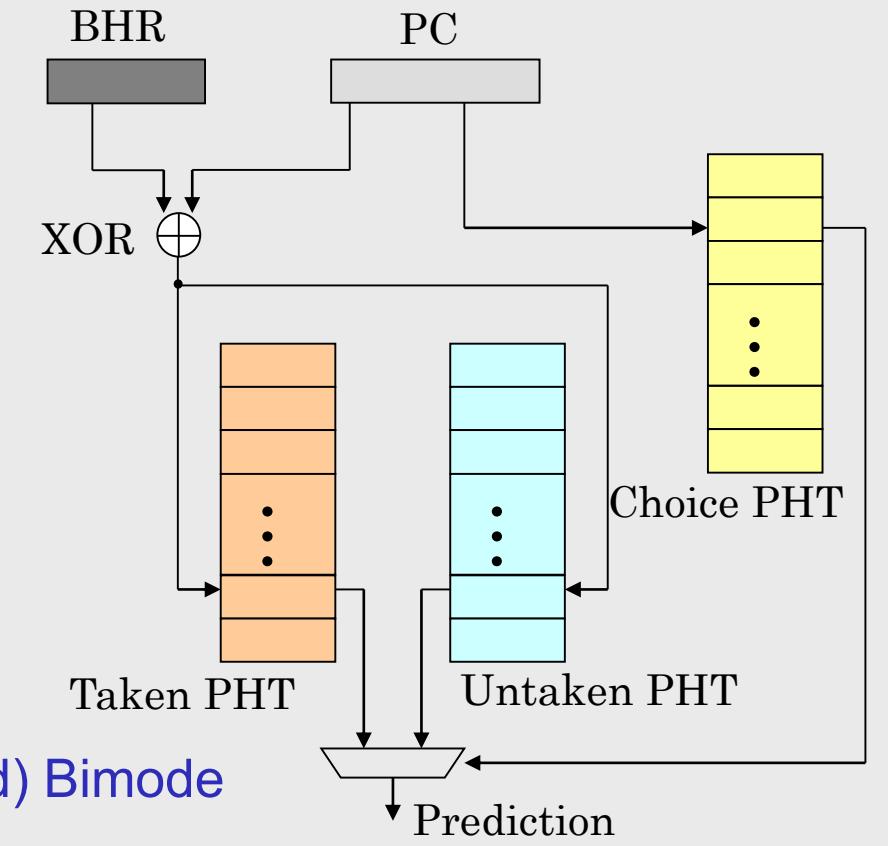
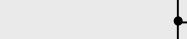
BHR



PC



XOR



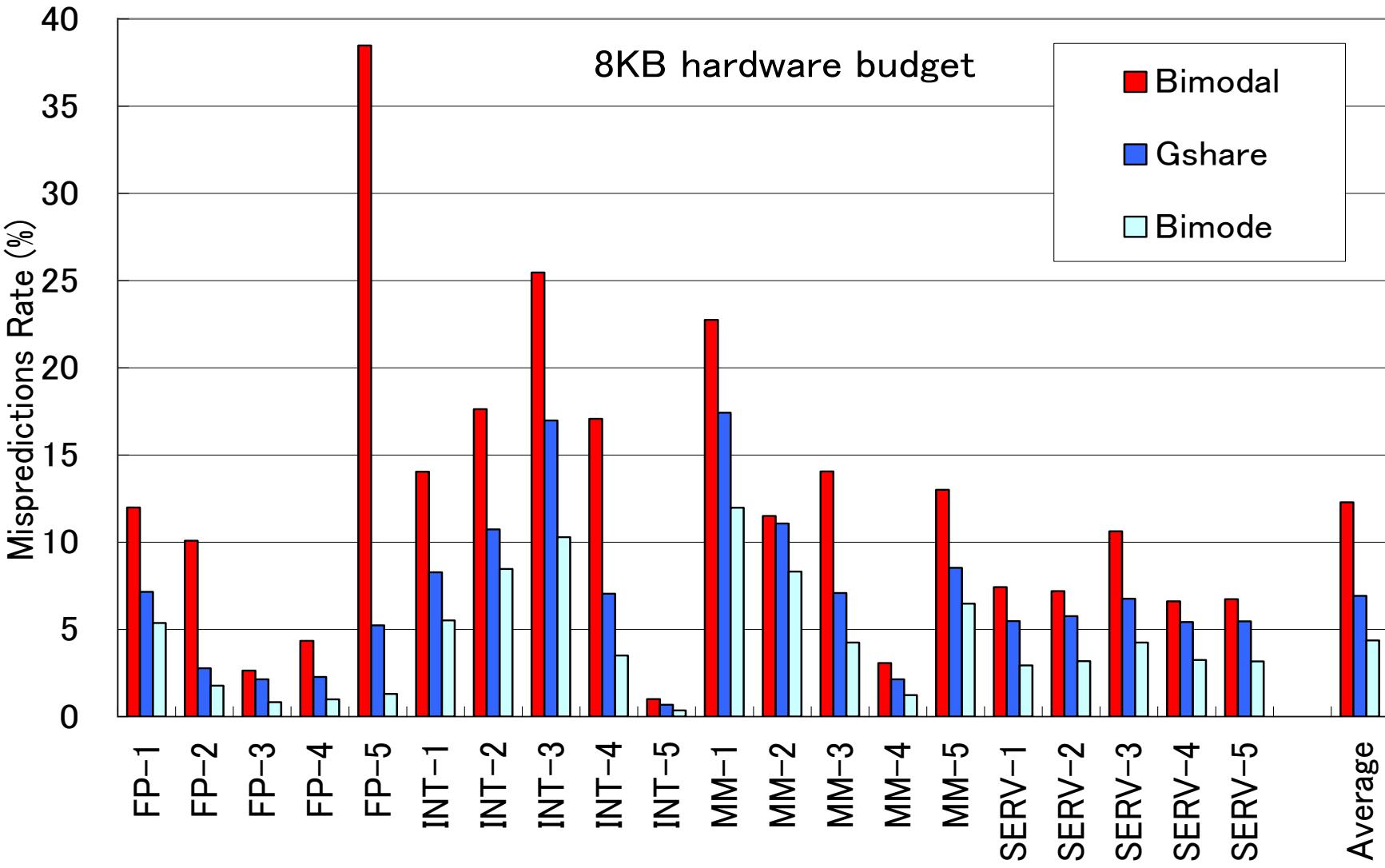
Taken PHT

Untaken PHT

Prediction

(d) Bimode

# 予測精度(予測ミス率)



Benchmark for CBP(2004) by Intel MRL and IEEE TC uARCH.

# 高いバンド幅の命令フェッチ



- パイプラインにバブルを生じさせないために、条件分岐命令をフェッチした時に次の3つを予測 (prediction) する。
  - フェッチしている命令が分岐命令か？
  - 分岐先アドレス(32bit)
  - 分岐方向
    - 分岐成立(Taken)か分岐不成立(Not taken, Untaken)か？
- Speculation assuming the prediction is true
  - No penalty by a branch instruction when the prediction hits
  - Recovery when the miss is detected

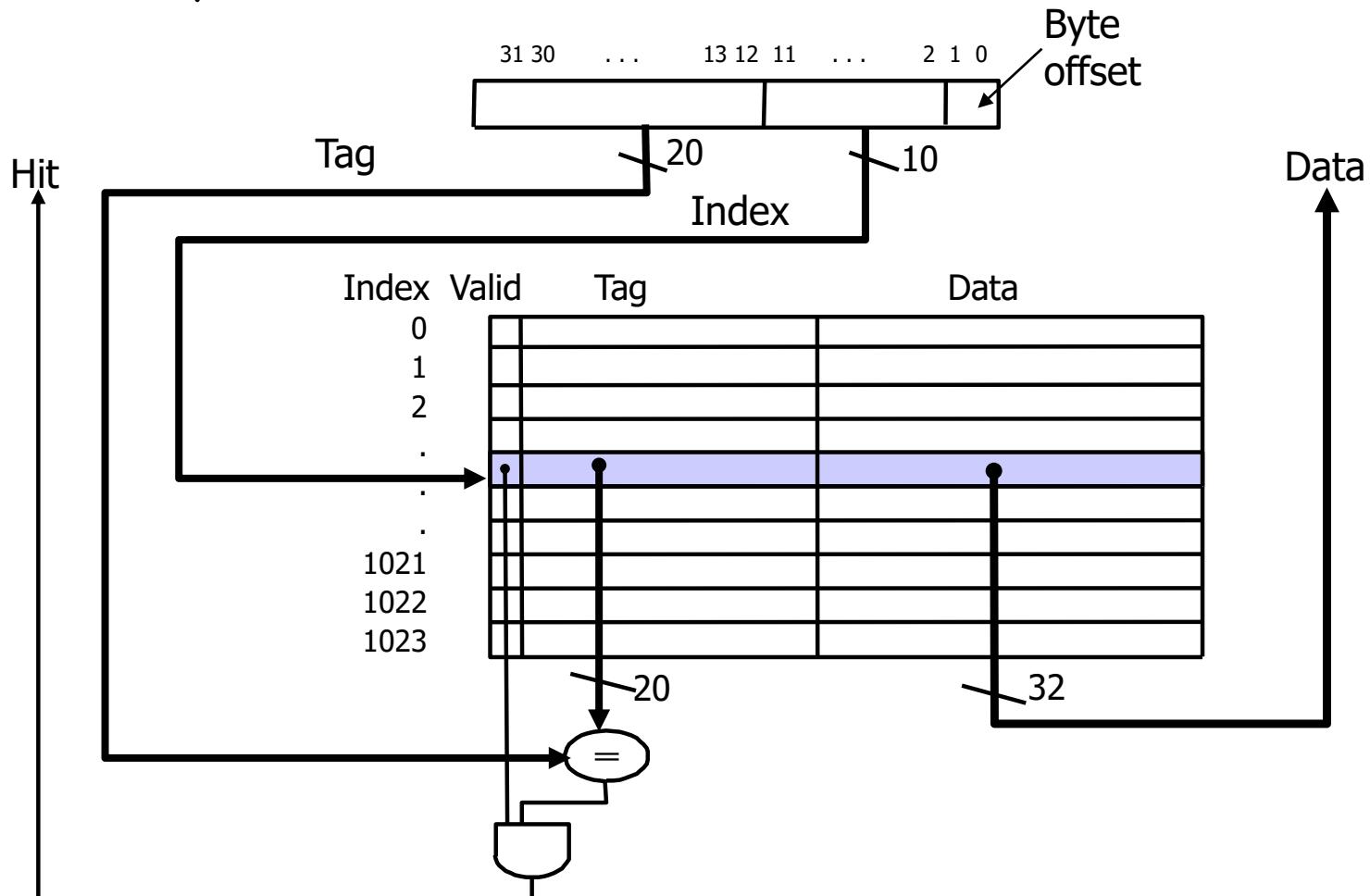


# BTB (Branch Target Buffer)

- フェッチしている命令が分岐命令か？
- 分岐先アドレス(32bit)

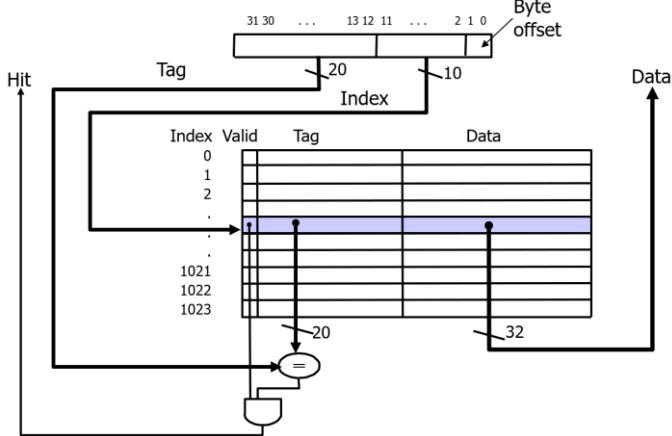
# Direct Mapped Cache Example

- One word/block, cache size = 1K words

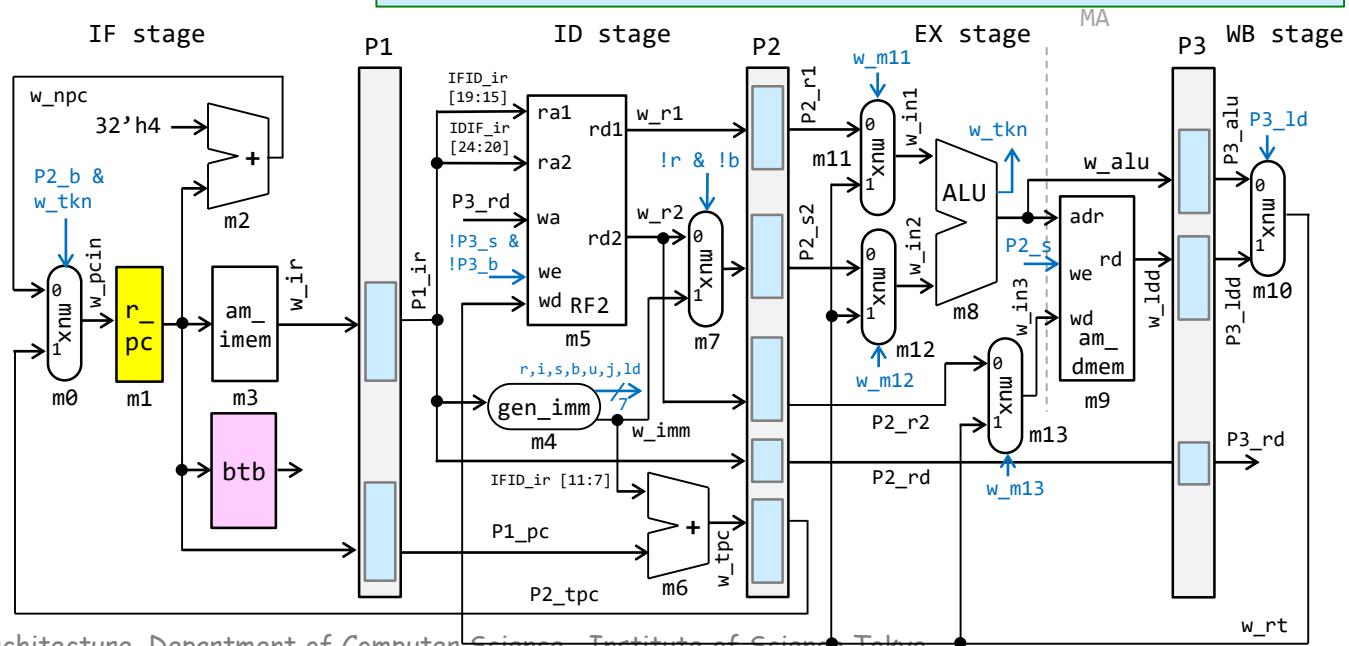


*What kind of locality are we taking advantage of?*

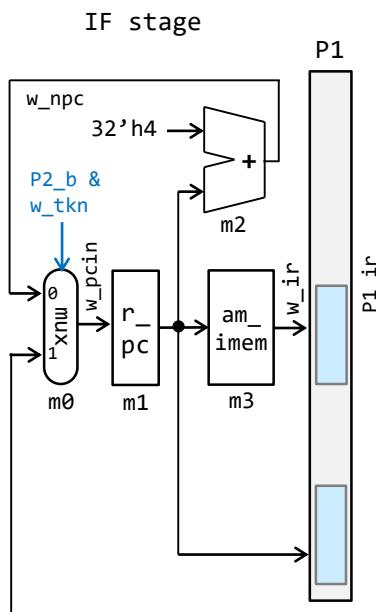
# Branch Target Buffer (BTB)



```
module m_btb (
    input wire      w_clk, w_we,
    input wire [31:0] w_adr,
    input wire [4:0] w_wadr,
    input wire [57:0] w_wd,
    output wire     w_hit,
    output wire [31:0] w_dout
);
    wire w_v;
    wire [24:0] w_tag;
    reg [57:0] mem [0:31];
    always @(posedge w_clk) if (w_we) mem[w_wadr] <= w_wd;
    assign {w_v, w_tag, w_dout} = mem[w_adr[6:2]];
    assign w_hit = w_v & (w_tag==w_adr[31:7]);
    integer i; initial for (i=0; i<32; i=i+1) mem[i] = 0;
endmodule
```



常に「分岐が成立しない」と予測。



```

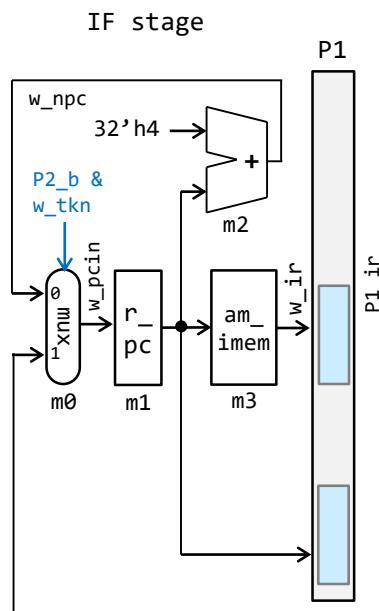
module m_proc8s(w_clk);
    input wire w_clk;
    reg [31:0] P1_ir=32'h13, P1_pc=0, P2_pc=0, P3_pc=0;
    reg [31:0] P2_r1=0, P2_s2=0, P2_r2=0, P2_tpc=0;
    reg [31:0] P3_alu=0, P3_ldd=0;
    reg P2_r=0, P2_s=0, P2_b=0, P2_ld=0, P3_s=0, P3_b=0, P3_ld=0;
    reg [4:0] P2_rd=0, P2_rs1=0, P2_rs2=0, P3_rd=0;
    reg P1_v=0, P2_v=0, P3_v=0;
    wire [31:0] w_npc, w_ir, w_imm, w_r1, w_r2, w_s2, w_rt;
    wire [31:0] w_alu, w_ldd, w_tpc, w_pcin, w_in1, w_in2, w_in3;
    wire w_r, w_i, w_s, w_b, w_u, w_j, w_ld, w_tkn;
    reg [31:0] r_pc = 0; // m1
    wire w_miss = P2_b & w_tkn & P2_v;
    assign w_pcin = (w_miss) ? P2_tpc : w_npc; // m0
    assign w_npc = r_pc + 32'h4; // m2
    m_am_imem m3 (r_pc, w_ir);
    m_gen_imm m4 (P1_ir, w_imm, w_r, w_i, w_s, w_b, w_u, w_j, w_ld);
    m_RF2 m5 (w_clk, P1_ir[19:15], P1_ir[24:20], w_r1, w_r2,
               P3_rd, !P3_s & !P3_b & P3_v, w_rt);
    assign w_tpc = P1_pc + w_imm; // m6
    assign w_s2 = (!w_r & !w_b) ? w_imm : w_r2; // m7
    always @(posedge w_clk) begin
        {P1_v, P2_v, P3_v} <= {!w_miss, !w_miss & P1_v, P2_v};
        {r_pc, P1_ir, P1_pc, P2_pc} <= {w_pcin, w_ir, r_pc, P1_pc};
        {P2_r1, P2_r2, P2_s2, P2_tpc} <= {w_r1, w_r2, w_s2, w_tpc};
        {P2_r, P2_s, P2_b, P2_ld} <= {w_r, w_s, w_b, w_ld};
        {P2_rs2, P2_rs1, P2_rd} <= {P1_ir[24:15], P1_ir[11:7]};
        {P3_pc, P3_ld} <= {P2_pc, P2_ld};
        {P3_alu, P3_ldd, P3_rd} <= {w_alu, w_ldd, P2_rd};
    end

```

|   |  |                         |
|---|--|-------------------------|
| 1 | `MM[0]={12'd0,5'd0,3'h0,5'd10,7'h13};          | // 00 addi x10,x0,0     |
| 2 | `MM[1]={12'd0,5'd0,3'h0,5'd3,7'h13};           | // 04 addi x3,x0,0      |
| 3 | `MM[2]={12'd101,5'd0,3'h0,5'd1,7'h13};         | // 08 addi x1,x0,101    |
| 4 | `MM[3]={7'd0,5'd3,5'd10,3'h0,5'd10,7'h33};     | // 0c L: add x10,x10,x3 |
| 5 | `MM[4]={12'd1,5'd3,3'h0,5'd3,7'h13};           | // 10 addi x3,x3,1      |
| 6 | `MM[5]={~12'd0,5'd1,5'd3,3'h1,5'b11001,7'h63}; | // 14 bne x3,x1,L       |
| 7 | `MM[6]=32'h00050f13;                           | // 18 HALT              |

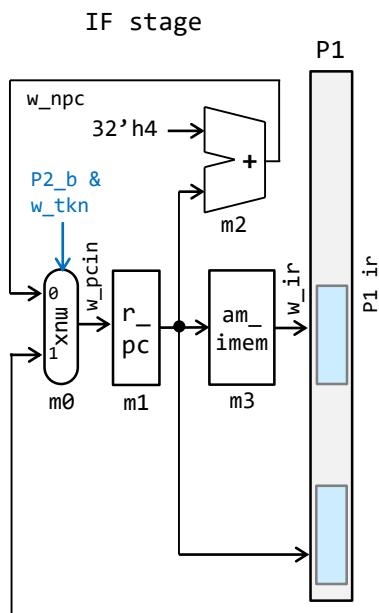
11  
13  
: P2\_s2; // m12

常に「分岐が成立する」と予測。



```
module m_proc8s(w_clk);
  input wire w_clk;
  reg [31:0] P1_ir=32'h13, P1_pc=0, P2_pc=0, P3_pc=0;
  reg [31:0] P2_r1=0, P2_s2=0, P2_r2=0, P2_tpc=0;
  reg [31:0] P3_alu=0, P3_ldd=0;
  reg P2_r=0, P2_s=0, P2_b=0, P2_ld=0, P3_s=0, P3_b=0, P3_ld=0;
  reg [4:0] P2_rd=0, P2_rs1=0, P2_rs2=0, P3_rd=0;
  reg P1_v=0, P2_v=0, P3_v=0;
  wire [31:0] w_npc, w_ir, w_imm, w_r1, w_r2, w_s2, w_rt;
  wire [31:0] w_alu, w_ldd, w_tpc, w_pcin, w_in1, w_in2, w_in3;
  wire w_r, w_i, w_s, w_b, w_u, w_j, w_ld, w_tkn;
  reg [31:0] r_pc = 0; // m1
  wire [31:0] w_ppc = 32'hc;
  wire w_bp_tkn = (r_pc==32'h14);
  wire [31:0] w_truepc = (P2_v & P2_b & w_tkn) ? P2_tpc : P2_pc+4;
  wire w_miss = P2_v & P2_b & P1_v & (P1_pc!=w_truepc);
  assign w_pcin = (w_miss) ? w_truepc : (w_bp_tkn) ? w_ppc : w_npc;
  assign w_npc = r_pc + 32'h4; // m2
  m_am_imem m3 (r_pc, w_ir);
  m_gen_imm m4 (P1_ir, w_imm, w_r, w_i, w_s, w_b, w_u, w_j, w_ld);
  m_RF2 m5 (w_clk, P1_ir[19:15], P1_ir[24:20], w_r1, w_r2,
             P3_rd, !P3_s & !P3_b & P3_v, w_rt);
  assign w_tpc = P1_pc + w_imm; // m6
  assign w_s2 = (!w_r & !w_b) ? w_imm : w_r2; // m7
  always @(posedge w_clk) begin
    {P1_v, P2_v, P3_v} <= {!w_miss, !w_miss & P1_v, P2_v};
    {r_pc, P1_ir, P1_pc, P2_pc} <= {w_pcin, w_ir, r_pc, P1_pc};
    {P2_r1, P2_r2, P2_s2, P2_tpc} <= {w_r1, w_r2, w_s2, w_tpc};
    {P2_r, P2_s, P2_b, P2_ld} <= {w_r, w_s, w_b, w_ld};
    {P2_rs2, P2_rs1, P2_rd} <= {P1_ir[24:15], P1_ir[11:7]};
  end
endmodule
```

- 1 `MM[0]={12'd0,5'd0,3'h0,5'd10,7'h13}; // 00 addi x10,x0,0
- 2 `MM[1]={12'd0,5'd0,3'h0,5'd3,7'h13}; // 04 addi x3,x0,0
- 3 `MM[2]={12'd101,5'd0,3'h0,5'd1,7'h13}; // 08 addi x1,x0,101
- 4 `MM[3]={7'd0,5'd3,5'd10,3'h0,5'd10,7'h33}; // 0c L:add x10,x10,x3
- 5 `MM[4]={12'd1,5'd3,3'h0,5'd3,7'h13}; // 10 addi x3,x3,1
- 6 `MM[5]={~12'd0,5'd1,5'd3,3'h1,5'b11001,7'h63}; // 14 bne x3,x1,L
- 7 `MM[6]=32'h00050f13; // 18 HALT



```

module m_proc8s(w_clk);
    input wire w_clk;
    reg [31:0] P1_ir=32'h13, P1_pc=0, P2_pc=0, P3_pc=0;
    reg [31:0] P2_r1=0, P2_s2=0, P2_r2=0, P2_tpc=0;
    reg [31:0] P3_alu=0, P3_ldd=0;
    reg P2_r=0, P2_s=0, P2_b=0, P2_ld=0, P3_s=0, P3_b=0, P3_ld=0;
    reg [4:0] P2_rd=0, P2_rs1=0, P2_rs2=0, P3_rd=0;
    reg P1_v=0, P2_v=0, P3_v=0;
    wire [31:0] w_npc, w_ir, w_imm, w_r1, w_r2, w_s2, w_rt;
    wire [31:0] w_alu, w_ldd, w_tpc, w_pcin, w_in1, w_in2, w_in3;
    wire w_r, w_i, w_s, w_b, w_u, w_j, w_ld, w_tkn;
    reg [31:0] r_pc = 0; // m1
    wire [31:0] w_ppc;
    wire [57:0] w_btb_wd = {1'b1, P2_pc[31:7], P2_tpc};
    wire w_hit;
    m_btb m14 (w_clk, r_pc, w_hit, w_ppc, P2_pc[6:2], P2_v & P2_b, w_btb_wd);
    wire w_bp_tkn = (r_pc==32'h14) & w_hit;
    wire [31:0] w_truepc = (P2_v & P2_b & w_tkn) ? P2_tpc : P2_pc+4;
    wire w_miss = P2_v & P2_b & P1_v & (P1_pc!=w_truepc);
    assign w_pcin = (w_miss) ? w_truepc : (w_bp_tkn) ? w_ppc : w_npc;
    assign w_npc = r_pc + 32'h4; // m2
    m_am_imem m3 (r_pc, w_ir);
    m_gen_imm m4 (P1_ir, w_imm, w_r, w_i, w_s, w_b, w_u, w_j, w_ld);
    m_RF2 m5 (w_clk, P1_ir[19:15], P1_ir[24:20], w_r1, w_r2,
                P3_rd, !P3_s & !P3_b & P3_v, w_rt);
    assign w_tpc = P1_pc + w_imm; // m6
    assign w_s2 = (!w_r & !w_b) ? w_imm : w_r2; // m7
    always @(posedge w_clk) begin
        {P1_v, P2_v, P3_v} <= {!w_miss, !w_miss & P1_v, P2_v};
        {r_pc, P1_ir, P1_pc, P2_pc} <= {w_pcin, w_ir, r_pc, P1_pc};
    end

```

- 1 `MM[0]={12'd0,5'd0,3'h0,5'd10,7'h13}; // 00 addi x10,x0,0
- 2 `MM[1]={12'd0,5'd0,3'h0,5'd3,7'h13}; // 04 addi x3,x0,0
- 3 `MM[2]={12'd101,5'd0,3'h0,5'd1,7'h13}; // 08 addi x1,x0,101
- 4 `MM[3]={7'd0,5'd3,5'd10,3'h0,5'd10,7'h33}; // 0c L:add x10,x10,x3
- 5 `MM[4]={12'd1,5'd3,3'h0,5'd3,7'h13}; // 10 addi x3,x3,1
- 6 `MM[5]={~12'd0,5'd1,5'd3,3'h1,5'b11001,7'h63}; // 14 bne x3,x1,L
- 7 `MM[6]=32'h00050f13; // 18 HALT

# Gshare branch predictor

```
1 module m_gshare (
2     input wire w_clk, w_we,
3     input wire [4:0] w_wadr, w_radr,
4     input wire w_tkn,
5     input wire [4:0] w_r_brh, w_w_bhr,
6     output wire w_pred
7 );
8     reg [1:0] mem [0:31];
9     wire [1:0] w_data = mem[w_radr ^ w_r_bhr];
10    assign w_pred = w_data[1];
11
12    wire [1:0] w_cnt = mem[w_wadr ^ w_w_bhr];
13    always @(posedge w_clk) if (w_we)
14        mem[w_wadr ^ w_w_bhr] <= (w_cnt < 3 & w_tkn) ? w_cnt + 1 :
15                                (w_cnt > 0 & !w_tkn) ? w_cnt - 1 : w_cnt;
16    integer i; initial for (i=0; i<32; i=i+1) mem[i] = 1;
17 endmodule
```

