

Branch Prediction and Caches

Computer Architecture
ECE 6913

Brandon Reagen

What we'll cover today

- 1) Logistics
- 2) Review and finish branch prediction
- 3) Introduce caches

Logistics

- 1) Homework 2 is released, solutions next week
- 2) Office hours Friday 11am-12pm
- 3) Course feedback (please fill out!)
- 4) Midterm Oct 22
 - 1) Will be a take-home test
 - 2) Should take ~2.5 hours
 - 3) Test will be released on Oct 22 at 11AM

There will **not be a lecture this day**, use this time to complete the midterm

To accommodate everyone, you'll have **25 hours to complete** the test

Issued: Oct 22, 11am

Due: Oct 23, 12PM (noon)

Midterm material

- 1) MIPS ISA
- 2) MIPS datapath and pipelining
- 3) Hazards & forwarding
- 4) Performance metrics
- 5) Branch Prediction

We'll save caches for the final.

Hazards!

Hazard (formal) := “when the next instruction cannot execute in the following clock cycle”

Hazard (me) := “you need something you don’t have”

There are three of them..

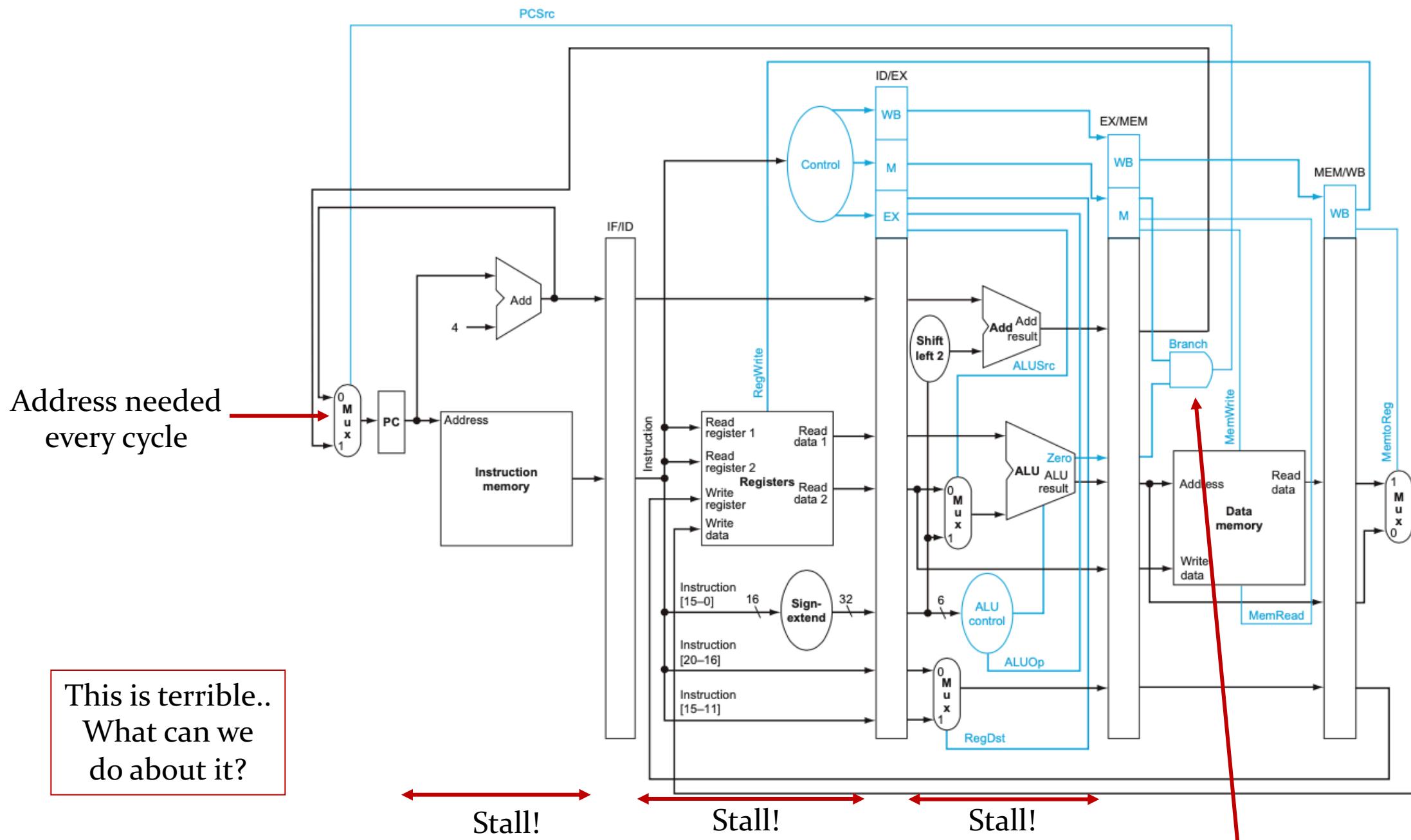
Think about what goes into a program and how they execute

- 1) Structural : hardware resources : SOLVED!
- 2) Data : need result from prior computation : Forwarding!
- 3) Control : need to know where to go next => Branch prediction (later)



NYU

TANDON SCHOOL
OF ENGINEERING



NYU

TANDON SCHOOL
OF ENGINEERING

Solution 1: Assume not taken

- Assume not taken doesn't really change the design
 - Fetch instructions each cycle from PC+4
 - Resolve branches in MEM stage
- Pro:
 - What happens to CPI when branch resolves NOT taken?
- Con:
 - What happens to CPI when branch IS taken? (Common case..)
 - How do we implement this?
 - Flush: Discard instructions due to misprediction.
 - Insert NOPs
 - How and where?



NYU

TANDON SCHOOL
OF ENGINEERING

Solution 2: Move branch resolution “up”

- Branches resolved in MEM stage... but why?
- We have all the information we need during ID!
 - Know it's a branch, which registers to compare, value to update PC
- What does this do?
 - Reduces the branch delay
 - This implies that when we flush, we waste fewer cycles
- MIPS branch conditions intentionally kept simple
 - E.g., beq and bne only need to compare values
 - We can do this by adding small circuits to the ID stage

Challenges with early branching (1)

- Computing branch target address
 - First, we already have the info we need in IF/ID pipeline reg:
PC and branch Immediate
 - Second, we can move the branch adder from EX stage to the ID stage

Challenges with early branching (2)

- Evaluating branch decision
 - New circuits late in the ID stage
 - After we read registers, have to compare values.
 - How can we implement comparators?
 - New bypass/forward logic!
 - Now forward data must also come into ID stage
 - Specifically, to the inputs of the equality unit
 - Can come from EX/MEM or MEM/WB pipeline regs
 - Requires new control logic
 - Cannot come from ID/EX regs, why?
 - Even with forwarding, we can stall. How?
 - If branch depends on previous ALU instruction: 1 cycle
 - If branch depends on previous MEM instruction: 2 cycles

Can only forward from the
beginning of cycle

What happens when pipelines get “deeper” and “wider”?

- Assume a pipeline resolves branches in stage 20 (not 2) and fetches 5 inst per cycle
- Also, 20% (1 in 5) instructions are branches, uniform distribution
- How long does it take to fetch 500 instructions?
 - 100% accuracy:
 - 100 cycles, no flushes => No wasted work!
 - 99% accuracy:
 - $100 \text{ (correct)} + 20 \text{ (penalty)} * 1 \text{ (wrong frequency * branches)} = \underline{120}$ cycles
 - 20% more cycles..
 - 95% accuracy
 - $100 \text{ (correct)} + 20 \text{ (penalty)} * 5 \text{ (wrong frequency * num branches)} = \underline{200}$ cycles
 - **100% extra instructions fetched.**



NYU

TANDON SCHOOL
OF ENGINEERING

Dynamic Hardware Branch Prediction

- Need to do better than always fetching not taken instructions
 - Performs poorly (think about loops)
- Branch prediction
 - **Guess** which direction to go!
- Dynamic
 - Use results of recently taken branches while program executes
- Hardware
 - Implement with special circuits for speed.
 - Would branch prediction be part of architecture or microarchitecture?
 - Could we do it in software?



NYU

TANDON SCHOOL
OF ENGINEERING

Why are Branches Predictable?

```
for (i=0; i<100; i++) {  
    ....  
}
```

```
addi r10, r0, 100  
add r1, r0, r0  
  
L1:  
....  
....  
addi r1, r1, 1  
bne r1, r10, L1  
....
```

```
if (aa==2)  
    aa = 0;  
if (bb==2)  
    bb = 0;  
if (aa!=bb)  
....
```

```
addi r2, r0, 2  
bne r10, r2, L_bb  
xor r10, r10, r10  
  
L_bb:  
bne r11, r2, L_xx  
xor r11, r11, r11  
  
L_xx:  
beq r10, r11, L_exit  
...  
Lexit:
```



NYU

TANDON SCHOOL
OF ENGINEERING

Branch predictors have two jobs

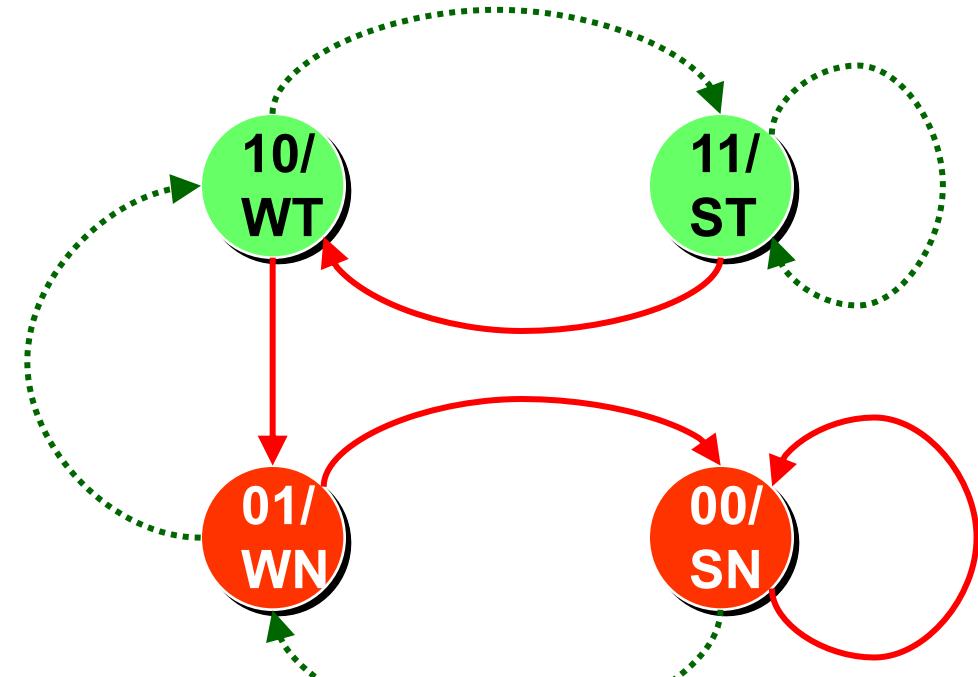
- In order to perform task, predictors must
 - Speculate whether the branch is to be **taken** or **not-taken**
 - Calculate the **target address**
 - Target address: Where to go when branch predicted taken
 - What about unconditional jumps?
- Past predicts the future
 - To make informed decisions, we'll track previous branch directions
 - As we'll see, this sounds simple but gets quite complex

This is what we'll focus on today.



2-bit Saturating Up/Down Counter Predictor

- Taken
- Not Taken
- Predict Not taken
- Predict taken



ST: Strongly Taken
WT: Weakly Taken
WN: Weakly Not Taken
SN: Strongly Not Taken

MSB tells us direction (taken/not taken)
LSB servers as hysteresis or inertia

The “Smith Predictor” (1981)

Rather than just using single 2b counter
for everything, create a **table of counters/FSMs**

Index this table using the branch address (PC)

Option 1: Use entire branch address

Too expensive!

Option 2: Use 1 entry/branch

Too expensive!

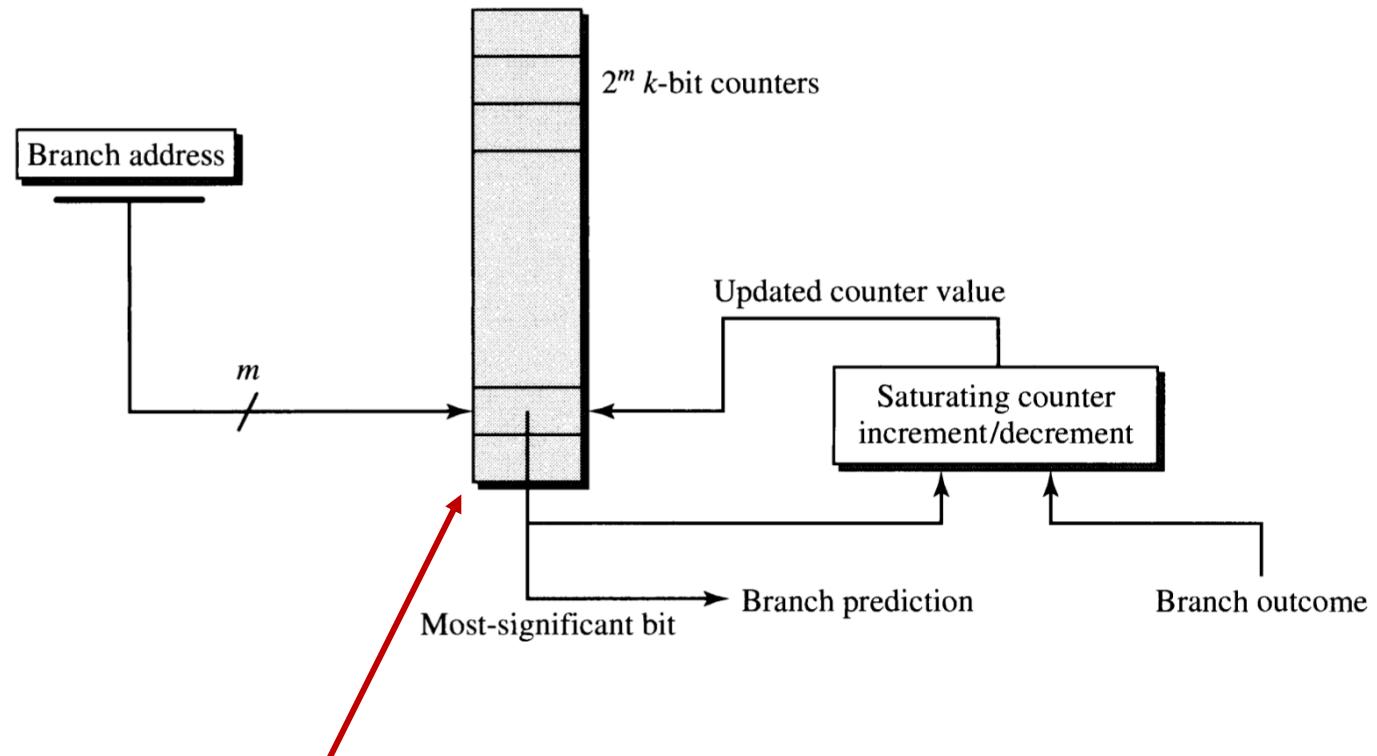
Option 3: Hash and hope collisions
aren't too bad

Option 3 is used today. Some uses actual hash
functions some take LSBs

General solution:

- m determines # entries (FSMs)
- k is the number of bits/entry

Total size of the predictor is: $k2^m$



Called the Pattern History Table (PHT)

Two –Level Prediction Tables

More advanced solution: Two-Level adaptive branch prediction

“Two-level” (Yeh and Patt, 91, 92, 93)

“Correlation branch prediction” (Pan, 92)

Key idea: Different branches tend to impact each other. We can make better predictions by using information from other recent branches (both locally and globally).

Smith’s predictor – Captures best branch decision

Two-level predictor – Also captures recent decision history

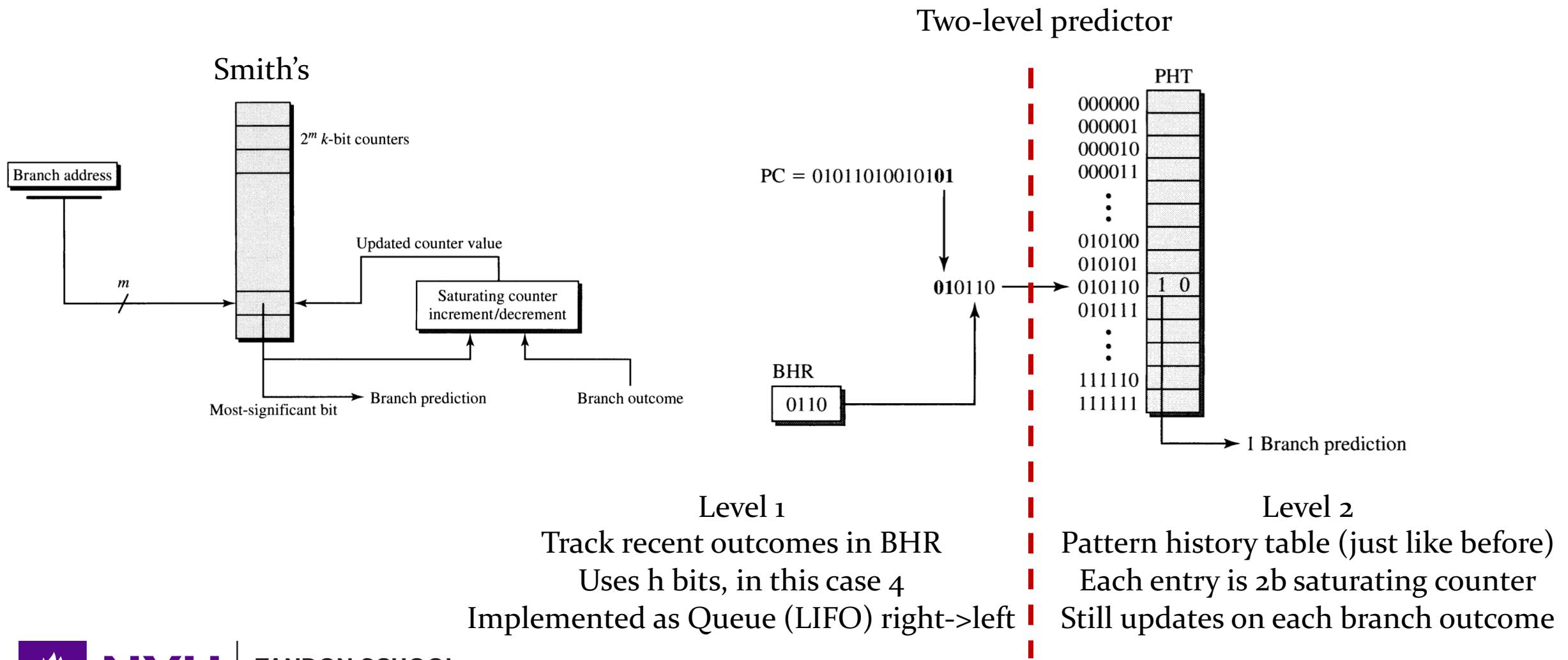


Key idea: Track global branch decisions

Branch history register (BHR)

- Use a single register to keep track of all recent branch outcomes
- Typically use a shift register
 - When branch resolves, shift in 1 for taken, 0 for not taken
- Use the BHR or BHR+PC to index PHT
 - Remember, PHT is just a bunch of FSMs
 - The BHR provides a different way of picking PHT entry

Example of a two-level predictor



NYU

TANDON SCHOOL
OF ENGINEERING

Indexing tradeoffs

How do we decide the PHT address?

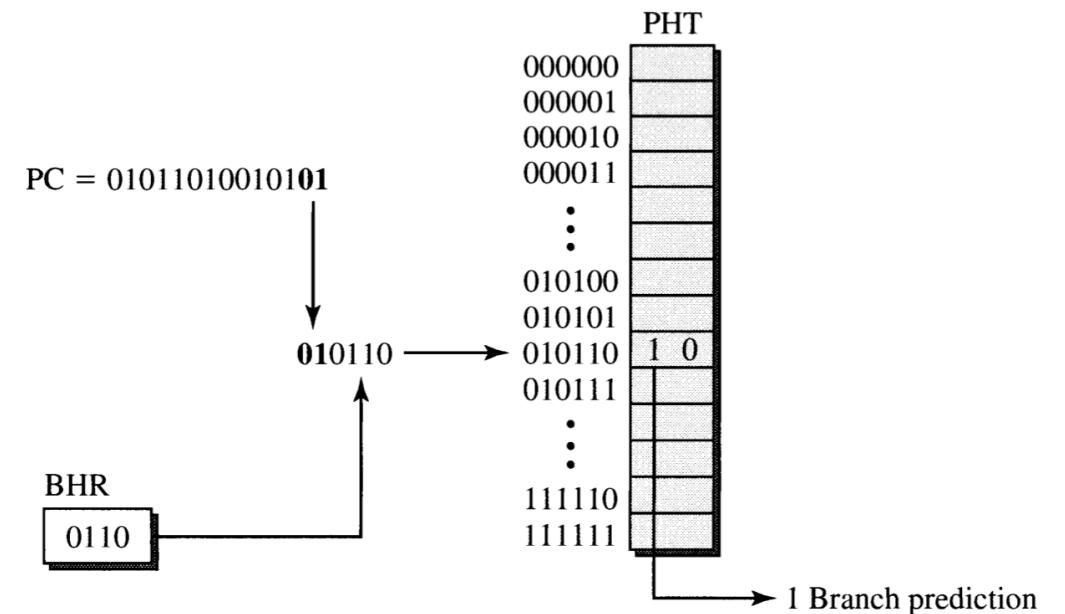
More PC bits implies fewer collisions

Collisions: different PC, same FSM

More BHR bits implies more correlation

Architects job is to balance these competing tradeoffs

What does Smith's predictor favor?



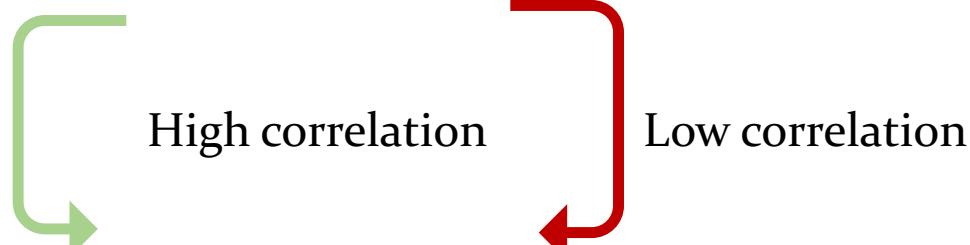
NYU

TANDON SCHOOL
OF ENGINEERING

Intuition

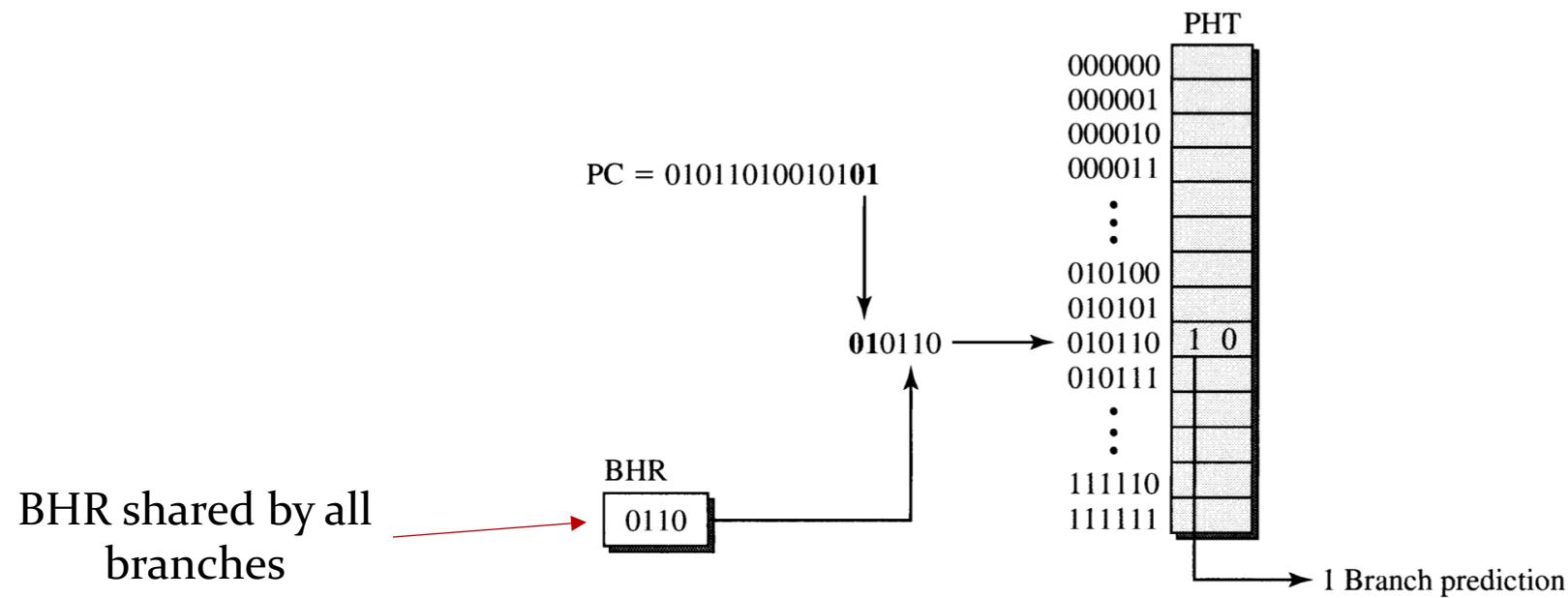
```
If (aa==2)  
    aa=0  
If (bb==2)  
    bb=0  
If (aa!=bb){  
    <do something>  
}
```

```
X = 0  
If (someCondition)  
    x=3  
If (someOtherCondition)  
    y+=19  
If (x <= 0){  
    <do something>  
}
```



BHR lets us capture these correlations!

“Global” two-level predictor



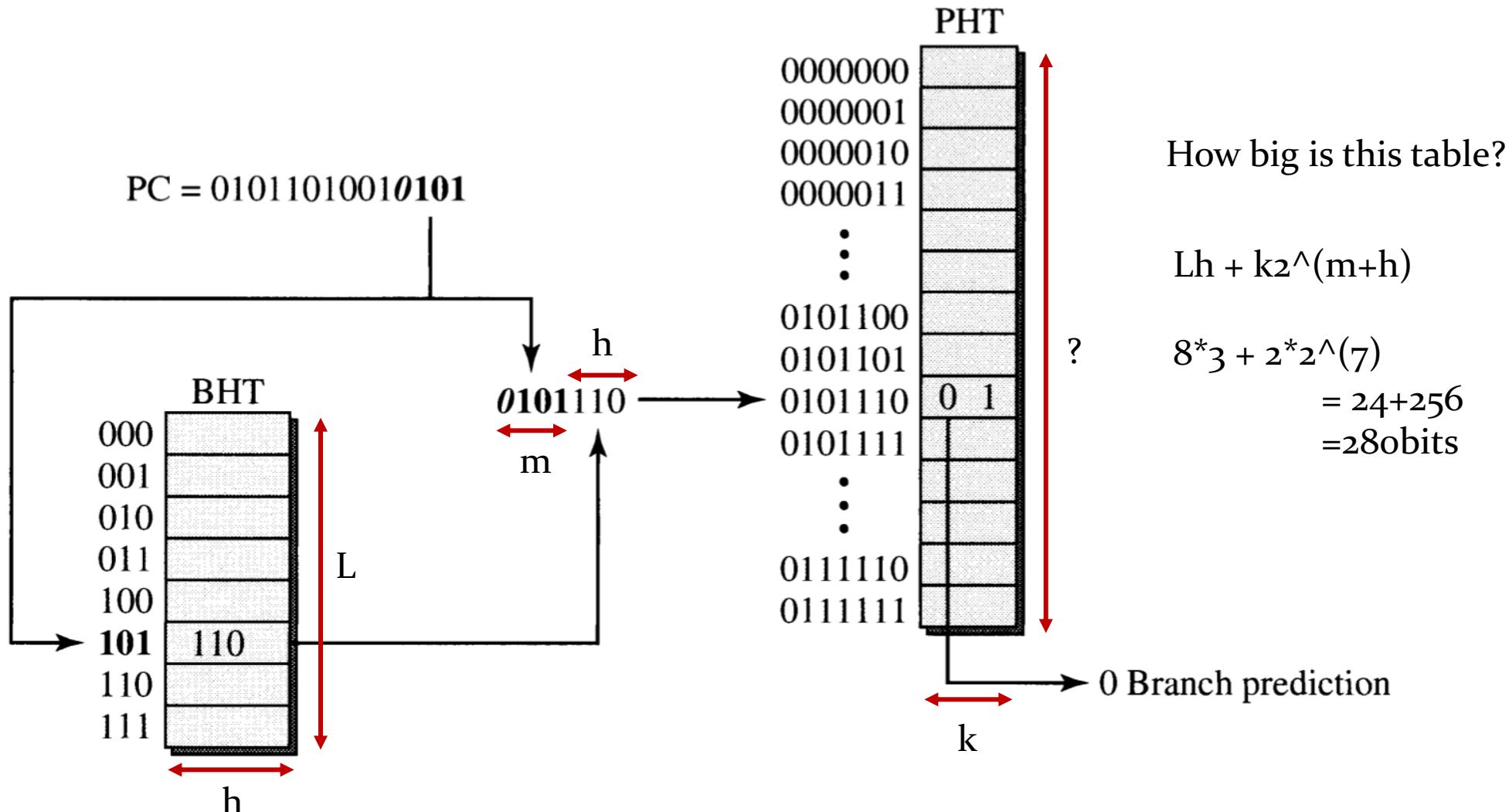
NYU

TANDON SCHOOL
OF ENGINEERING

“Local-history” two-level predictor

- Global BHR tracks history of all recent branches
- Local tracks outcomes of last several encounters of current branch
 - Generalization of the Global two-level predictor
 - Extend the single BHR to an array of BHRs
 - Called the **Branch History Table (BHT)**
 - Index the BHT using the LSBs of the PC
 - Each BHR tracks the branch history for that branch
 - Think of it as:
 - BHT index = look up the history pattern for this branch (PC)
 - BHT entry = the last few times this branch occurred, here's what happened
 - PHT index = Given I saw that pattern for that branch, what should I do?
 - PHT entry = I'll know what to do looking at the MSB of the 2b counter!

Local-History Two-Level Predictor



NYU

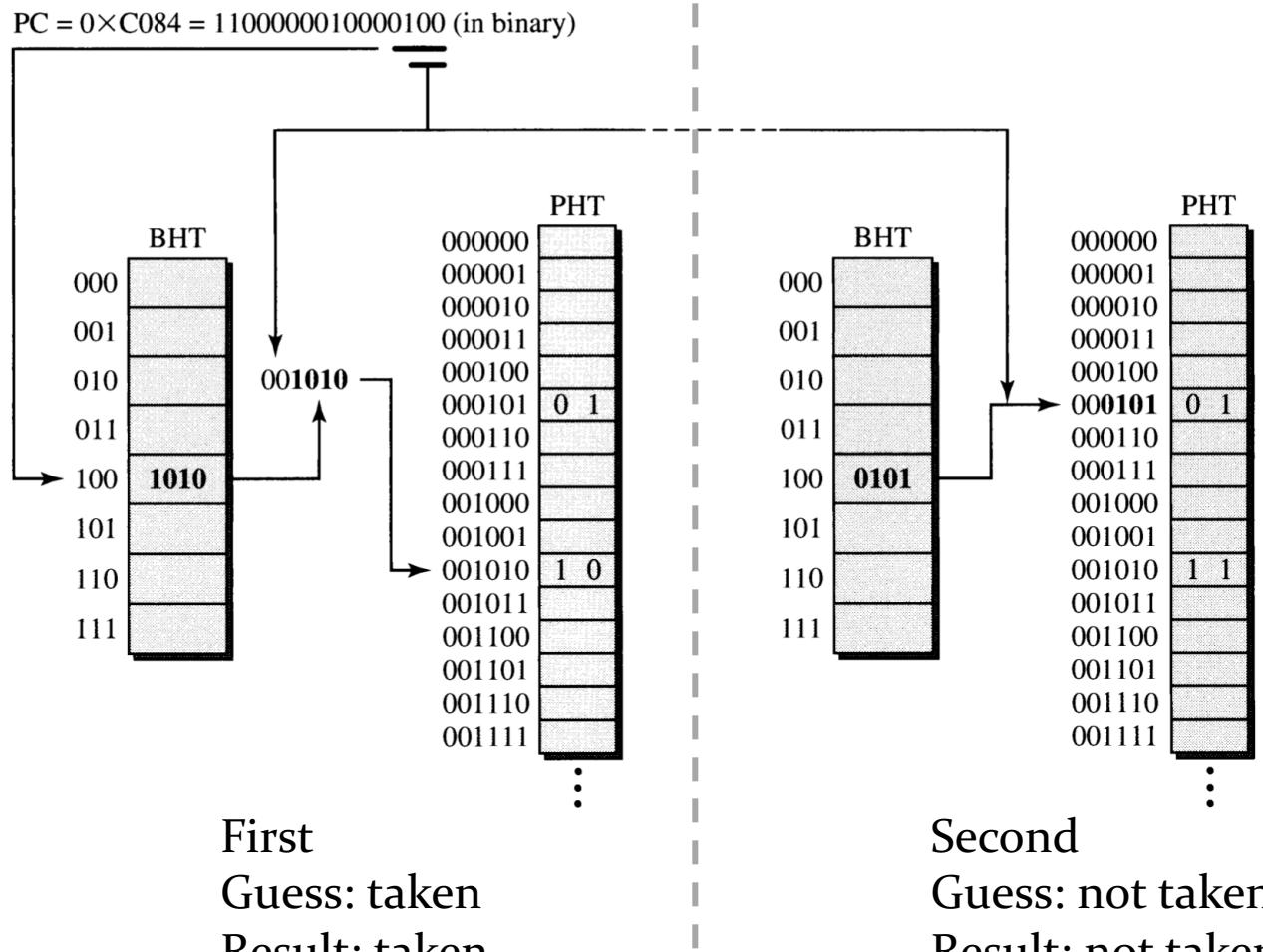
TANDON SCHOOL
OF ENGINEERING

Example 1: T, N, T, N, ...

- Simple pattern, is it simple to predict?
- What happens with 1b predictor initialized to not taken?
 - 0% accuracy, wrong every time!
- What about 2b saturating predictor initialized to Weakly T?
 - 50%
 - What if it's initialized to Strongly NT?
- What about a Locally-History Two-level predictor?

Example 1: Local-History 2-level predictor

T, NT, T, NT



What's going on here?

- 1) Always point to same BHR
- 2) PHT changes based on local branch history (key!)
- 3) Different FSMs tell us what to do depending on what we last saw

Intuition: Use both PC and pattern to make prediction.



NYU

TANDON SCHOOL
OF ENGINEERING

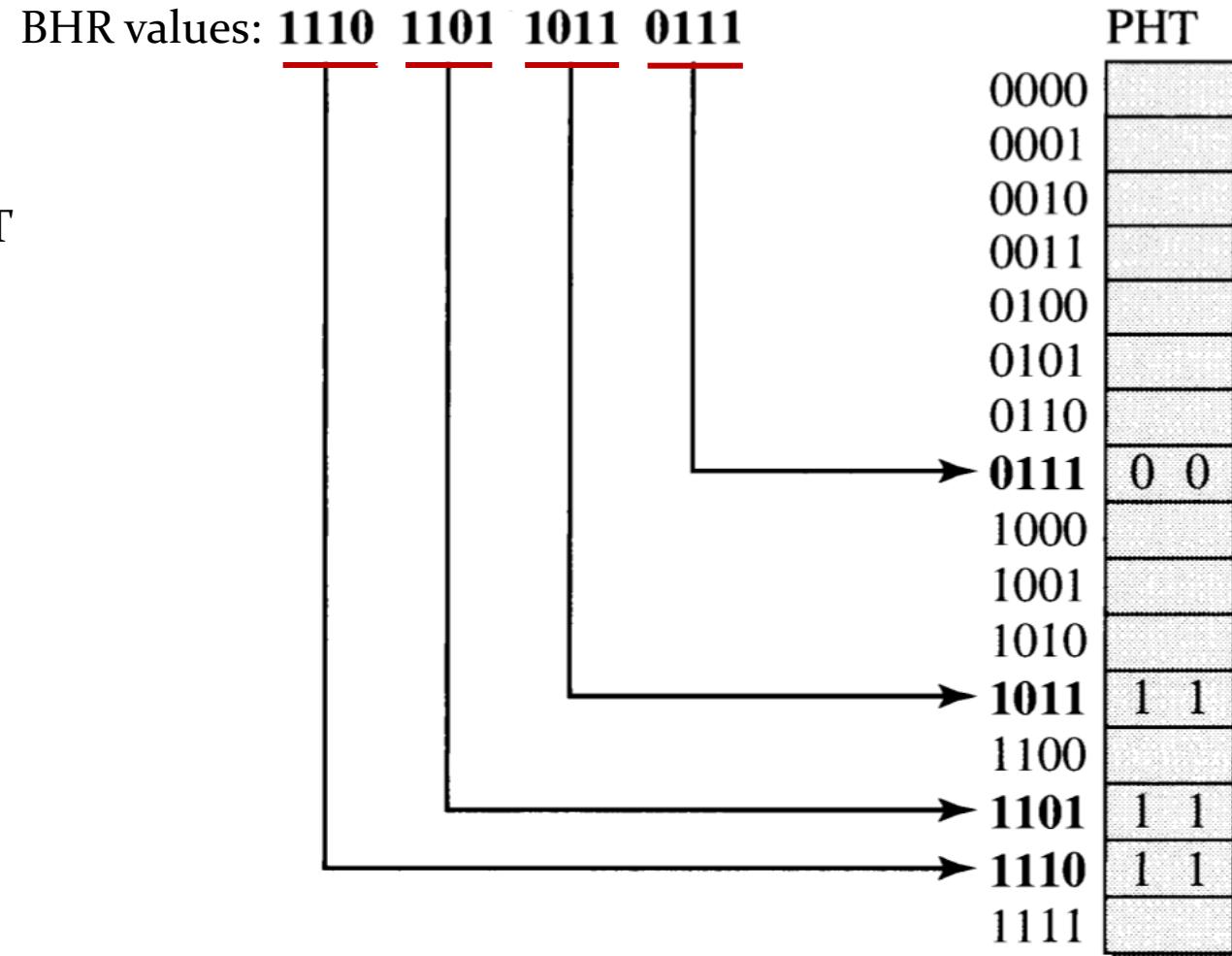
Example 2: 1110111011101

Say we only use BHR to index PHT
(no PC bits)

The branch pattern repeats every
four branches

If we have a PHT size 16, what
happens?

We can learn this pattern exactly!



NYU

TANDON SCHOOL
OF ENGINEERING

Taxonomy of branch predictors:

“Alternative Implementations of Two-Level Adaptive Branch Prediction”

Three types of Branch History Tables (BHTs)

- 1) Global (G): BHT is a single BHR shared by all branches
- 2) Set (S): Index BHR using hash function
- 3) Per address (P): Index BHR using PC address

“Global”

“Local”

Three types of Pattern History Tables (PHTs)

- 1) Global (g): 1 shared PHT, indexed by BHR
- 2) Set (s): PHT selected with branch hash, FSM selected w BHR
- 3) Per address (p): PHT selected with branch addr, FSM w BHR

“Local” and
“Global”

All 9 combos specified by $\langle X \rangle A \langle Y \rangle$ where $X = \{G, S, P\}$ and $Y = \{g, s, p\}$

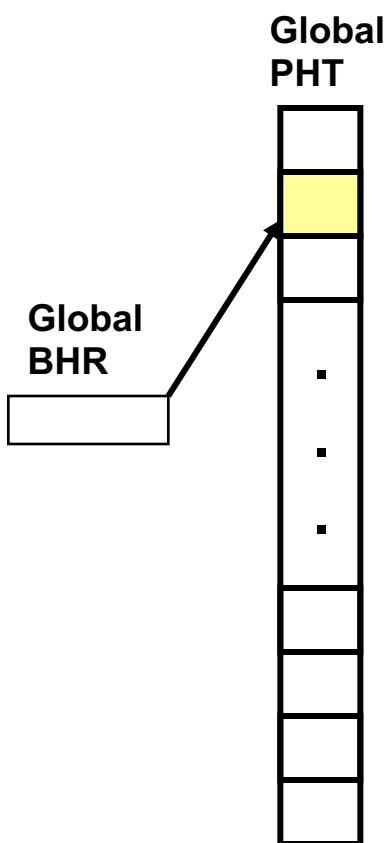


NYU

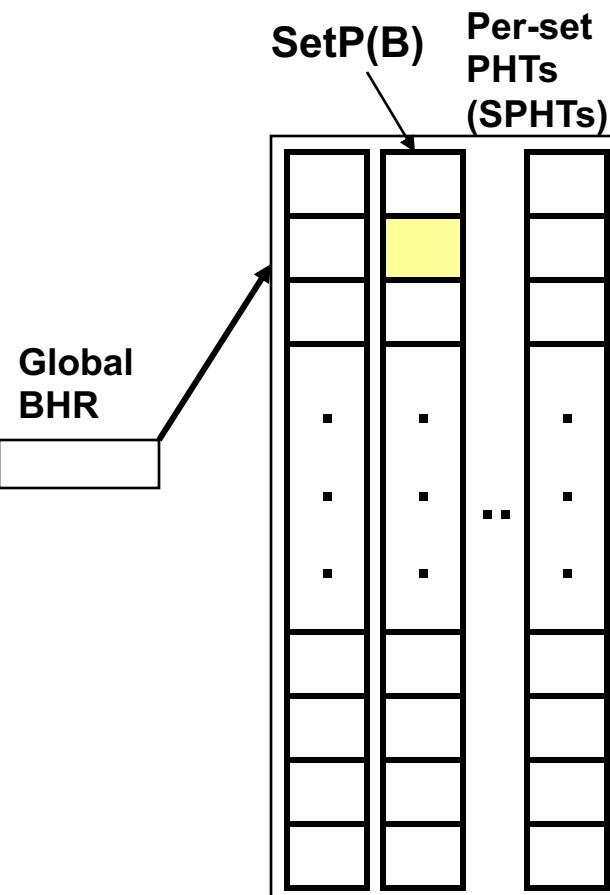
TANDON SCHOOL
OF ENGINEERING

Global History Schemes

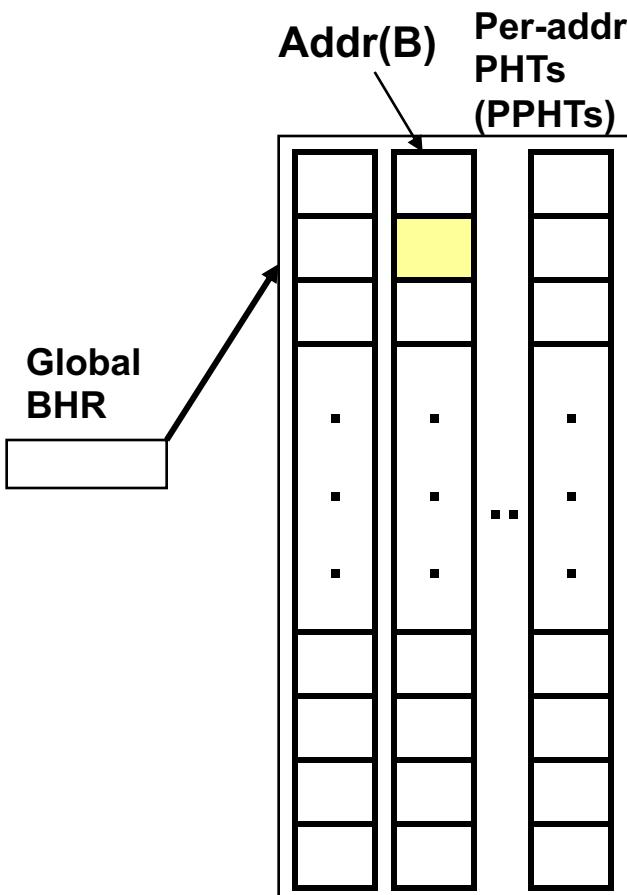
GAg



GAs



GAp



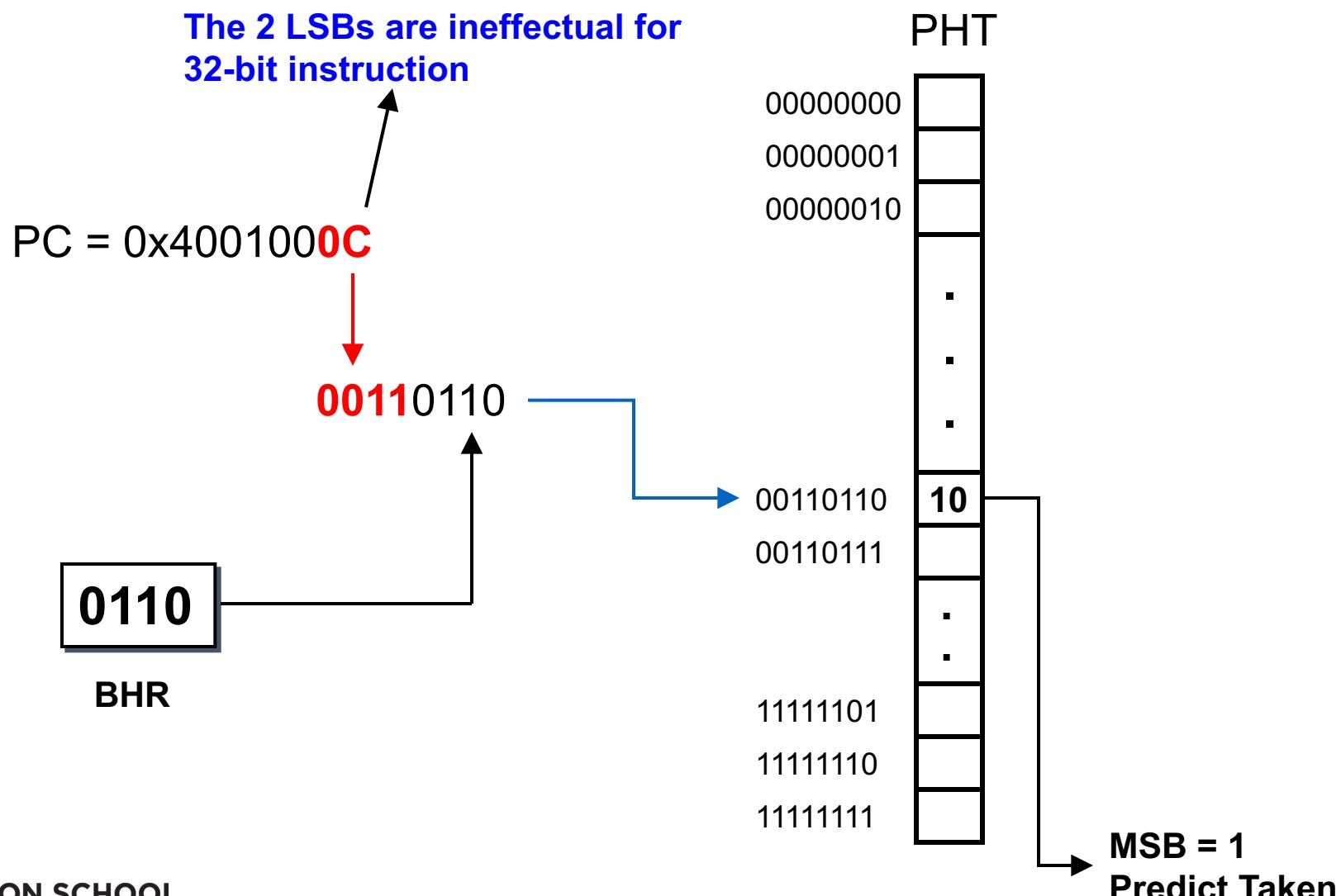
Set can be determined by branch
opcode, compiler classification,
or branch PC address.



NYU

TANDON SCHOOL
OF ENGINEERING

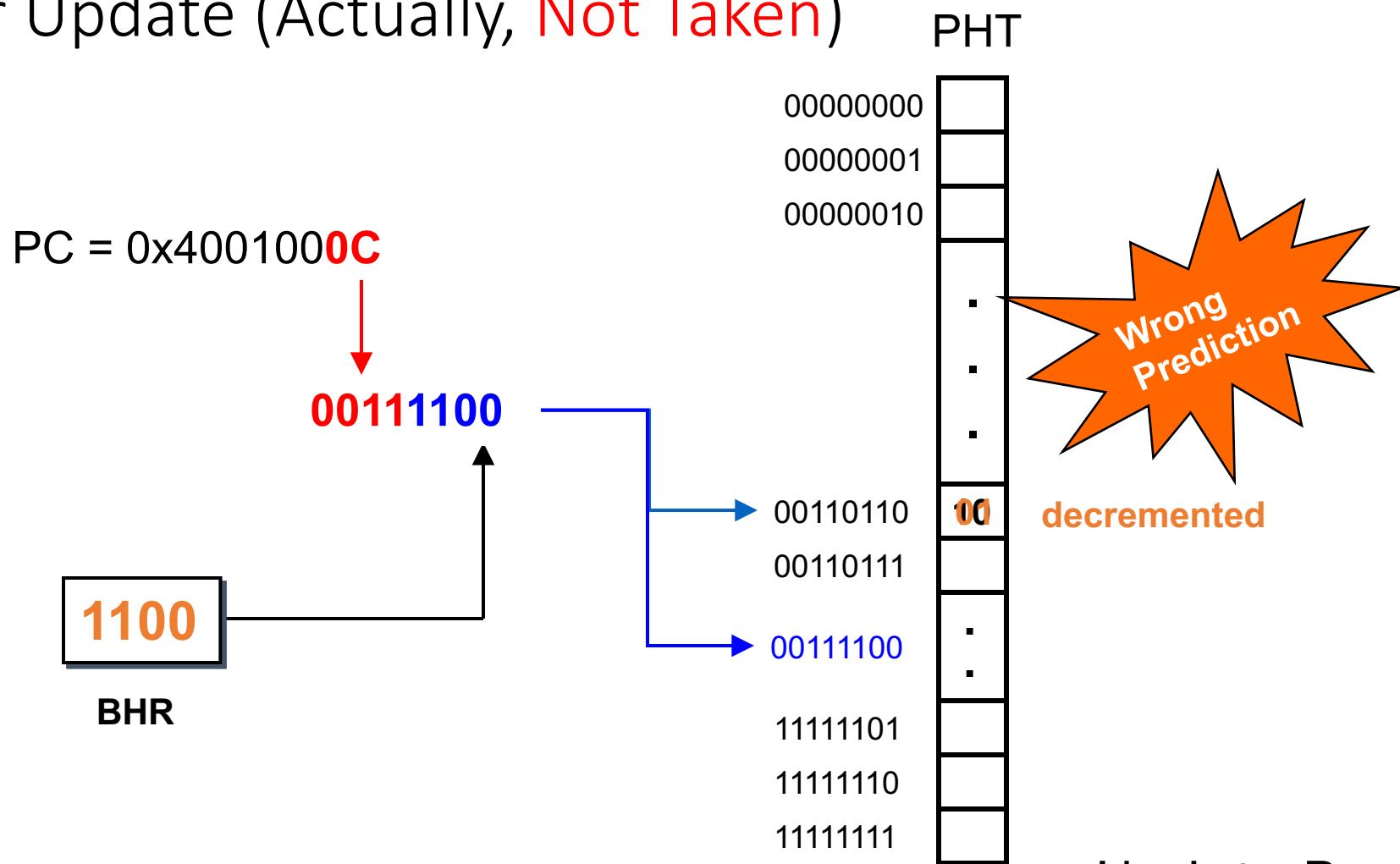
GAs Two-Level Branch Prediction



NYU

TANDON SCHOOL
OF ENGINEERING

Predictor Update (Actually, Not Taken)



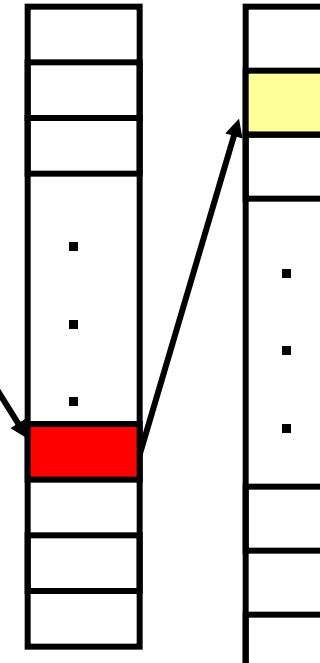
Update Predictor after
branch is resolved

Per-Address History Schemes

PAg

Per-addr
BHT (PBHT)

Addr(B)



Alpha 21264's
local predictor

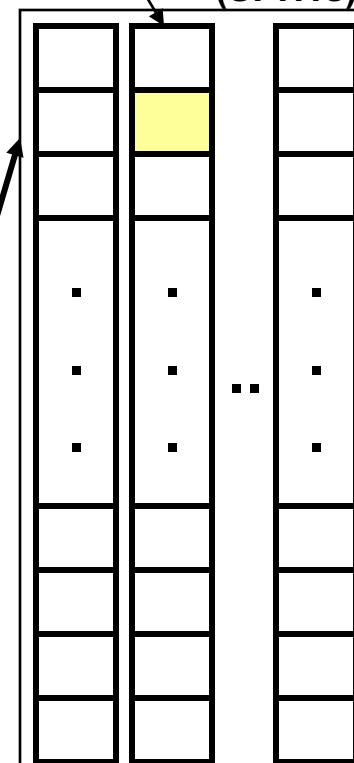
PAs

Global
PHT

Per-addr
BHT (PBHT)

Addr(B)

SetP(B)
Per-set
PHTs
(SPHTs)

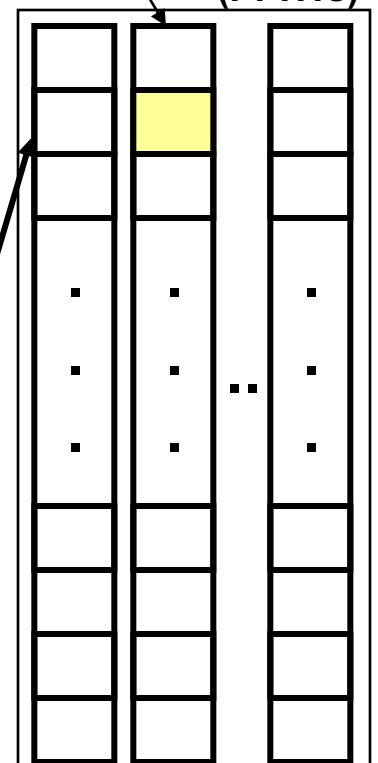


PAp

Per-set
BHT (PBHT)

Addr(B)

Addr(B)
Per-addr
PHTs
(PPHTs)



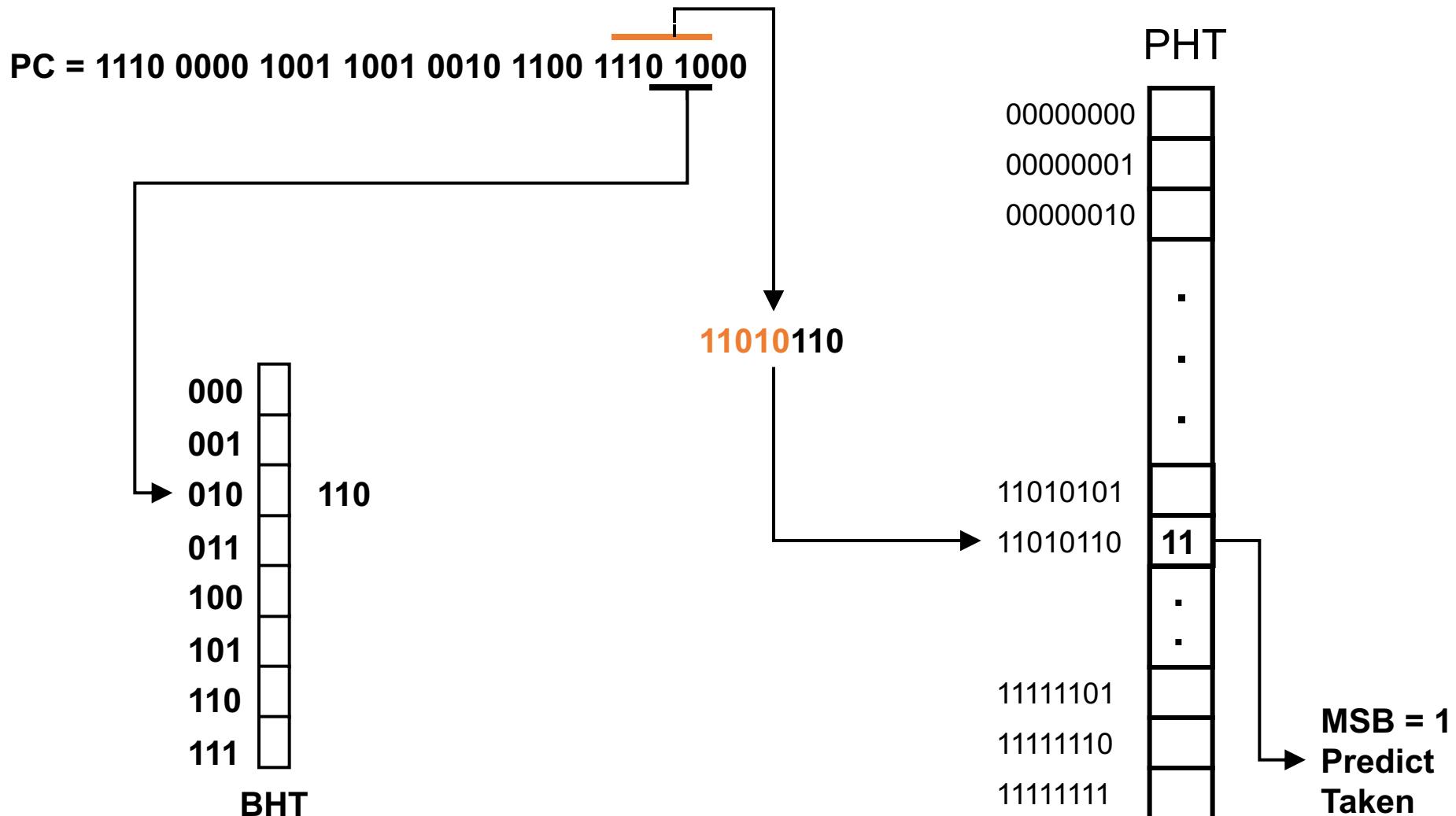
Ex: P6, Itanium



NYU

TANDON SCHOOL
OF ENGINEERING

PAP Two-Level Branch Predictor



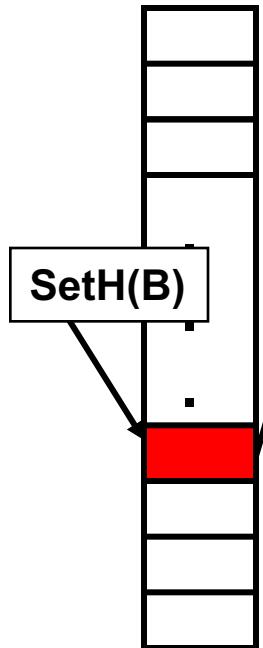
NYU

TANDON SCHOOL
OF ENGINEERING

Per-Set History Schemes

SAg

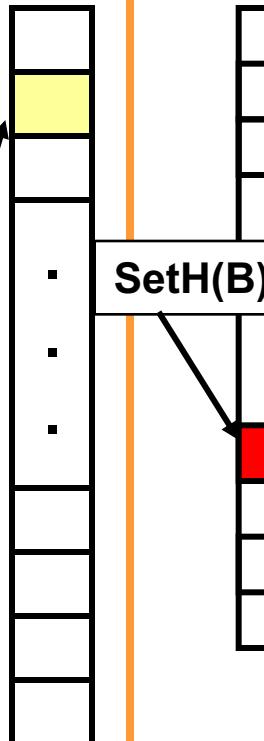
Per-set
BHT (SBHT)



SAs

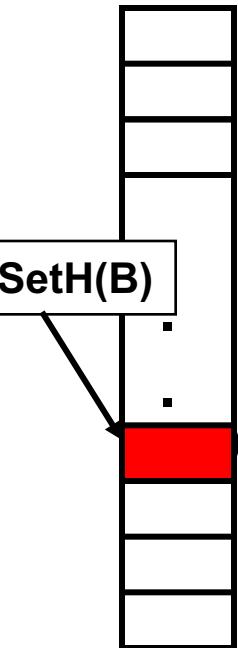
Global
PHT

Per-set
BHT (SBHT)

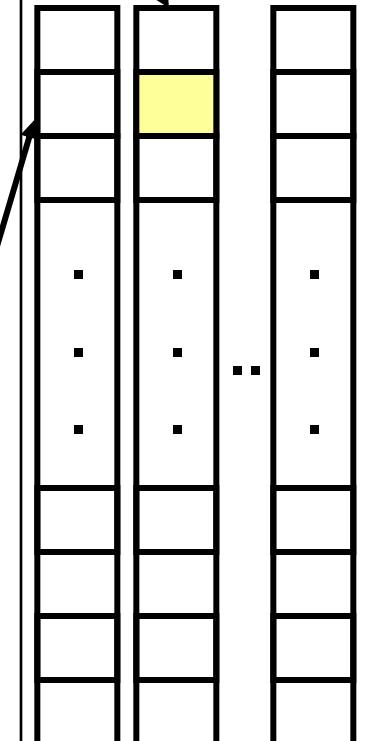


SAp

Per-set
BHT (SBHT)



Addr(B)
Per-addr
PHTs
(PPHTs)



NYU

TANDON SCHOOL
OF ENGINEERING

More tradeoffs

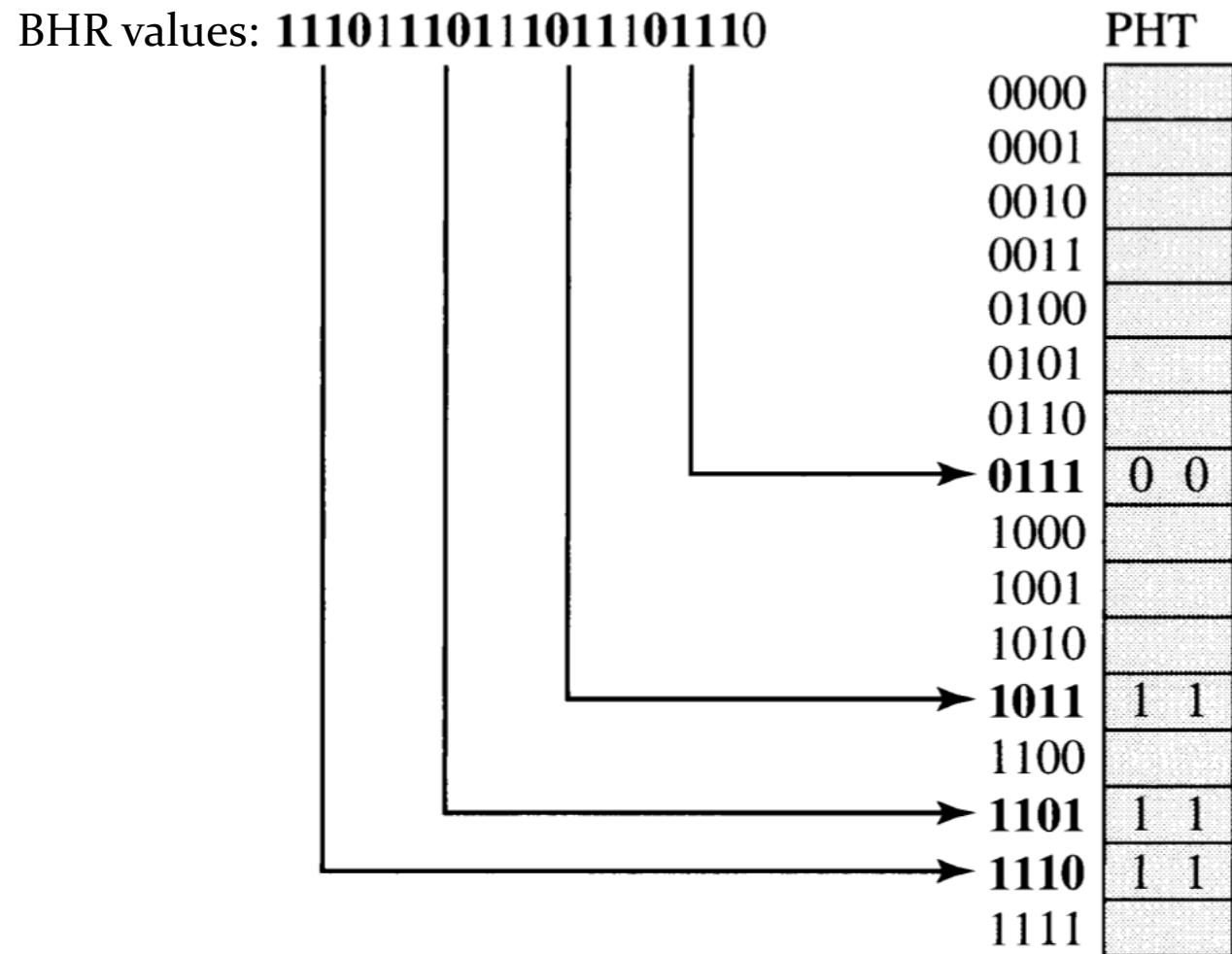
- How to index PHT?
 - Using more BHR bits enables better correlation
 - Fewer PC bits causes more PHT collisions/aliasing
- BHR-heavy indexing performs well but:
 - Causes aliasing
 - Results in poor PHT utilization
 - Why?
 - Normally only have a few common branch patterns, rest unused

Recall example

Only using 4 of 16 slots..
Poor utilization!
We can do better!

Q: What happens if another branch
aliases to the same PHT entries?

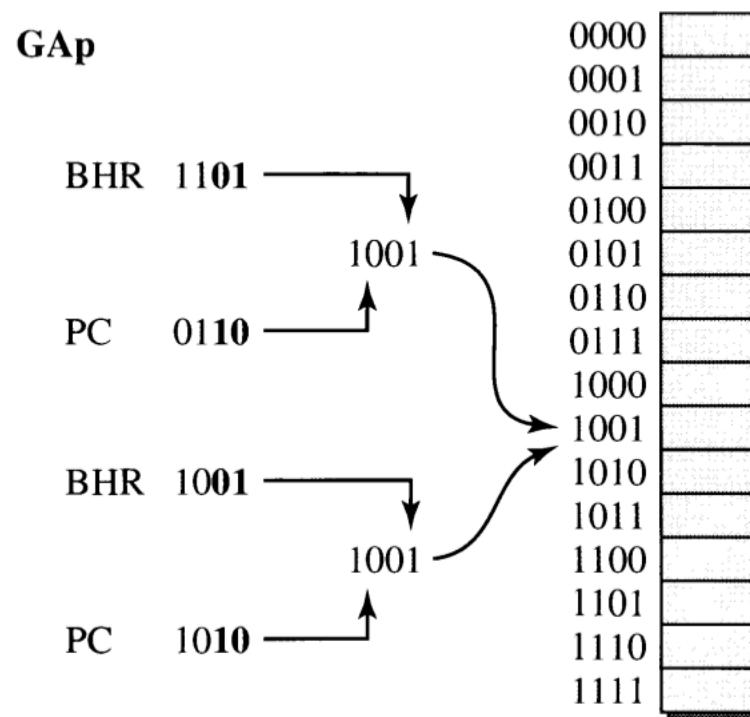
Everything works!
“Neutral Interference”



NYU

TANDON SCHOOL
OF ENGINEERING

Index-sharing: the gshare predictor



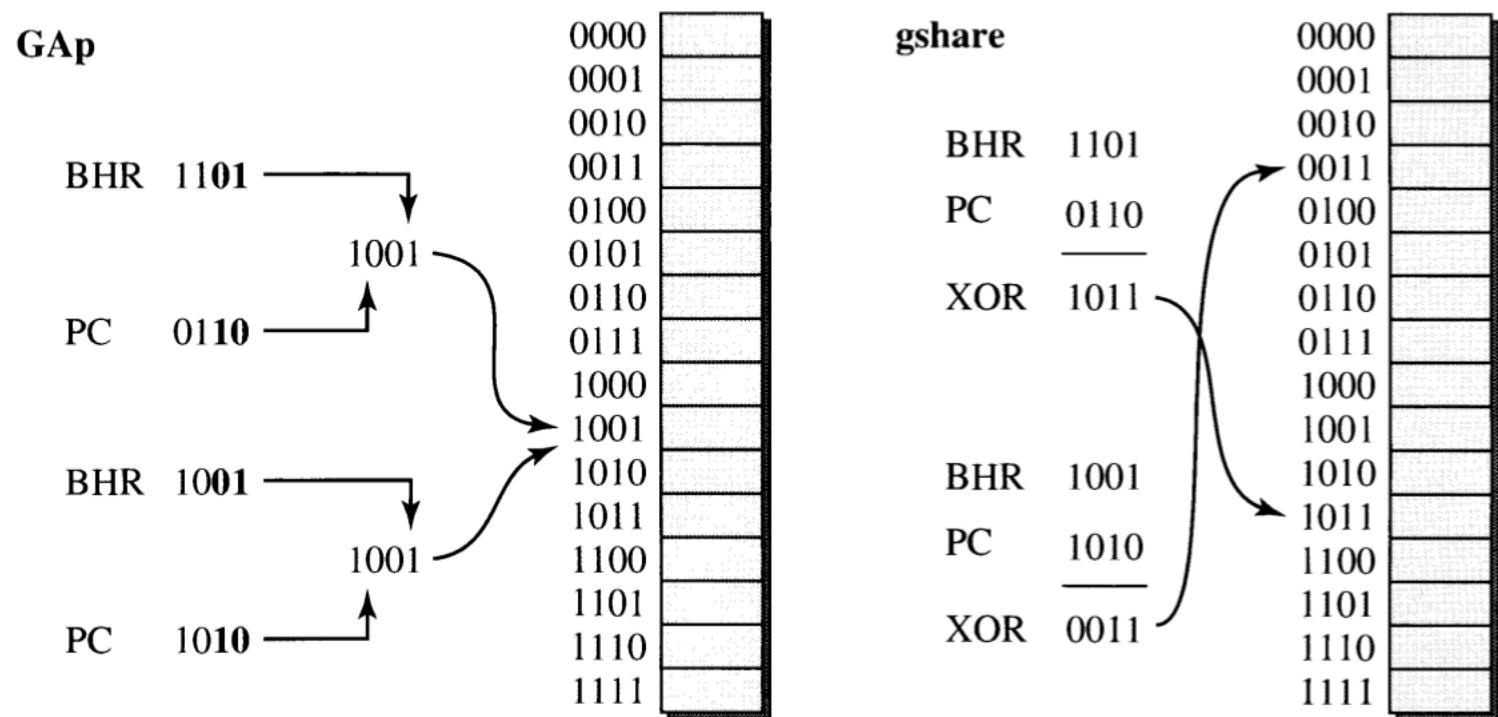
Index-sharing: the gshare predictor

Pros

- + Better utilize PHT
- + Resolves conflicts

Cons

- Can cause new conflicts
- Minor increase to critical path

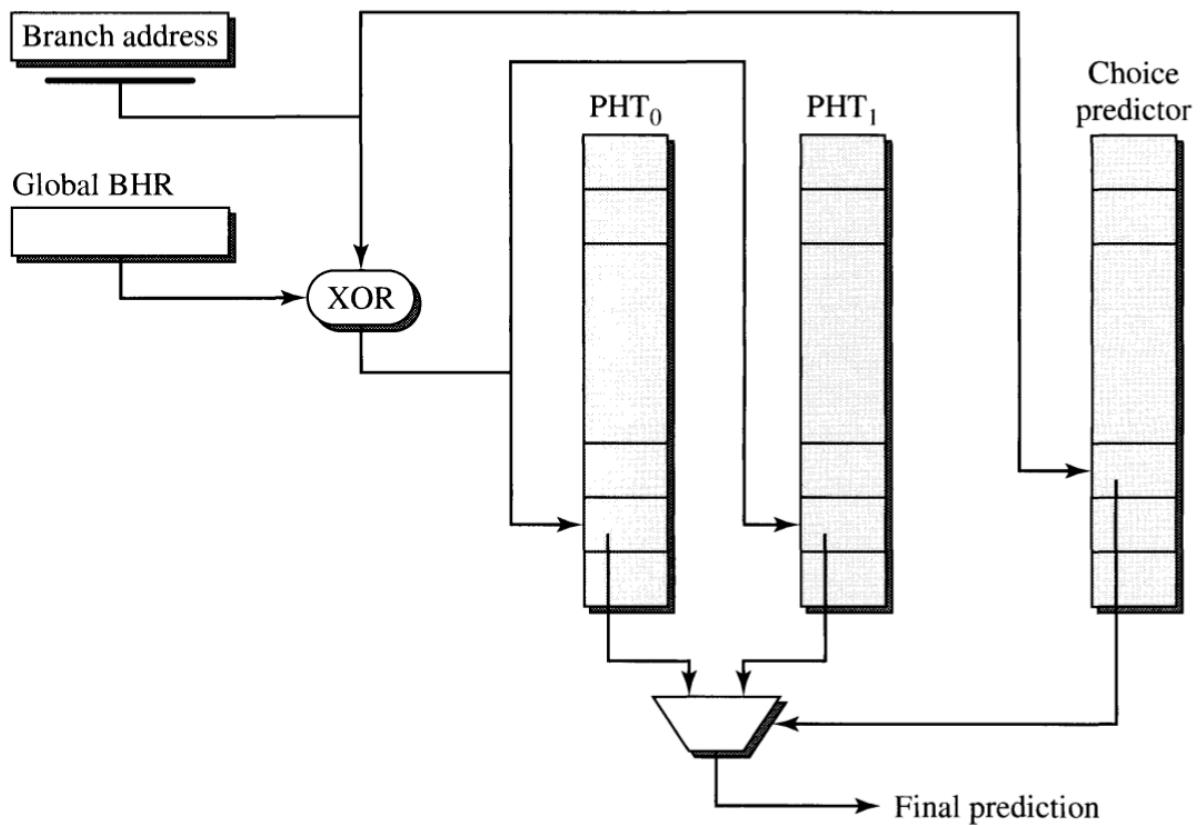


NYU

TANDON SCHOOL
OF ENGINEERING

Bi-Mode predictor

- Most branches are biased toward T/NT
 - PHTs learn bias
 - Choice predictor learns branch bias
- Turns negative conflicts into neutral
 - If collide in the same PHT,
Choice made it happen, so probably good
- Updating
 - Selected PHT always updated
 - Not selected bank not updated
 - Choice predictor
 - Always updated using branch direction
 - Unless Choice Pred. outcome was wrong
but selected PHT's guess was correct



NYU

TANDON SCHOOL
OF ENGINEERING

The gskew predictor

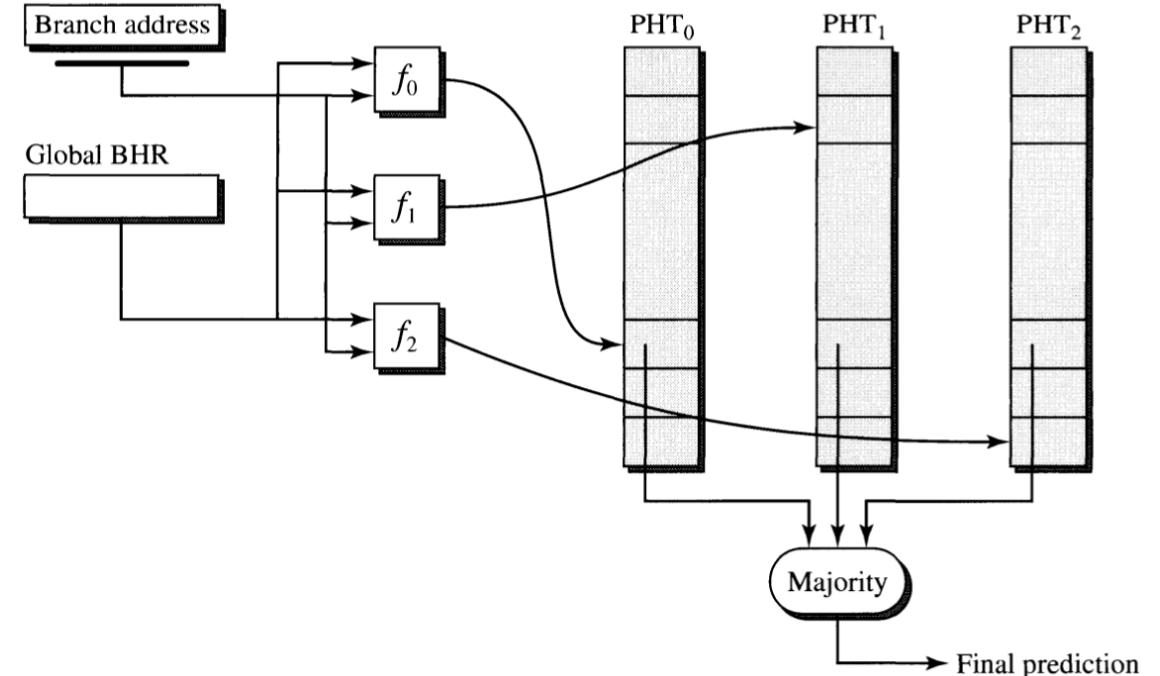
- Create 3+ PHT “banks”
- Hash Address-BHR pair
 - How does this compare to Bi-Mode?
- Read predictions from each banks, take majority
- Intuition:
 - Even if you have a conflict, the majority vote will mitigate negative impact
- Clever trick: f_0, f_1, f_2 can be created such that:
if $f_0(x) == f_1(x)$, then $f_1(x) != f_2(x), f_2(x) != f_0(x)$

What does this mean? If you have multiple PHT conflicts they're from independent branches

Two update policies:

- 1) Total update: update all PHTs after branch resolves
- 2) Partial update: If prediction correct but PHT wrong, don't update.

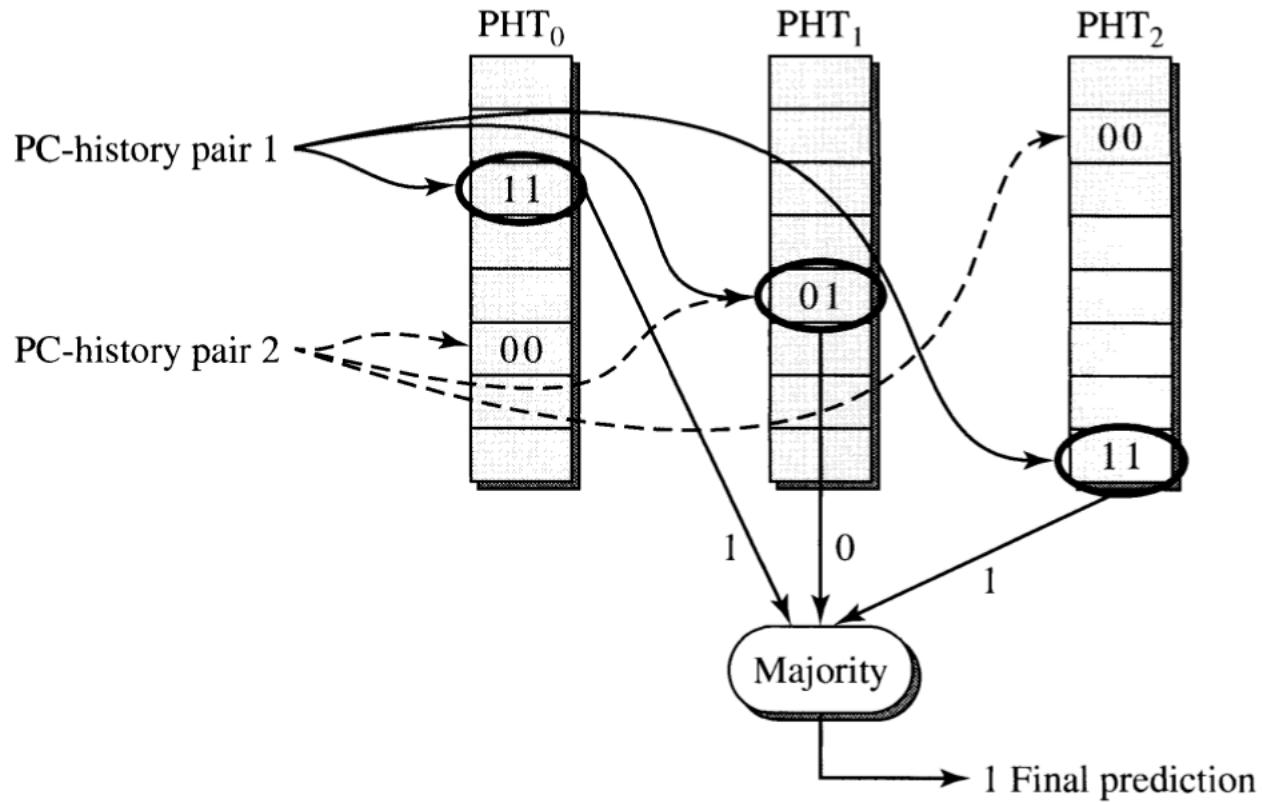
Partial tends to have best performance.



NYU

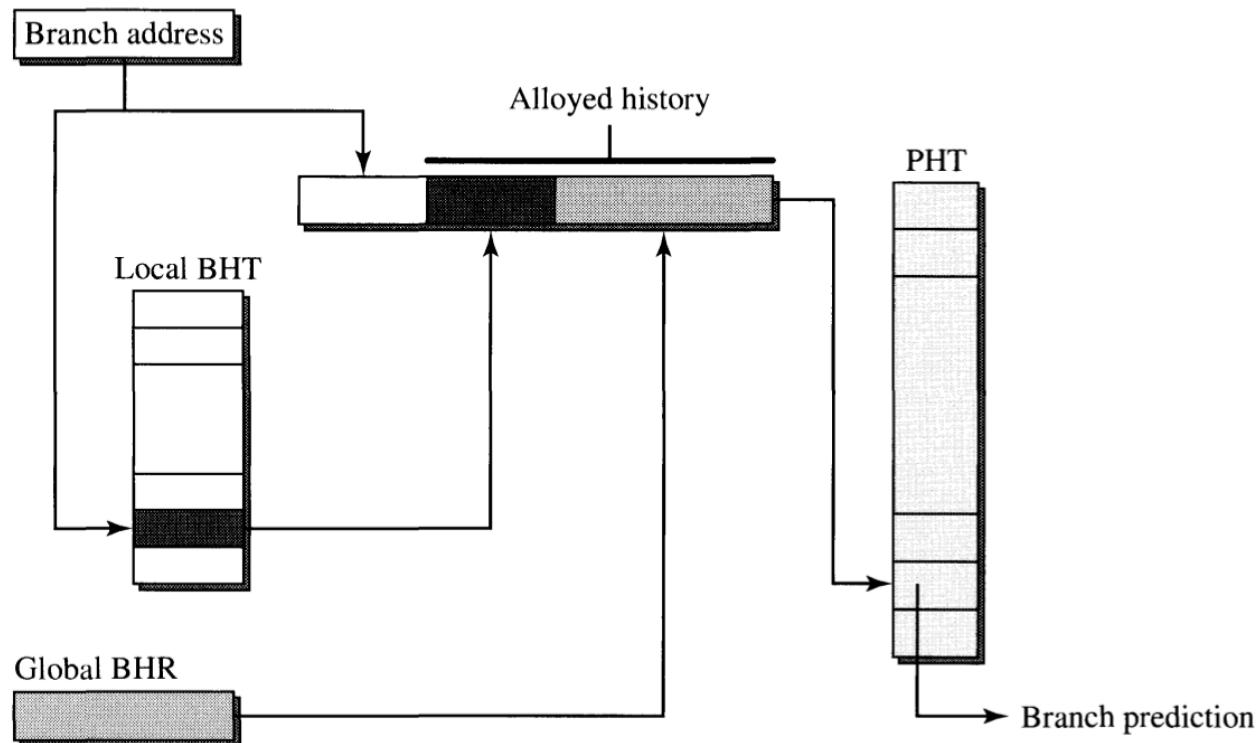
TANDON SCHOOL
OF ENGINEERING

Gskew example



The Alloyed History Predictor

- Some branches need global history (1 BHR)
some need local history (BHT)
What to do?
 - Combine them!
- Alloyed predictor tracks both global and local branch history
- Three benefits:
 - Cheap!
 - Best of both scopes
 - Can correlate even more, since some branches need context from both

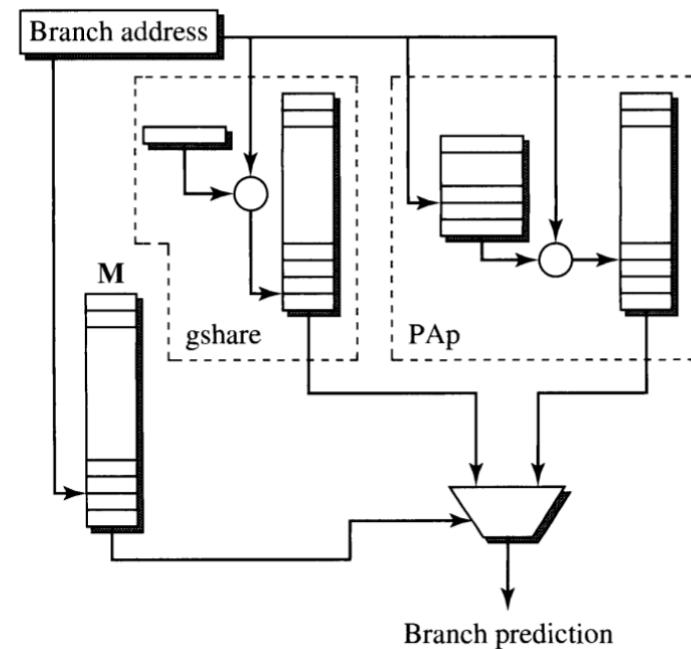
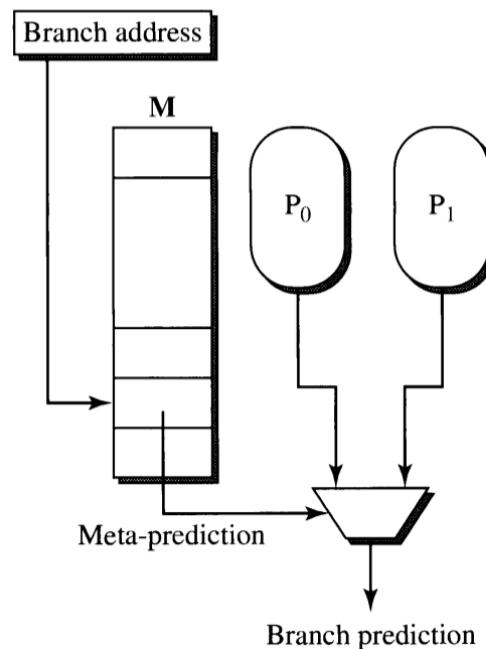


NYU

TANDON SCHOOL
OF ENGINEERING

Tournament predictors

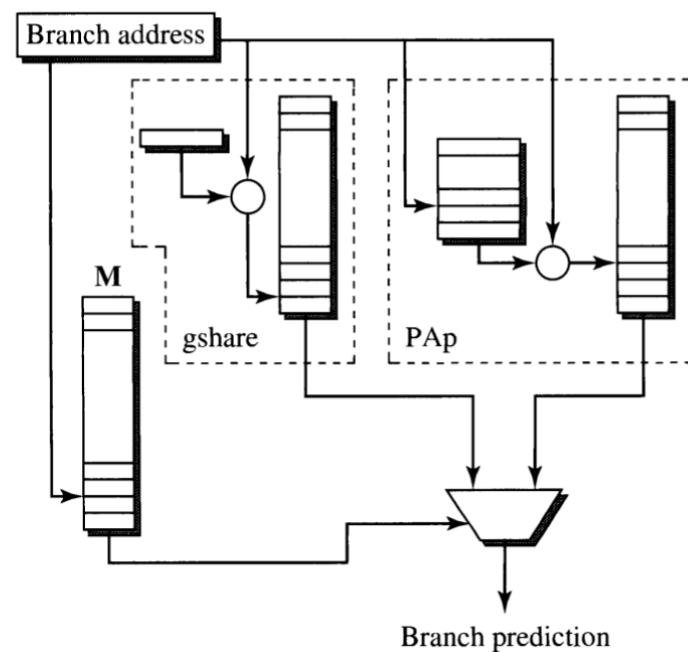
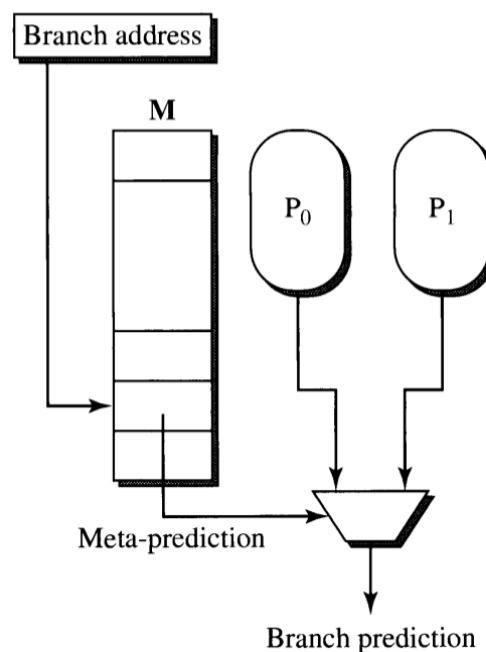
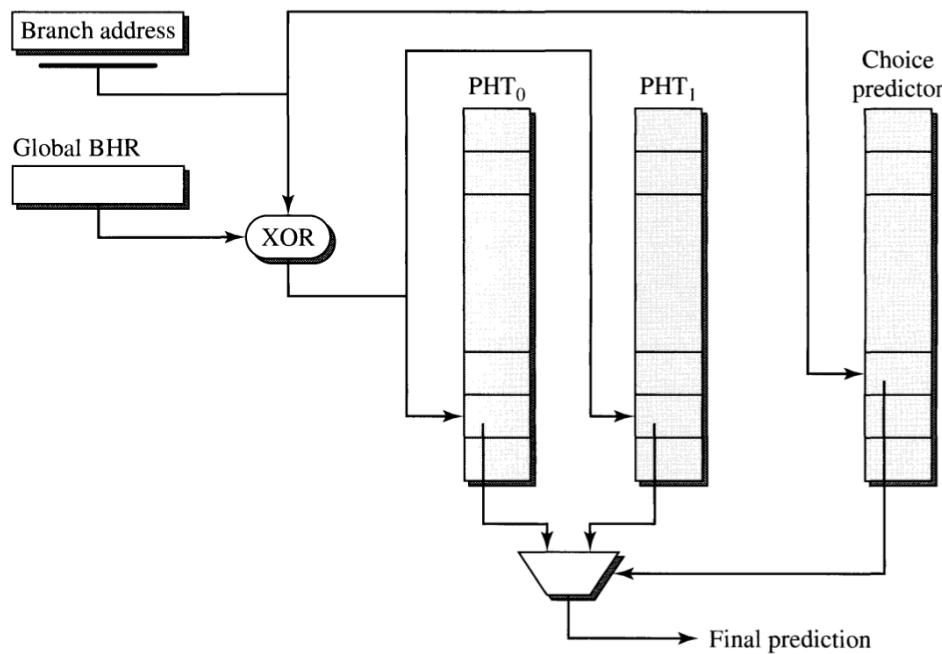
- Combine multiple BPs
- Use any 2 XAy BPs
- M is 2b counters for MUX select
- After branch resolved, update both P_0 , P_1
- Update M based on predictor mismatch, favoring correct
 - Increment if P_0 right, P_1 wrong
 - Decrement if P_0 wrong, P_1 right
 - If both right/wrong do nothing
- How is this different from Bi-Mode?



NYU

TANDON SCHOOL
OF ENGINEERING

Bi-mode and Tournament



What's happening today

Current state-of-the-art (rumored to be what Intel is using)

“**TAGE**: TAgged GEometric predictor”, (Seznac, JILP 2006)

What's research looking into?

Neural predictors!

- 1) Daniel Jimenez (Texas A&M) did lots of pioneering work
- 2) Yale Patt(?) has neural prediction paper at MICRO this month:

“BranchNet: A Convolutional Neural Network to Predict Hard-To-Predict Branches”

This sounds interesting.. I want to learn more

- 1) Highly suggest reading
Chapter 9 (Advanced Instruction Flow Techniques) in
“Modern Processor Design” by Mikko Lipasti
- 2) Read “TAGE” and work backwards through references
- 3) Try out the branch prediction championship:
<https://www.jilp.org/cbp2016/>
- 4) If you can beat the SOTA in branch prediction, I'll give you a PhD.

“Branch Prediction Is Not A Solved Problem: Measurements, Opportunities, and Future Directions” (IISWC 2019)

55% performance hit from miss prediction at 4x scaling

This performance benefit on par with process technology node

New foundries (chip making plants) cost tens of billions of dollars!

~\$20B by some estimates

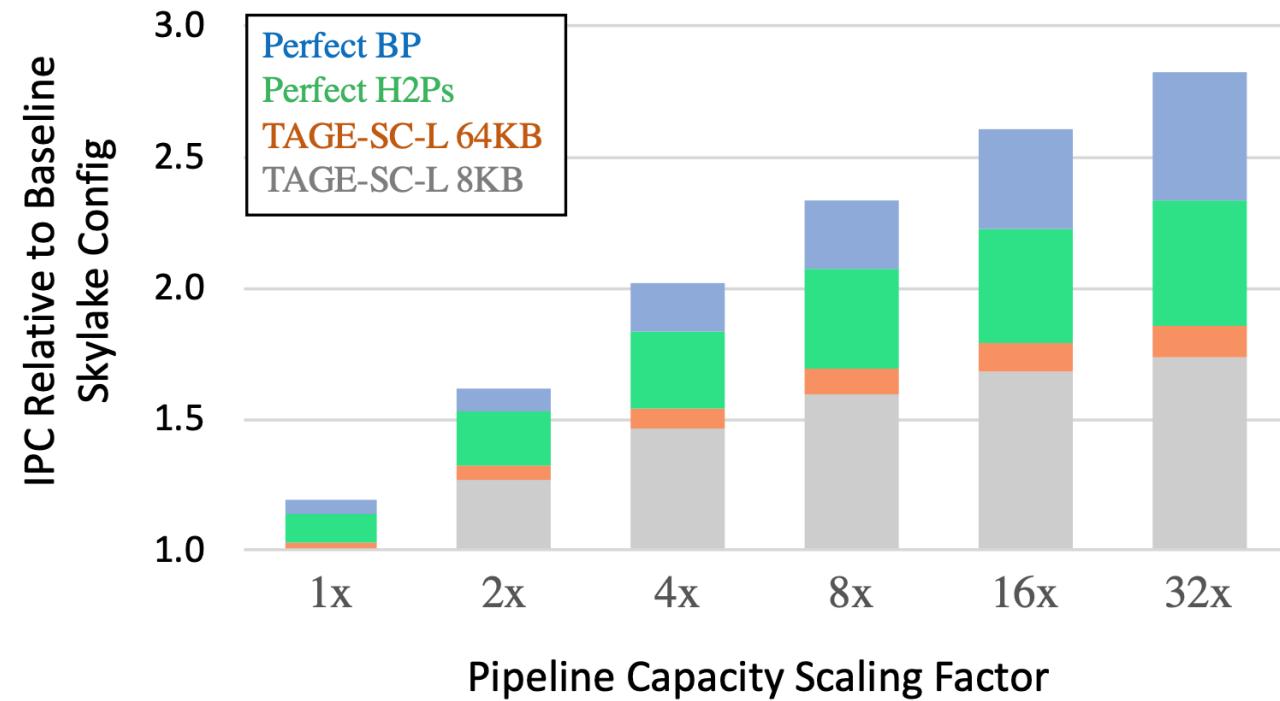
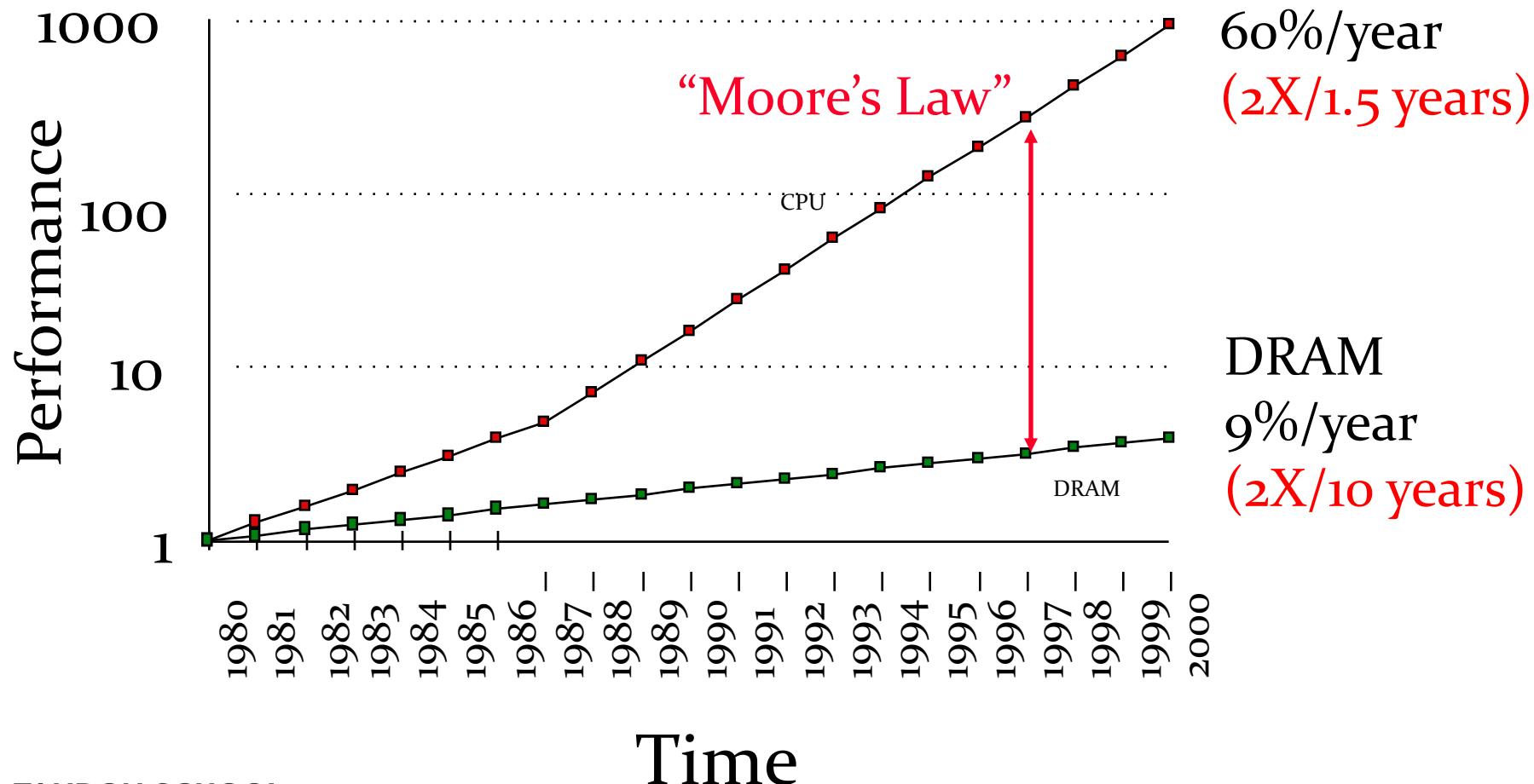


Fig. 1: Without better branch prediction, scaling the pipeline capacity of an Intel Skylake configuration will produce diminishing returns in single-threaded IPC for SPECint 2017.

Caches

Why Care About Memory Hierarchy?

Processor-DRAM Performance Gap grows 50% / year



Processor
60%/year
(2X/1.5 years)

DRAM
9%/year
(2X/10 years)

Memory Issues

Latency

- Time to move through the longest circuit path
- Time to read a word (from the start of request to the response)

Bandwidth

- Number of bits transported at once
- E.g., Gbs/second

Capacity

- Size of memory (32b vs. 64b max verses how much you actually have)

Power/Energy

- Cost of accessing memory (to read and write)
- Computers spend most their energy moving things around



NYU

TANDON SCHOOL
OF ENGINEERING

What can we do about it?

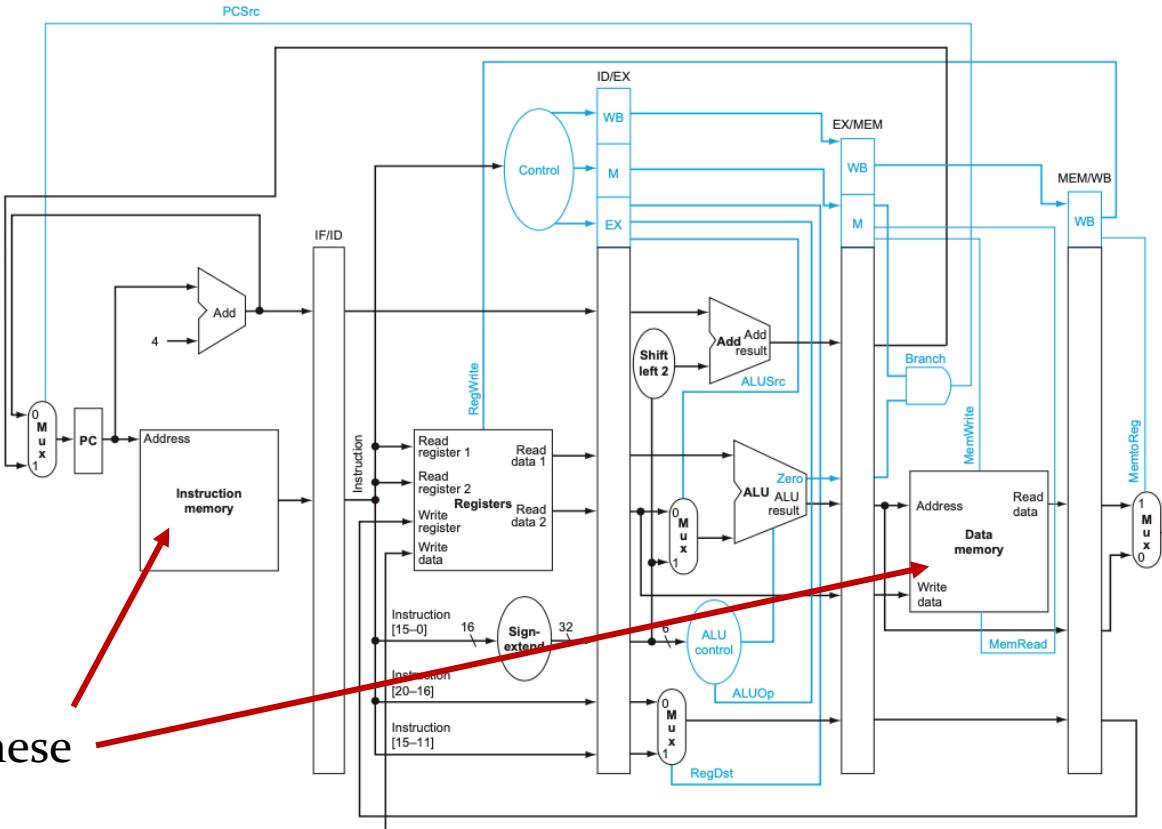
Tennent of architecture: “Smaller is faster.”

How would you solve this?

“Temporary” buffers between the register file and memory.

microArchitecture or Architecture technique?

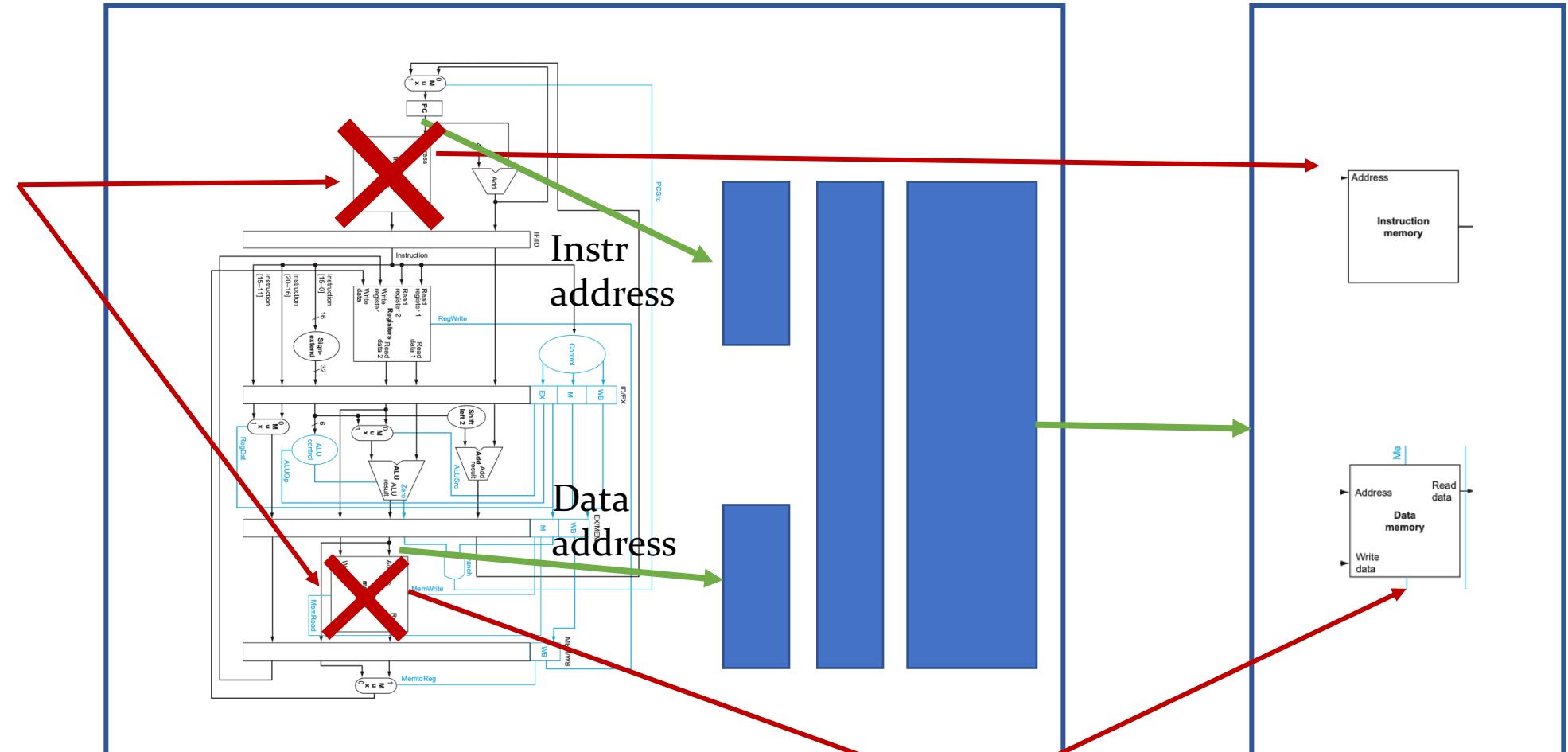
What can we do about it..



Need to replace these

What can we do about it..

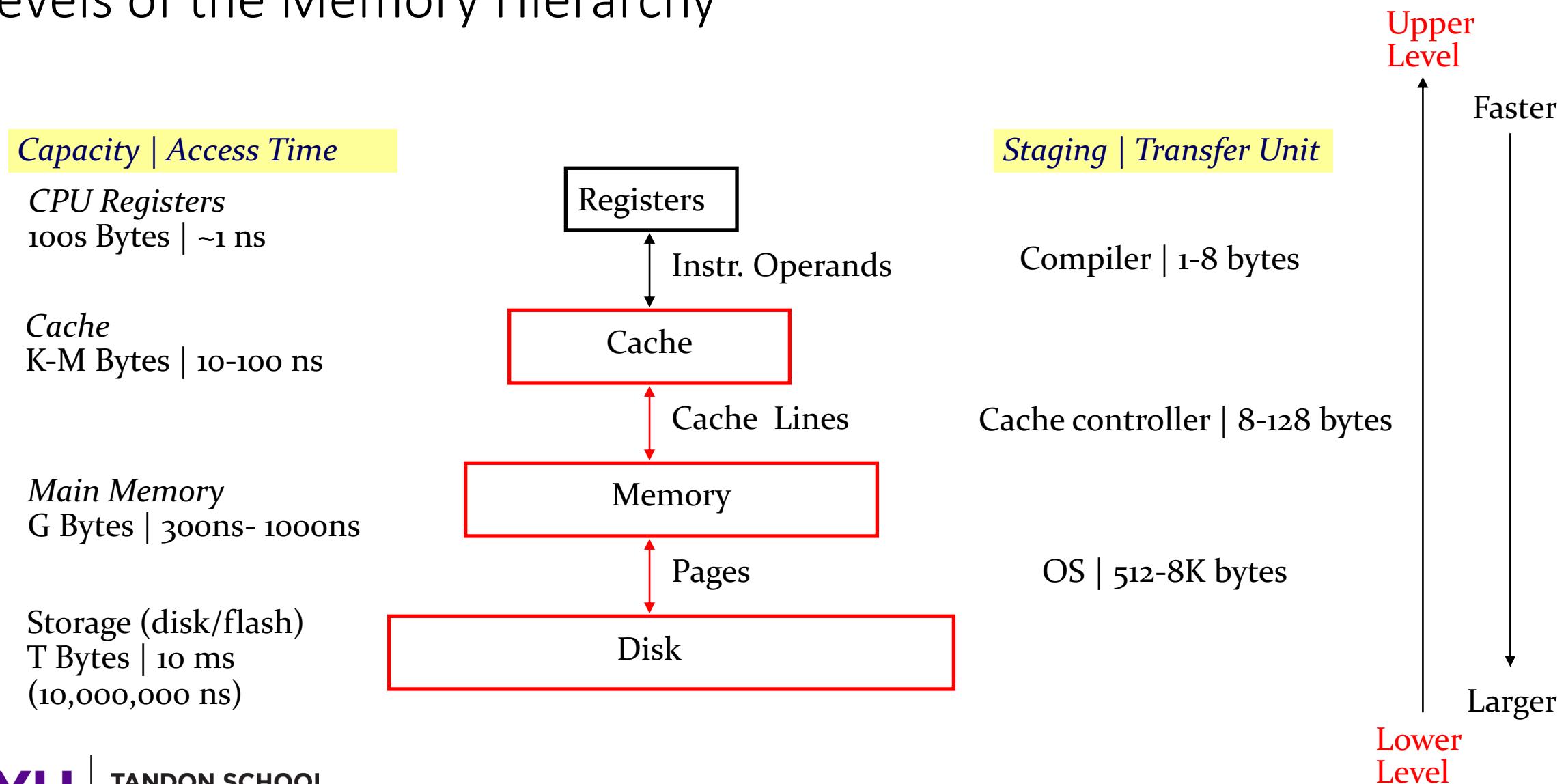
Need to replace



NYU

TANDON SCHOOL
OF ENGINEERING

Levels of the Memory Hierarchy



NYU

TANDON SCHOOL
OF ENGINEERING

Cache Organization

Caches stores data in **blocks**

- A block is a collection of contiguous bytes
 - Can be as small as 1 byte, but can be larger (example: 128 bytes)
 - What's the size of blocks today?
- Request to a cache operate on an entire cache block

Key questions we need to answer:

- **Placement**: Where does a block go when it is fetched into cache?
- **Identification**: How do we know if a block already exists in the cache?
- **Replacement**: Which block should we kick out if there isn't enough room?

Why do Caches Work: Principle of Locality

Programs access a relatively small portion of address space during different phases of execution

Two Types of Locality:

- Temporal Locality (Locality in Time): If an address is referenced, it tends to be referenced again
 - e.g., loops, reuse
- Spatial Locality (Locality in Space): If an address is referenced, neighboring addresses tend to be referenced
 - e.g., straightline code, array access

Traditionally, HW has relied on locality for speed

Locality is a program property that is exploited in machine design.



NYU

TANDON SCHOOL
OF ENGINEERING

Example of Locality

```
int A[100], B[100], C[100], D;  
for (i=0; i<100; i++) {  
    C[i] = A[i] * B[i] + D;  
}
```

2D view of
memory

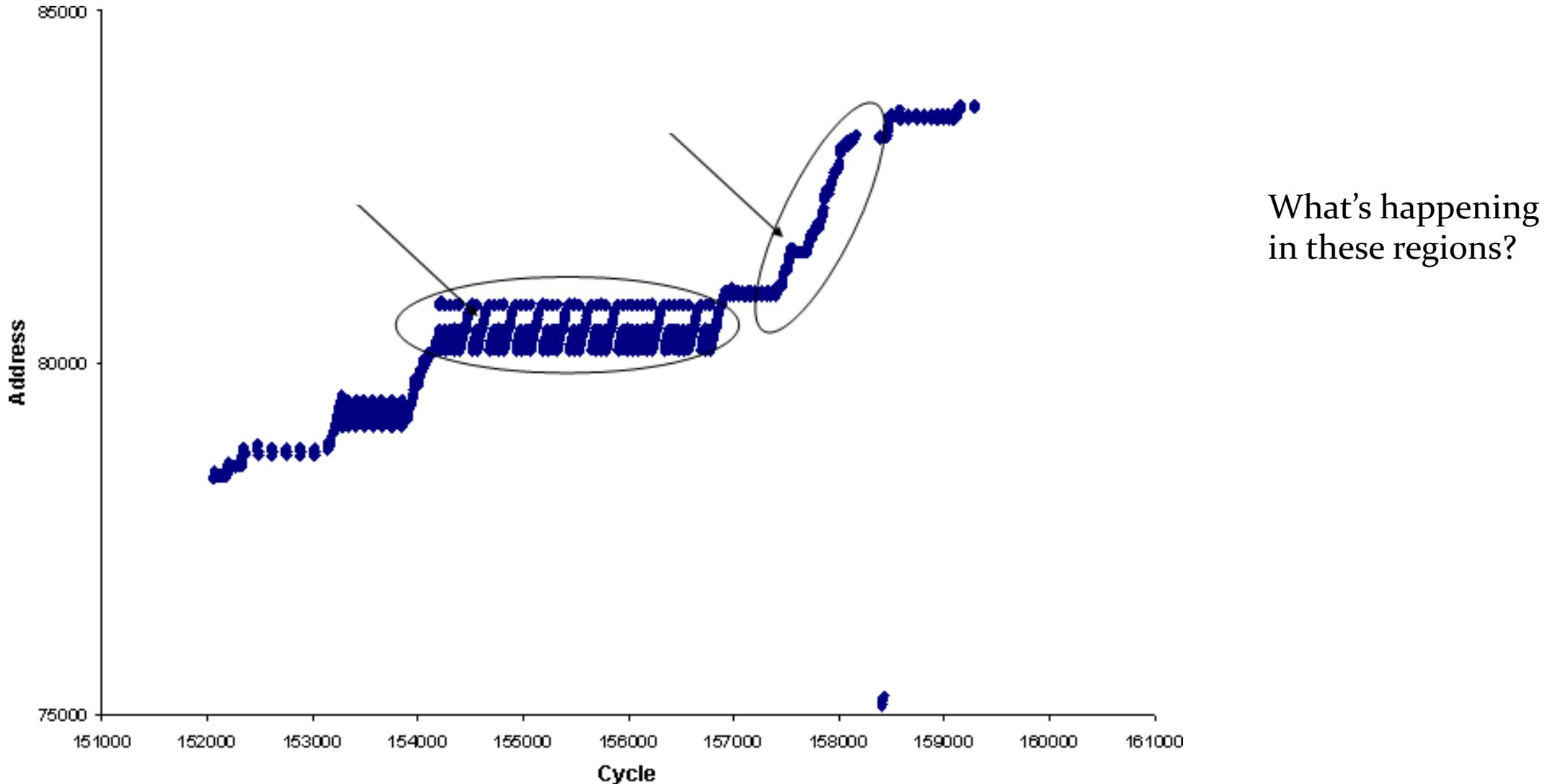
			D	C[99]	C[98]	C[97]	C[96]
.....
C[7]	C[6]	C[5]	C[4]	C[3]	C[2]	C[1]	C[0]
.....
B[11]	B[10]	B[9]	B[8]	B[7]	B[6]	B[5]	B[4]
B[3]	B[2]	B[1]	B[0]	A[99]	A[98]	A[97]	A[96]
.....
A[7]	A[6]	A[5]	A[4]	A[3]	A[2]	A[1]	A[0]

↔ A Cache Line (One fetch)

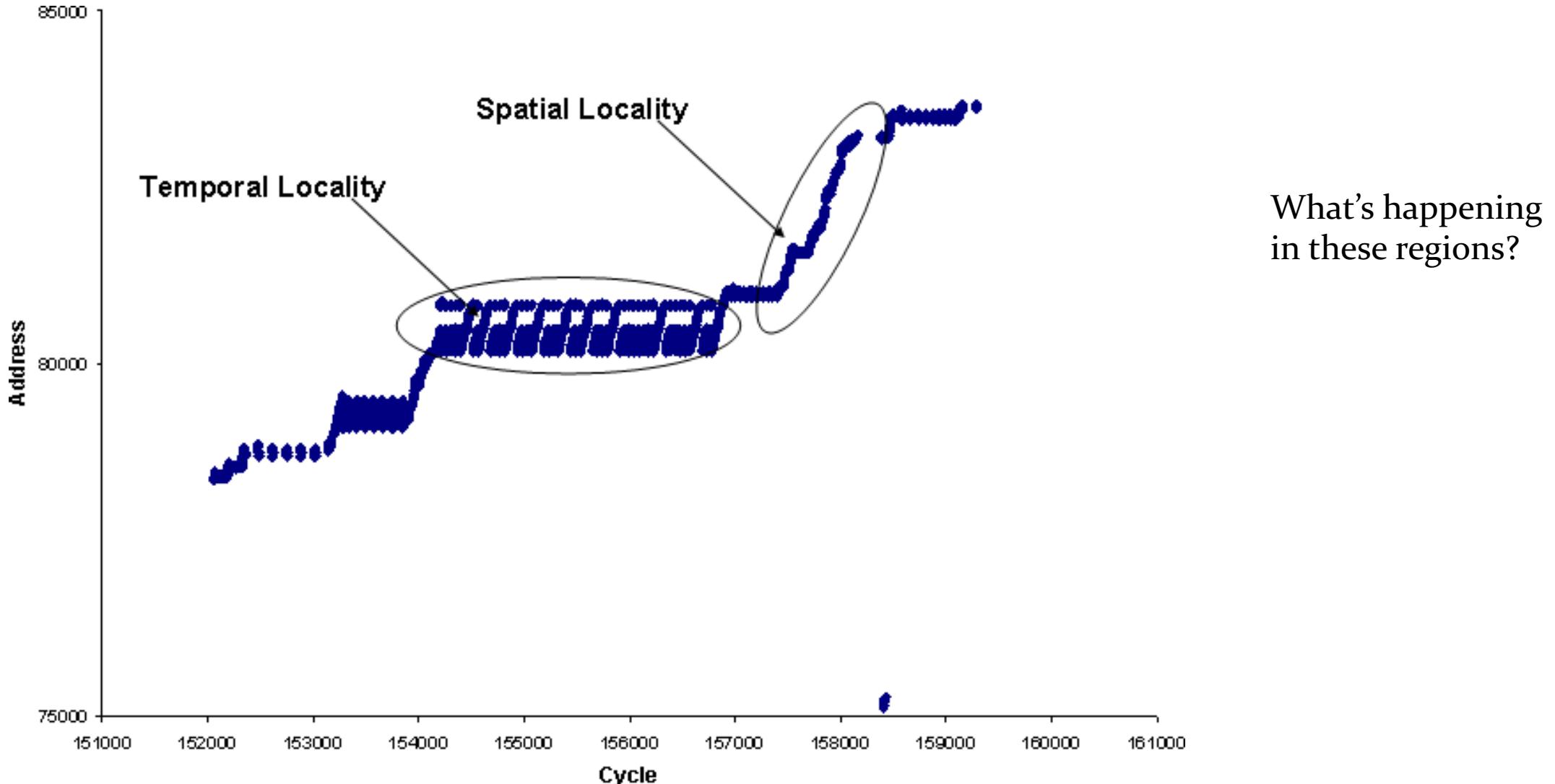
What's temporal?

Where's spatial?

Access Pattern Exhibiting Locality



Access Pattern Exhibiting Locality

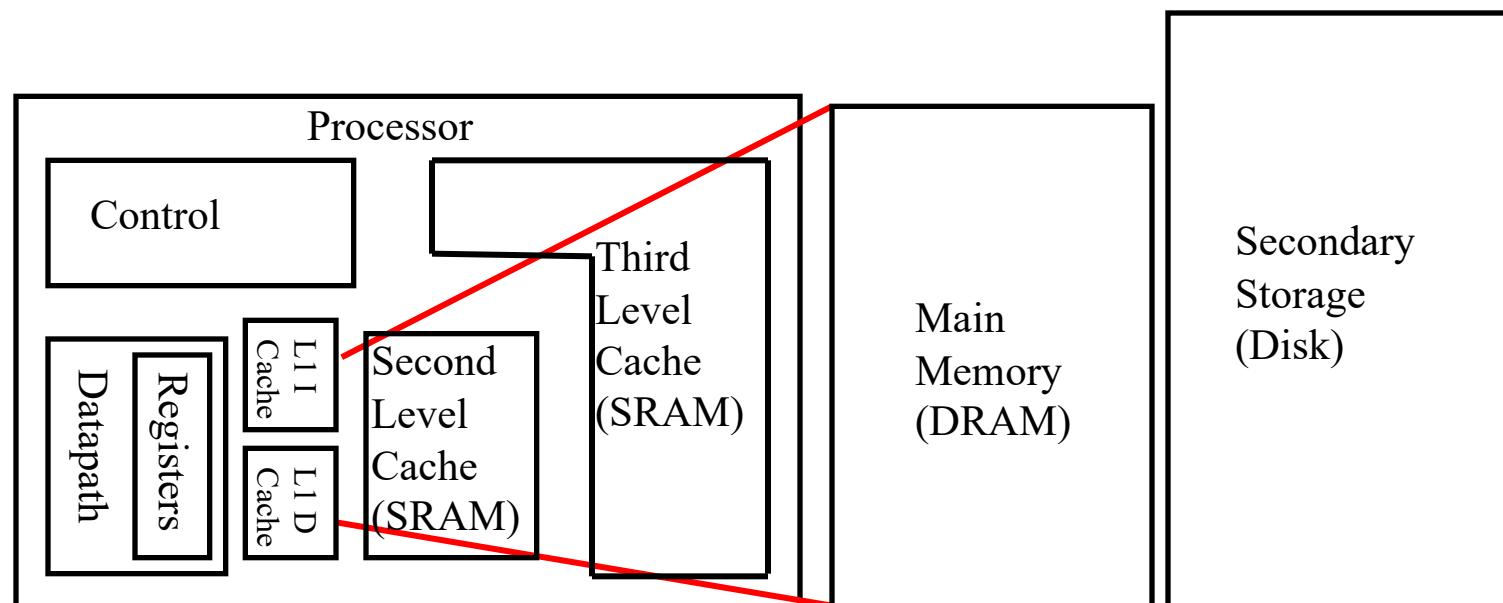


“Memory Hierarchy” or “memory subsystem”

This is how lw and sw are implemented ($L_1, L_2, L_3\dots$)

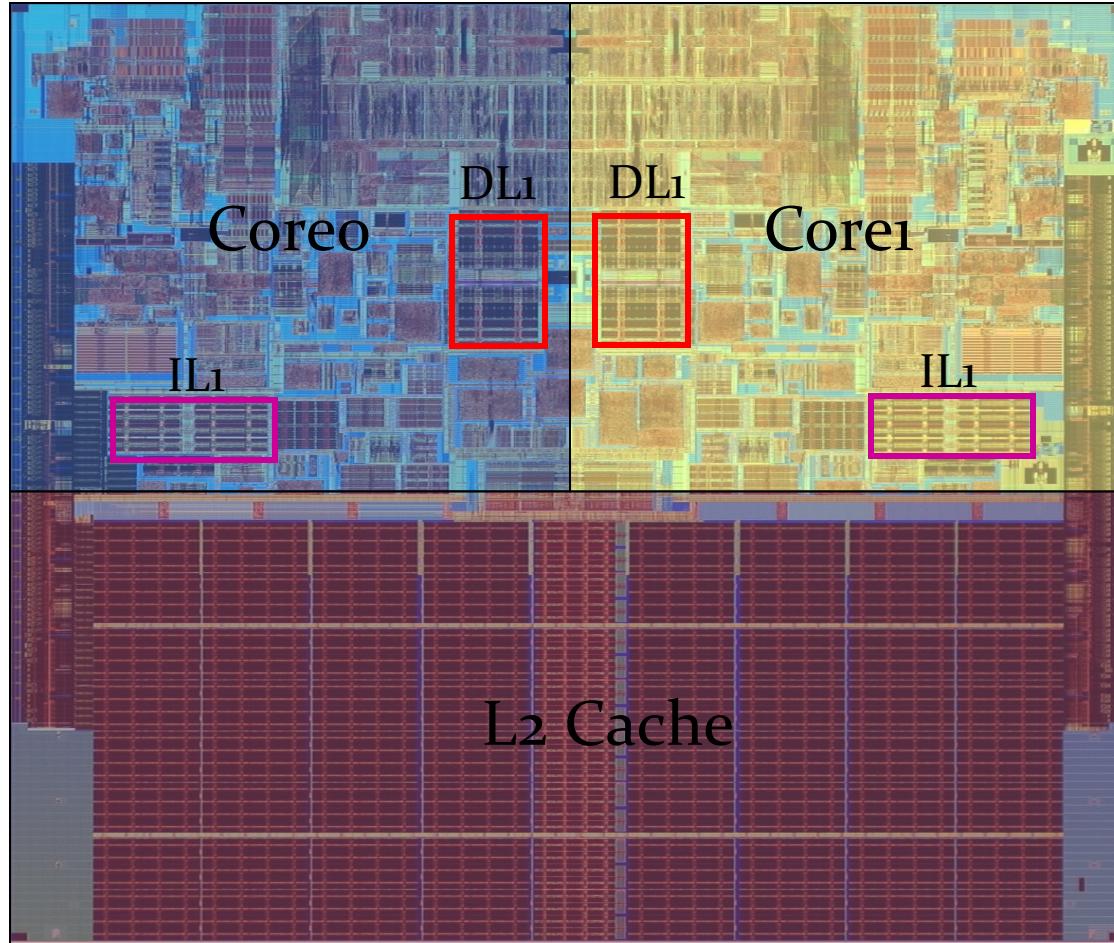
Taking advantage of locality gives the best of both worlds:

- Present the user with as much memory as is available in the cheapest technology
- Provide access at the speed offered by the fastest technology



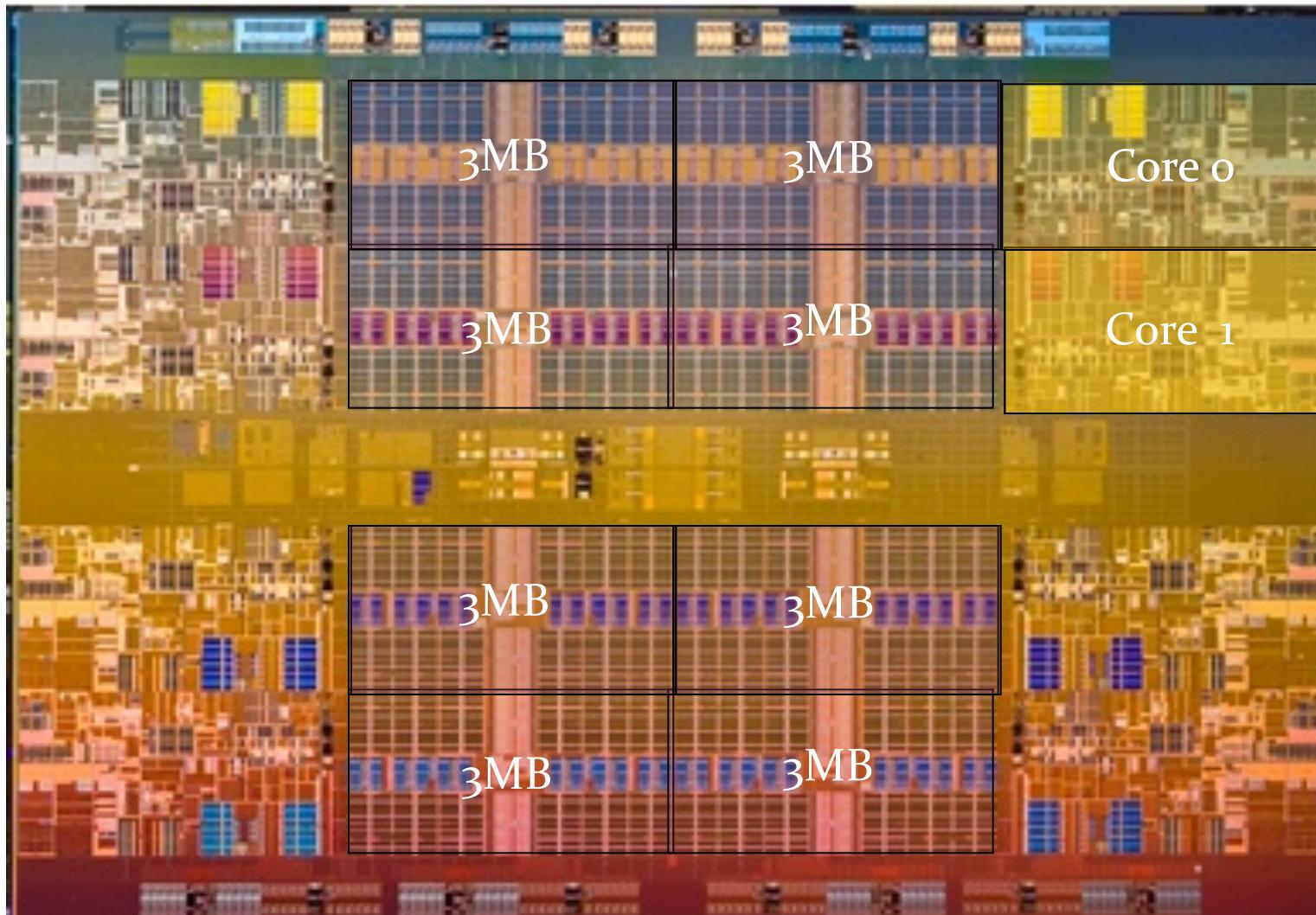
Example: Intel Core2 Duo

L1	32 KB, 8-Way, 64 Byte/Line, LRU, WB 3 Cycle Latency
L2	4.0 MB, 16-Way, 64 Byte/Line, LRU, WB 14 Cycle Latency



Source: <http://www.sandpile.org>

Intel Nehalem



24MB L3!

What's the
trend?

More area
to cache!



NYU

TANDON SCHOOL
OF ENGINEERING

Cache Terminology

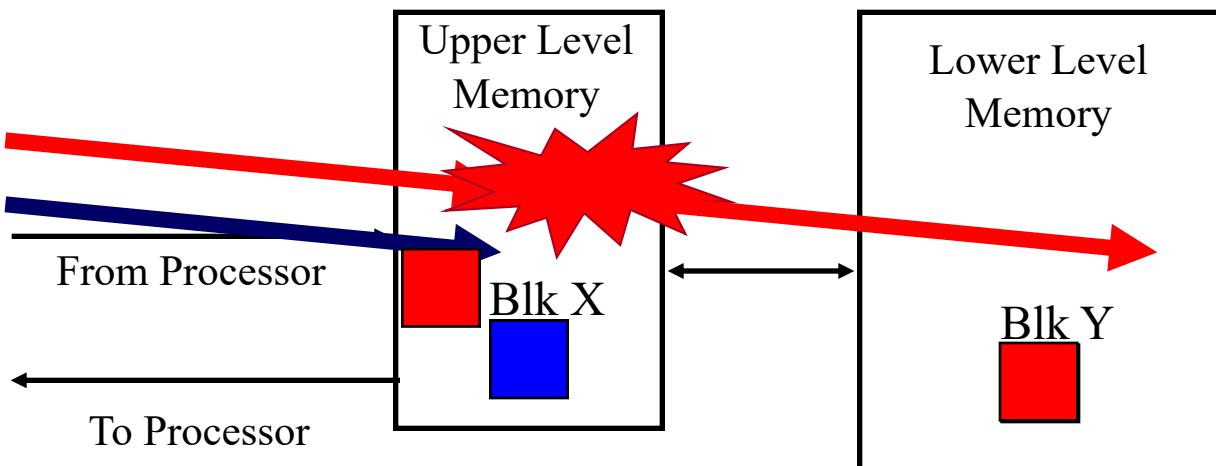
Hit: data is resident in the cache

- **Hit Rate:** the fraction of memory accesses that hit in a level
- **Hit Time:** Time to access the level (access time + time to determine if hit)

Miss: data needs to be retrieved from a block in the lower level

- **Miss Rate** = $1 - \text{Hit Rate}$
- **Miss Penalty:** Time to replace a block in the upper level +
Time to deliver the block to the processor

$\text{Hit Time} \ll \text{Miss Penalty}$



NYU

TANDON SCHOOL
OF ENGINEERING

Average Memory Access Time (AMAT)

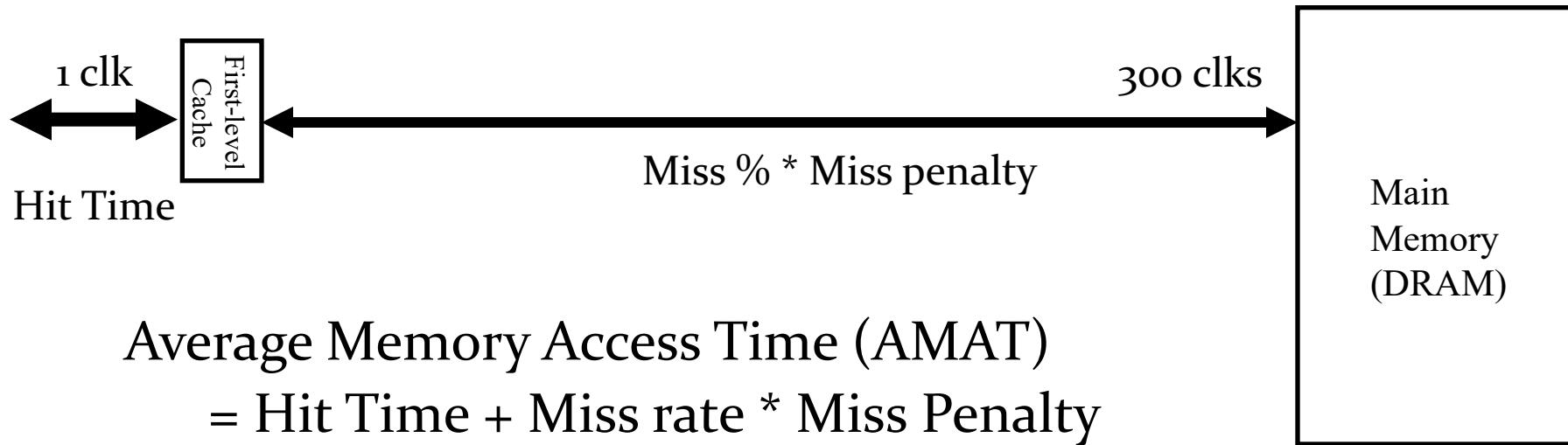
Average memory-access time

- How long does it get to read/write?
- AMAT = Hit time + Miss rate x Miss penalty

Miss penalty: time spent fetching a block from lower memory level

- *access time*: function of latency
- *transfer time*: function of bandwidth b/w levels
 - Transfer one “cache line/block” at a time
 - Transfer at the size of the memory-bus width

Memory Hierarchy Performance



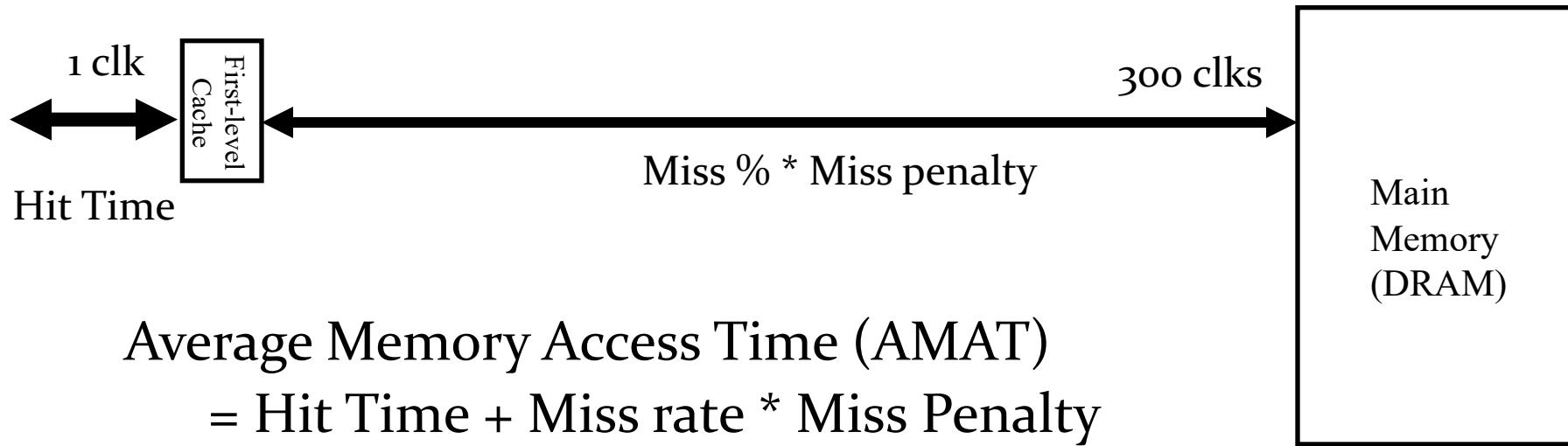
Average Memory Access Time (AMAT)

$$\begin{aligned} &= \text{Hit Time} + \text{Miss rate} * \text{Miss Penalty} \\ &= T_{\text{hit}}(L_1) + \text{Miss\%}(L_1) * T(\text{memory}) \end{aligned}$$

Example:

- Cache Hit = 1 cycle
- Miss rate = 10% = 0.1
- Miss penalty = 300 cycles
- AMAT = ?

Memory Hierarchy Performance



Average Memory Access Time (AMAT)

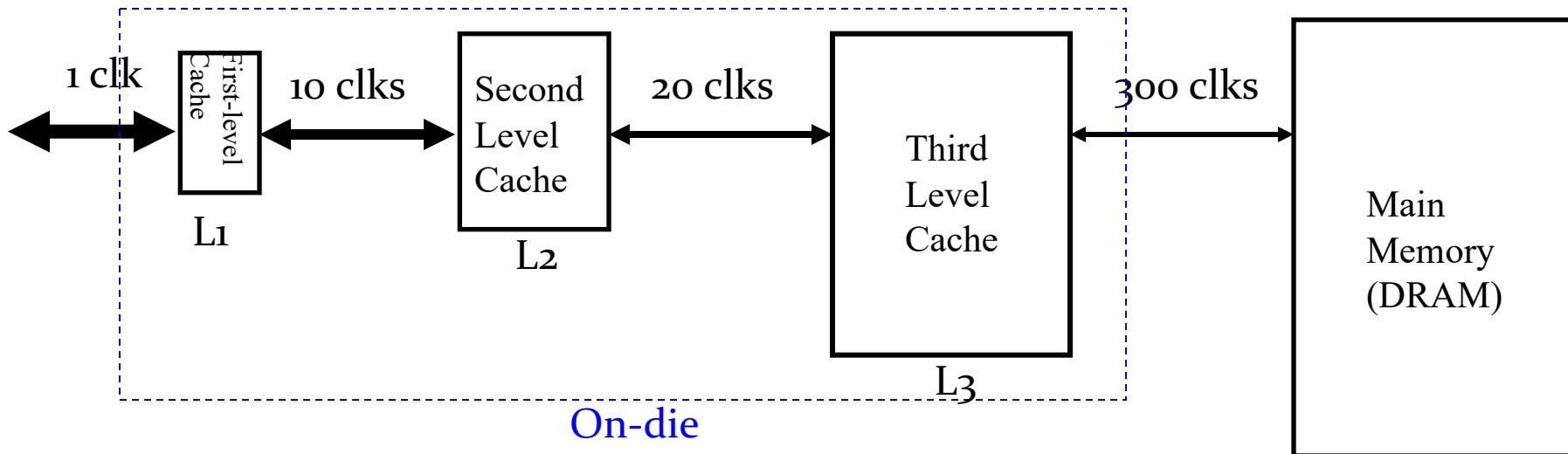
$$\begin{aligned} &= \text{Hit Time} + \text{Miss rate} * \text{Miss Penalty} \\ &= T_{\text{hit}}(L_1) + \text{Miss\%}(L_1) * T(\text{memory}) \end{aligned}$$

Example:

- Cache Hit = 1 cycle
- Miss rate = 10% = 0.1
- Miss penalty = 300 cycles
- AMAT = $1 + 0.1 * 300 = 31$ cycles

Can we improve it?

Reducing Penalty: Multi-Level Cache!



Average Memory Access Time (AMAT)

$$\begin{aligned} &= T_{\text{hit}}(L_1) + \\ &\quad \text{MissRate}(L_1) * [T_{\text{hit}}(L_2) + \\ &\quad \quad \text{MissRate}(L_2) * \{ T_{\text{hit}}(L_3) + \\ &\quad \quad \quad \text{MissRate}(L_3) * T(\text{memory}) \}] \end{aligned}$$

AMAT Example

Example:

- Miss rate L₁=10%, T_{hit}(L₁) = 1 cycle
- Miss rate L₂=5%, T_{hit}(L₂) = 10 cycles
- Miss rate L₃=1%, T_{hit}(L₃) = 20 cycles
- T(memory) = 300 cycles

AMAT = ?

- No cache: 300
- L₁ only: 31 [9.68 speedup!]
- All levels: 2.115 [14.7x speedup!]

141.8x total speedup.

Caches work very well!

Average Memory Access Time (AMAT)

$$\begin{aligned} &= T_{hit}(L_1) + \\ &\text{MissRate}(L_1) * [T_{hit}(L_2) + \\ &\quad \text{MissRate}(L_2) * \{ T_{hit}(L_3) + \\ &\quad \quad \text{MissRate}(L_3) * T(\text{memory}) \}] \end{aligned}$$



NYU

TANDON SCHOOL
OF ENGINEERING

Types of Caches

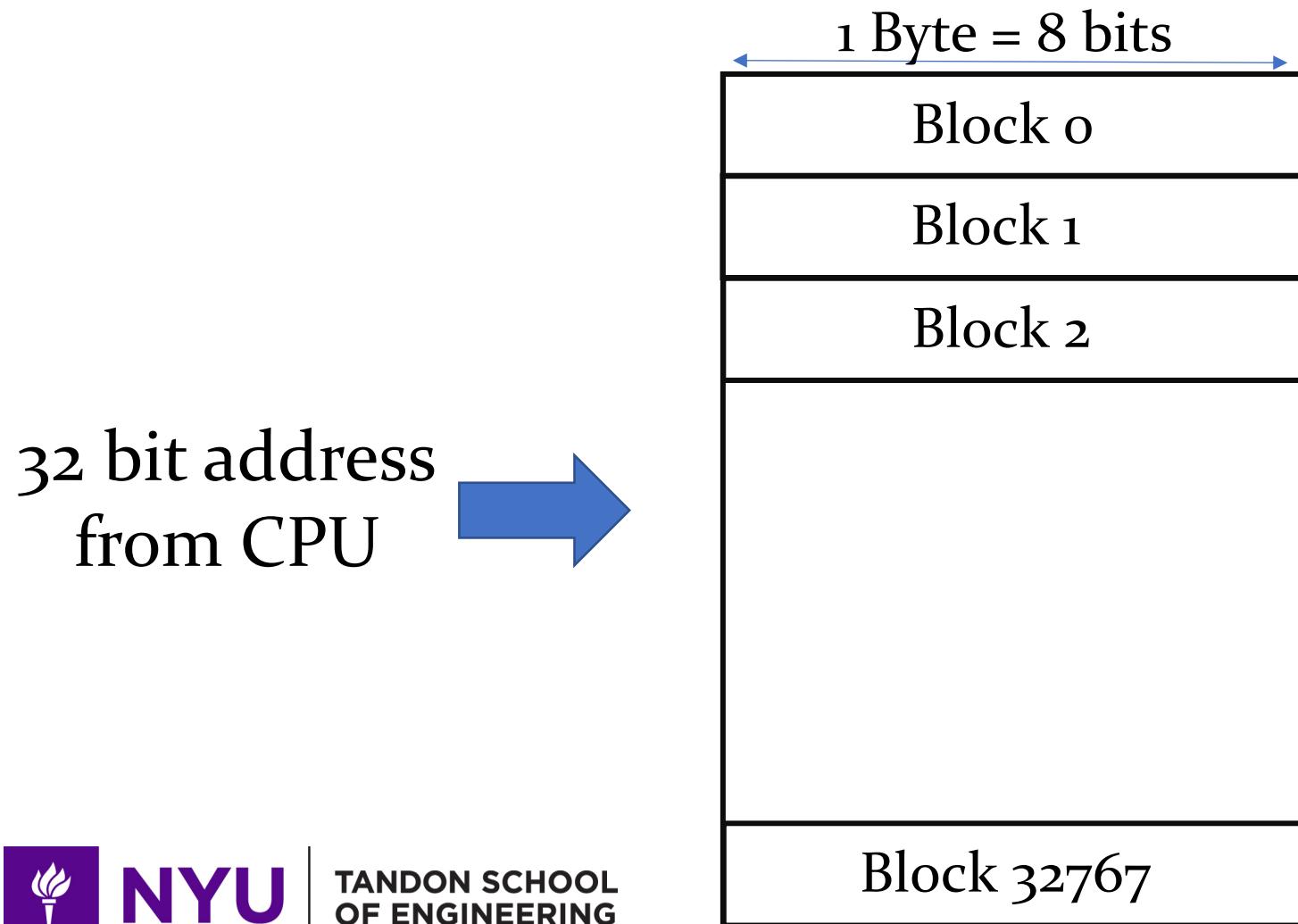
Type of cache	Placement: Mapping of data from memory to cache	Identification: Complexity of searching the cache
Direct mapped (DM)	<ul style="list-style-type: none">• DM and FA can be thought as special cases of SA<ul style="list-style-type: none">• DM → 1-way SA• FA → All-way SA	Fast indexing mechanism
Set-associative (SA)	A memory value can be placed in any of a set of locations in the cache	Slightly more involved search mechanism
Fully-associative (FA)	A memory value can be placed in any location in the cache	Extensive hardware resources required to search (CAM)



Direct Mapped Cache example

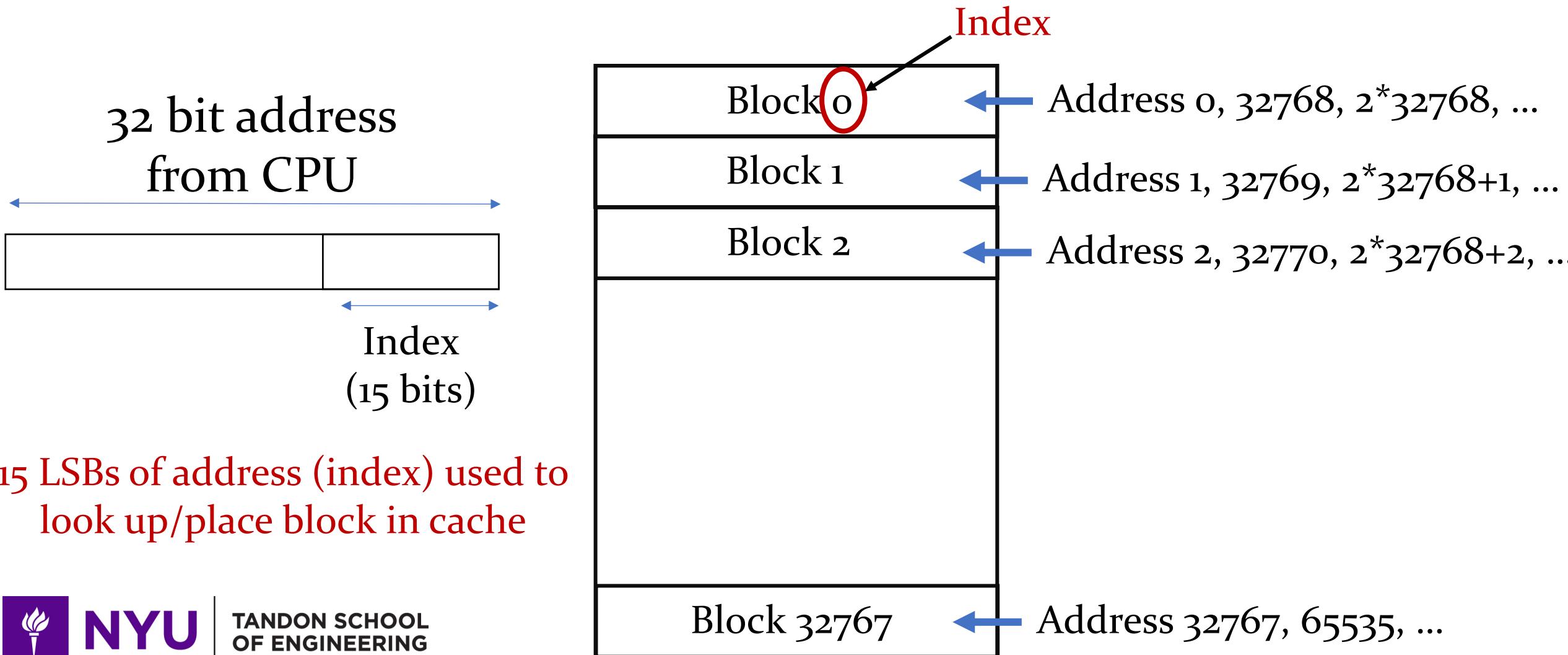
Assume 32 KByte cache with a 1 Byte cache block

- Note that $32\text{KB} = 2^{15}$ Bytes = 32768 Bytes

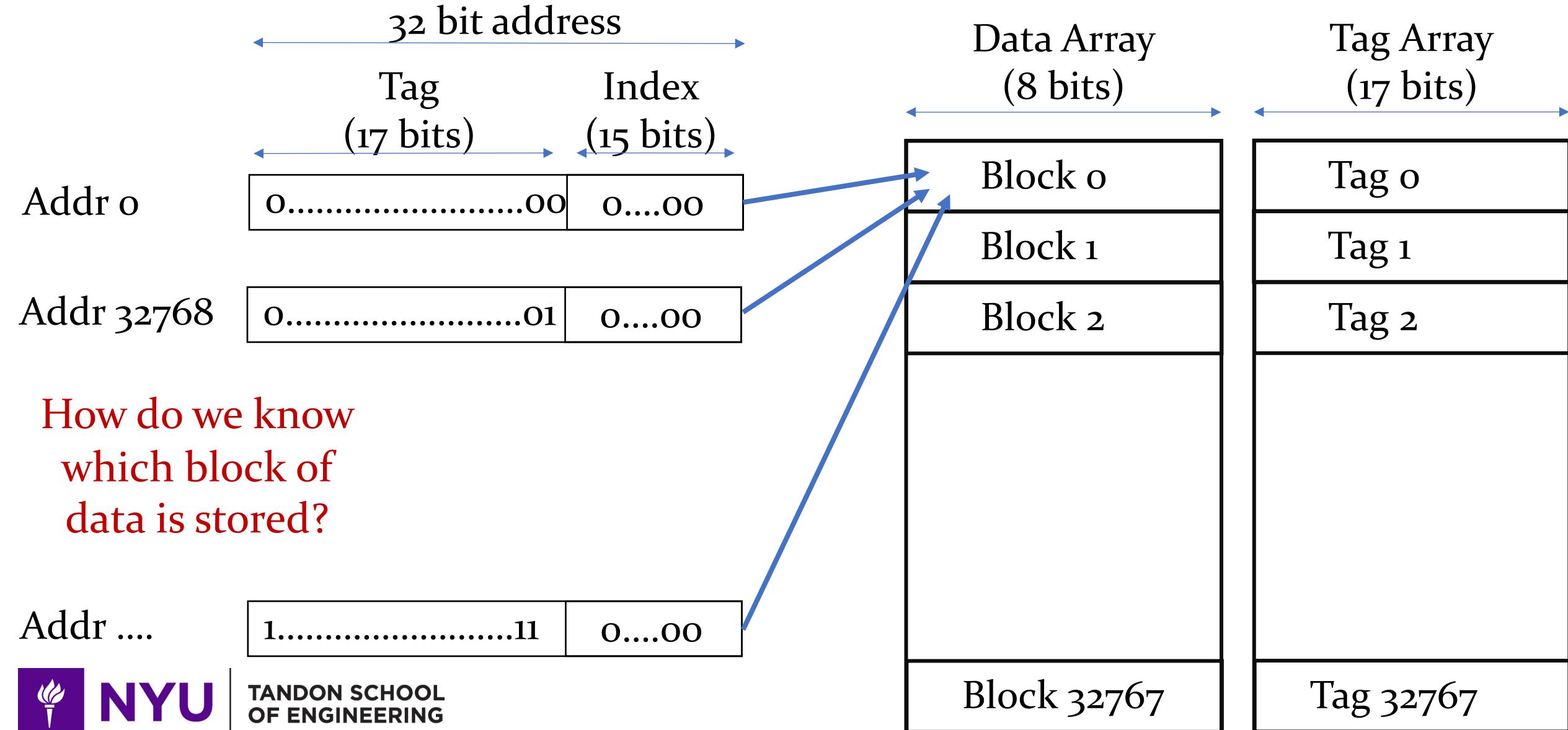


Direct Mapped Cache: Placement

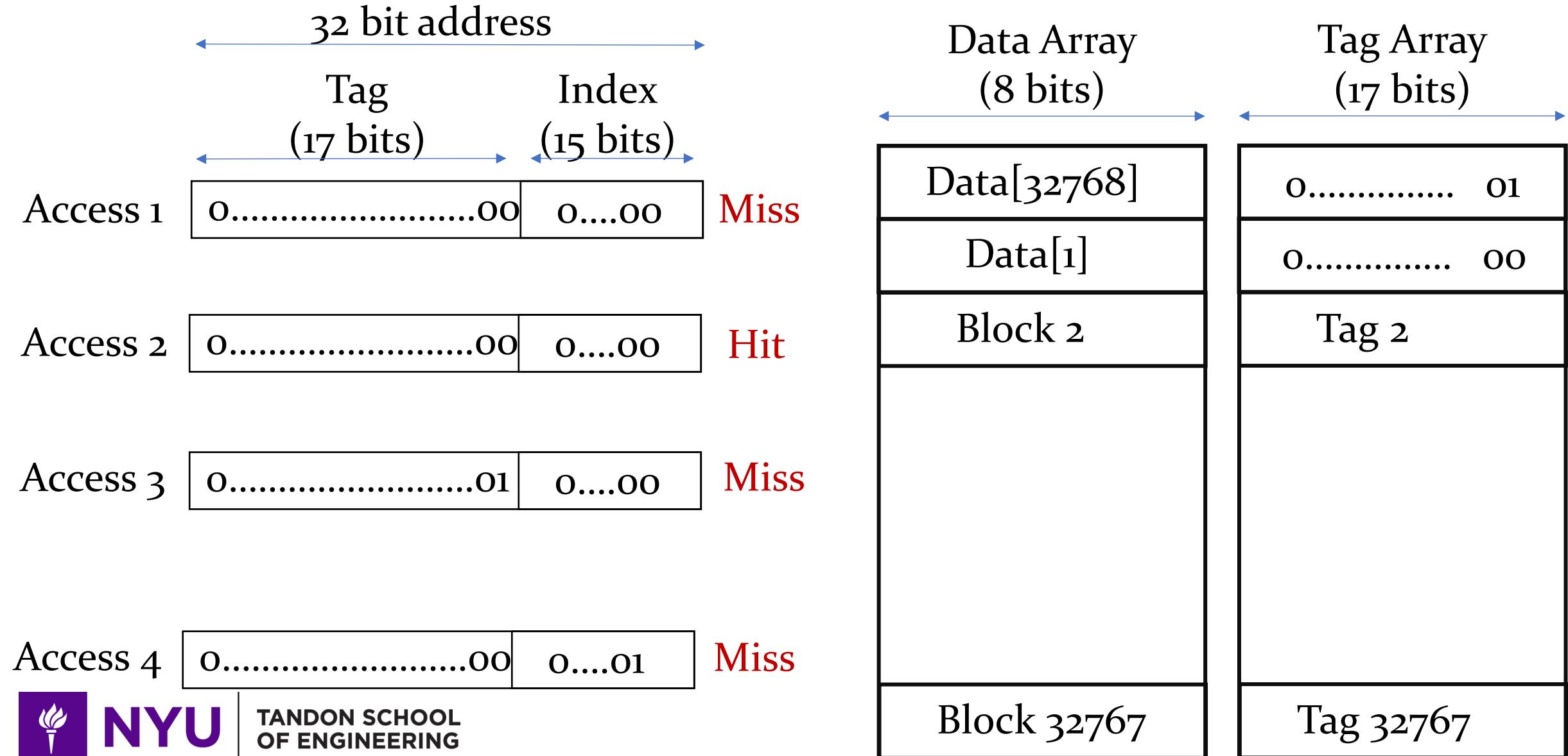
- Each block in the cache has a $\log_2(32768) = 15$ bit “index”



Direct Mapped Cache: Identification



Direct Mapped Cache: Replacement

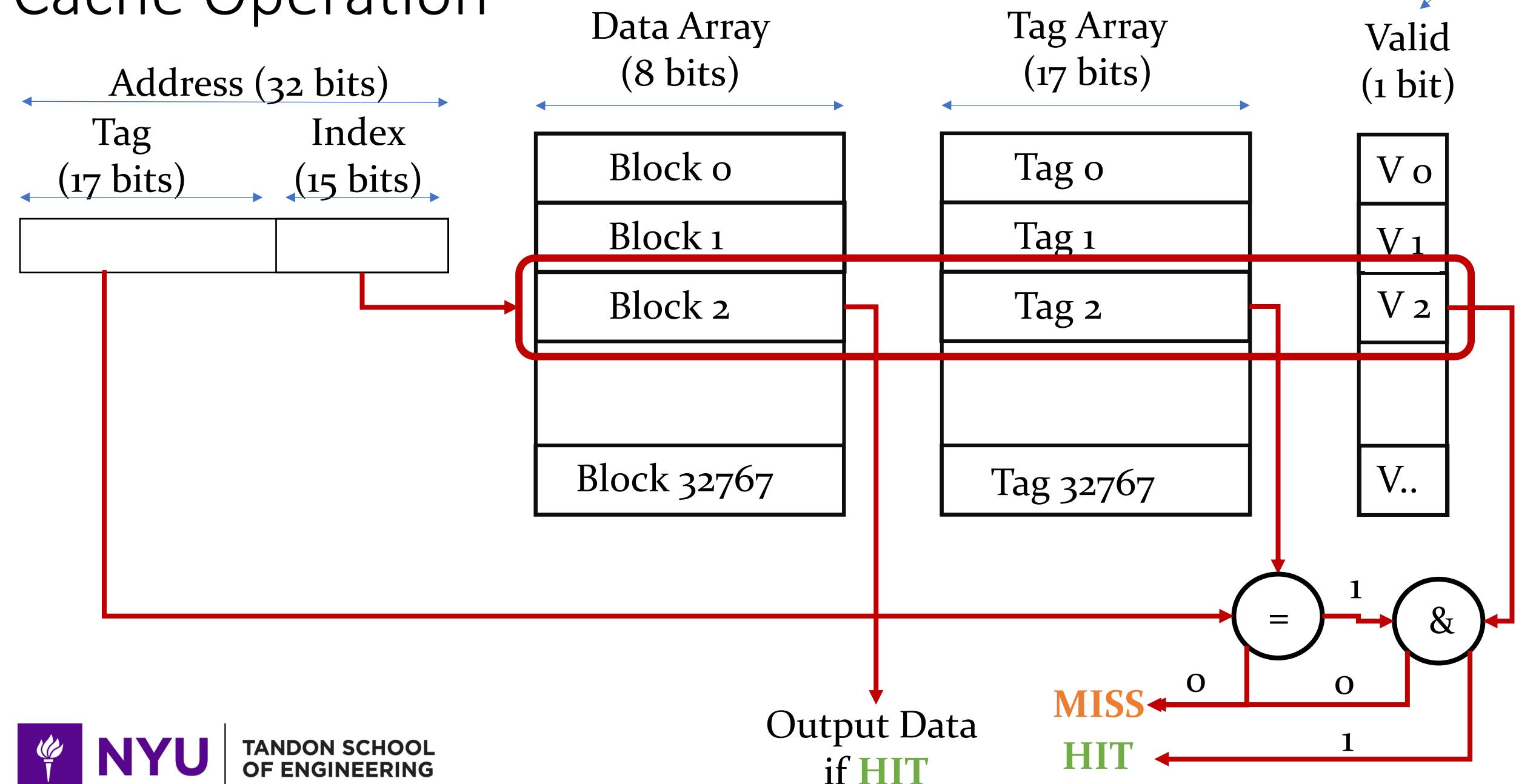


NYU

TANDON SCHOOL
OF ENGINEERING

(details later)

Cache Operation



NYU

TANDON SCHOOL
OF ENGINEERING

Alternative DM Placement

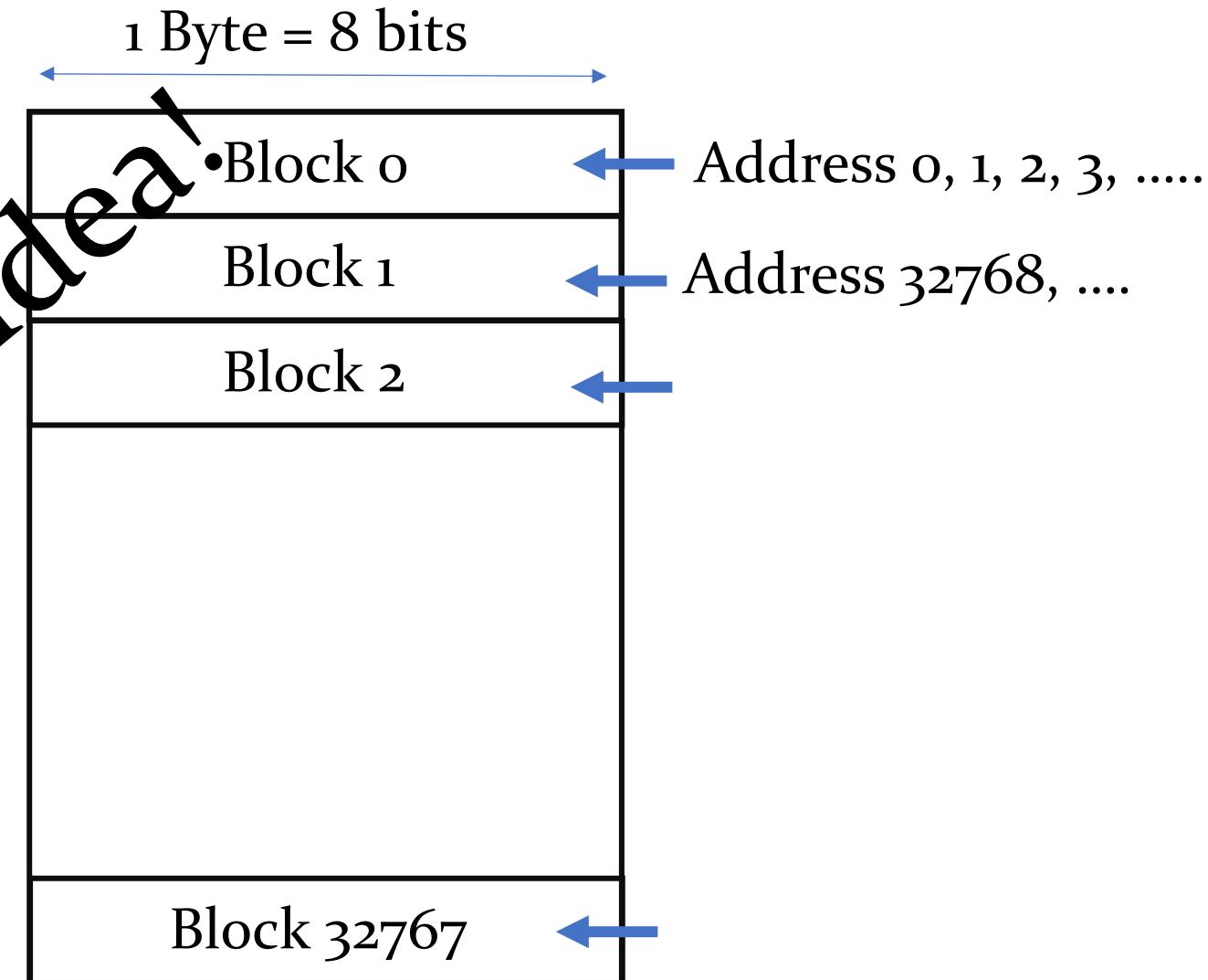
Assume 32 KByte cache with a 1 Byte cache block

- Use 15 **MSBs** of address to determine index

- Consider following stream of accesses:

Addr 0, Addr 1, Addr 0, Addr 1,...

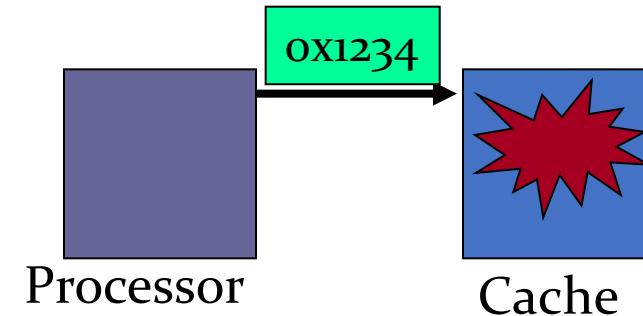
Bad Idea



Three Cs (Cache Miss Terms)

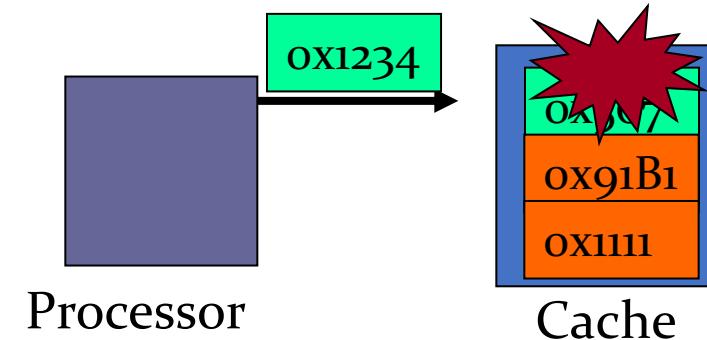
Compulsory Misses:

- “Cold start” misses
- Caches do not have valid data at the start of the program



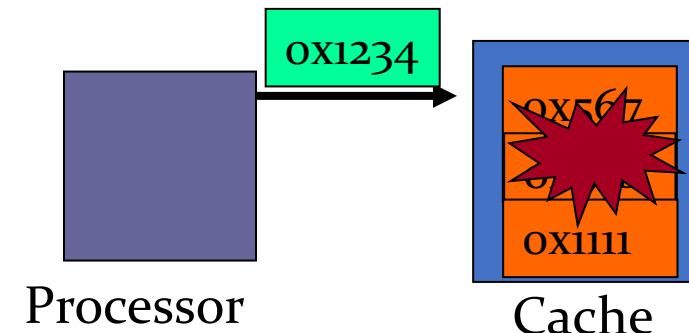
Conflict Misses:

- Increase cache size or associativity (next)
- What we saw with DM caches
- Associative caches reduce conflict misses



Capacity Misses:

- Increase cache size



NYU

TANDON SCHOOL
OF ENGINEERING

Set Associative Caches

Direct mapped caches can have high miss rates due to **conflicts**

- Each address maps to a unique location in the cache

Assume addresses A and B with same index bits

- Sequence: A, B, A, B, A..... results in 100% cache miss rate!

Set-associative cache: each address can map to N different locations in the cache!

- “N-way” set associative cache
- 2-way set associate cache has
~0% cache miss rate for sequence A, B, A, B ...

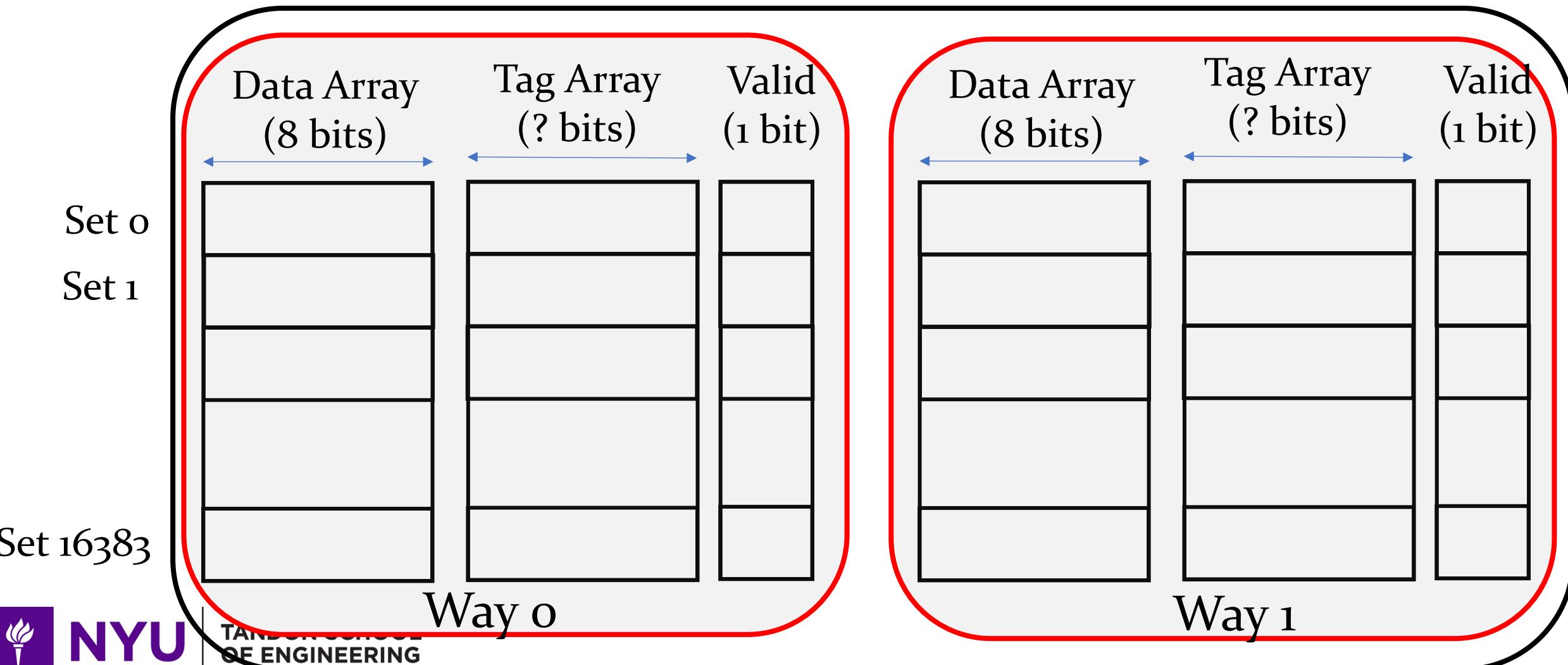


NYU

TANDON SCHOOL
OF ENGINEERING

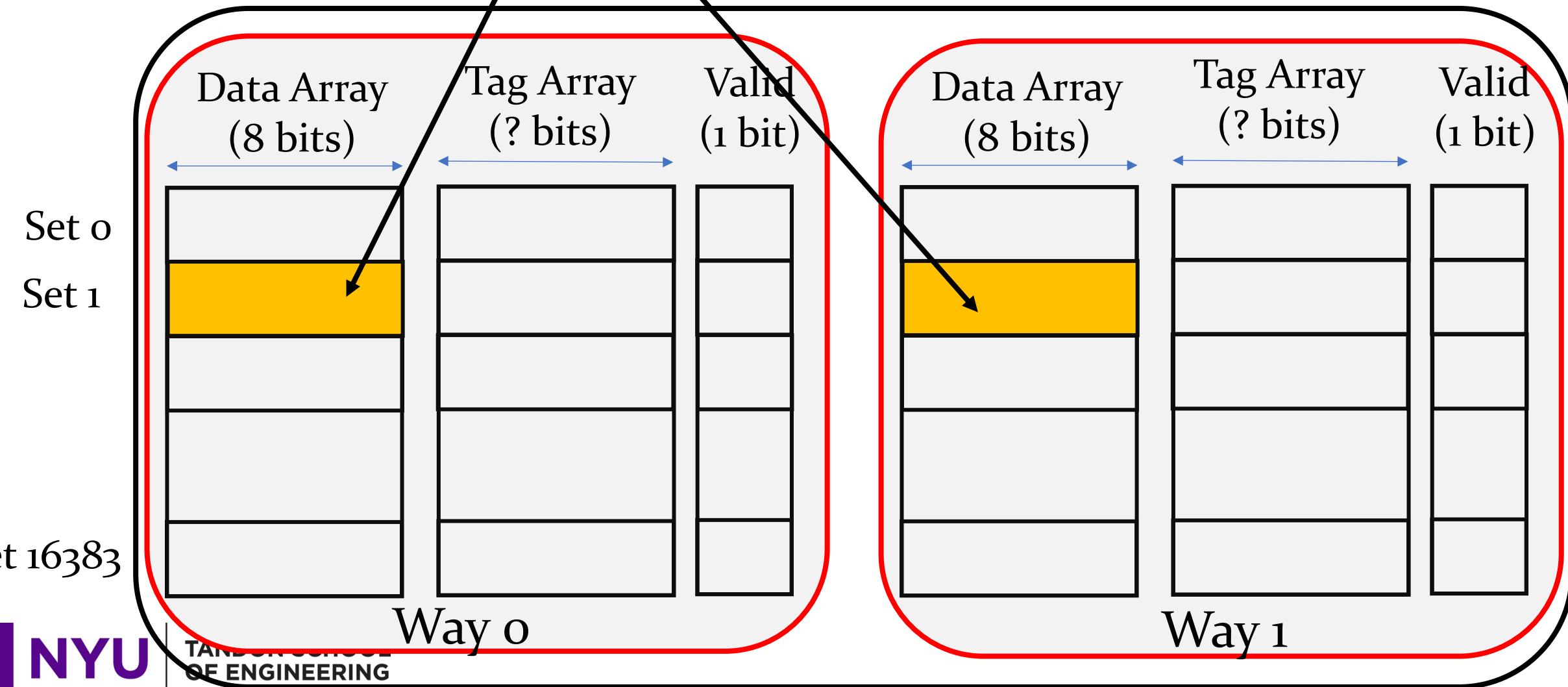
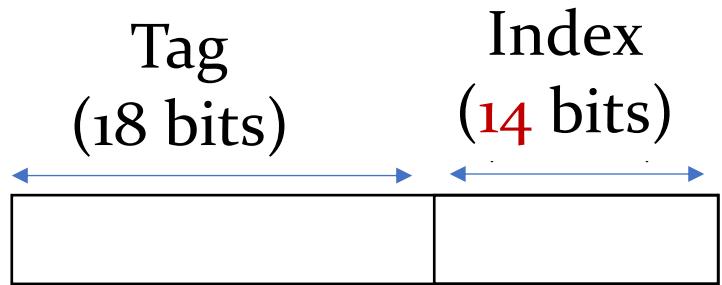
2-Way Set Associative Cache

2-way set-associative 32 KB cache with 1 Byte blocks

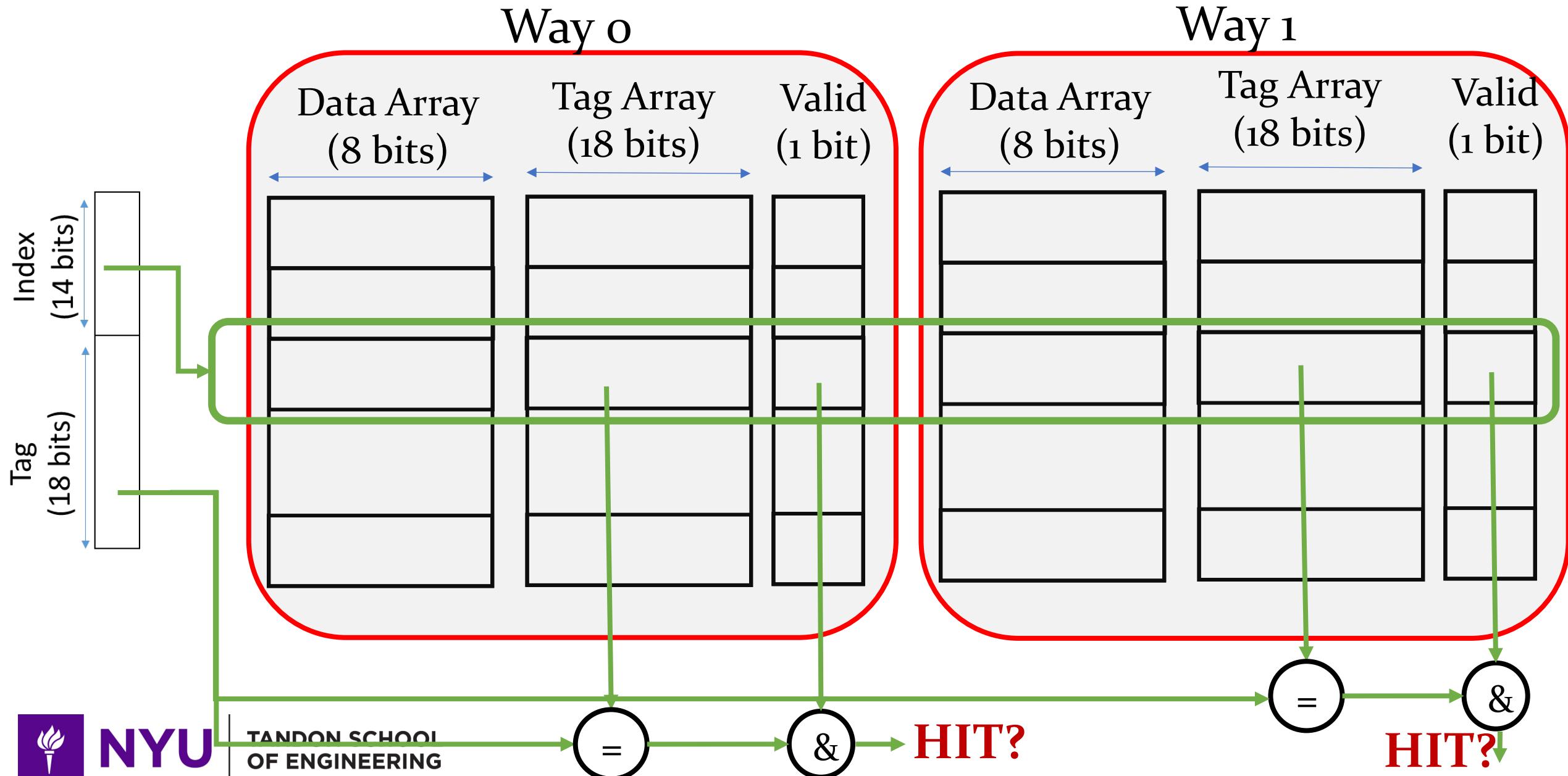


Placement

Data mapped to a unique set but can be placed in either way



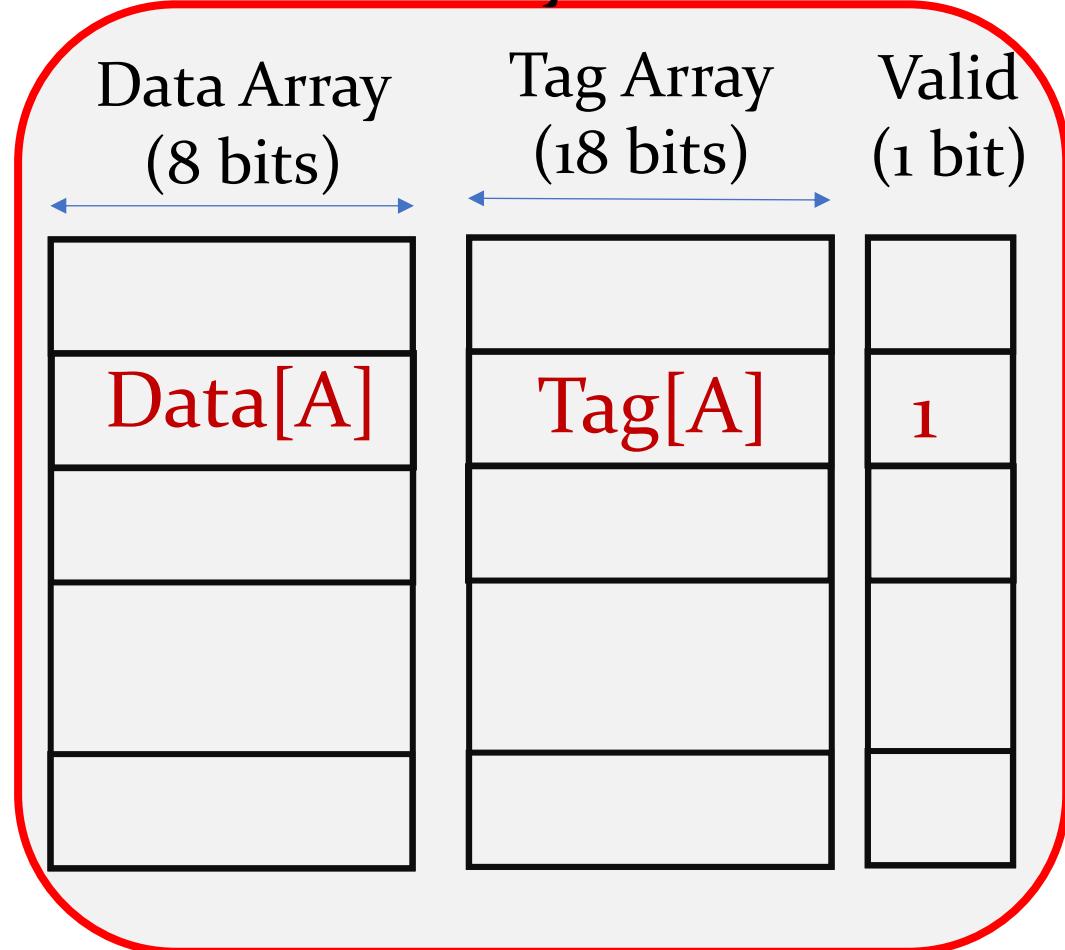
Identification



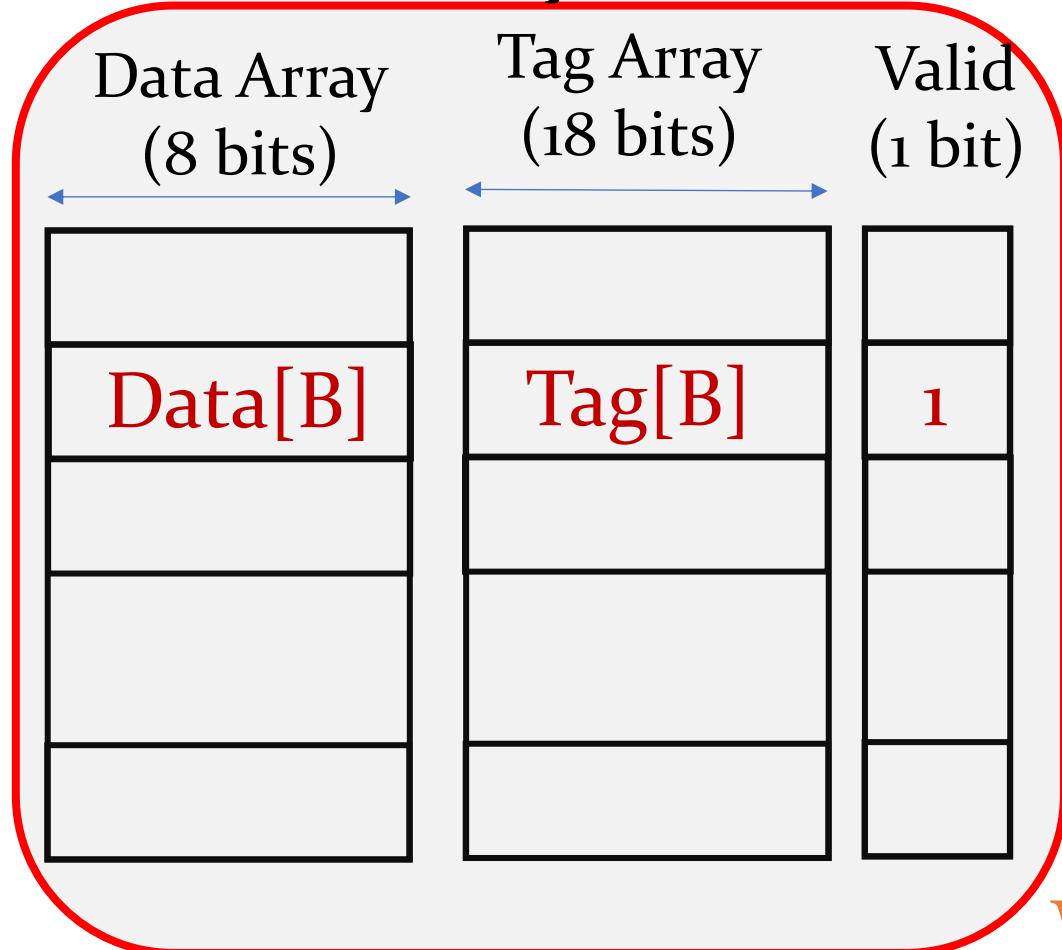
Replacement

(Addresses A, B, C map to same set)

Way 0



Way 1



Sequence: A B B A A A A C
 M M H H H H H M

Where
should
C go?

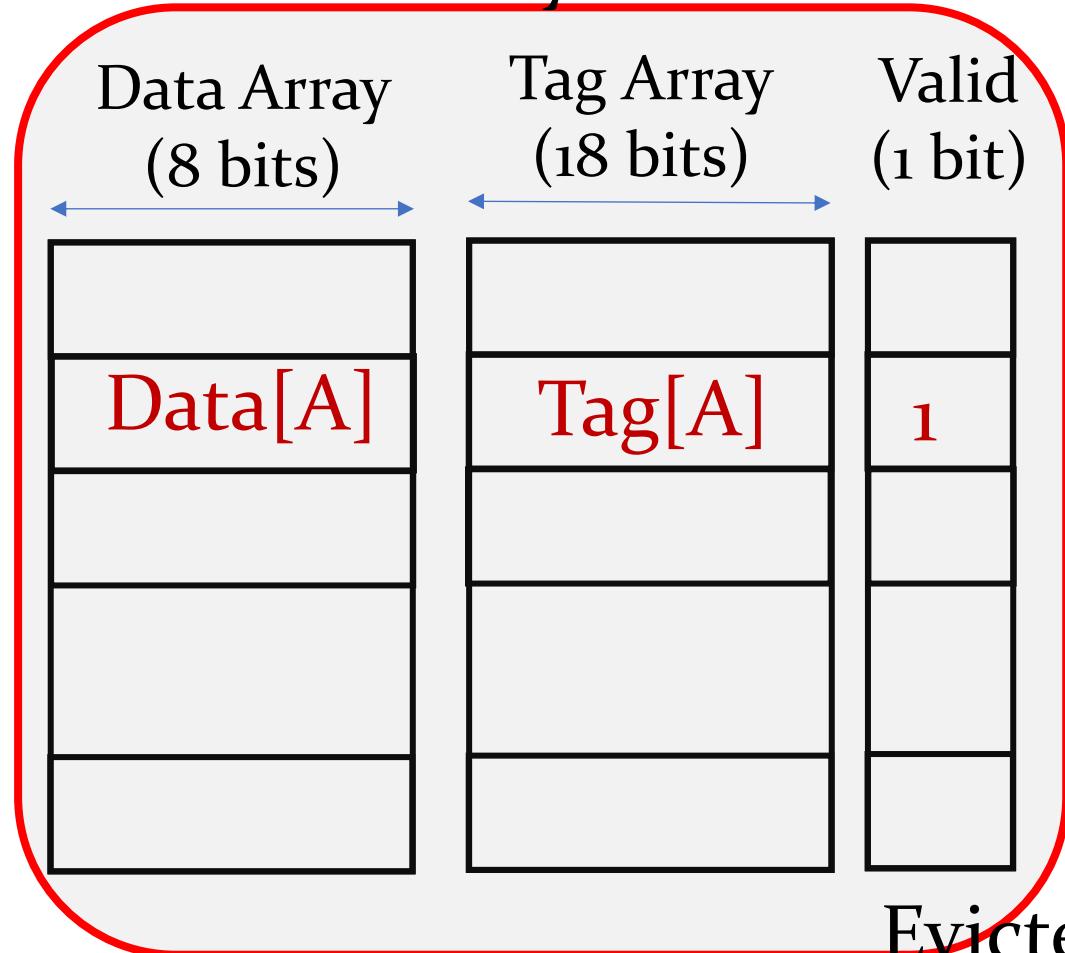


NYU

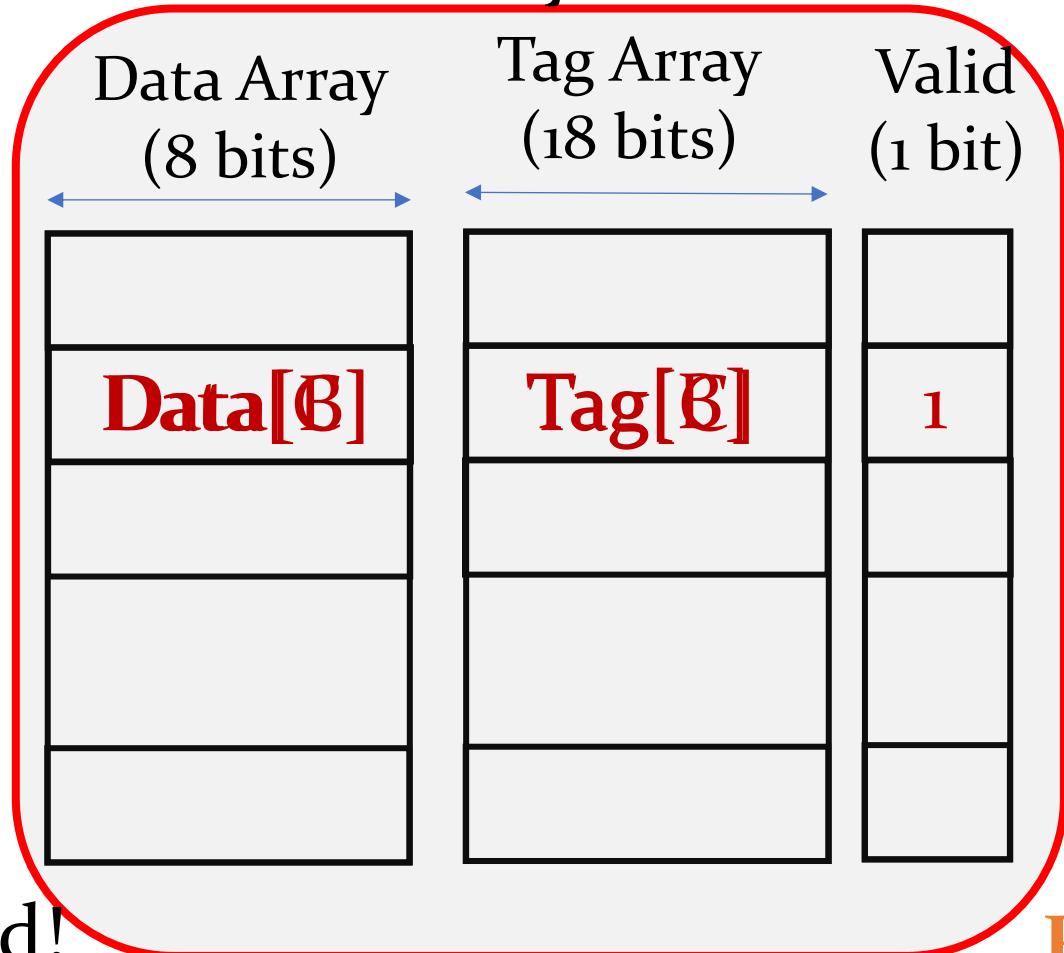
TANDON SCHOOL
OF ENGINEERING

Least Recently Used (LRU)

Way 0



Way 1



Evicted!

Sequence:

A B B
M M H

A A A
H H H

C
A C
H M

Replace least recently used block

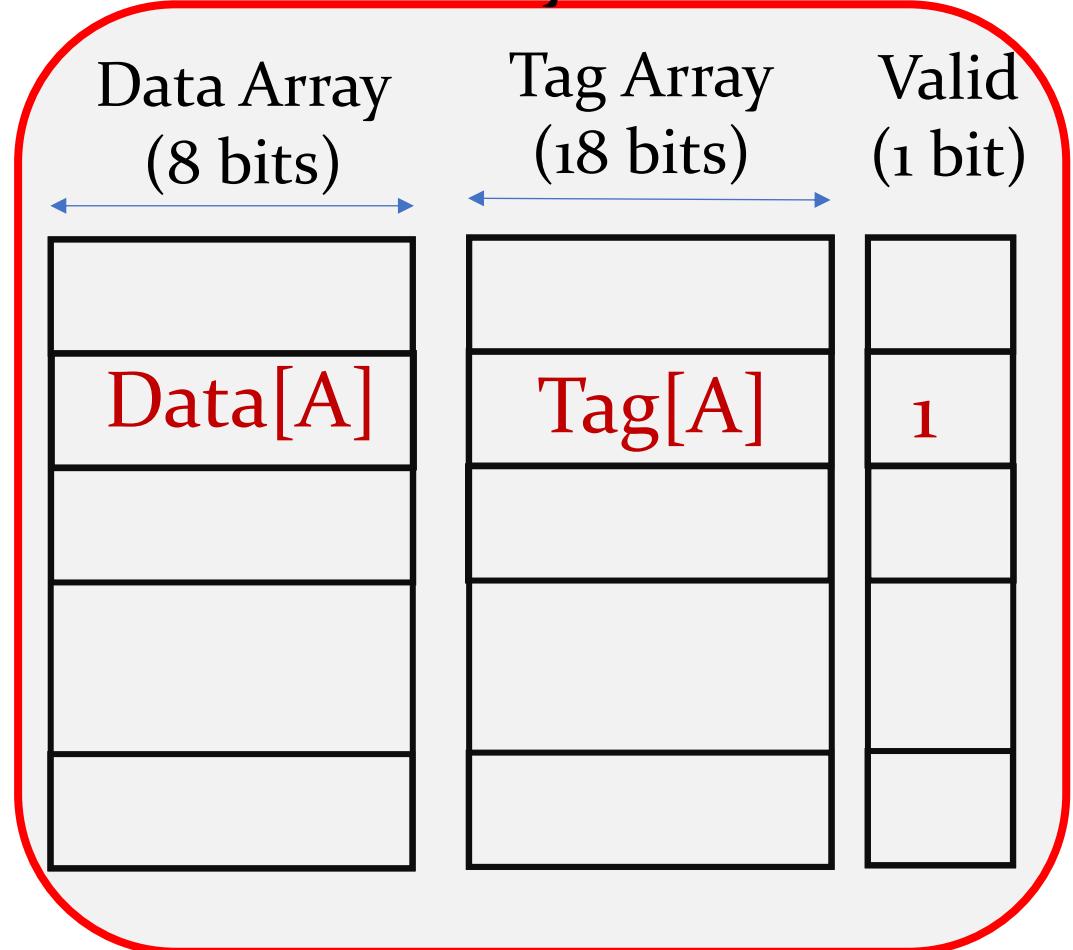


NYU

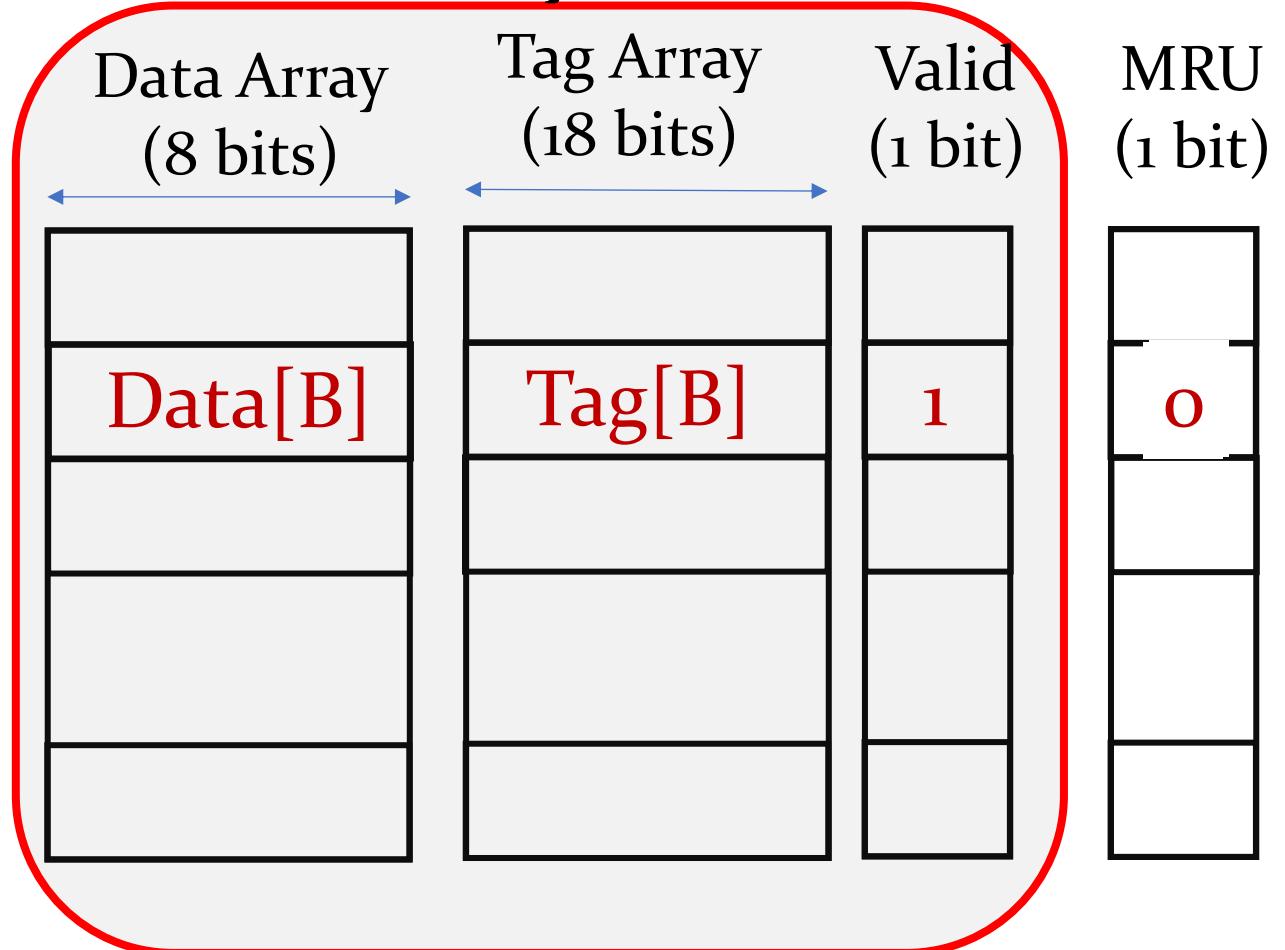
TANDON SCHOOL
OF ENGINEERING

Implementing LRU

Way 0



Way 1



Sequence:

A M M H H H H H C M

**LRU =
NOT MRU**



NYU

TANDON SCHOOL
OF ENGINEERING

Cache Replacement Policy

Random

- Replace a randomly chosen line

FIFO

- Replace the oldest line

LRU (Least Recently Used)

- Replace the least recently used line

NRU (Not Recently Used)

- Replace one of the lines that is not recently used
- Commonly implemented



NYU

TANDON SCHOOL
OF ENGINEERING

Practical LRU Implementation

LRU is hard to implement in hardware when $N > 2$

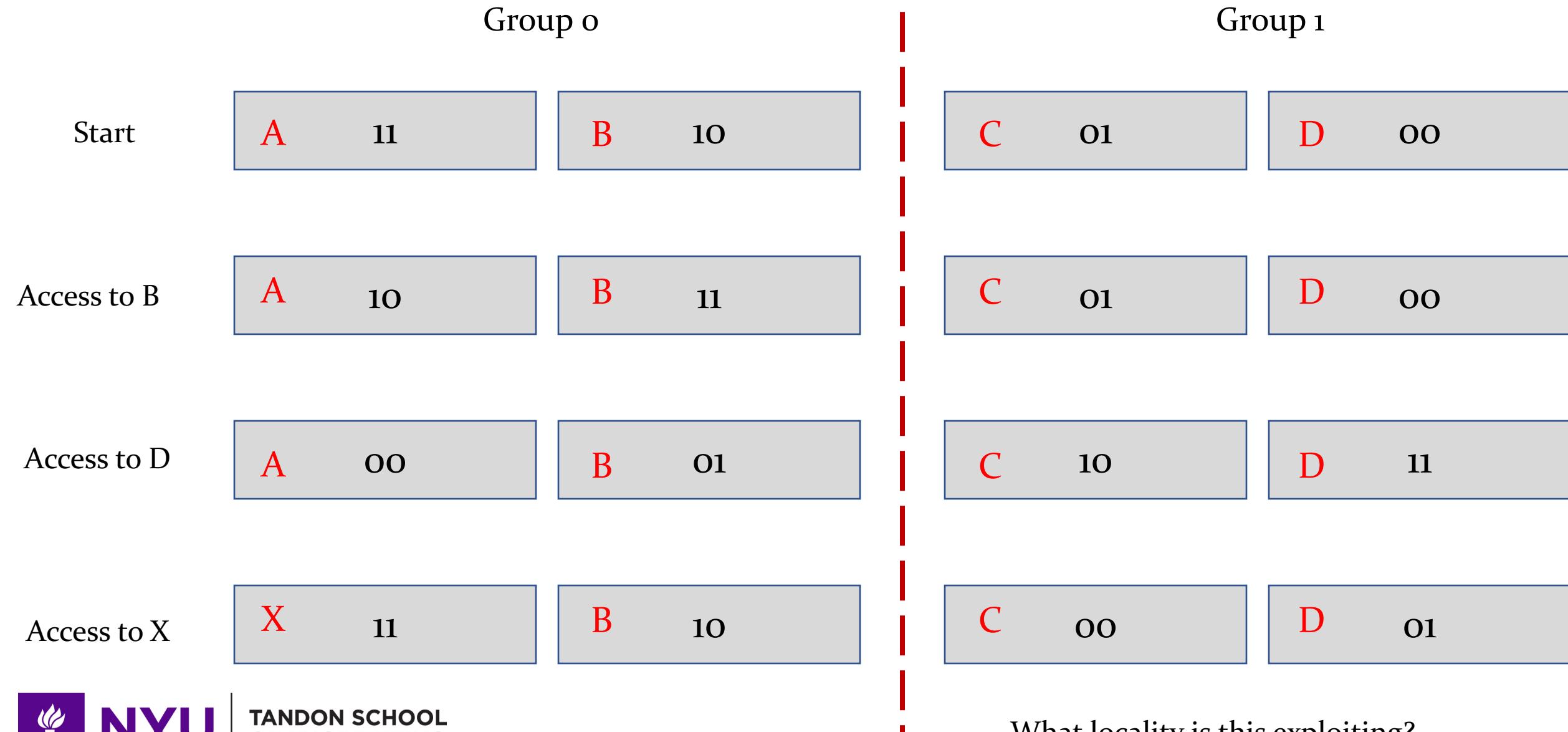
- Keep track of all possible $N!$ orderings of N ways
- A linked list in which the head points to MRU and tail points to LRU
- Example: $N=4$; 4×2 bits = 8 bits per cache set and extra logic to update list on every access

Intuition: LRU is an approximation anyways..

Alternative policies that are more hardware friendly

- **NOT MRU**: same as LRU for $N=2$, requires only $\log(N)$ bits, easy update
- **Hierarchical**: for $N=4$, divide ways into 2 groups of 2 ways

Hierarchical LRU example



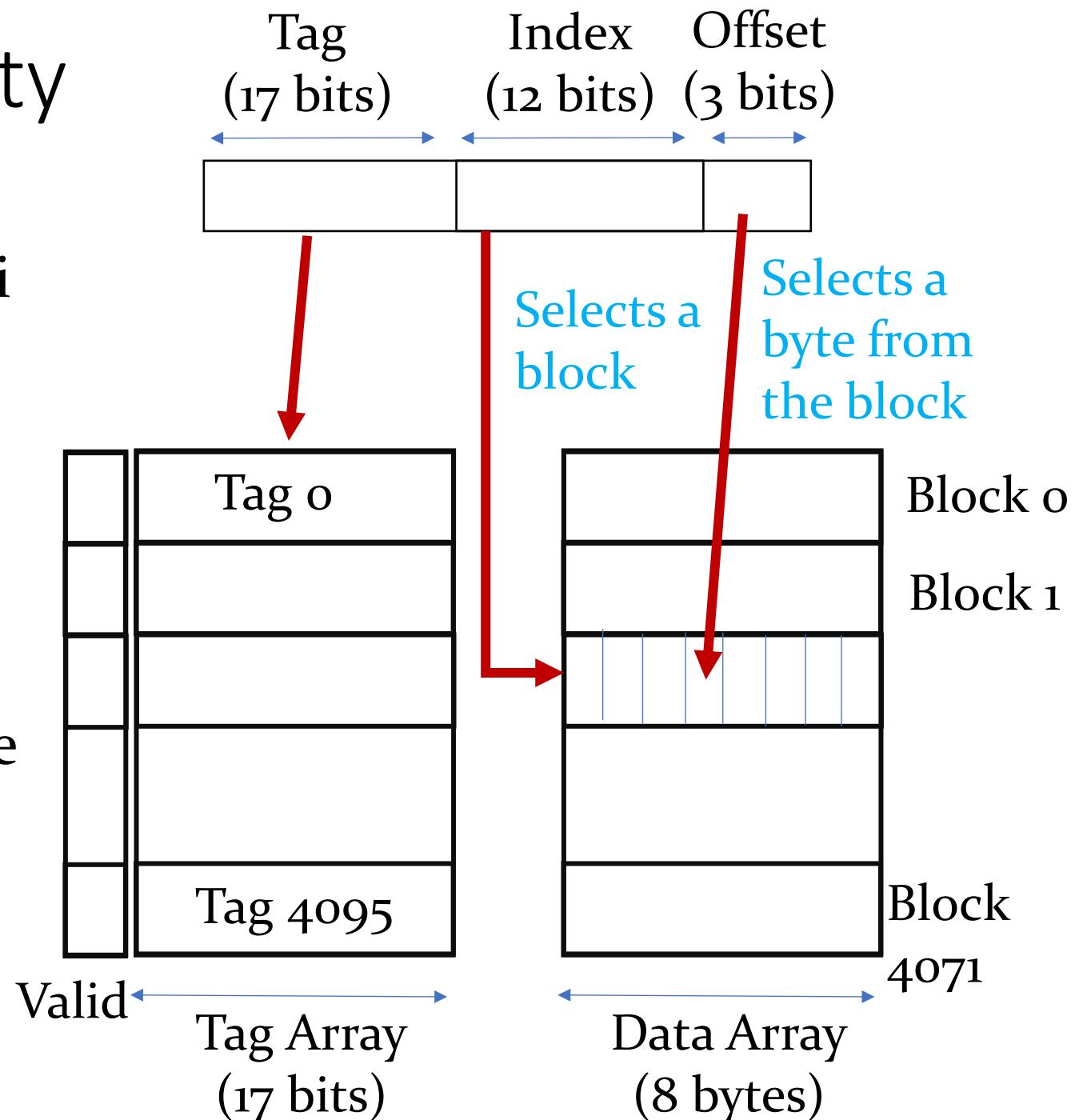
Exploiting Spatial Locality

Recall that if the byte from address i is accessed, then byte from address $i+1$ is likely to be accessed

- Pull in *multiple* contiguous bytes of data in each access
- Use larger block size!

Example: 32KB direct mapped cache with 8 Byte (64 bit) blocks

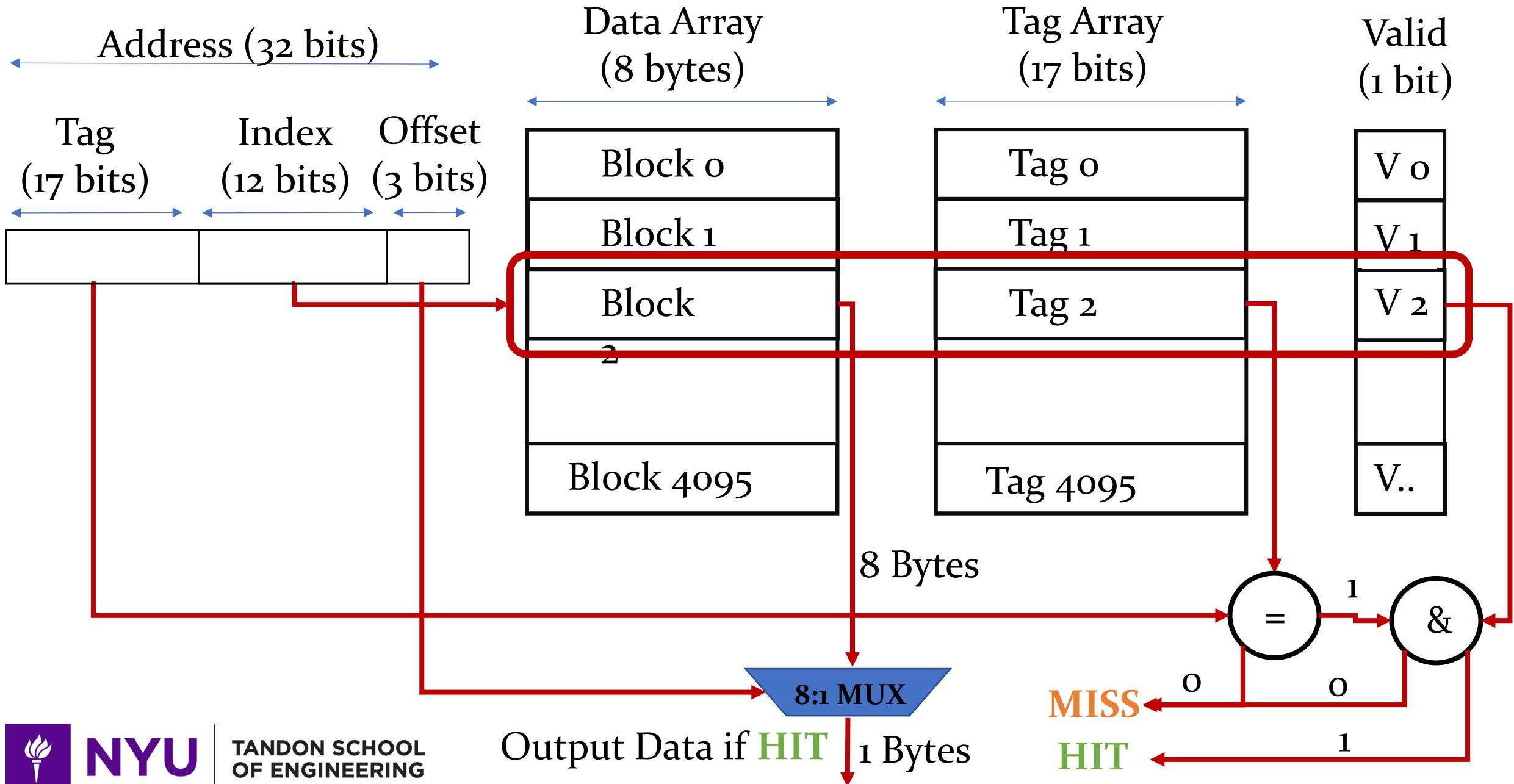
- i.e., cache has 4096 Byte blocks



NYU

TANDON SCHOOL
OF ENGINEERING

Cache Operation



NYU

TANDON SCHOOL
OF ENGINEERING

Impact of Block Size

Sequence of addresses: A, A+1, A+2, A+3 ...

- 4 consecutive misses for 1 byte block size
- 1 miss and 3 hits for 4 byte block size

Small block sizes don't exploit any spatial locality

What happens if the block size increases for the same cache size

- Fewer number of larger blocks

