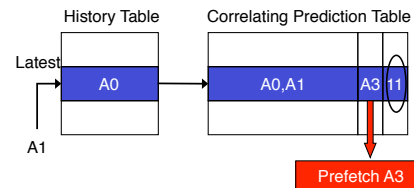


## Lecture 7 Prefetching

Spring 2016  
Pejman Lotfi-Kamran



Adapted from slides originally developed by Profs. Falsafi, Hill, Hoe, Lipasti, Shen, Smith, Sohi, and Vijaykumar of Carnegie Mellon University, EPFL, Purdue University, and University of Wisconsin.

Lecture 7  
Slide 1

## Where Are We?

Fr	Sa	Su	Mo	Tu
	27-Shahrivar		29-Shahrivar	
	3-Mehr		5-Mehr	
	10-Mehr		12-Mehr	
	17-Mehr		19-Mehr	
	24-Mehr		26-Mehr	
	1-Aban		3-Aban	
	8-Aban		10-Aban	
	15-Aban		17-Aban	
	22-Aban		24-Aban	
	29-Aban		1-Azar	
	6-Azar		8-Azar	
	13-Azar		15-Azar	
	20-Azar		22-Azar	
	27-Azar		29-Azar	
	4-Dey		6-Dey	

This Lecture  
▢ Prefetching

Next Lecture:  
▢ Virtual Memory

Lecture 7  
Slide 2

## Why Prefetch?

Cache miss is either expensive or very expensive

If we can foretell which address the program will reference in the future then we can ensure the location is in the cache ahead of time ⇒ *No cache miss!!*

Prefetching takes advantage of regular/repeatable program behavior (in this case the memory reference pattern)

Prefetching is quite safe -- only involves prediction but no “speculative execution”

*Prefetching the wrong location can only affect processor performance (more misses) but not correctness*

Lecture 7  
Slide 3

## What is Prefetching?

- Fetch memory ahead of time
- Targets compulsory, capacity, & coherence misses

Big challenges:

1. knowing “what” to fetch
  - Fetching useless info wastes valuable resources
2. “when” to fetch it
  - Fetching too early clutters storage
  - Fetching too late defeats the purpose of “pre”-fetching

Lecture 7  
Slide 4

## Software Prefetching

Compiler/programmer places prefetch instructions

- ▢ requires ISA support
- ▢ why not use regular loads?
- ▢ found in ISA's such as SPARC V-9

Prefetch into

- ▢ register (binding)
- ▢ caches (non-binding): preferred in multiprocessors

Lecture 7  
Slide 5

## Software Prefetching (Cont.)

e.g.,

```
for (I = 1; I < rows; I++)  
  for (J = 1; J < columns; J++)  
  {  
    prefetch(&x[I+1,J]);  
    sum = sum + x[I,J];  
  }
```

Lecture 7  
Slide 6

## Software Prefetching Support

PowerPC Data Cache Block Touch Instruction (**dcbt EA**)

*"a hint that performance will probably be improved if the block containing the byte addressed by **EA** is fetched into the data cache"*

A correct implementation of **dcbt** is to do nothing

Or, as a load instruction with no destination register  
except it should not trigger page or protection faults

Where should compilers insert **dcbt**?

- ▢ in front of every load: *wastes I-cache and D-cache bandwidth*
- ▢ where are loads likely to miss
  - When traversing large data sets (arrays in scientific code)
- ▢ where load misses would really hurt performance
  - pointer arguments to functions
  - linked-list traversal - find loads whose data address is itself the result of a previous load

Lecture 7  
Slide 7

## Hardware Prefetching

What to prefetch?

- ▢ one block spatially ahead?
- ▢ use address predictors → work well for regular patterns (e.g., x, x+8, x+16,...)

When to prefetch?

- ▢ on every reference
- ▢ on every miss
- ▢ when prior prefetched data is referenced

Where to put prefetched data?

- ▢ auxiliary buffers
- ▢ caches

Lecture 7  
Slide 8

## Spatial Locality and Sequential Prefetching

Works well for I-cache

- Instruction fetching tend to access memory sequentially

Doesn't work very well for D-cache

- More irregular access pattern
- regular patterns may have non-unit stride (e.g. matrix code)

Relatively easy to implement

- Large cache block size already have the effect of prefetching
- After loading one-cache line, start loading the next line automatically if the line is not in cache and the bus is not busy

What if you fetch at the wrong time ....

*Imagine if you started sequential prefetching of a long cache line and so happens you get a load miss to the middle of that line?*

*A critical-word-first reload triggered by the load miss itself may actually have restarted computation sooner!!*

Lecture 7  
Slide 9

## Stride Prefetchers

Access pattern for a particular static load is more predictable

Reference Prediction Table

Load Inst PC	Load Inst.	Last Address	Last	Flags
	PC (tag)	Referenced	Stride	
	.....	.....	.....	

Remembers previously executed loads, their PC, the last address referenced, stride between the last two references

When executing a load, look up in RPT and compute the distance between the current data addr and the last addr

- if the new distance matches the old stride  
⇒ found a pattern, go ahead and prefetch “current addr+stride”
- update “last addr” and “last stride” for next lookup

Lecture 7  
Slide 10

## Stream Buffers

Each stream buffer holds one stream of sequentially prefetched cache lines

*No cache pollution*

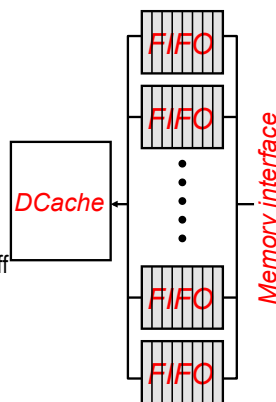
On a load miss check the head of all stream buffers for an address match

- if hit, pop the entry from FIFO, update the cache with data
- if not, allocate a new stream buffer to the new miss address (may have to recycle a stream buffer following LRU policy)

Stream buffer FIFOs are continuously topped-off with subsequent cache lines whenever there is room and the bus is not busy

Stream buffers can incorporate stride prediction mechanisms to support non-unit-stride streams

*Indirect array accesses (e.g., A[B[i]])?*



Lecture 7  
Slide 11

## Generalized Access Pattern Prefetchers

How do you prefetch

1. Heap data structures?
2. Indirect array accesses?
3. Generalized memory access patterns?

Current proposals:

- Precomputation prefetchers
- Address correlating prefetchers

Lecture 7  
Slide 12

## Runahead Prefetchers

Proposed for I/O prefetching first (Gibson et al.)

Duplicate the program

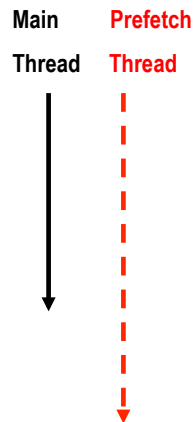
- Only execute the address generating stream
- Let it run ahead

May run as a thread on

- A separate processor
- The same multithreaded processor

Or custom address generation logic

Many names: slipstream, precomp., runahead, ...



Lecture 7  
Slide 13

## Runahead Prefetcher

To get ahead:

- Must avoid waiting
- Must compute less

Predict

1. Control flow thru branch prediction
2. Data flow thru value prediction
3. Address generation computation only

- + Prefetch any pattern (need not be repetitive)
- Prediction only as good as branch + value prediction

How much prefetch lookahead?

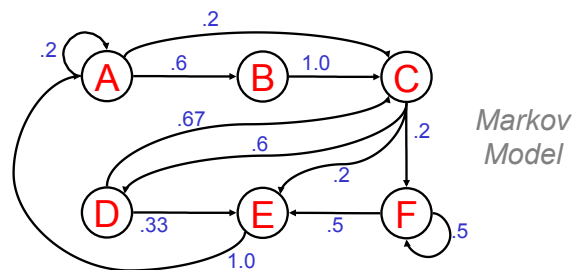
Lecture 7  
Slide 14

## Correlation-Based Prefetching

Consider the following history of Load addresses emitted by a processor

A, B, C, D, C, E, A, C, F, F, E, A, A, B, C, D, E, A, B, C, D, C

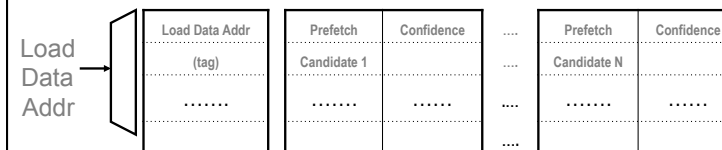
After referencing a particular address (say A or E), are some addresses more likely to be referenced next



Markov  
Model

Lecture 7  
Slide 15

## Correlation-Based Prefetching



Track the likely next addresses after seeing a particular addr.

Prefetch accuracy is generally low so prefetch up to N next addresses to increase coverage (but this wastes bandwidth)

Prefetch accuracy can be improved by using longer history

- Decide which address to prefetch next by looking at the last K load addresses instead of just the current one
- e.g. index with the XOR of the data addresses from the last K loads
- Using history of a couple loads can increase accuracy dramatically

This technique can also be applied to just the load miss stream

Lecture 7  
Slide 16

## Example Address Correlating: Markov Prefetchers

Markov Prefetchers (Joseph & Grunwald, ISCA'97)

- Correlate subsequent cache misses
- Trigger prefetch on miss
- Predict & prefetch 4 candidates: predicting 1 results in low coverage!
- Prefetch into a buffer

Lecture 7  
Slide 17

## Insufficient Lookahead in Markov

Distance between two misses is usually small  
Smaller in out-of-order cores

load/store A1 (miss)  
load/store A1 (hit)  
...  
load/store C3 (miss)  lookahead  
load/store A3 (miss) 

Lecture 7  
Slide 18



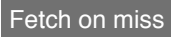

## Prefetch on Last Touch

Predict & fetch on last touch

Evict **dead block**

Enhance fetch lookahead

Fetch directly into L1

load/store A1 (miss)	load/store A1 (miss)
load/store A1 (hit)	load/store A1 (hit)
...	...
load/store C3 (miss) 	load/store C3 (miss)
load/store A3 (miss) 	load/store A3 (miss)
	

Lecture 7  
Slide 19

## Dead-Block Prediction

Correlate a **trace** of memory accesses to a block  
Uniquely identify different dead-times

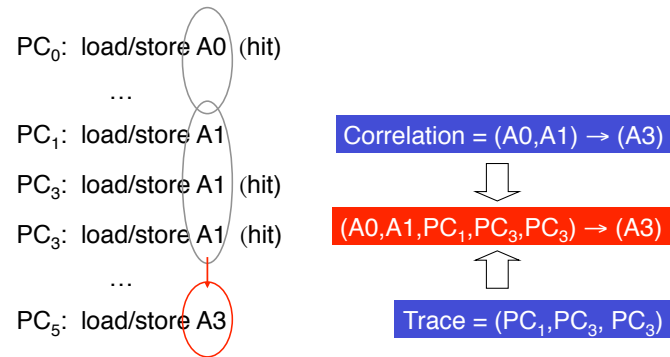
Access stream to a block frame

PC <sub>0</sub> : load/store A0 (hit)	
PC <sub>1</sub> : load/store A1 (miss)	<b>First touch</b> 
PC <sub>3</sub> : load/store A1 (hit)	
PC <sub>3</sub> : load/store A1 (hit)	<b>Last touch</b> 
PC <sub>5</sub> : load/store A3 (miss)	

Lecture 7  
Slide 20

## Miss-Address Prediction

Correlate last 2 misses within a cache block frame



Lecture 7  
Slide 21

## Improving Cache Performance: Summary

### Miss rate

- ▢ large block size
- ▢ higher associativity
- ▢ victim caches
- ▢ skewed-/pseudo-associativity
- ▢ hardware/software prefetching
- ▢ compiler optimizations

### Miss penalty

- ▢ give priority to read misses over writes/writebacks
- ▢ subblock placement
- ▢ early restart and critical word first
- ▢ non-blocking caches
- ▢ multi-level caches

### Hit time (difficult?)

- ▢ small and simple caches
- ▢ avoiding translation during L1 indexing (later)
- ▢ subblock placement for fast write hits in write through caches

Lecture 7  
Slide 22