

# Not Quite My Tempo<sup>1</sup>

Matching Prefetches to Memory Access Times

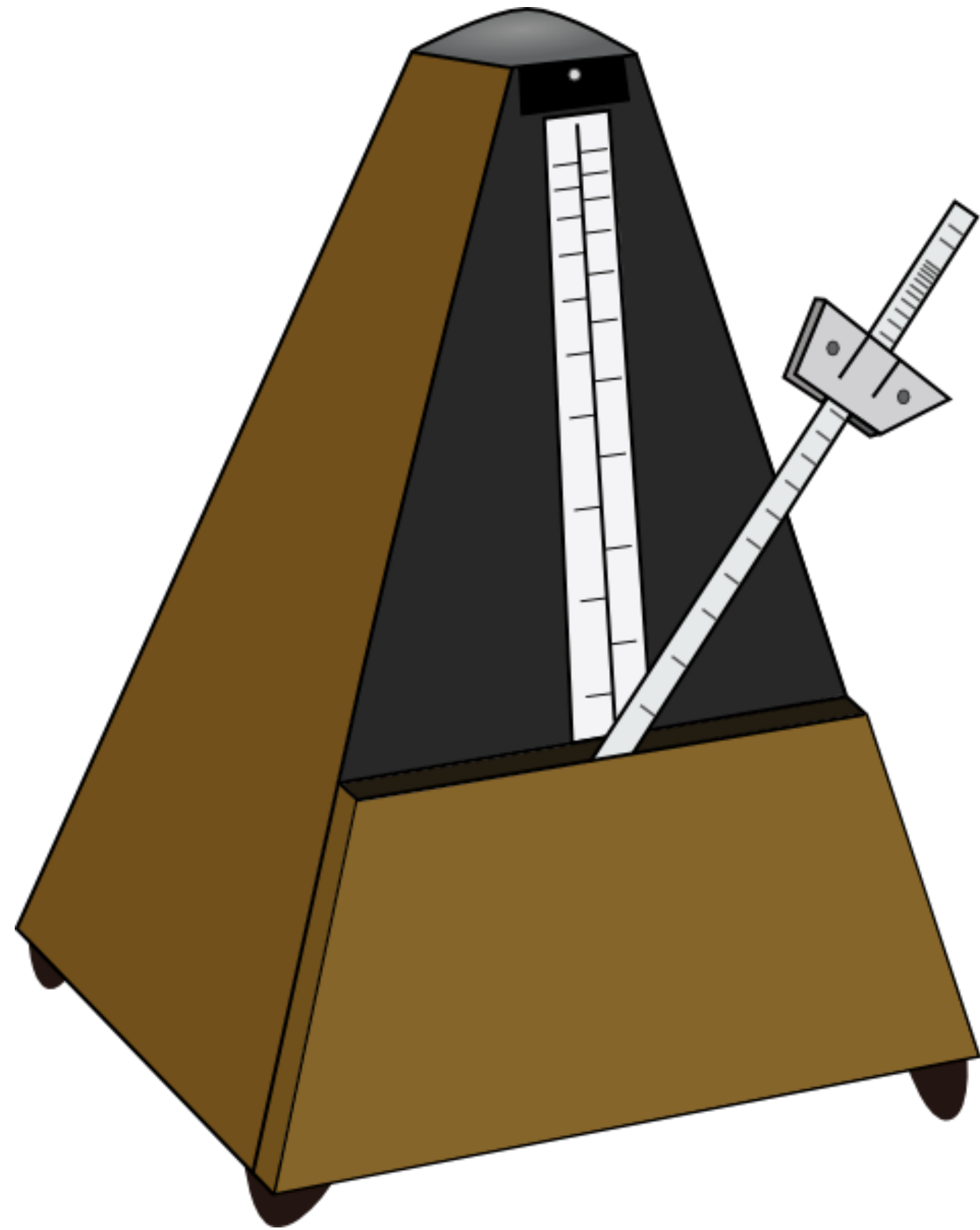
**Mark Sutherland**

Ajaykumar Kannan

Natalie Enright Jerger

{suther68,kannanaj,enright}@ece.utoronto.ca

1. This work was inspired by Chazelle, Simmons, and Teller's foundational work in this area: *Whiplash*.

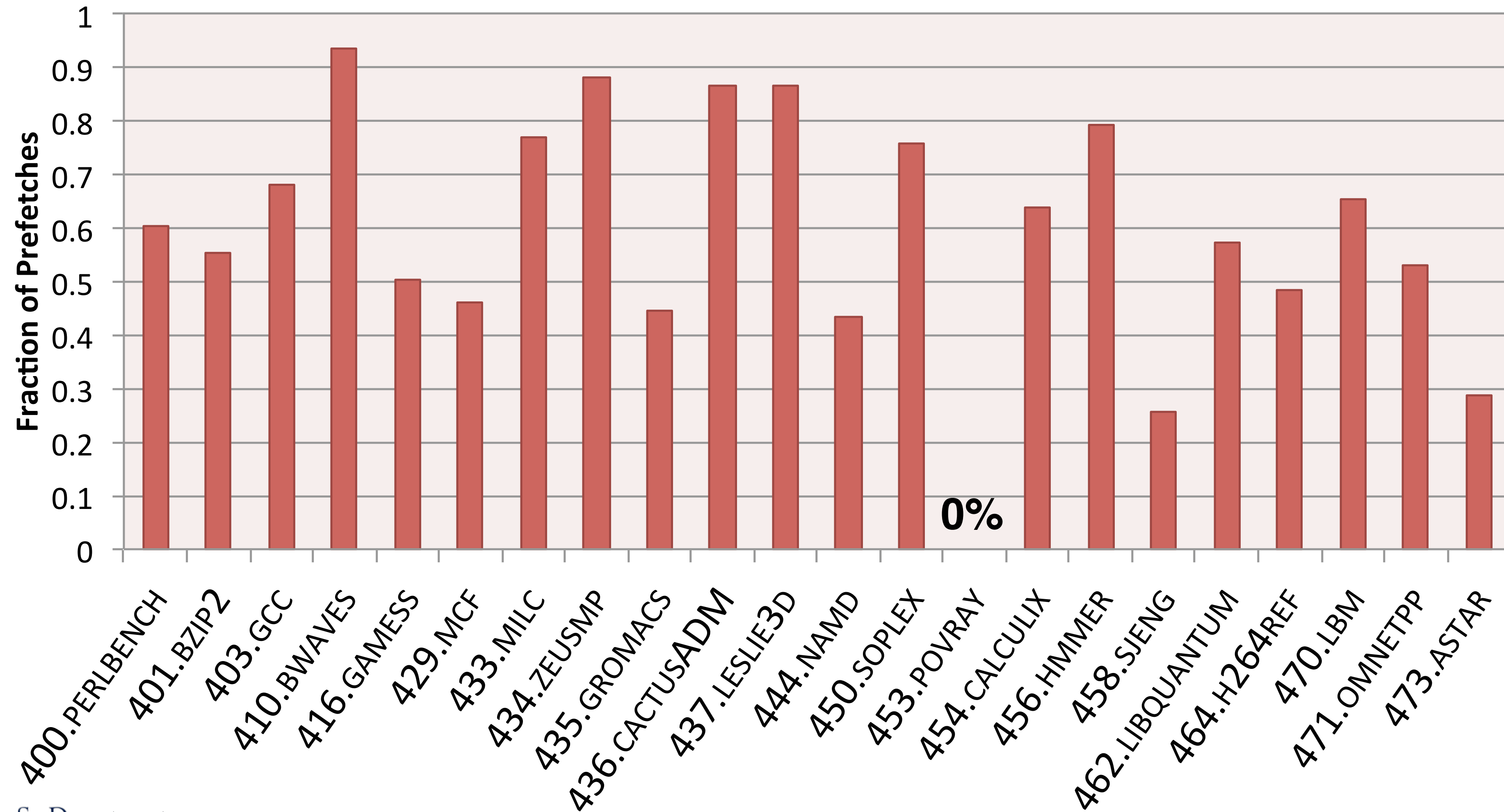


Map  
History Access  
Next-Line  
Stride Feedback  
Pattern SMS  
Stream Markov  
Local  
GHB

**Goal:** Add temporal information to an SMS prefetcher.

# Motivation

## Useless Prefetches in SMS



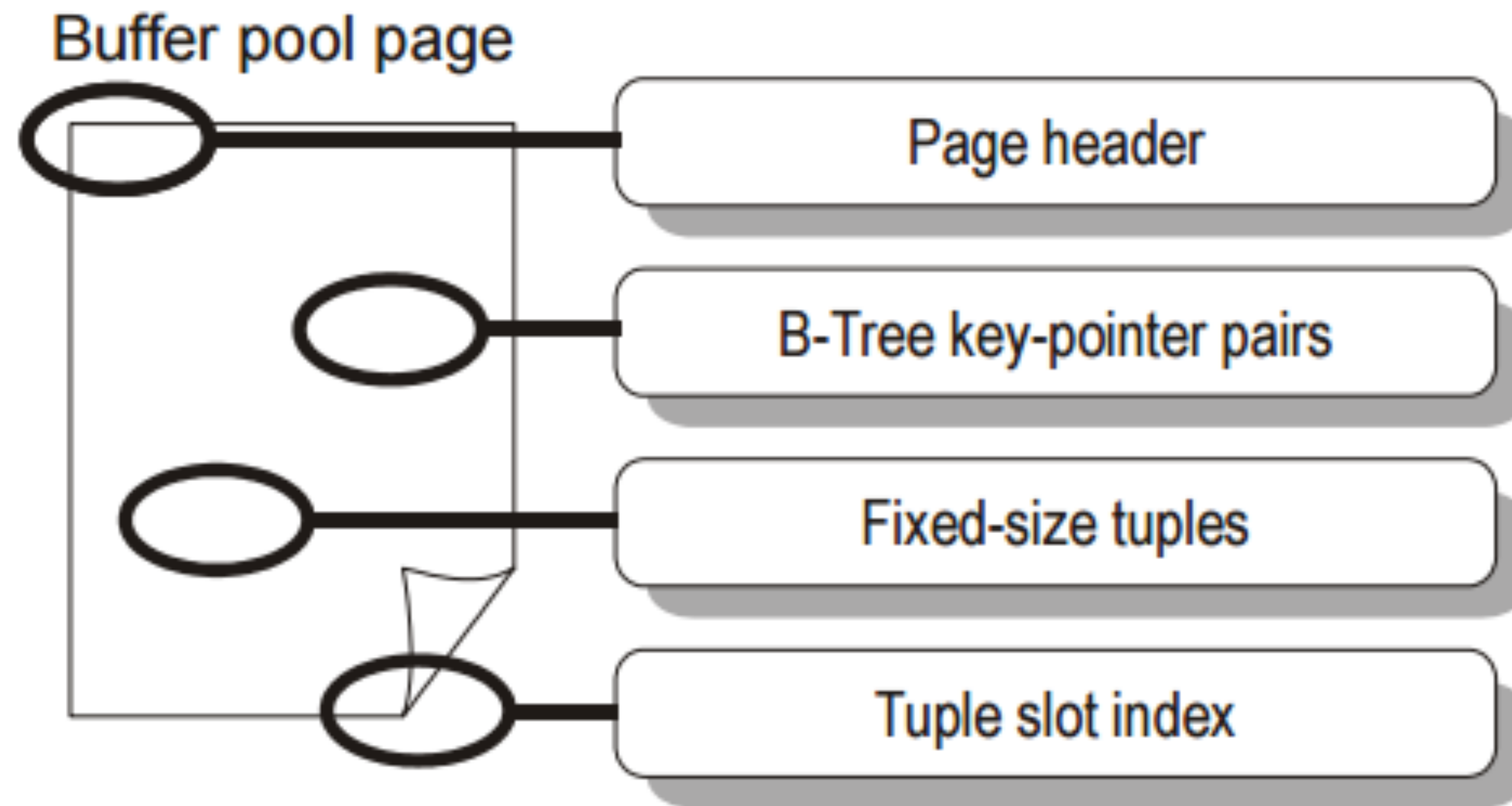
# Our Contributions

1. Propose a framework for timely prefetching in SMS, which utilizes *temporal* information to filter prefetches in the Pattern History Table.
2. Demonstrate practical improvements over SMS with a 32kB storage budget.



# SMS Prefetcher Background

## Examples of spatially-correlated elements in DBMSs



Source: S. Somogyi, T. F. Wenisch, A. Ailamaki, B. Falsafi, and A. Moshovos.  
*Spatial Memory Streaming. SIGARCH Comput. Archit. News*, 2006.

# SMS Prefetcher Background

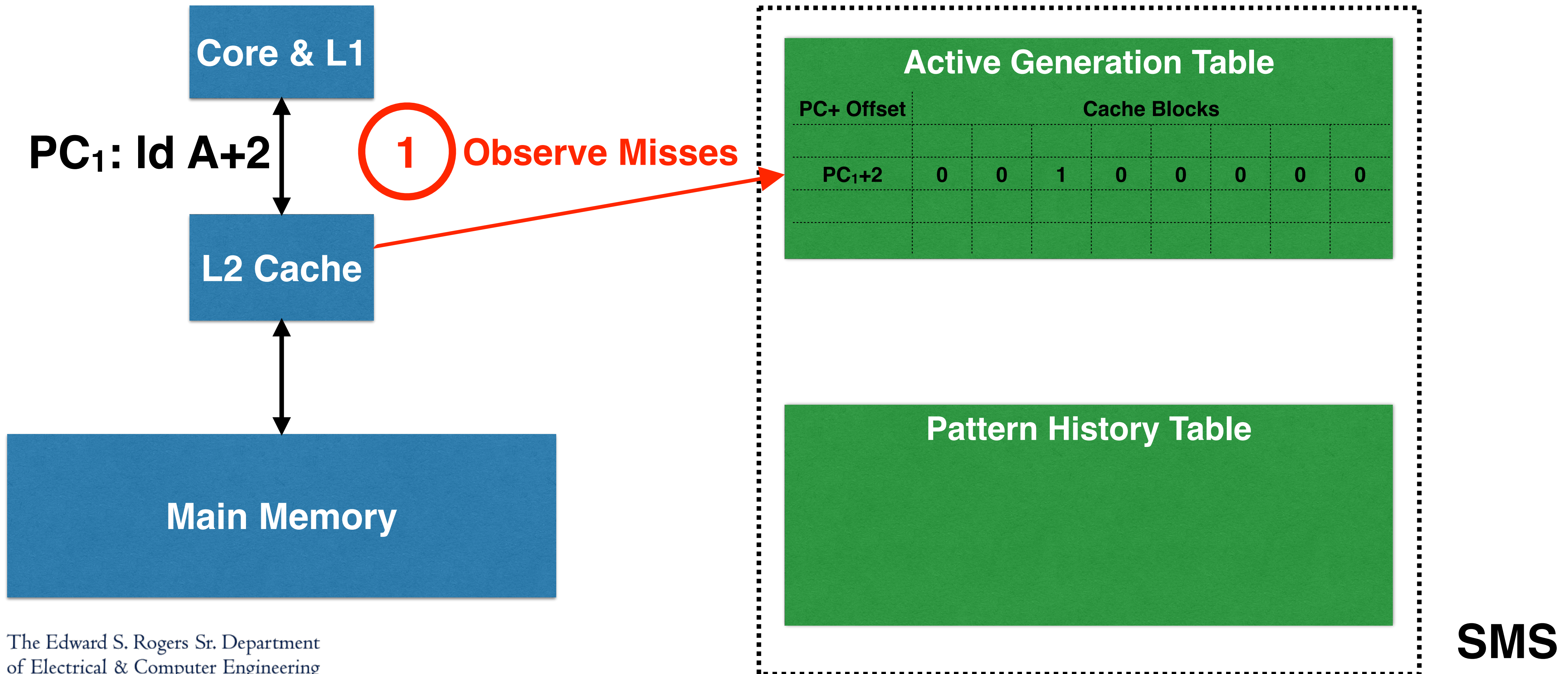
## Representation of a Memory Region

Page Bits		Accessed Cache Blocks = 1															
100110100	...	1	1	1	1	0	0	0	1	1	0	1	0	1	0	1	1
110100101	...	1	0	1	0	0	1	0	1	0	0	1	0	0	1	1	0
100000010	...	0	1	1	0	0	0	0	1	0	0	0	0	0	0	0	0

- Memory regions (most often the size of a physical page) are bit vectors, where a set bit indicates that the program accessed this block.

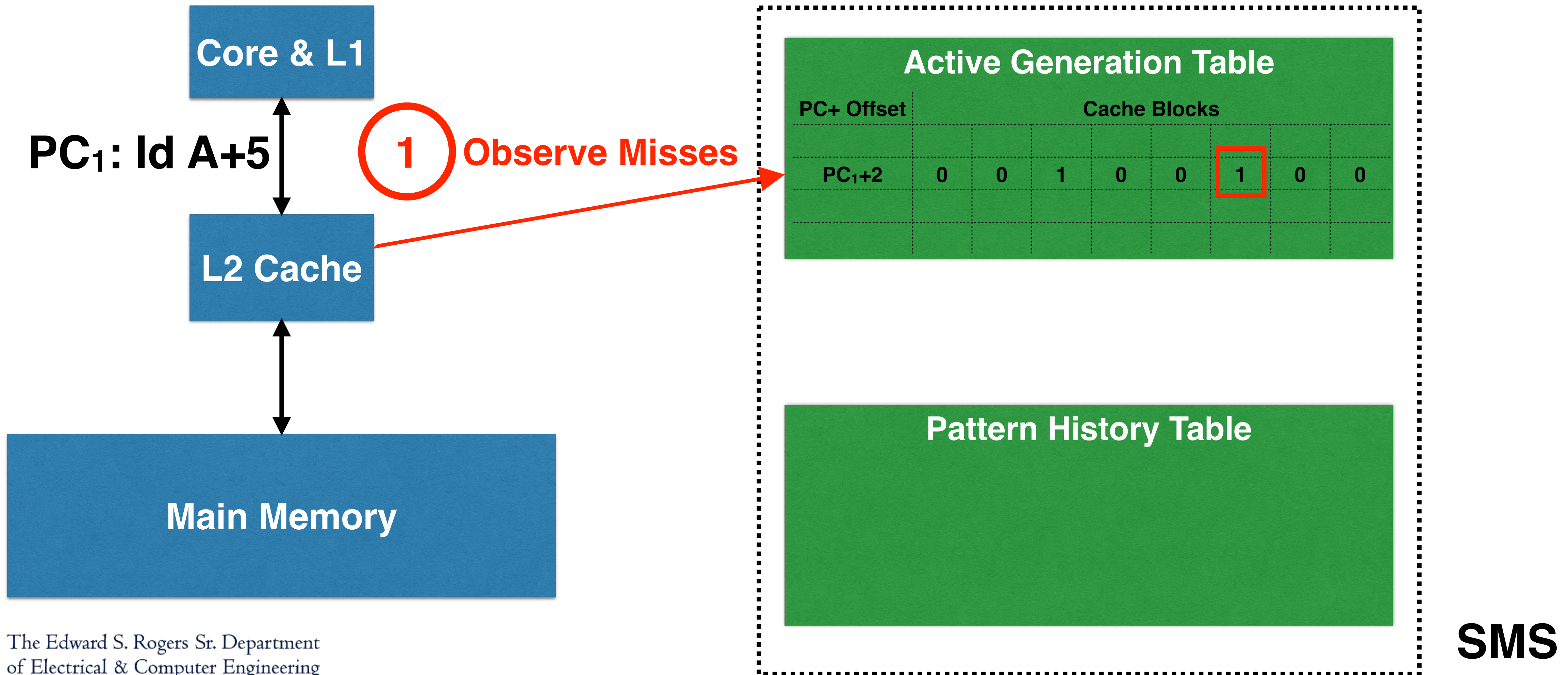


# SMS Prefetcher Background



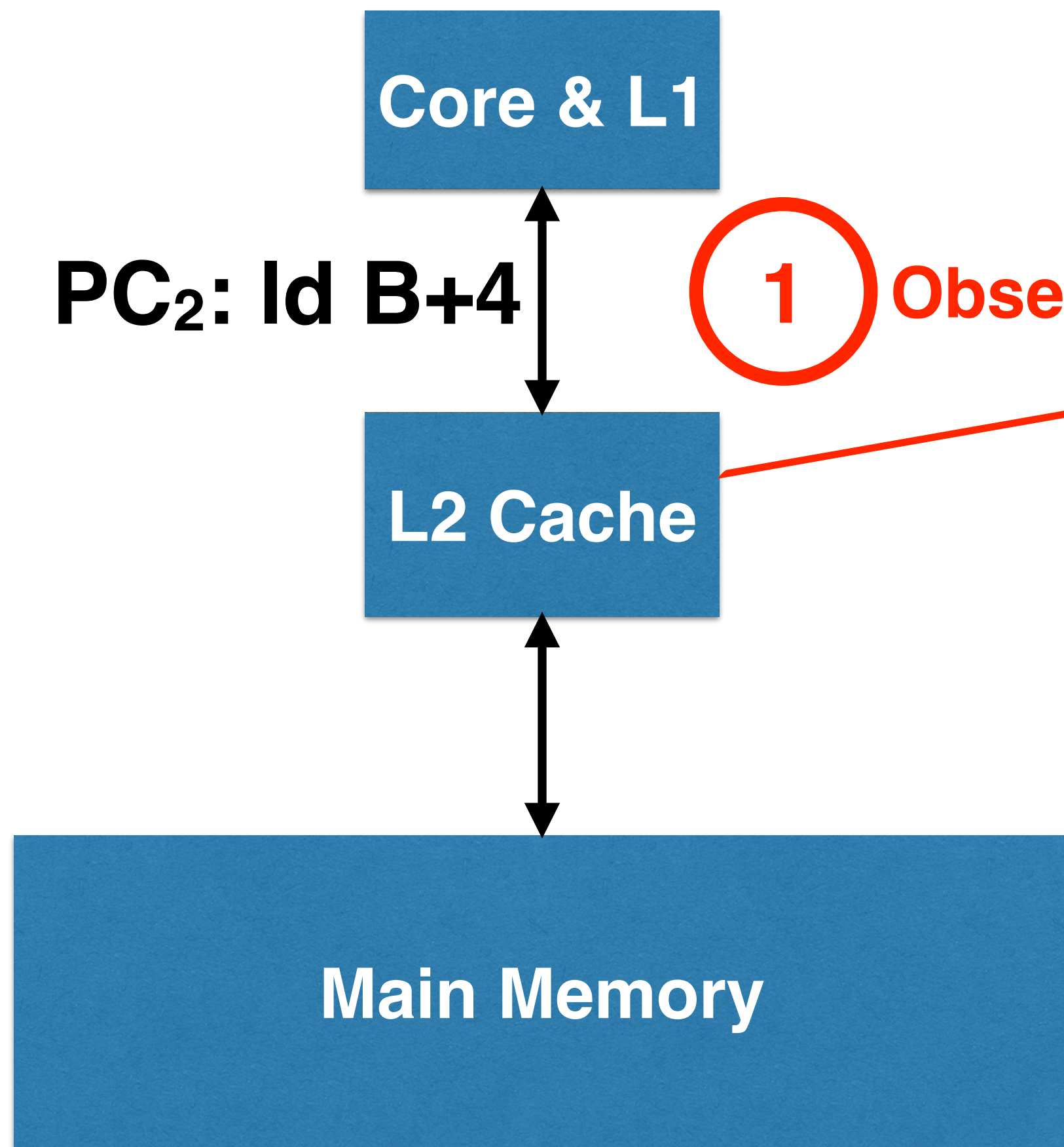


# SMS Prefetcher Background





# SMS Prefetcher Background



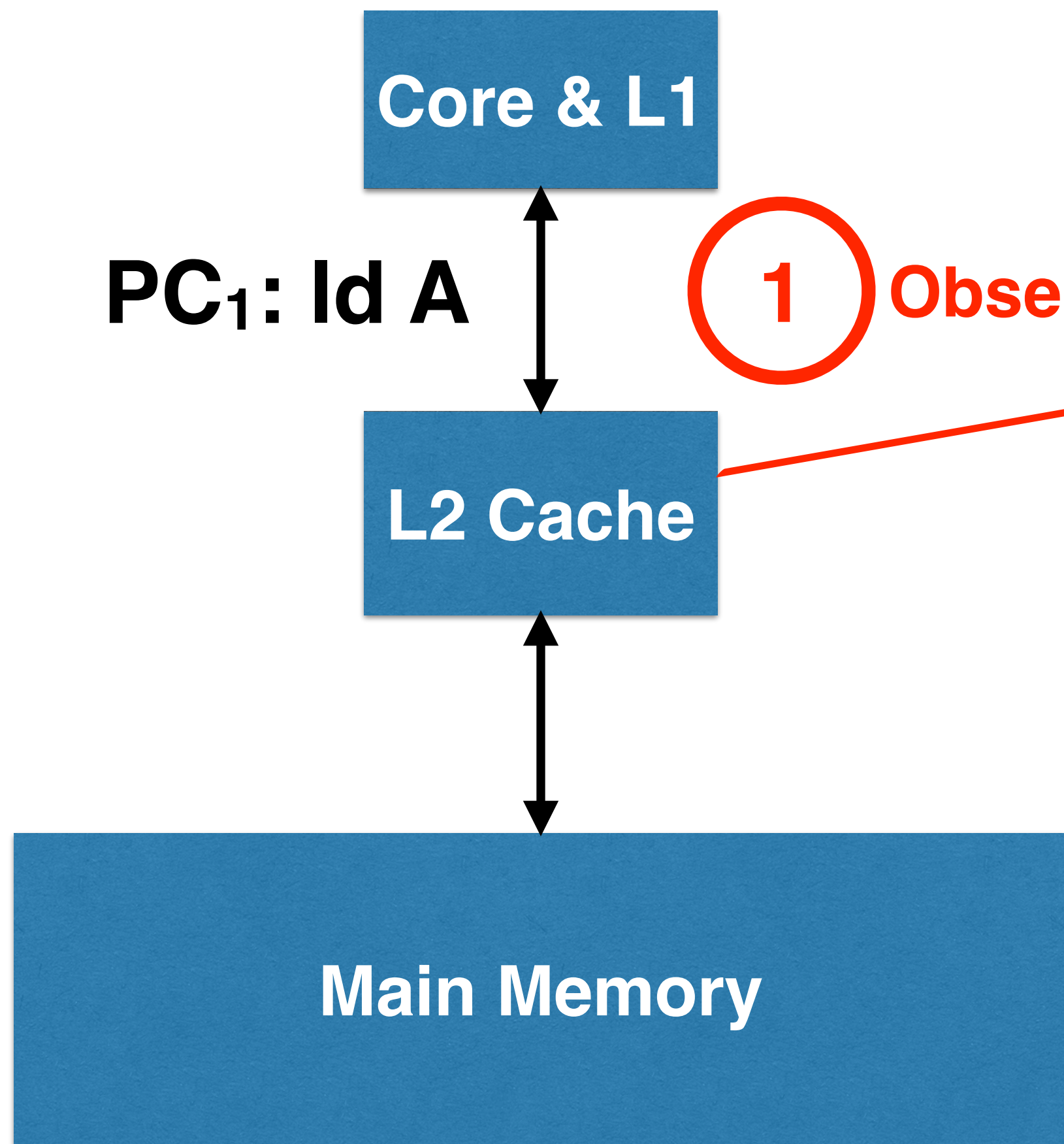
Active Generation Table									
PC+ Offset	Cache Blocks								
PC <sub>1</sub> +2	0	0	1	0	0	1	0	0	
PC <sub>2</sub> +4	0	0	0	0	1	0	0	0	

Pattern History Table									

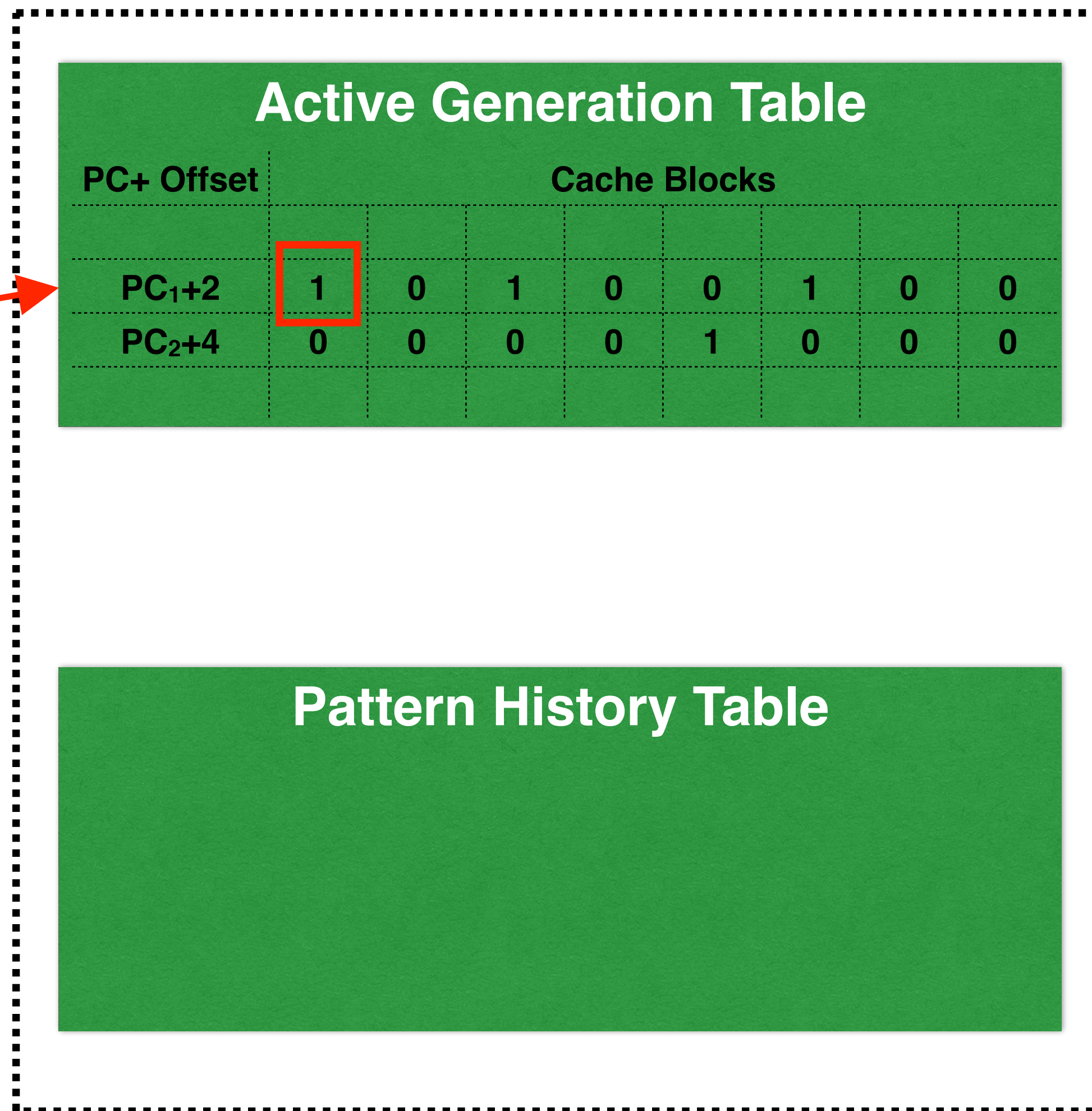
**SMS**



# SMS Prefetcher Background



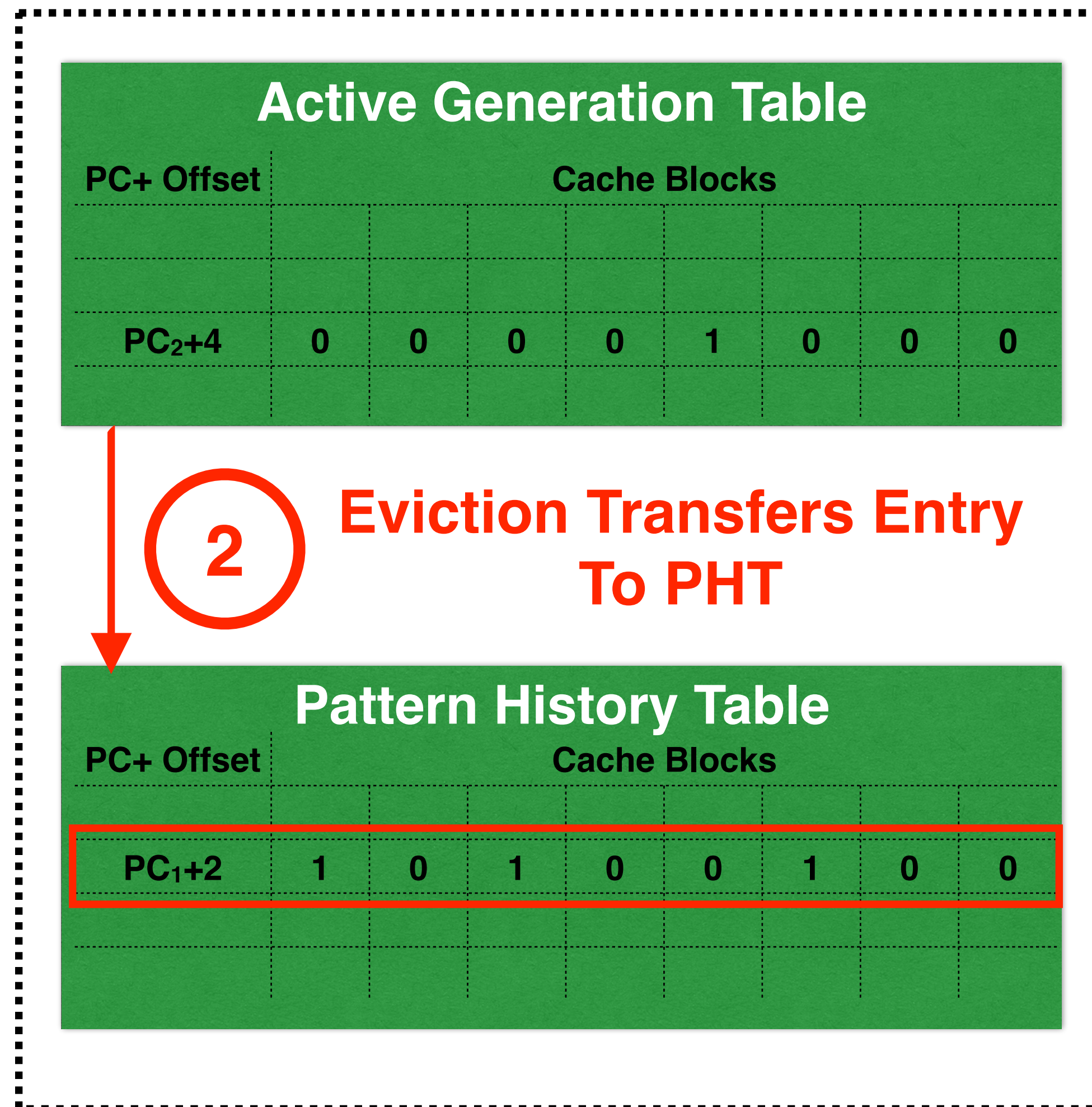
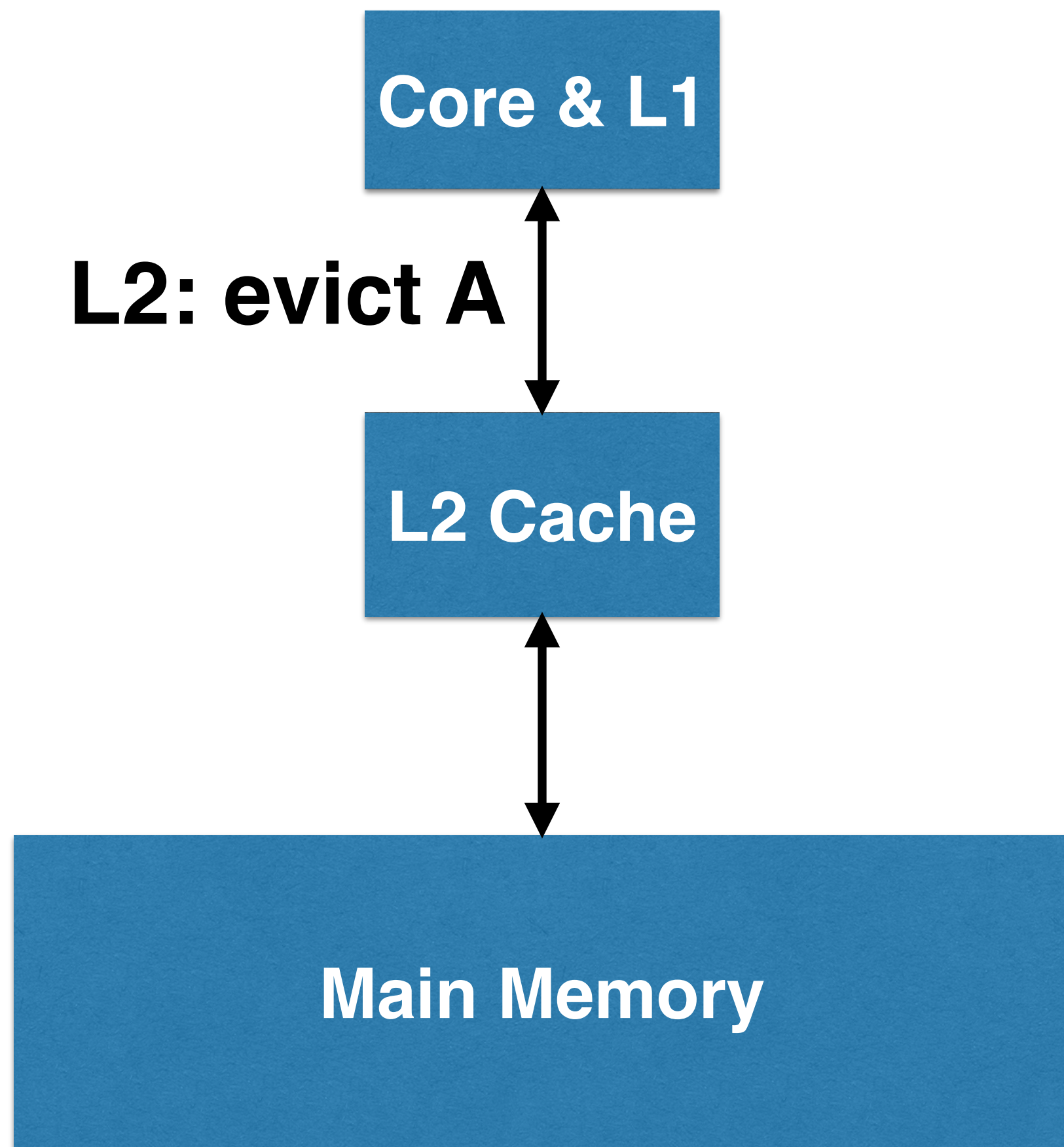
**1** Observe Misses



**SMS**



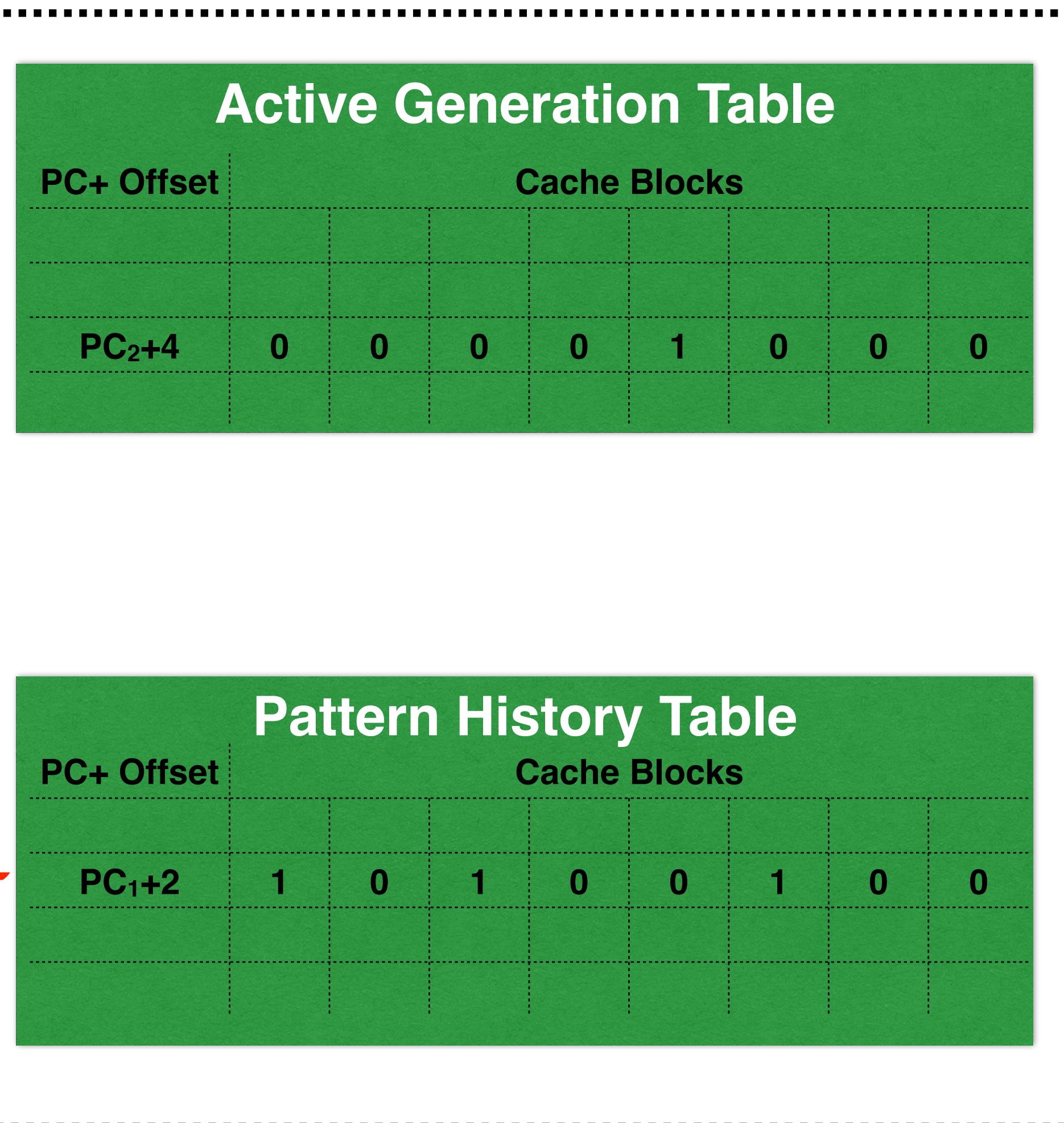
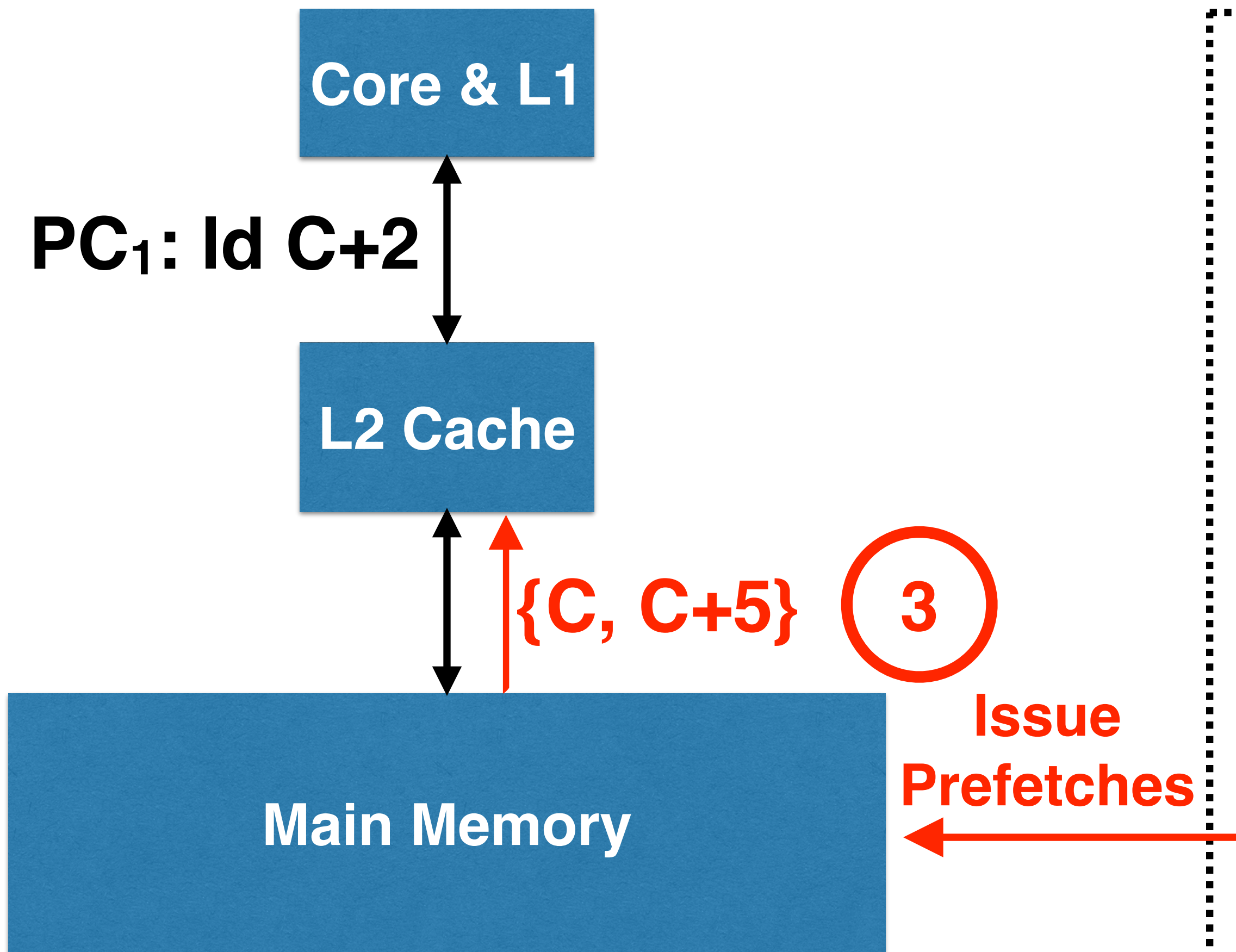
# SMS Prefetcher Background



**SMS**



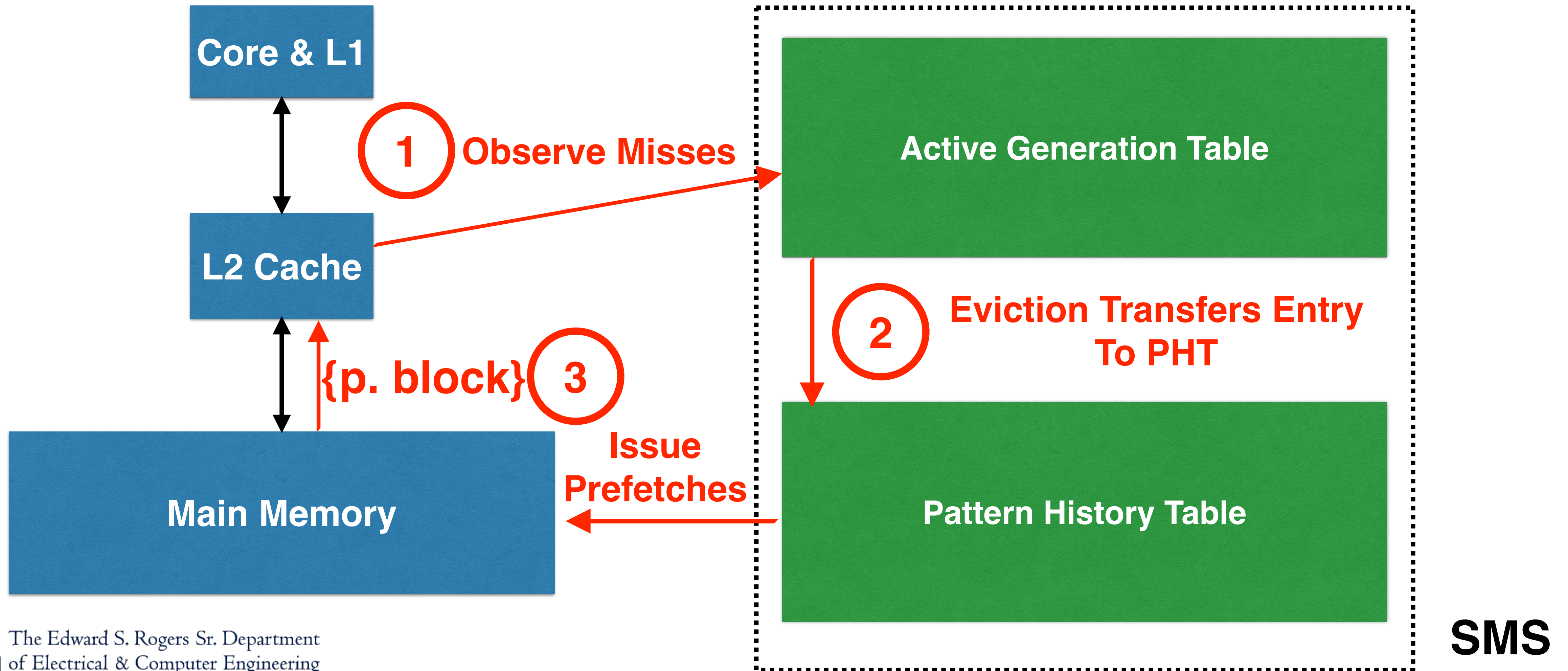
# SMS Prefetcher Background



**SMS**



# SMS Prefetcher Background

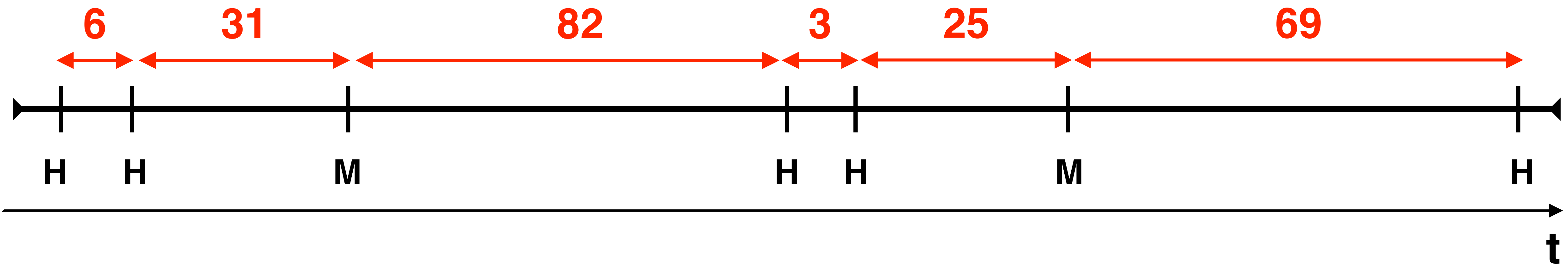


# Adding Temporal Information

- **Problem:** In the bit vectors stored by SMS, we do not have any information about the *order* of the accesses.
- **Goal:** A mechanism to distinguish various accesses within a spatial region, preferably one that allows us to prefetch blocks in a more timely fashion.

# Adding Temporal Information

- Observe that many of our applications see the same PC generating memory requests with different time deltas between each request.





# Adding Temporal Information

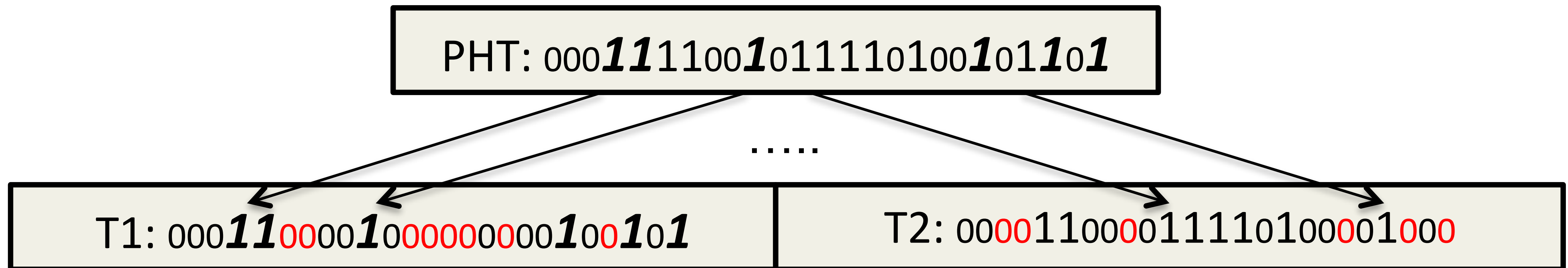
- Add a global counter to the L2 that ticks on every cache miss.
- Define **access tempo**: Delta between two cache accesses.





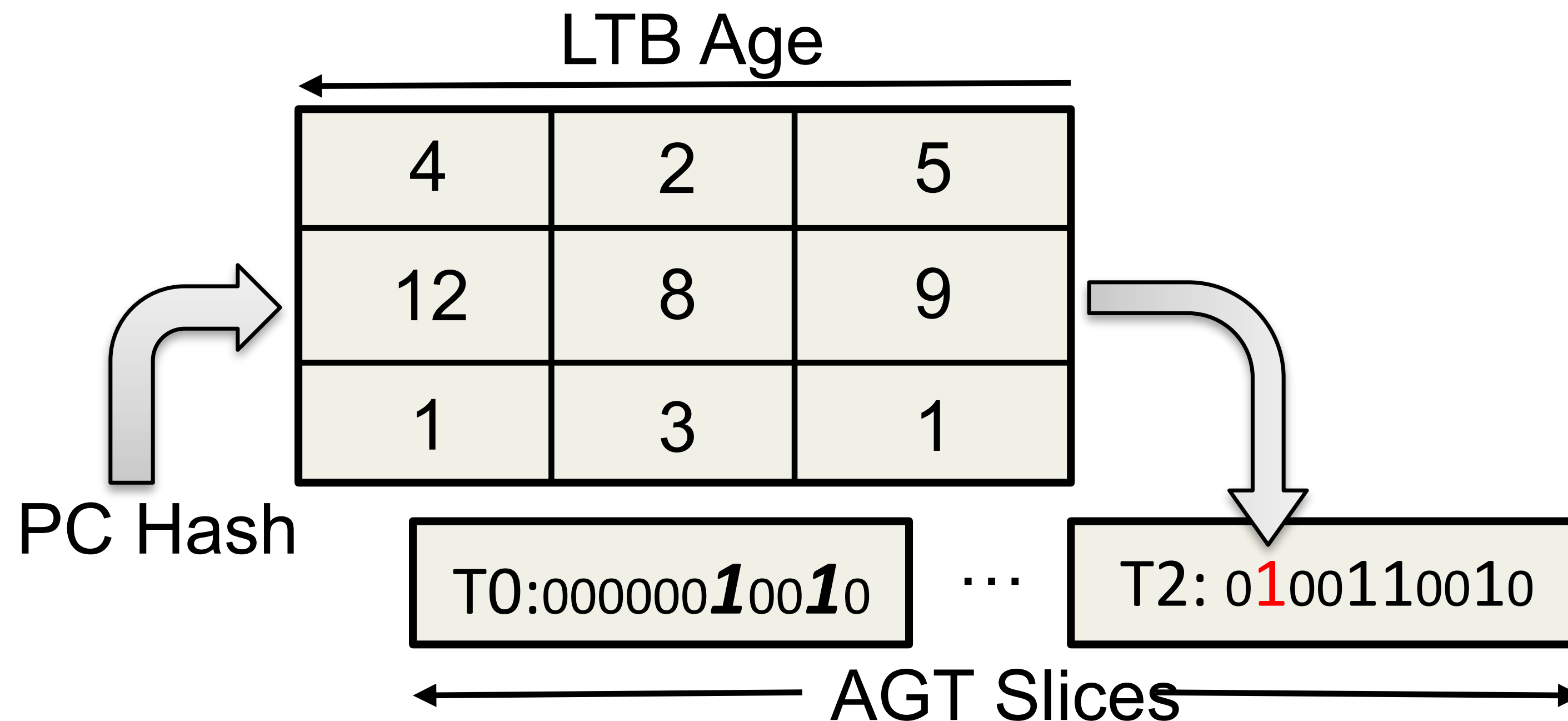
# Adding Temporal Information

- Decompose all PHT entries into multiple **access tempos**.
- Store each disjoint tempo in its own “slice” of SMS.



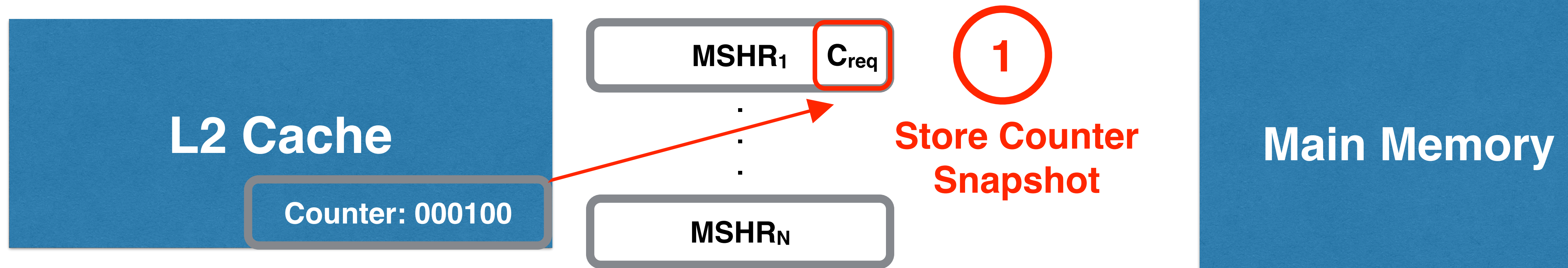
# Training Based on Tempo

- Store the access tempos in on-chip Localized Tempo Buffers (LTBs).
- During training, only set the bits of the vector corresponding to the **currently measured tempo**.



# Prefetching Based on Tempo

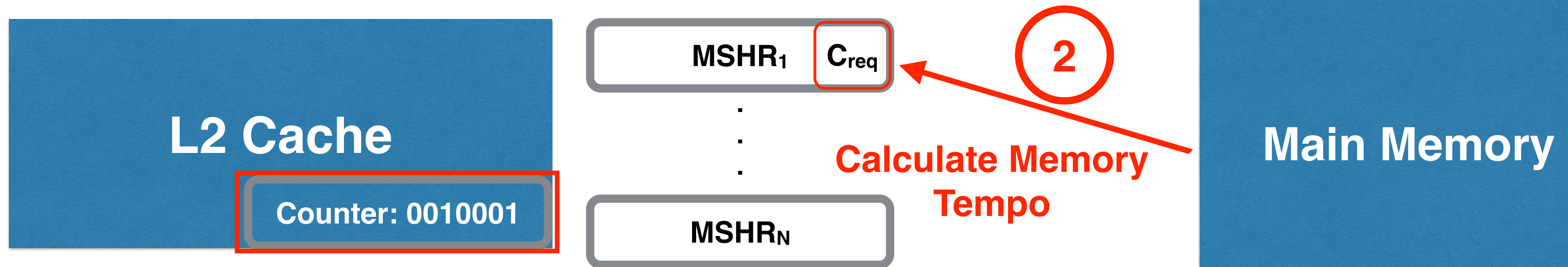
- Simply knowing our current *access tempo* is not enough to ensure prefetch timeliness.
- Define **memory tempo**: Delta between when a memory request is issued, and when the block is filled into the cache.





# Prefetching Based on Tempo

- Simply knowing our current *access* tempo is not enough to ensure prefetch timeliness.
- Define **memory tempo**: Delta between when a memory request is issued, and when the block is filled into the cache.



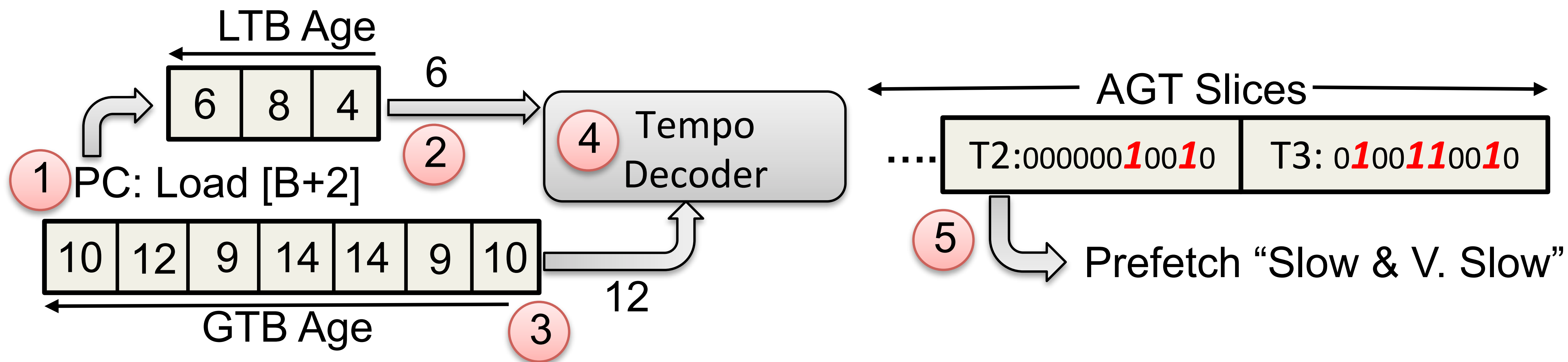
# Rushing or Dragging?

- Can compare the memory's tempo against the current PC's access tempo at prefetch time.

PC > Memory (Rushing)	Memory > PC (Dragging)
<ul style="list-style-type: none"><li>• Generating requests faster than memory is returning them.</li><li>• Prefetch only <b>slower</b> tempos.</li></ul>	<ul style="list-style-type: none"><li>• Memory is fast enough to keep up with our requests.</li><li>• Prefetch <b>faster</b> tempos.</li></ul>

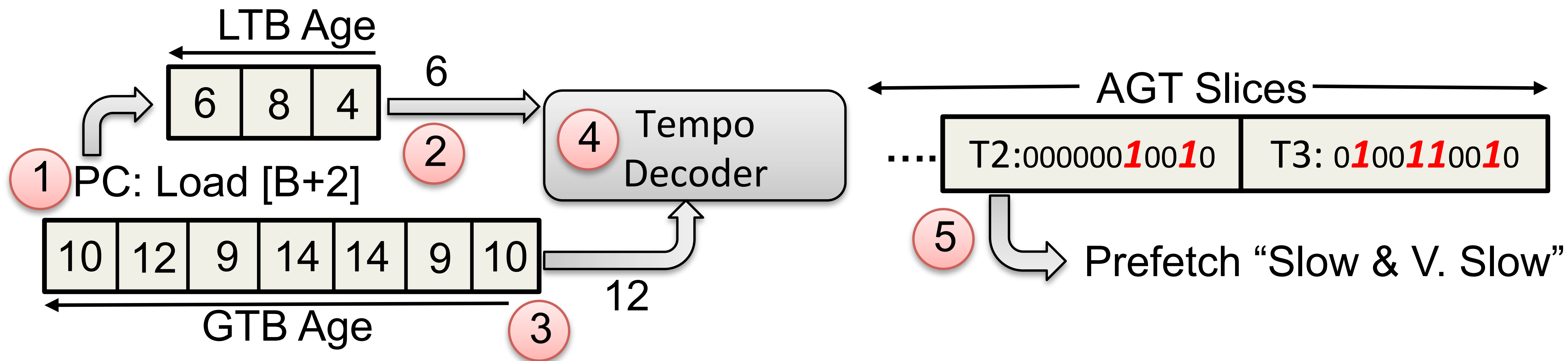


# Walkthrough Example



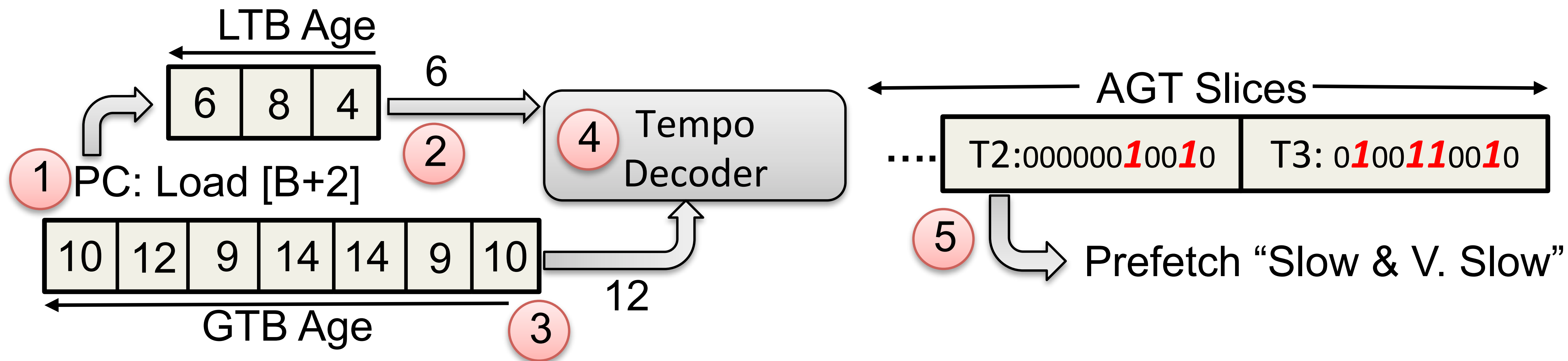
## 1. PC Observes Demand Request to Address $[B+2]$

# Walkthrough Example



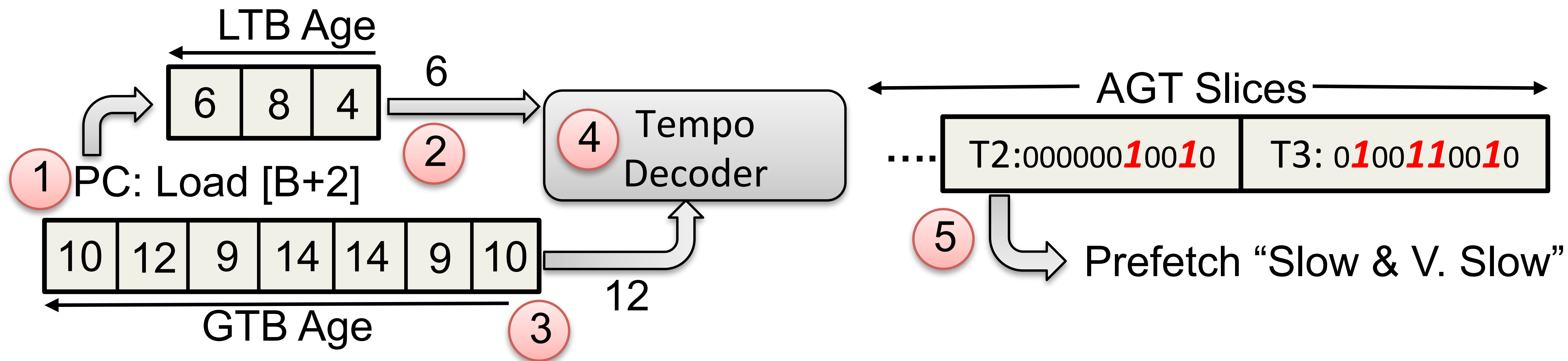
## 2. Lookup Current Local Access Tempo (6)

# Walkthrough Example



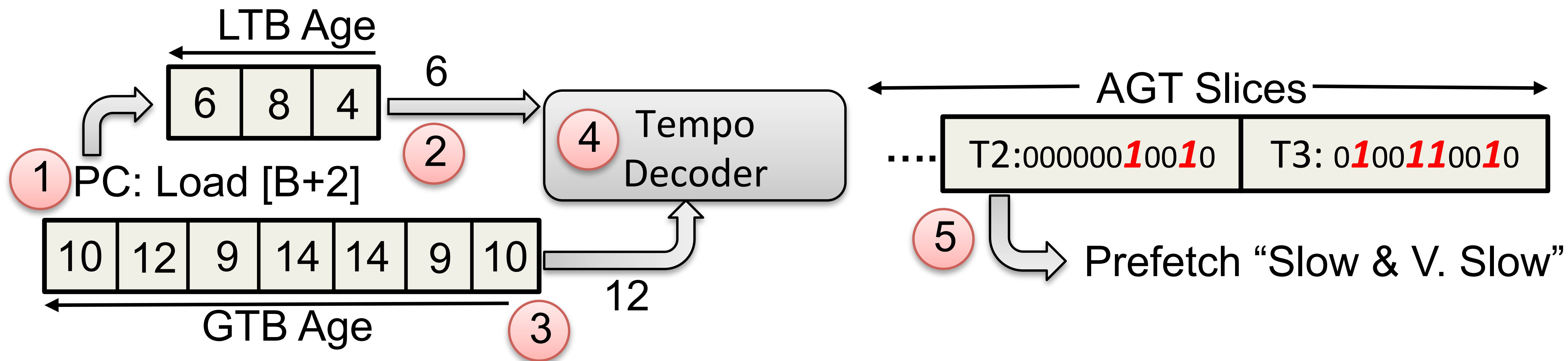
## 3. Average All Elements of Global Tempo Buffer (12)

# Walkthrough Example



**4. Decode "Rushing" or "Dragging".**

# Walkthrough Example



## 5. Prefetch Matching Tempos

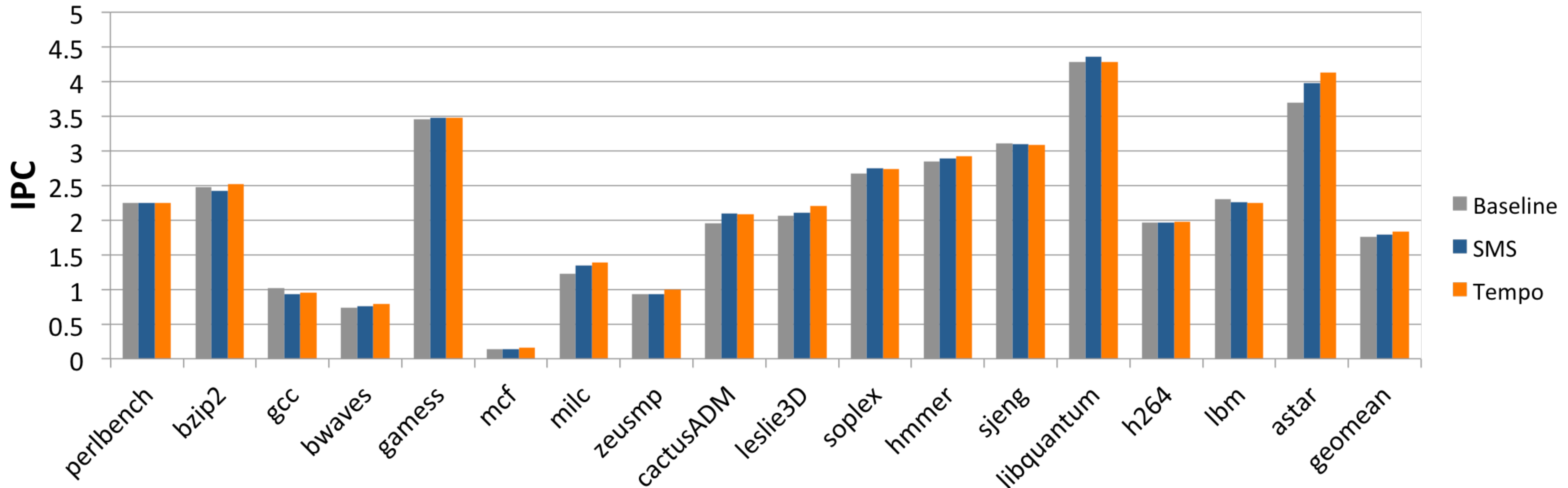


# Evaluation

Core (OoO)	Memory
6-wide, 256 entry instruction window	16kB L1, 128kB L2, Fully inclusive L3
<b>No fetch hazards.</b>	32 entry L2 request queue
Issue 2 loads & 1 store / cycle.	Open Row FR-FCFS Mem. Scheduling
	$t_{l2} = 10$ cycles
	<b>Config 1:</b> 1MB L3, 12.8 GB/s memory
	<b>Config 2:</b> 256kB L3, 12.8 GB/s memory
	<b>Config 3:</b> 1 MB L3, 3.2 GB/s memory

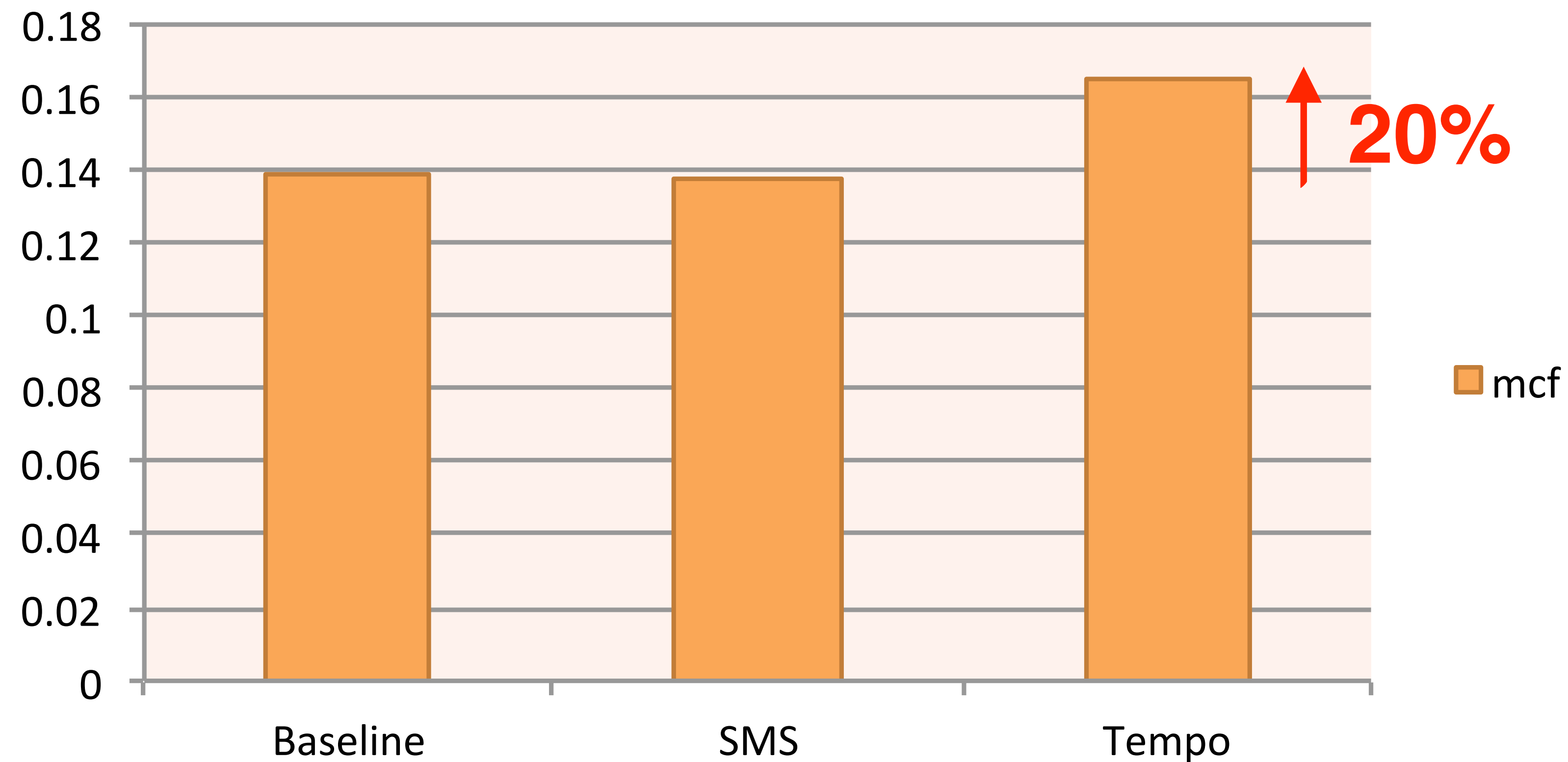
# Absolute IPC Results

## Normal Configuration

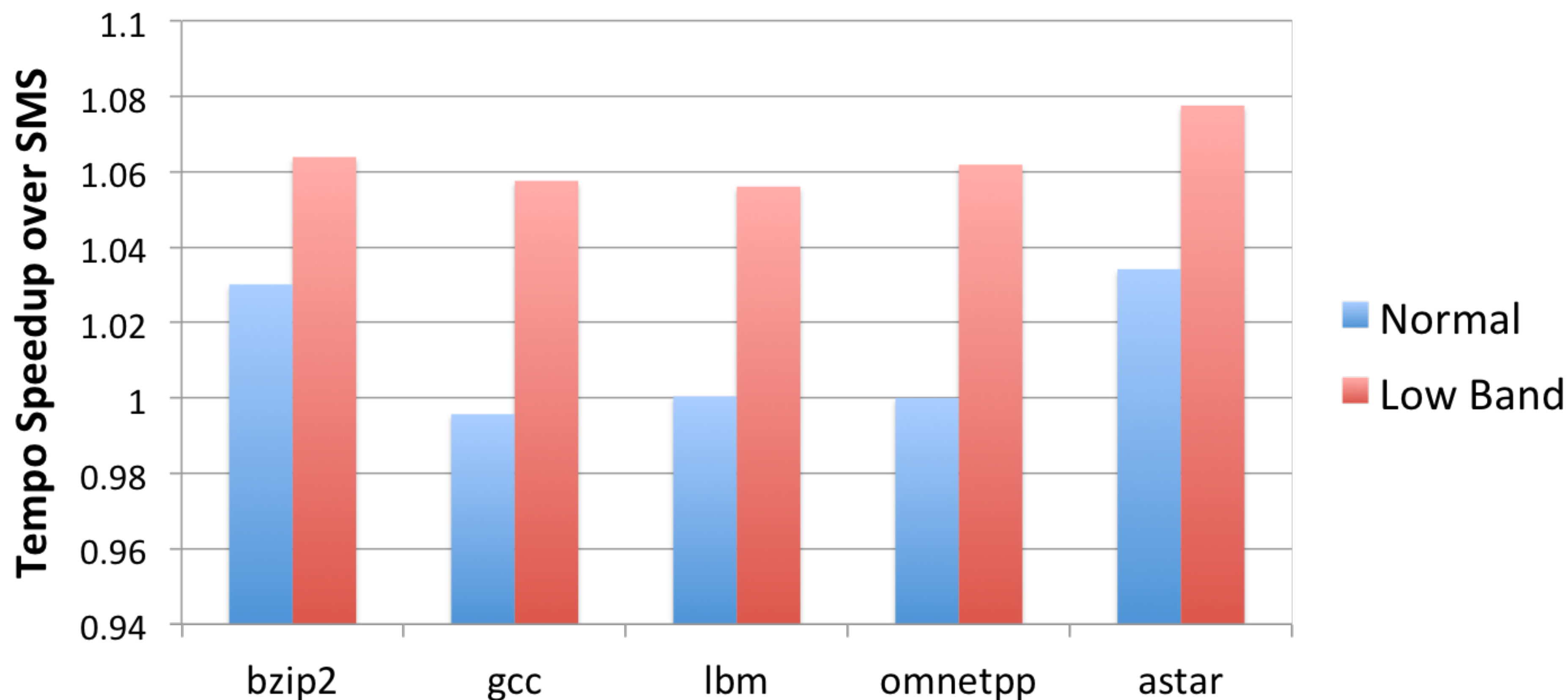


# Absolute IPC Results

## IPC for *mcf*



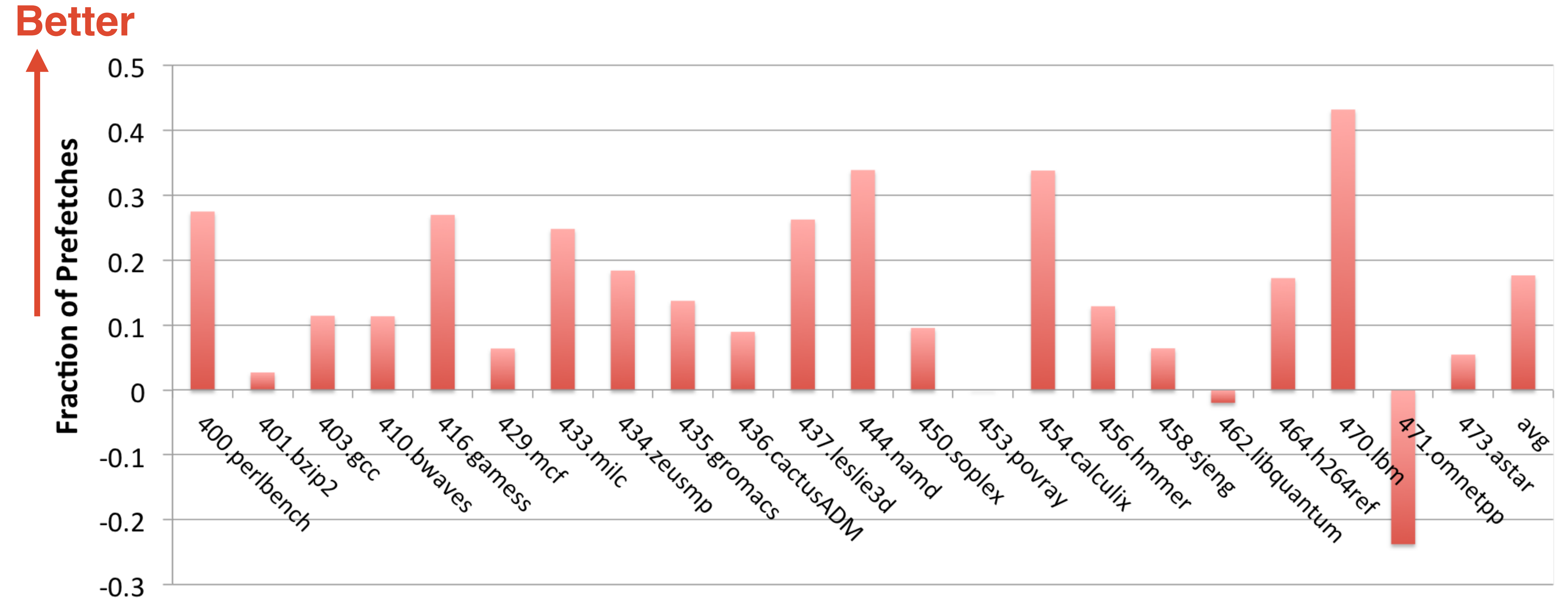
# Increased Low B/W Performance



- Decreasing memory B/W increases Tempo's gains.



# Reduction in Useless Prefetches



# Conclusion

1. Presented a “sliced” version of SMS that filters prefetching decisions based on repeated PC access times.
2. Achieved **1.45%** and **2.57%** IPC improvement on Normal and Low Bandwidth configurations.
3. Reduced the number of useless prefetches by **17.6%**.

?