

ENTANGLING PREFETCHERS

Alberto Ros

University of Murcia, Spain

Oct 22, 2021

PREFETCHING

- Processors need to fetch instructions and data from memory
- High-performance processors would need a very large memory with a low access latency
- This is not possible due to technology limitations
- Computer architects came with a solution to this problem:
prefetching

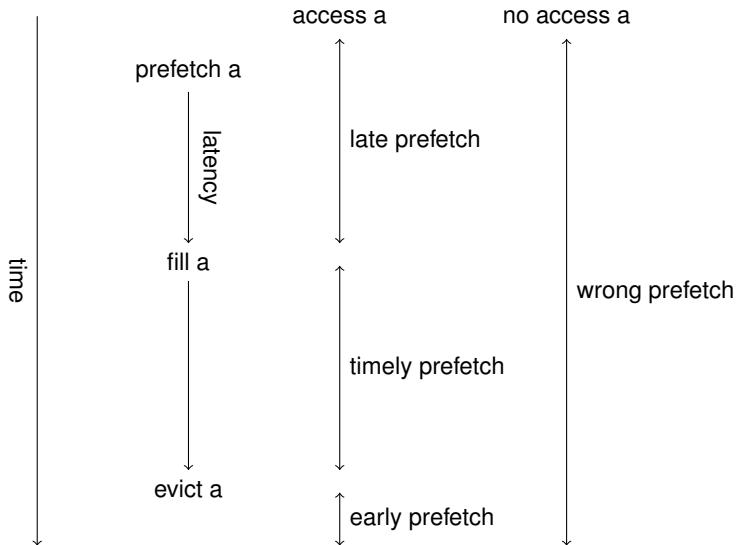
PREFETCHING

- Processors need to fetch instructions and data from memory
- High-performance processors would need a very large memory with a low access latency
- This is not possible due to technology limitations
- Computer architects came with a solution to this problem: **prefetching**

PREFETCHING

Predict **which** memory addresses will be accessed by the processor and fetch them **before** the processor requests them

TIMELINESS



PREFETCHING METRICS

- **COVERAGE**

- Fraction of cache misses covered by the prefetch
- Indicator of performance benefits

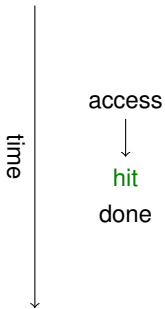
$$Coverage = \frac{timely\ prefetches}{timely\ prefetches + cache\ misses}$$

- **ACCURACY**

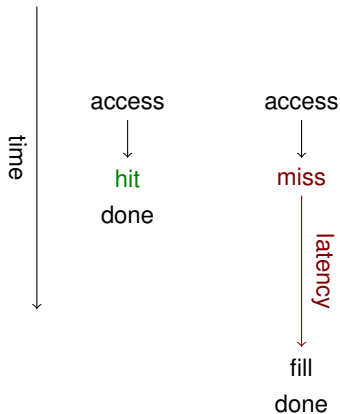
- Fraction of useful prefetches
- Indicator of efficiency

$$Accuracy = \frac{timely\ prefetches + late\ prefetches}{prefetches\ to\ next\ cache\ level}$$

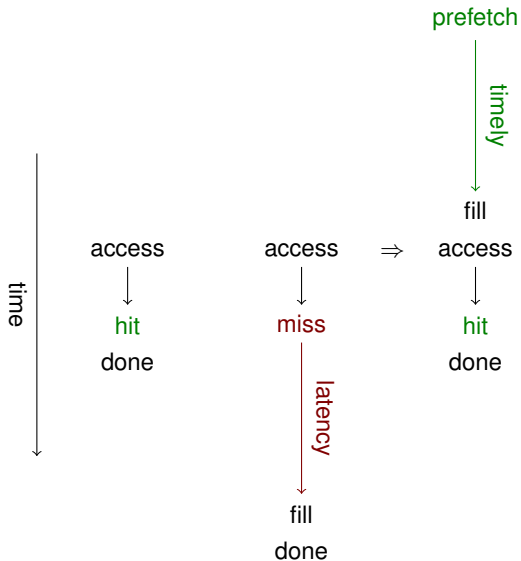
MOTIVATION: TIMELINESS



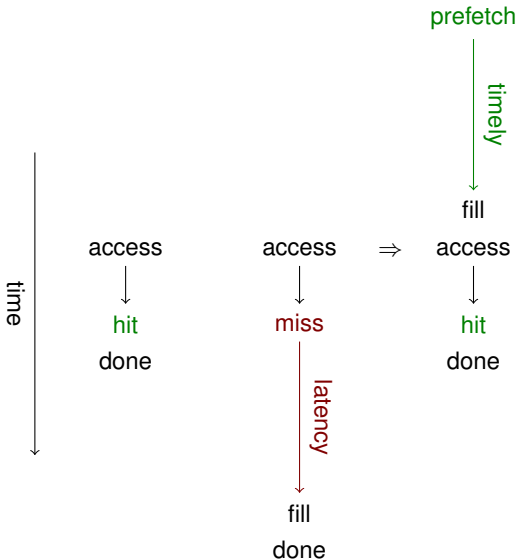
MOTIVATION: TIMELINESS



MOTIVATION: TIMELINESS



MOTIVATION: TIMELINESS



Timely prefetches
for all misses:
Coverage 100%

And only for misses:
Accuracy 100%

OVERVIEW: INSTRUCTION PREFETCHER

- Server and cloud apps getting larger, far from fitting in L1
 ⇒ stalls processor front-end, performance degradation

OVERVIEW: INSTRUCTION PREFETCHER

- Server and cloud apps getting larger, far from fitting in L1
 - ⇒ stalls processor front-end, performance degradation

- Prefetching instructions is fundamental for performance
 - Even when a decoupled front-end is implemented

OVERVIEW: INSTRUCTION PREFETCHER

- Server and cloud apps getting larger, far from fitting in L1
 - ⇒ stalls processor front-end, performance degradation

- Prefetching instructions is fundamental for performance
 - Even when a decoupled front-end is implemented

- Our contribution: An **ENTANGLING** prefetcher
 - **ENTANGLING**: adaptive correlation based on latency
 - **Winner** of the 1st Instruction Prefetching Championship
 - A **cost-effective** prefetcher
 - Prefetcher code is **available**¹

¹ <https://github.com/alberto-ros/EntanglingInstructionPrefetcher>

CONCEPT OF ENTANGLED ACCESSSES



CONCEPT OF ENTANGLED ACCESSSES



CONCEPT OF ENTANGLED ACCESSSES



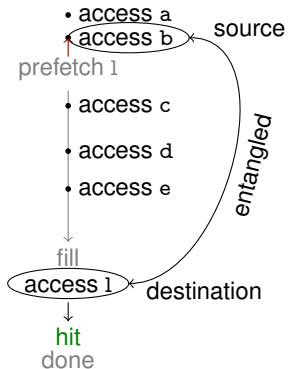
CONCEPT OF ENTANGLED ACCESSSES



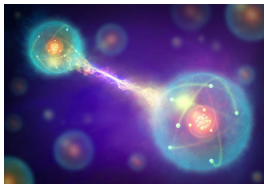
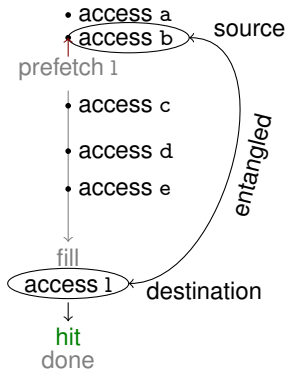
CONCEPT OF ENTANGLED ACCESSSES



CONCEPT OF ENTANGLED ACCESSSES



CONCEPT OF ENTANGLED ACCESSSES

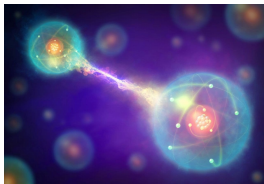
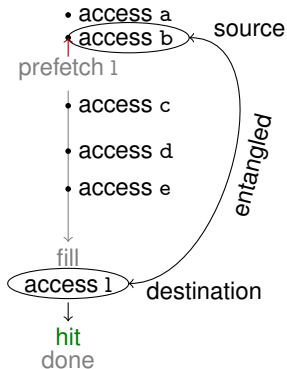


Quantum entanglement

(Image: © MARK GARLICK/SCIENCE

PHOTO LIBRARY/Getty)

CONCEPT OF ENTANGLED ACCESSSES



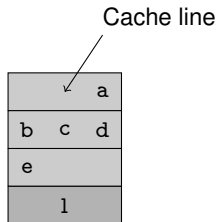
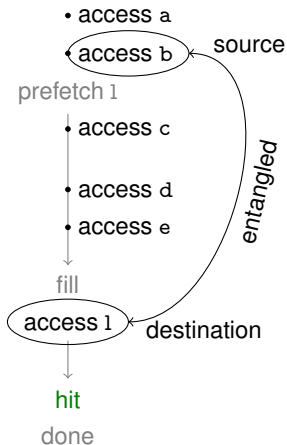
Quantum entanglement

(Image: © MARK GARLICK/SCIENCE

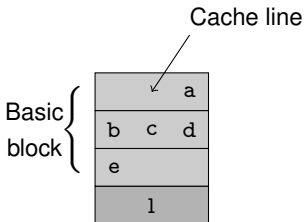
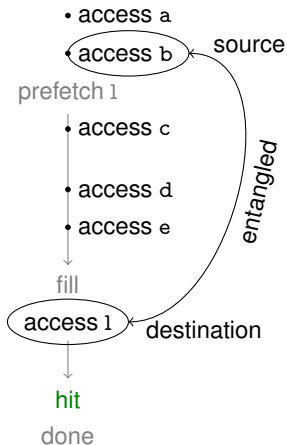
PHOTO LIBRARY/Getty)

THE ENTANGLING PREFETCHER FOR INSTRUCTIONS

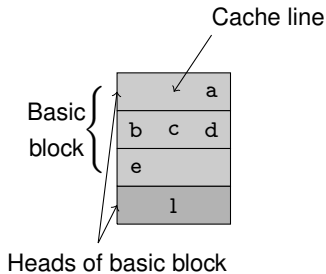
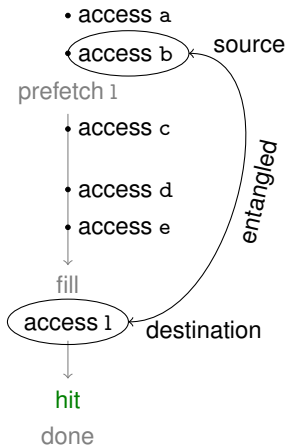
ENTANGLING CACHE LINES HEAD OF BASIC BLOCKS



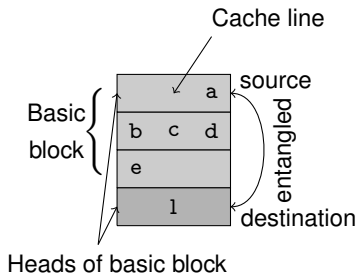
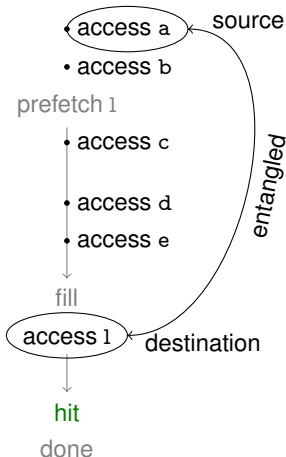
ENTANGLING CACHE LINES HEAD OF BASIC BLOCKS



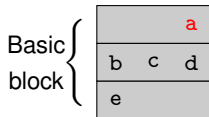
ENTANGLING CACHE LINES HEAD OF BASIC BLOCKS



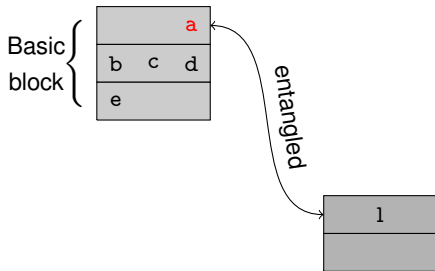
ENTANGLING CACHE LINES HEAD OF BASIC BLOCKS



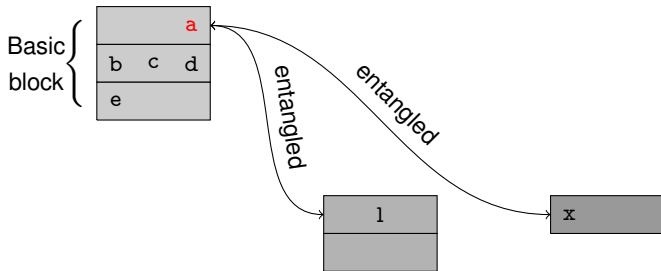
WHAT TO PREFETCH ON AN ACCESS TO a?



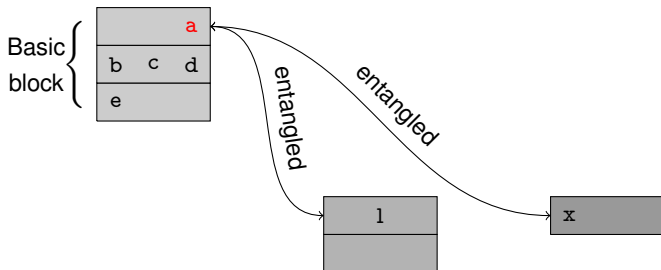
WHAT TO PREFETCH ON AN ACCESS TO a?



WHAT TO PREFETCH ON AN ACCESS TO a?

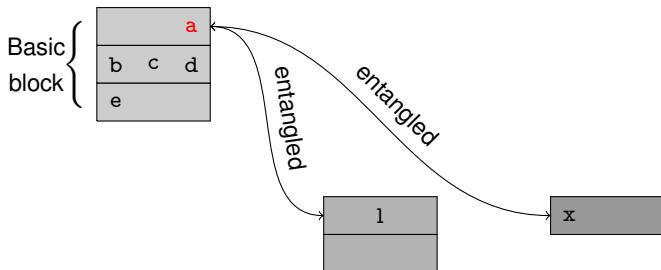


WHAT TO PREFETCH ON AN ACCESS TO a?



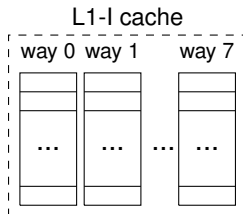
- Too much? (Max entangled = 6, Max BB size = 64)

WHAT TO PREFETCH ON AN ACCESS TO a?



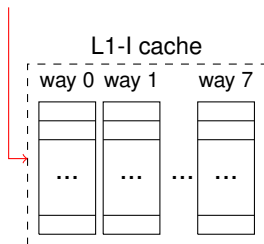
- Too much? (Max entangled = 6, Max BB size = 64)
 - Most of the time no prefetches are issued (no head of basic block)
 - Average number of prefetches per access to **head of basic block** ranging from ≈ 9 to ≈ 17

DESIGN OF THE ENTANGLING PREFETCHER



DESIGN OF THE ENTANGLING PREFETCHER

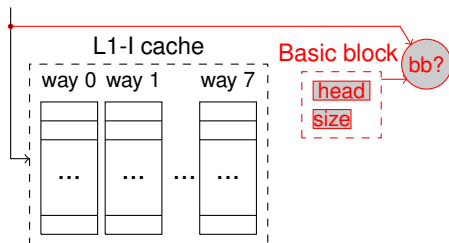
L1-I access



DESIGN OF THE ENTANGLING PREFETCHER - FINDING BASIC BLOCKS

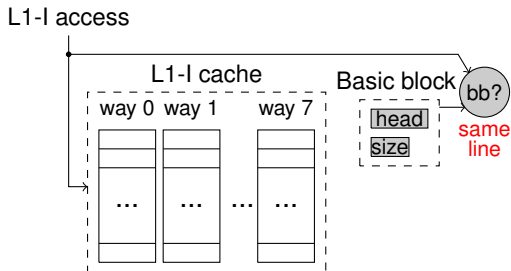
0	a
1	b c d
2	e
0	1

L1-I access

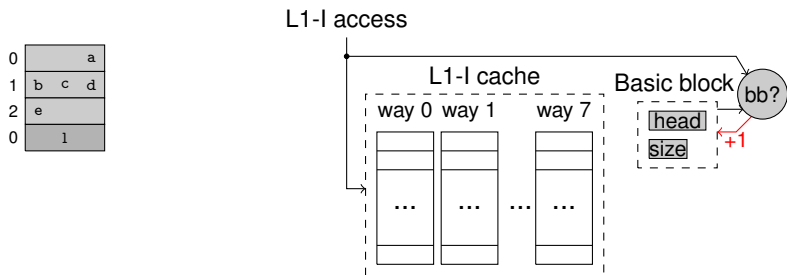


DESIGN OF THE ENTANGLING PREFETCHER - FINDING BASIC BLOCKS

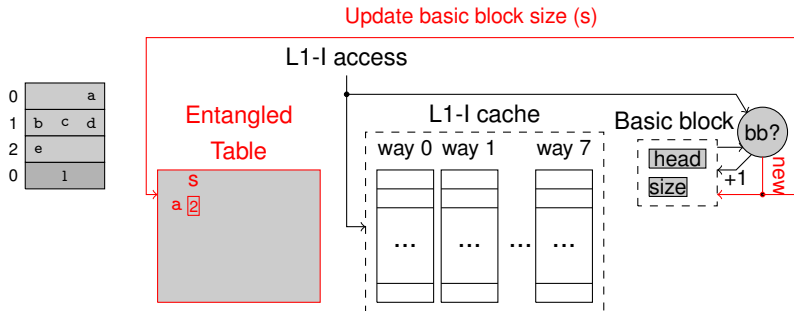
0	a
1	b c d
2	e
0	1



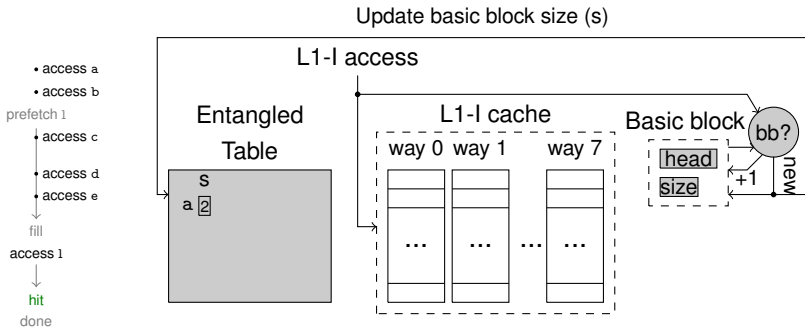
DESIGN OF THE ENTANGLING PREFETCHER - FINDING BASIC BLOCKS



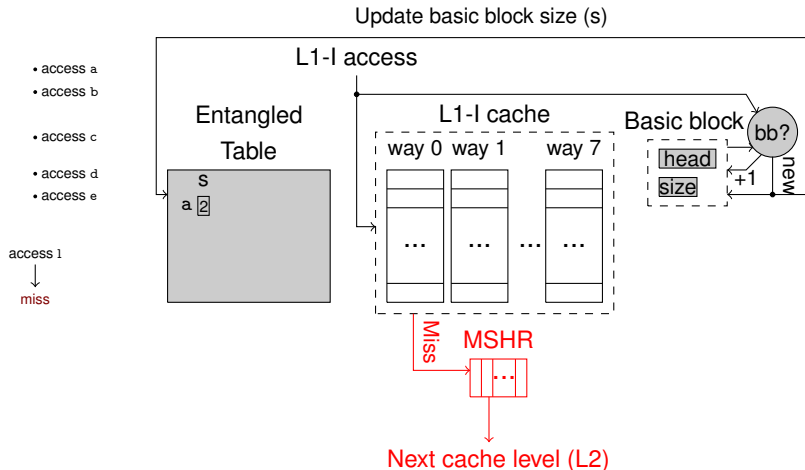
DESIGN OF THE ENTANGLING PREFETCHER - FINDING BASIC BLOCKS



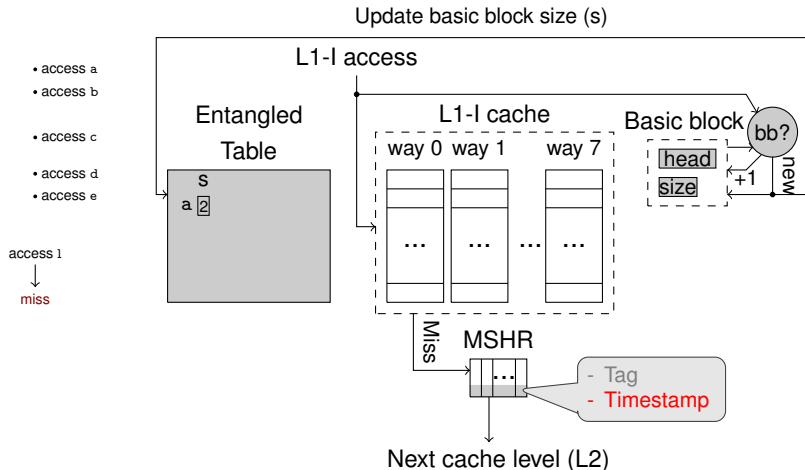
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



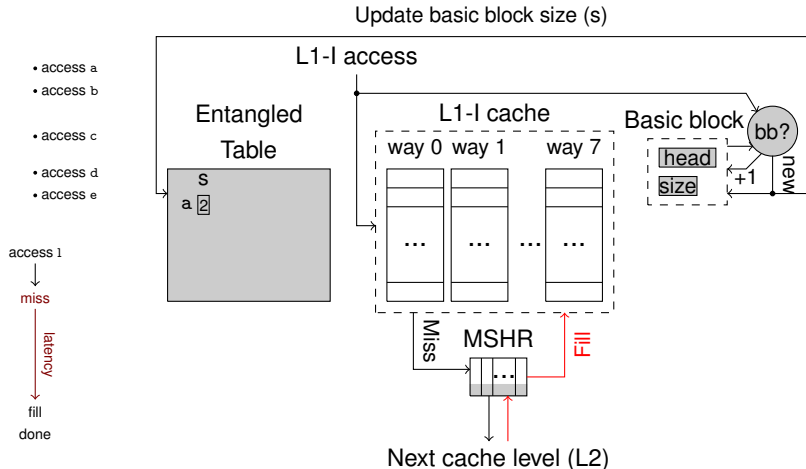
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



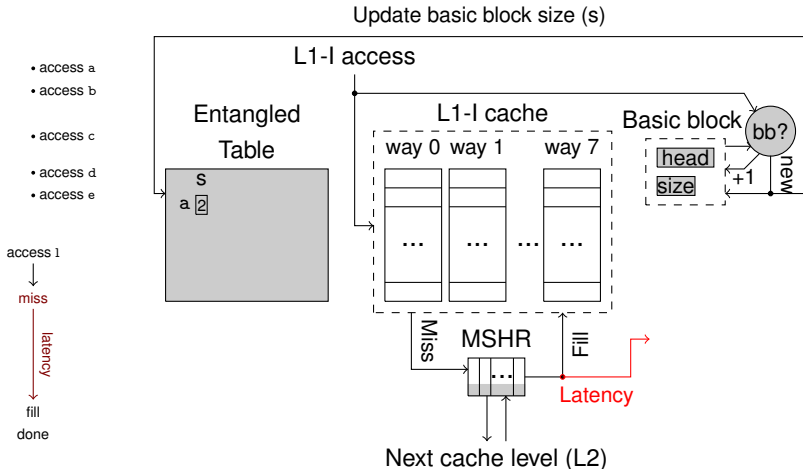
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



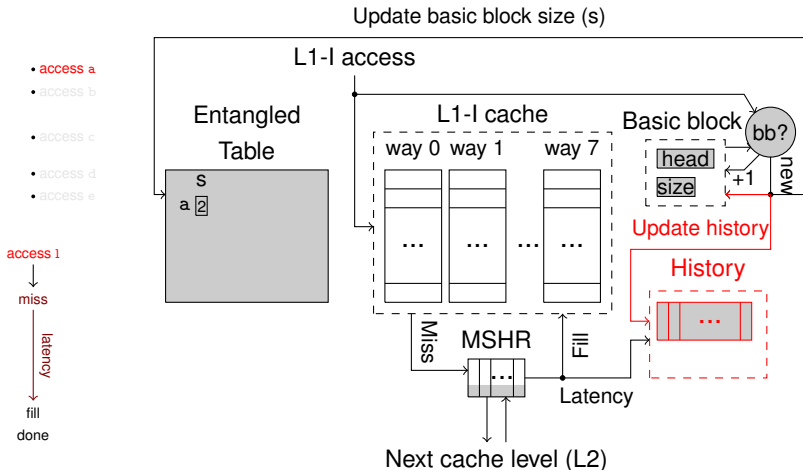
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



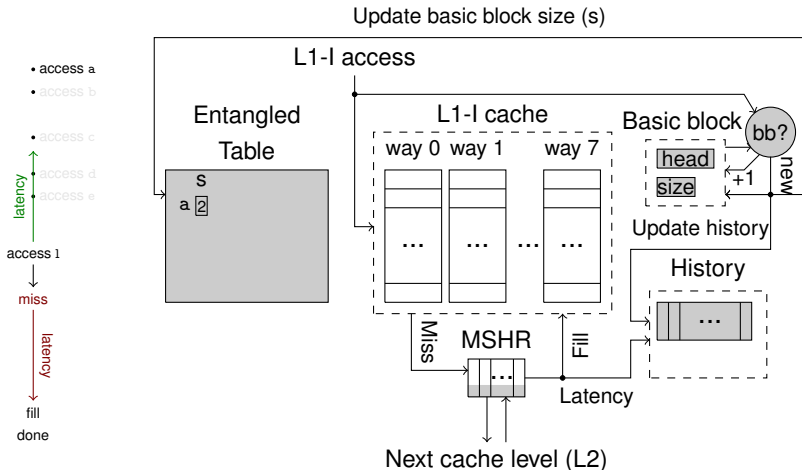
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



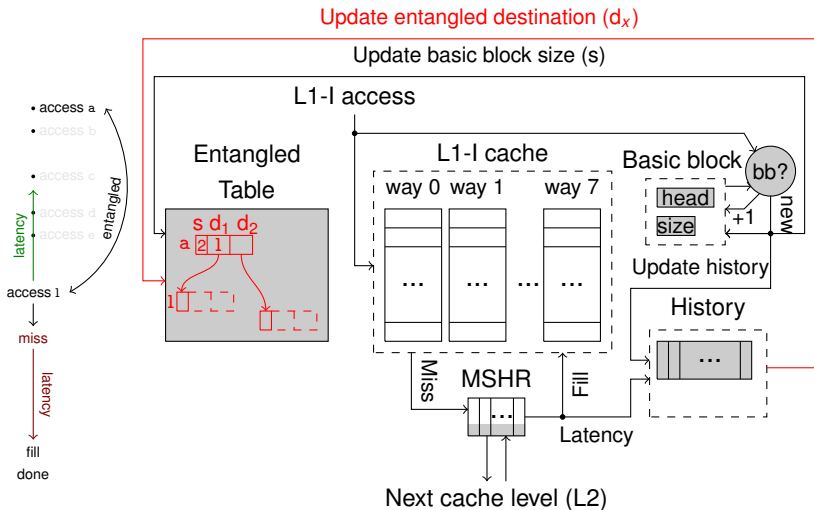
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



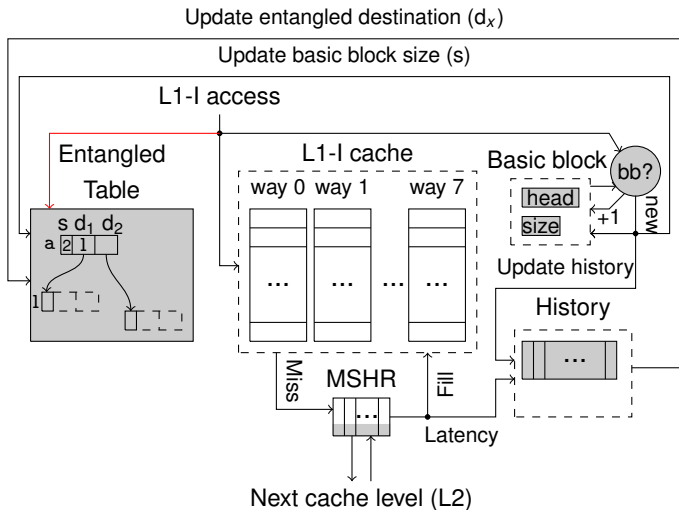
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



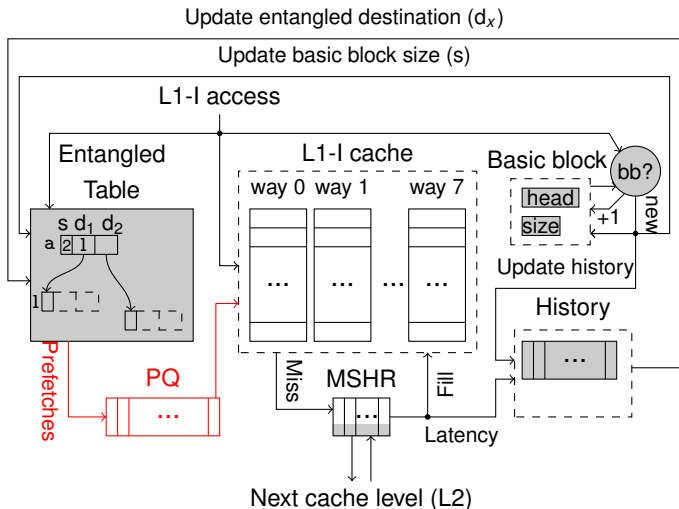
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



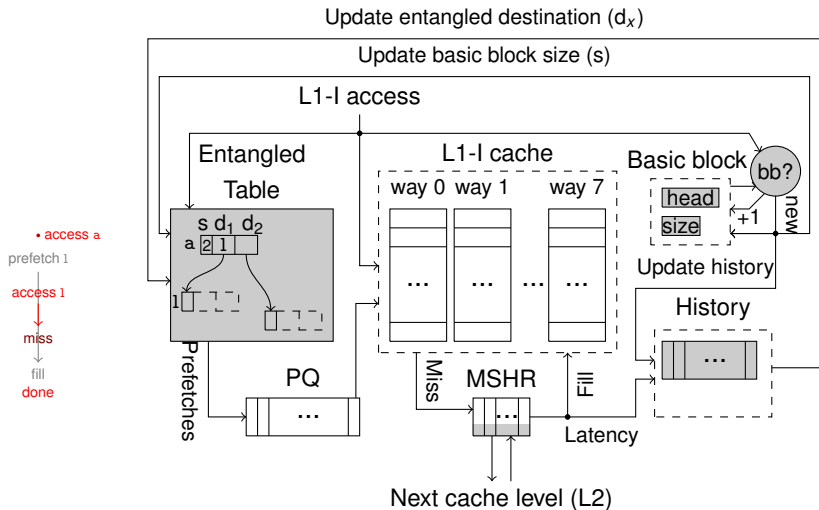
DESIGN OF THE ENTANGLING PREFETCHER - ISSUING PREFETCHES



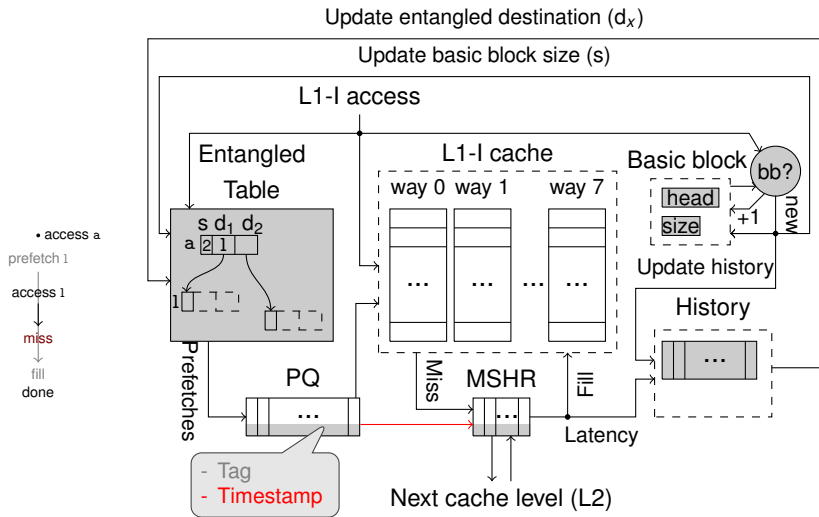
DESIGN OF THE ENTANGLING PREFETCHER - ISSUING PREFETCHES



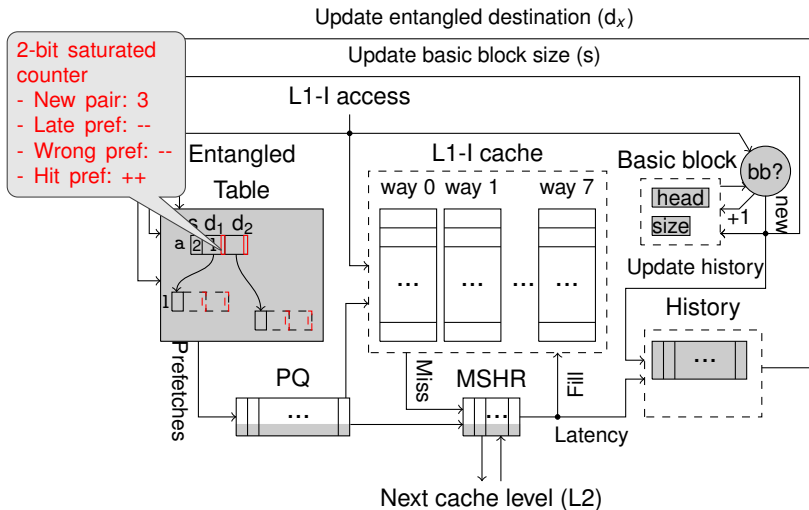
DESIGN OF THE ENTANGLING PREFETCHER - FIXING LATE PREFETCHES



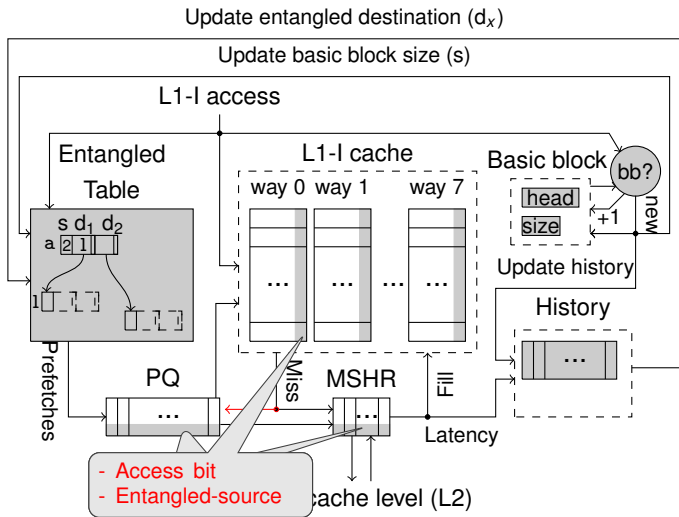
DESIGN OF THE ENTANGLING PREFETCHER - FIXING LATE PREFETCHES



DESIGN OF THE ENTANGLING PREFETCHER - CONFIDENCE FOR ENTANGLED PAIRS

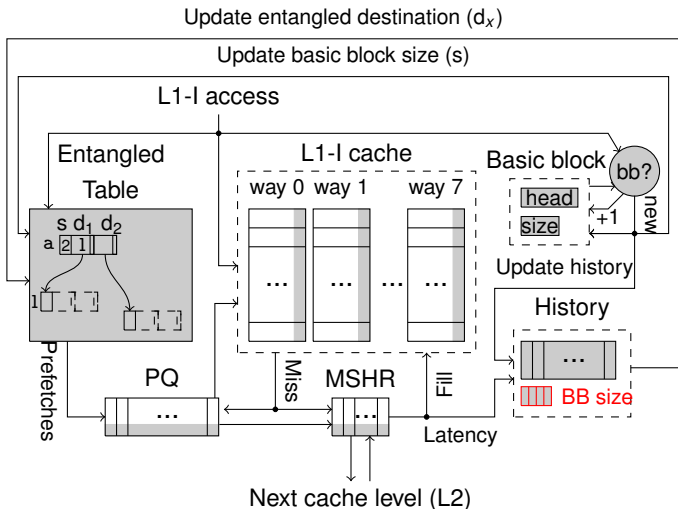


DESIGN OF THE ENTANGLING PREFETCHER - CONFIDENCE FOR ENTANGLED PAIRS

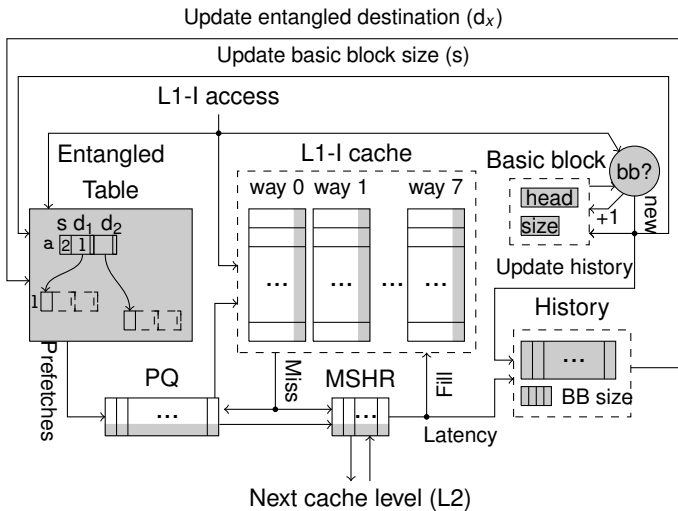


DESIGN OF THE ENTANGLING PREFETCHER - MERGING

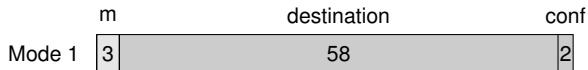
BASIC BLOCKS



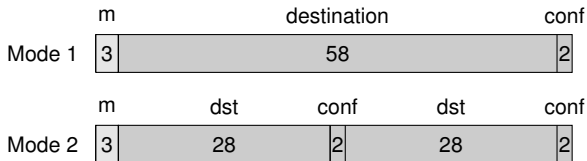
DESIGN OF THE ENTANGLING PREFETCHER



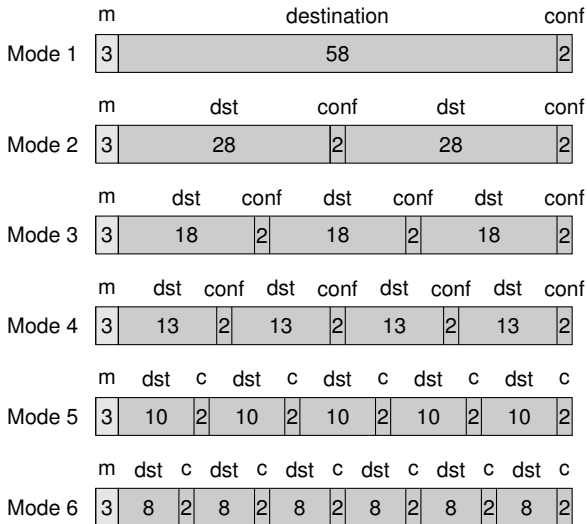
COMPRESSING DESTINATIONS



COMPRESSING DESTINATIONS



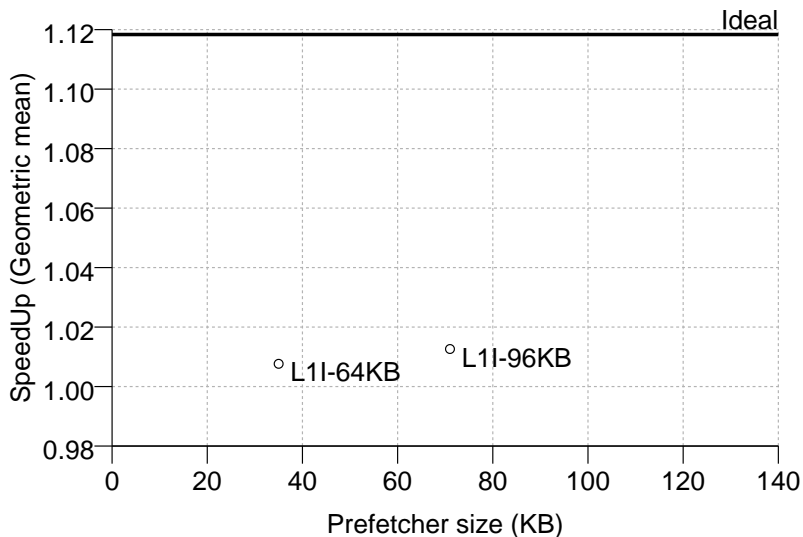
COMPRESSING DESTINATIONS



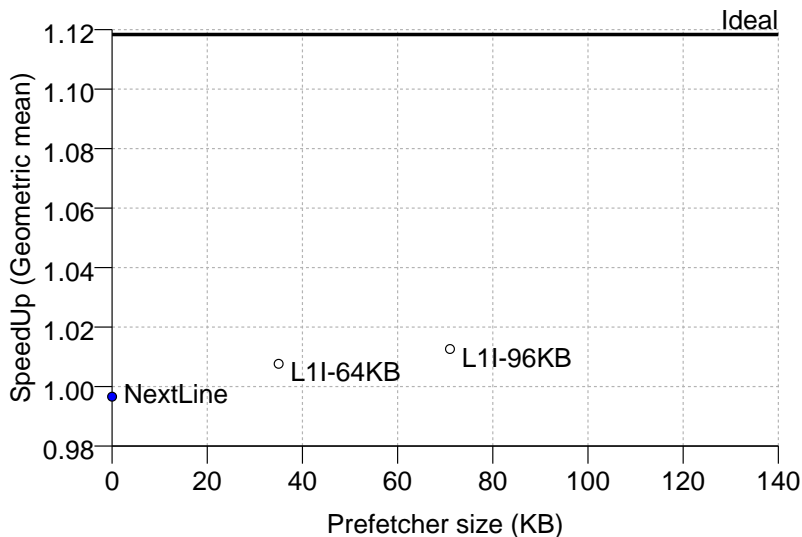
METHODOLOGY

- **ChampSim** develop branch (nov 2020)
- **Baseline:**
 - Sunny Cove-like system
 - Decoupled front-end (64-entry fetch queue)
 - 32KB L1I
- **ENTANGLED:**
 - *History buffer:* 16 entries
 - *Entangled table:* 2K, 4K and 8K entries
- **Applications**
 - 959 traces from the Championship Value Prediction (provided by Qualcomm)
 - Cloud Suite
- **Analysis** both for virtual and physical prefetching

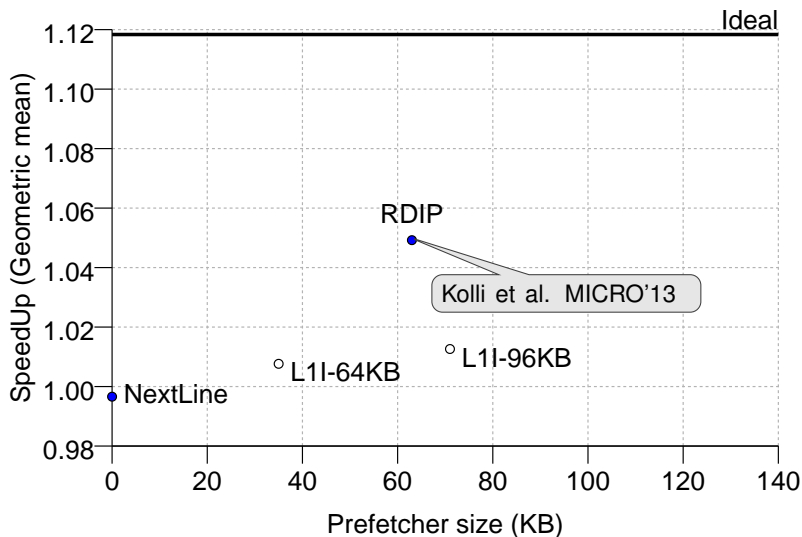
RESULTS: IPC VS MEMORY OVERHEAD



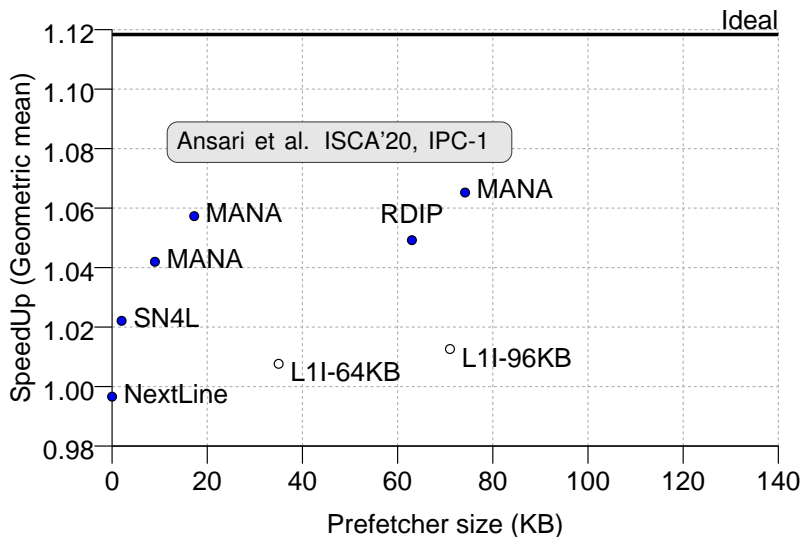
RESULTS: IPC VS MEMORY OVERHEAD



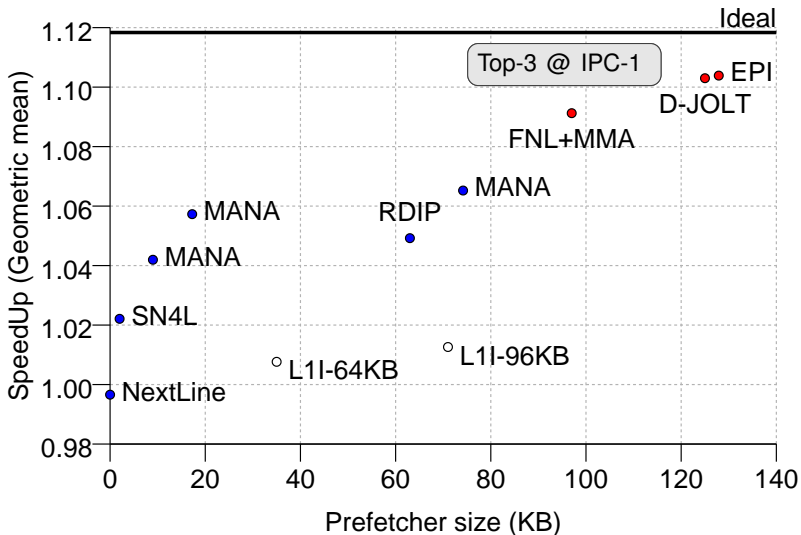
RESULTS: IPC VS MEMORY OVERHEAD



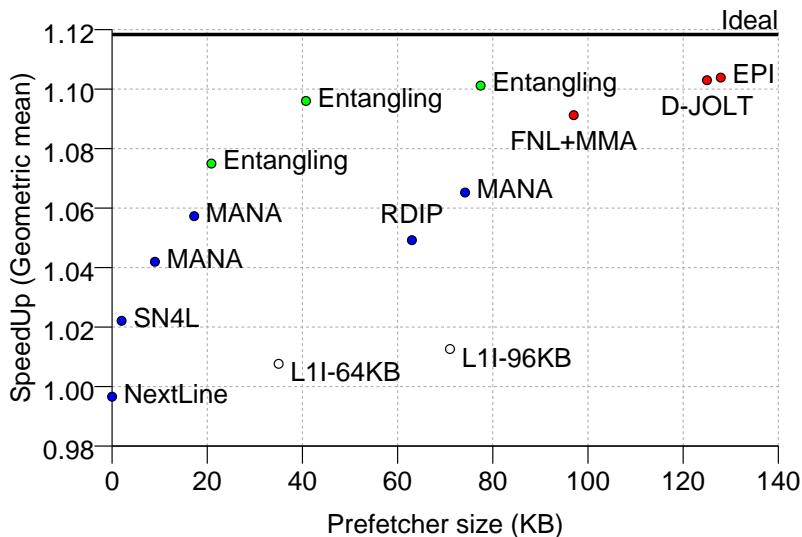
RESULTS: IPC VS MEMORY OVERHEAD



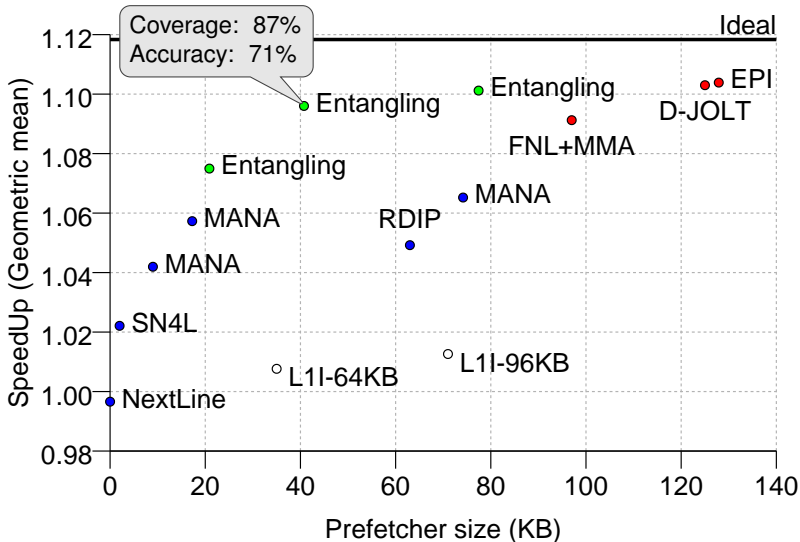
RESULTS: IPC VS MEMORY OVERHEAD



RESULTS: IPC VS MEMORY OVERHEAD



RESULTS: IPC VS MEMORY OVERHEAD



CONCLUDING REMARKS

- **Timeliness** as a key property
- **Entangles** heads of basic blocks to trigger timely prefetches
- Near **ideal** performance with just 40KB

ENTANGLING PREFETCHERS

Alberto Ros

aros@ditec.um.es

Thank you!



ECHO, ERC Consolidator Grant (No 819134)