

BOOSTING DATA CENTERS PERFORMANCE WITH THE ENTANGLING INSTRUCTION PREFETCHER

Alberto Ros

University of Murcia, Spain

Dec 2, 2021

DATA CENTERS

- Data centers serve most devices
 - Internet of things, smartphones, self-driving cars, ...
 - Energy costs expected to reach **8% of the global consumption** by 2030¹

¹ Andrae et al. *On Global Electricity Usage of Communication Technology: Trends to 2030*, Callenges 2015.

DATA CENTERS

- Data centers serve most devices
 - Internet of things, smartphones, self-driving cars, ...
 - Energy costs expected to reach **8% of the global consumption** by 2030¹
- They run increasingly complex applications
 - Deep software stacks

¹ Andrae et al. *On Global Electricity Usage of Communication Technology: Trends to 2030*, Callenges 2015.

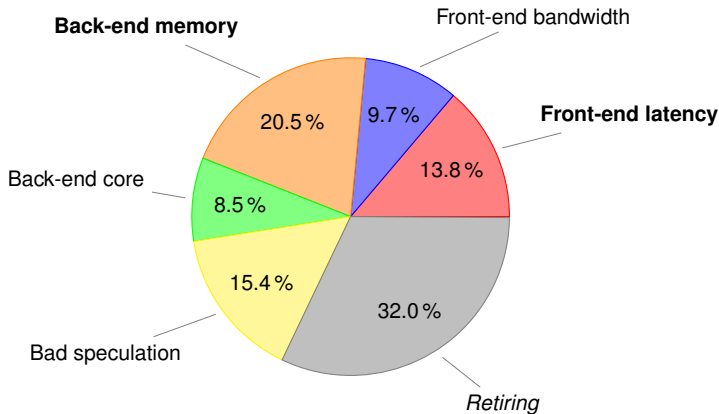
DATA CENTERS

- Data centers serve most devices
 - Internet of things, smartphones, self-driving cars, ...
 - Energy costs expected to reach **8% of the global consumption** by 2030¹
- They run increasingly complex applications
 - Deep software stacks
- **Instruction footprint** constantly growing
 - Far from fitting in small instruction caches (L1I)
 - And **growing by 20% per year!**²

¹ Andrae et al. *On Global Electricity Usage of Communication Technology: Trends to 2030*, Challenges 2015.

² Kanev et al. *Profiling a warehouse-scale computer*, ISCA 2015.

DATA CENTERS BOTTLENECKS³



³ Ayers et al. *AsmDB: Understanding and Mitigating Front-End Stalls in Warehouse-Scale Computers*, ISCA 2019.

DATA CENTERS BOTTLENECKS

FRONT-END LATENCY (13.8%)

- Dominated by **instruction cache (L1I) misses**
 - Hitting in the second level cache (L2) or last level cache (LLC)
- Latency more important than bandwidth
- Critical as processors need to keep the pipeline full

BACK-END MEMORY (20.5%)

- Due to **data cache (L1D) misses**
 - Many of them reaching main memory
- Cause significant stalls and late detection of **BAD SPECULATION (15.4%)**

PREFETCHING TO THE RESCUE

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

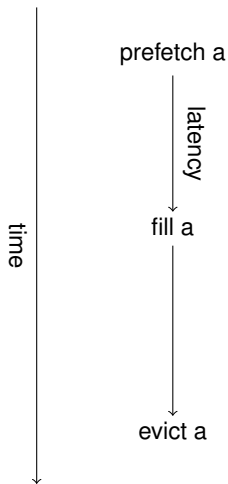
PREFETCHING TO THE RESCUE

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

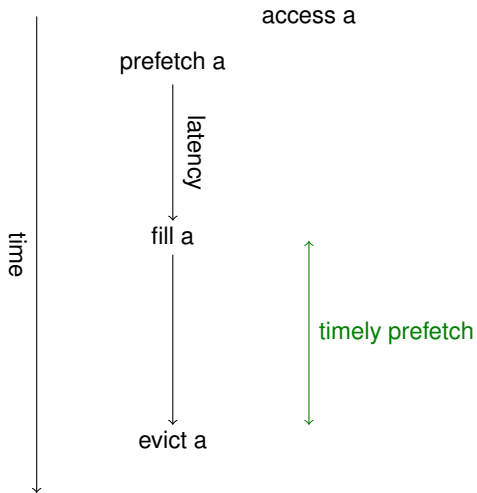
PREFETCHING

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

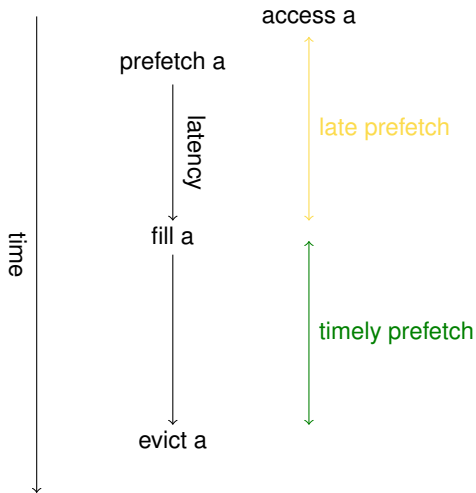
TIMELINESS IS THE KEY PROPERTY



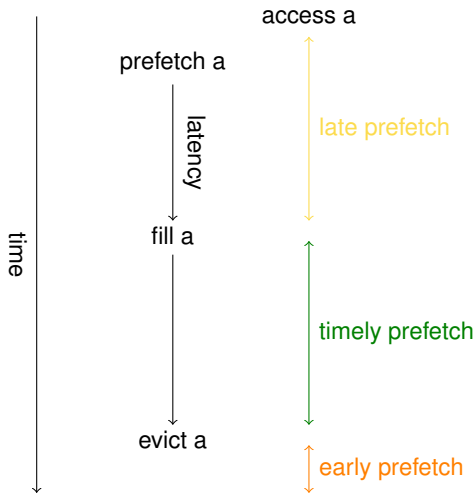
TIMELINESS IS THE KEY PROPERTY



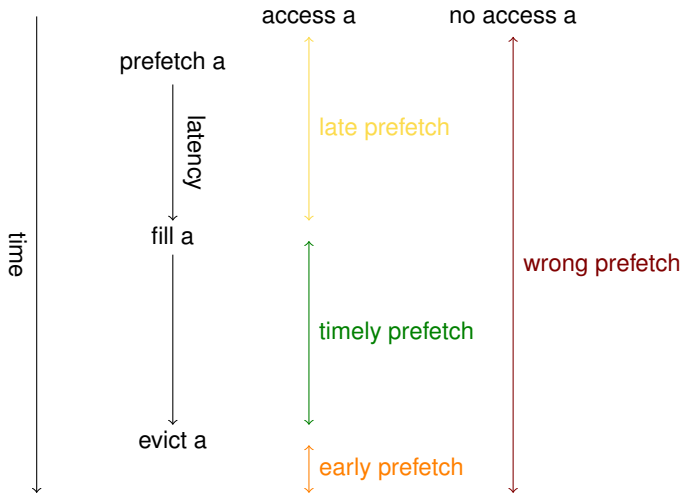
TIMELINESS IS THE KEY PROPERTY



TIMELINESS IS THE KEY PROPERTY



TIMELINESS IS THE KEY PROPERTY



TIMELINESS IS THE KEY PROPERTY

- Two prefetchers with a strong focus on timeliness

TIMELINESS IS THE KEY PROPERTY

- Two prefetchers with a strong focus on timeliness
 - ① The BLUE Data Prefetcher⁴
 - An LLC prefetcher
 - **Winner** of the 1st ML-based Data Prefetching Competition
 - Organized by Google
 - With a non-ML solution!

⁴ Ros, *BLUE: A Timely, IP-based Data Prefetcher*, ML-DPC-1 2021

TIMELINESS IS THE KEY PROPERTY

- Two prefetchers with a strong focus on timeliness
 - 1 The **BLUE** Data Prefetcher⁴
 - An LLC prefetcher
 - **Winner** of the 1st ML-based Data Prefetching Competition
 - Organized by Google
 - With a non-ML solution!
 - 2 The **ENTANGLING** Instruction Prefetcher⁵
 - An L1I prefetcher
 - **Winner** of the 1st Instruction Prefetching Championship
 - Organized by Intel
 - Follow up paper published at ISCA'21

⁴ Ros, *BLUE: A Timely, IP-based Data Prefetcher*, ML-DPC-1 2021

⁵ Ros and Jimborean, *The Entangling Instruction Prefetcher*, IPC-1 2020

THE ENTANGLING INSTRUCTION PREFETCHER

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

THE ENTANGLING 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

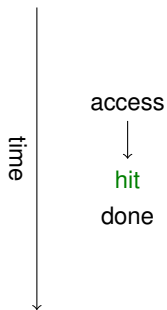
THE ENTANGLING INSTRUCTION PREFETCHER

- 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
- Solution: The **ENTANGLING** instruction prefetcher⁶
 - **ENTANGLING**: adaptive correlation based on latency
 - A **cost-effective** prefetcher
 - Prefetcher code is **available**⁷

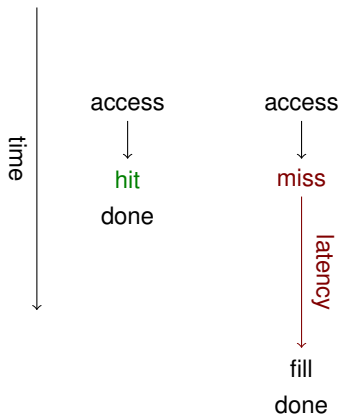
⁶ Ros and Jimborean, *A Cost-Effective Entangling Prefetcher for Instructions*, ISCA 2021

⁷ <https://github.com/alberto-ros/EntanglingInstructionPrefetcher>

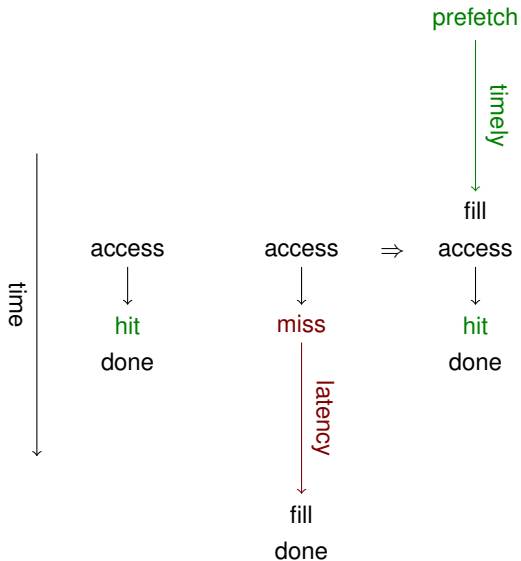
MOTIVATION: TIMELINESS



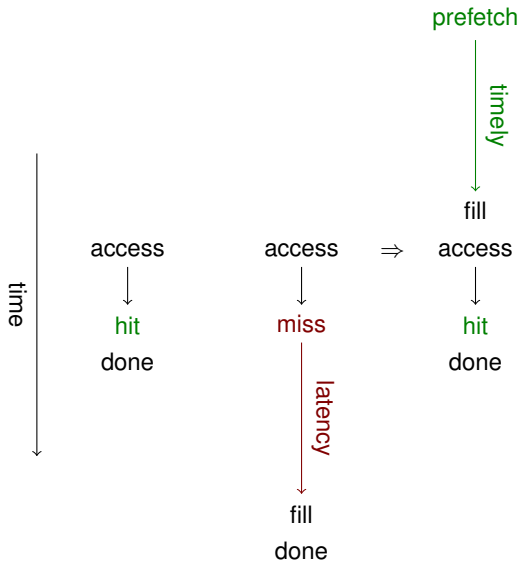
MOTIVATION: TIMELINESS



MOTIVATION: TIMELINESS



MOTIVATION: TIMELINESS



Timely prefetches
for all misses:
Coverage 100%

And only for misses:
Accuracy 100%

CONCEPT OF ENTANGLED ACCESSSES



CONCEPT OF ENTANGLED ACCESSSES

prefetch 1



access 1



miss

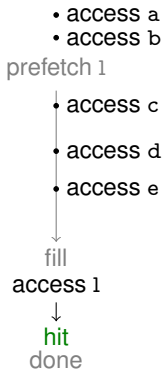


fill
done

CONCEPT OF ENTANGLED ACCESSSES



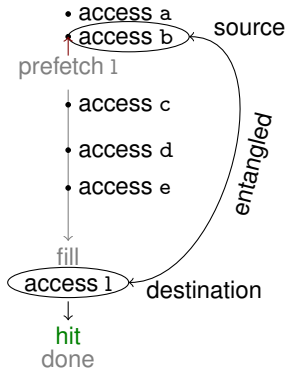
CONCEPT OF ENTANGLED ACCESSSES



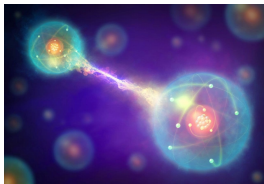
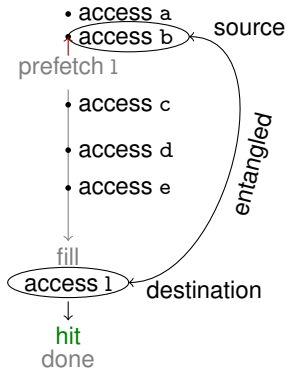
CONCEPT OF ENTANGLED ACCESSSES



CONCEPT OF ENTANGLED ACCESSSES



CONCEPT OF ENTANGLED ACCESSSES

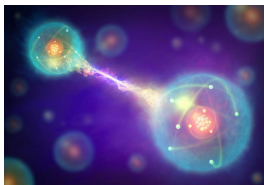
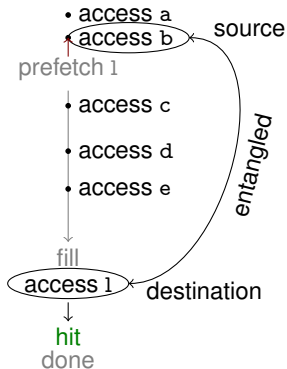


Quantum entanglement

(Image: © MARK GARLICK/SCIENCE

PHOTO LIBRARY/Getty)

CONCEPT OF ENTANGLED ACCESSSES

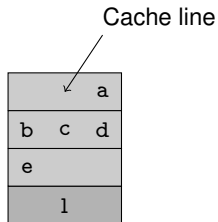
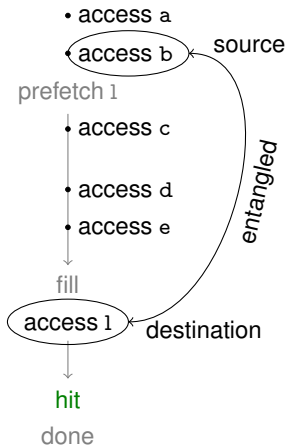


Quantum entanglement

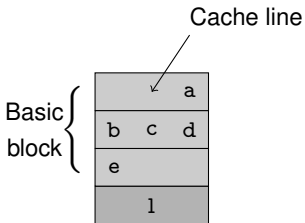
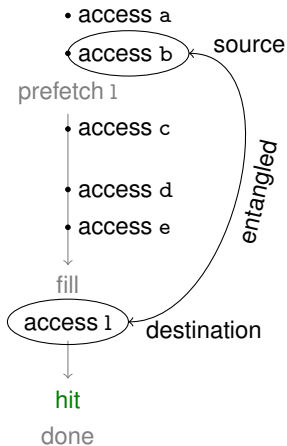
(Image: © MARK GARLICK/SCIENCE
PHOTO LIBRARY/Getty)

THE ENTANGLING PREFETCHER FOR INSTRUCTIONS

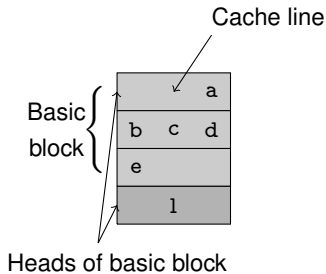
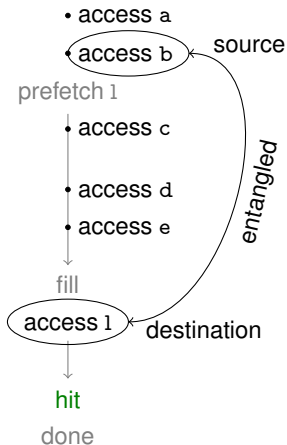
ENTANGLING CACHE LINES HEAD OF BASIC BLOCKS



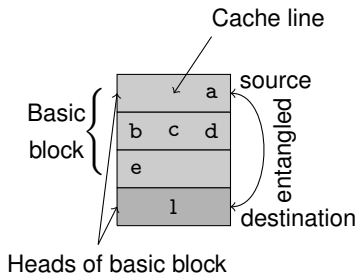
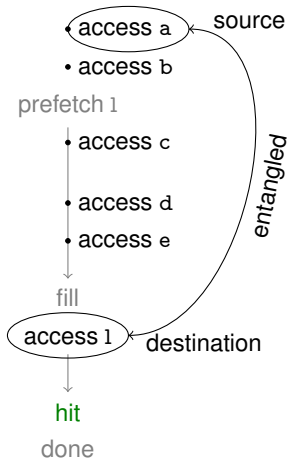
ENTANGLING CACHE LINES HEAD OF BASIC BLOCKS



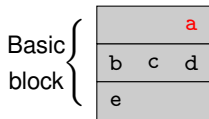
ENTANGLING CACHE LINES HEAD OF BASIC BLOCKS



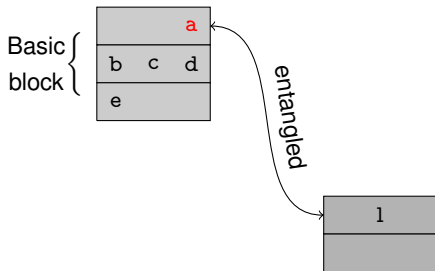
ENTANGLING CACHE LINES HEAD OF BASIC BLOCKS



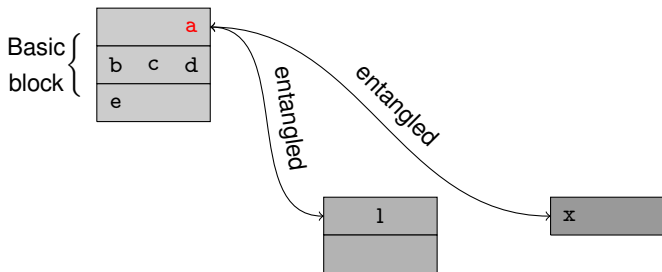
WHAT TO PREFETCH ON AN ACCESS TO a?



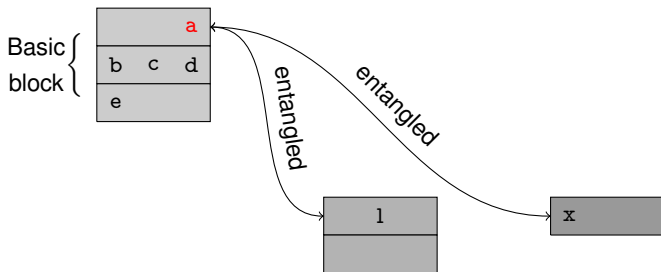
WHAT TO PREFETCH ON AN ACCESS TO a?



WHAT TO PREFETCH ON AN ACCESS TO a?

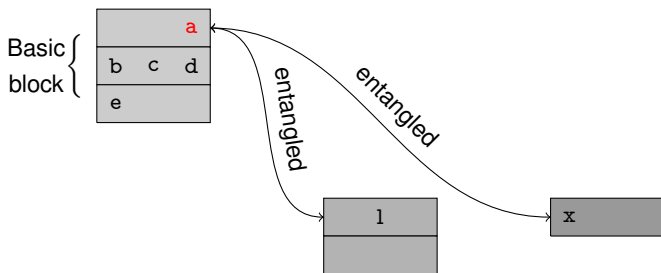


WHAT TO PREFETCH ON AN ACCESS TO a?



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

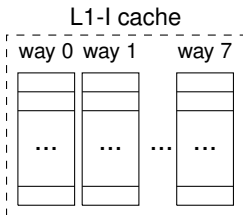
WHAT TO PREFETCH ON AN ACCESS TO a?



- 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
 - Remember: Front-end latency more important than bandwidth⁸

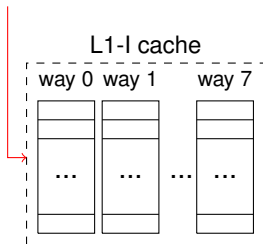
⁸ Kanev et al. *Profiling a warehouse-scale computer*, ISCA 2015.

DESIGN OF THE ENTANGLING PREFETCHER



DESIGN OF THE ENTANGLING PREFETCHER

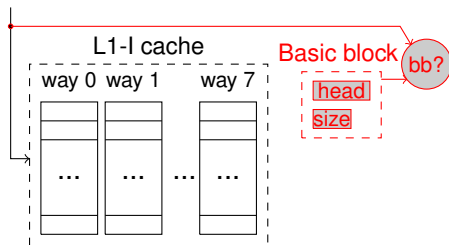
L1-I access



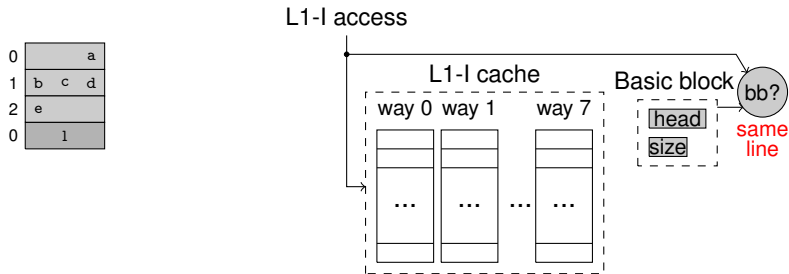
DESIGN OF THE ENTANGLING PREFETCHER - FINDING BASIC BLOCKS

0	a
1	b c d
2	e
0	1

L1-I access



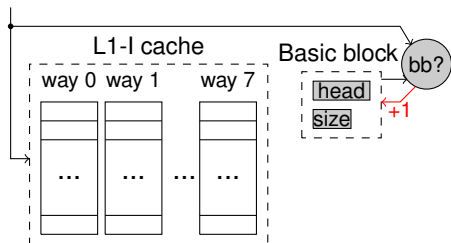
DESIGN OF THE ENTANGLING PREFETCHER - FINDING BASIC BLOCKS



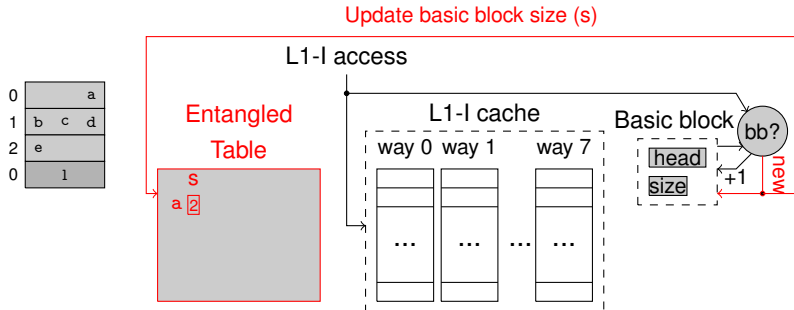
DESIGN OF THE ENTANGLING PREFETCHER - FINDING BASIC BLOCKS

0	a
1	b c d
2	e
0	1

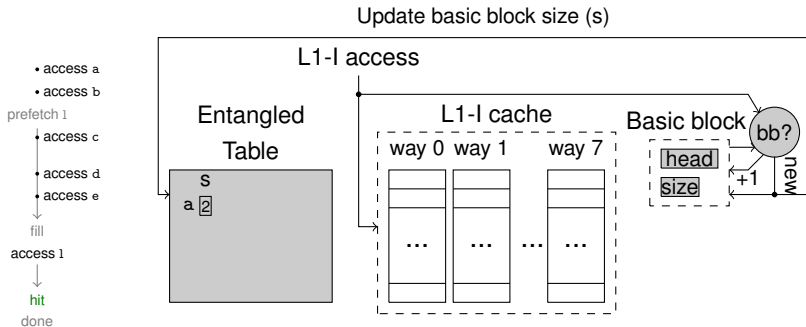
L1-I access



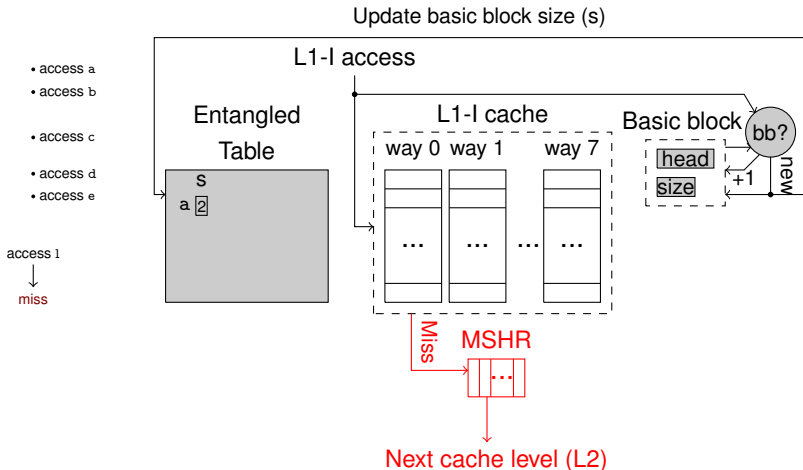
DESIGN OF THE ENTANGLING PREFETCHER - FINDING BASIC BLOCKS



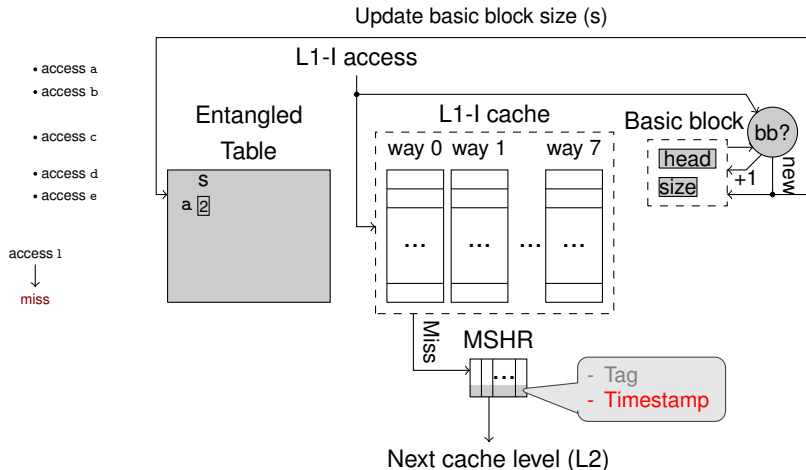
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



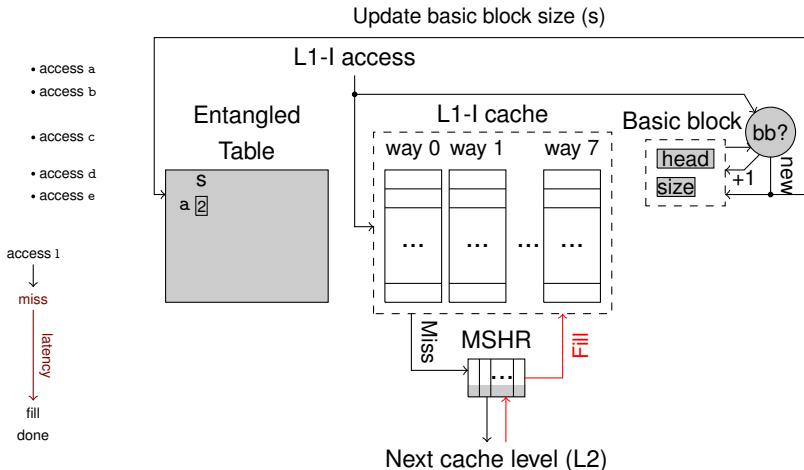
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



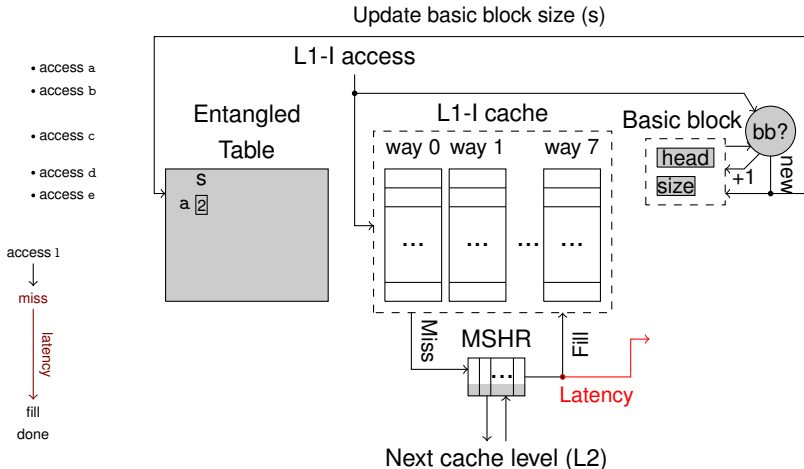
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



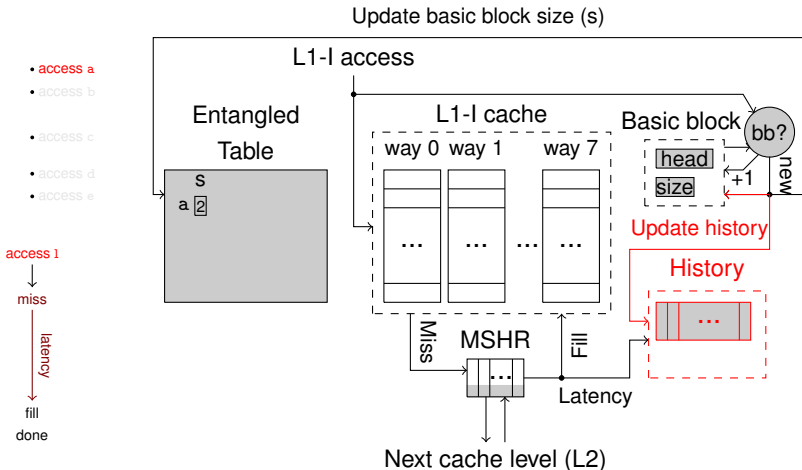
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



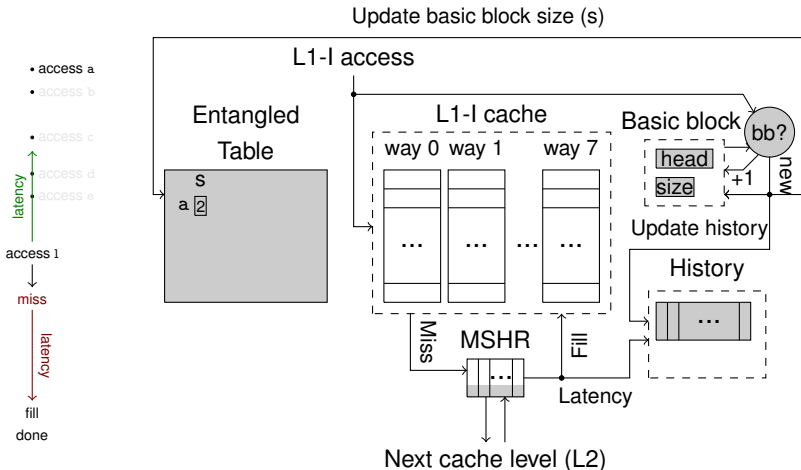
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



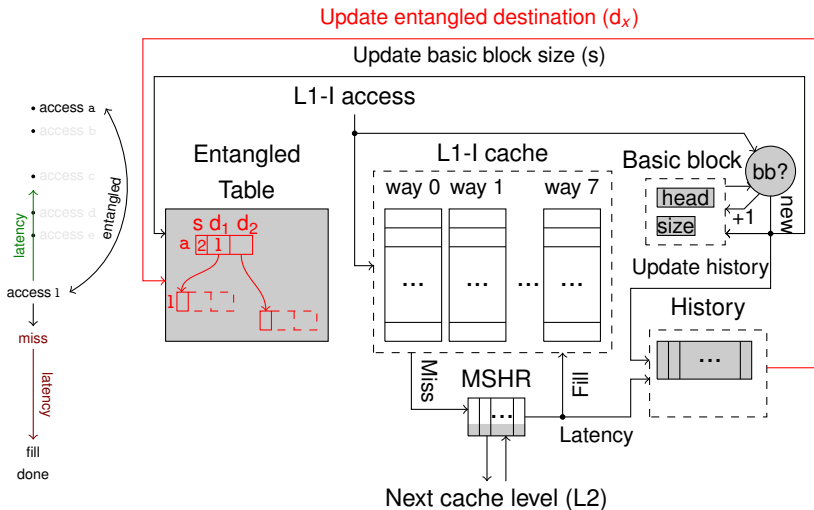
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



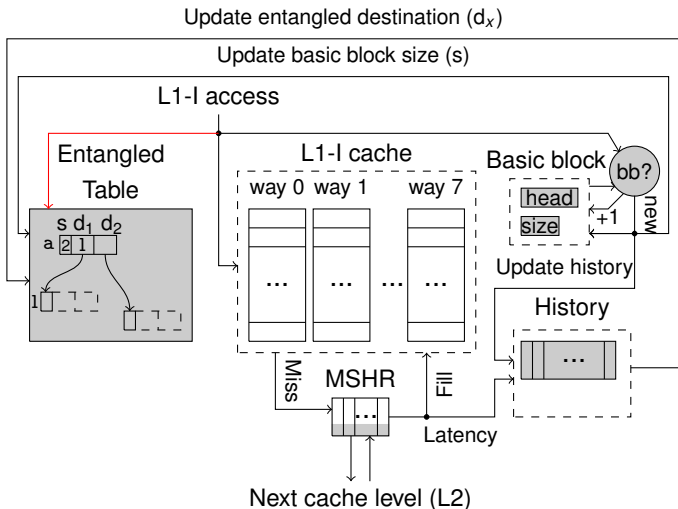
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



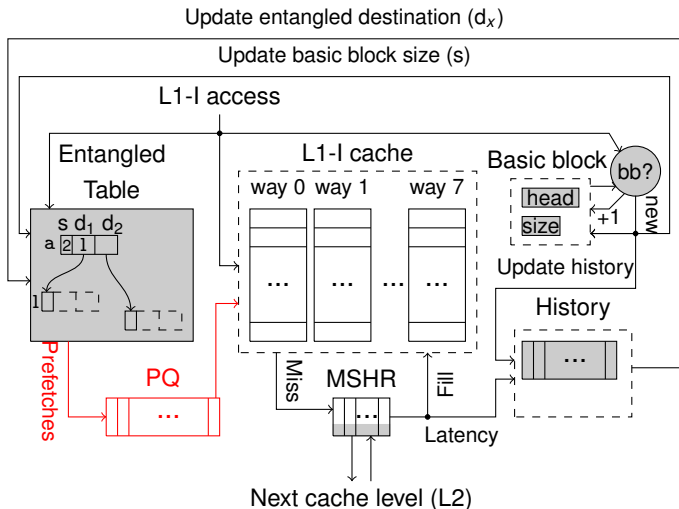
DESIGN OF THE ENTANGLING PREFETCHER - ENTANGLING CACHE LINES



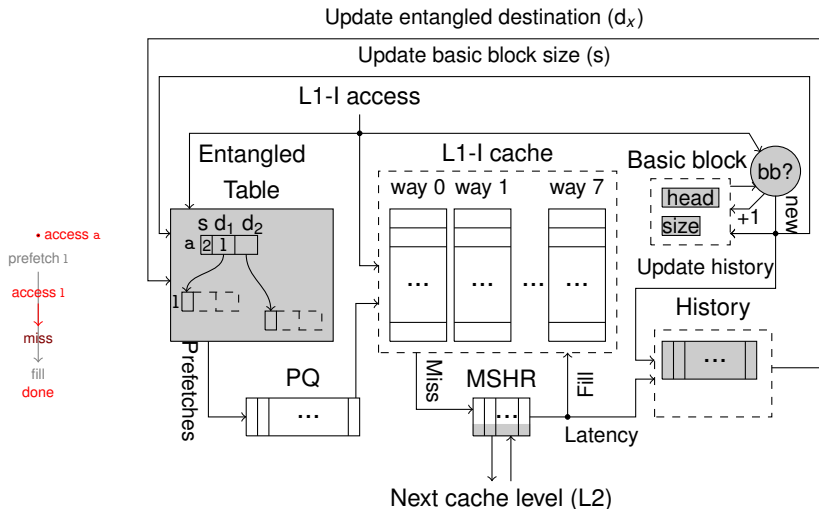
DESIGN OF THE ENTANGLING PREFETCHER - ISSUING PREFETCHES



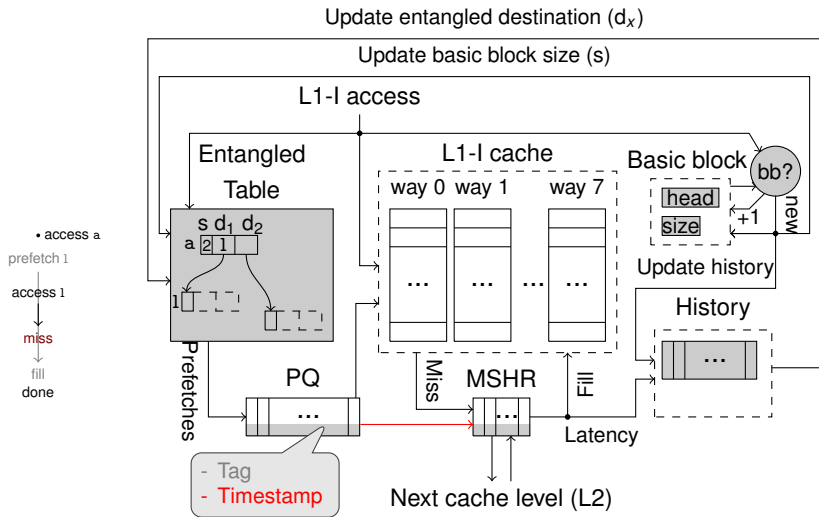
DESIGN OF THE ENTANGLING PREFETCHER - ISSUING PREFETCHES



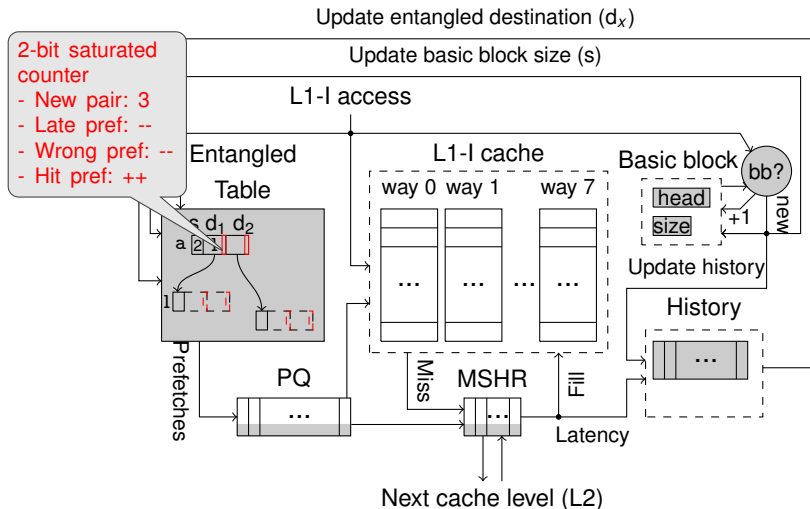
DESIGN OF THE ENTANGLING PREFETCHER - FIXING LATE PREFETCHES



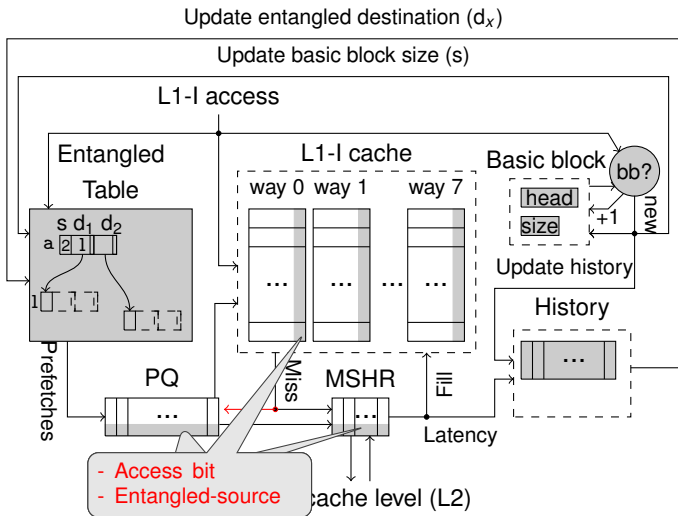
DESIGN OF THE ENTANGLING PREFETCHER - FIXING LATE PREFETCHES



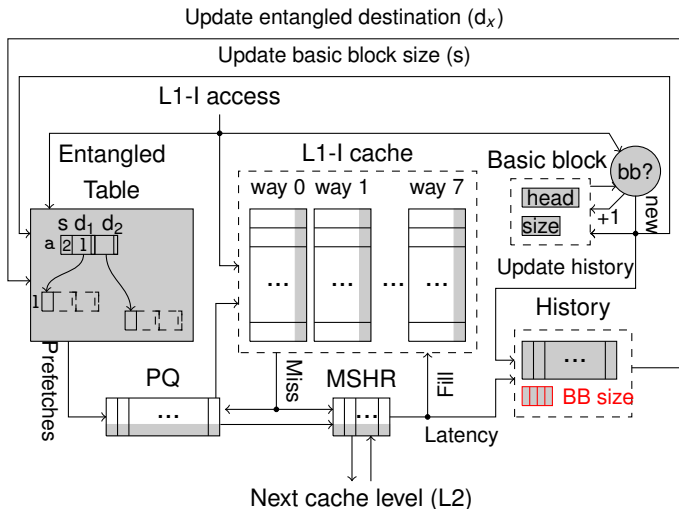
DESIGN OF THE ENTANGLING PREFETCHER - CONFIDENCE FOR ENTANGLED PAIRS



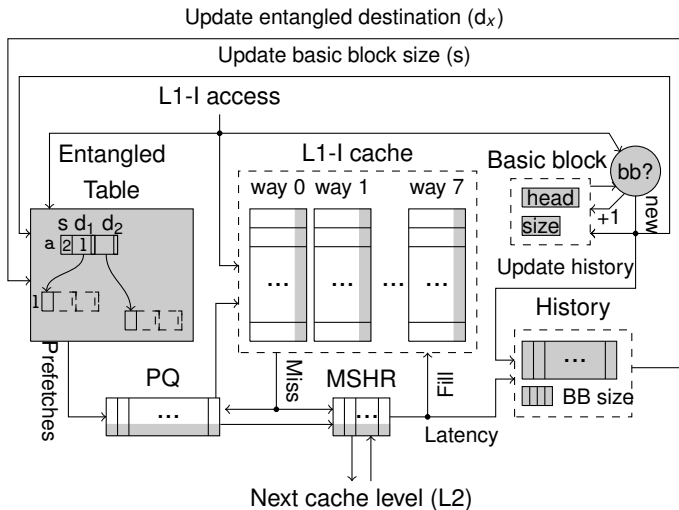
DESIGN OF THE ENTANGLING PREFETCHER - CONFIDENCE FOR ENTANGLED PAIRS



DESIGN OF THE ENTANGLING PREFETCHER - MERGING BASIC BLOCKS



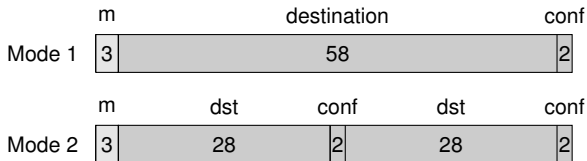
DESIGN OF THE ENTANGLING PREFETCHER



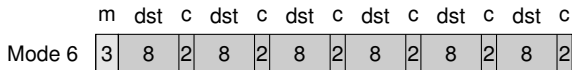
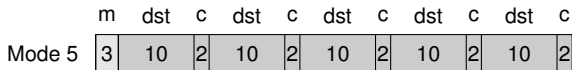
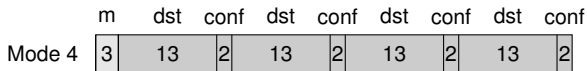
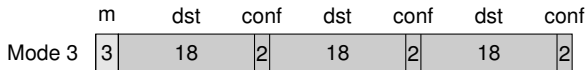
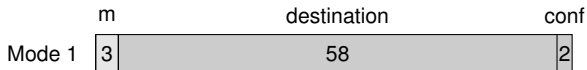
COMPRESSING DESTINATIONS

	m	destination	conf
Mode 1	3	58	2

COMPRESSING DESTINATIONS



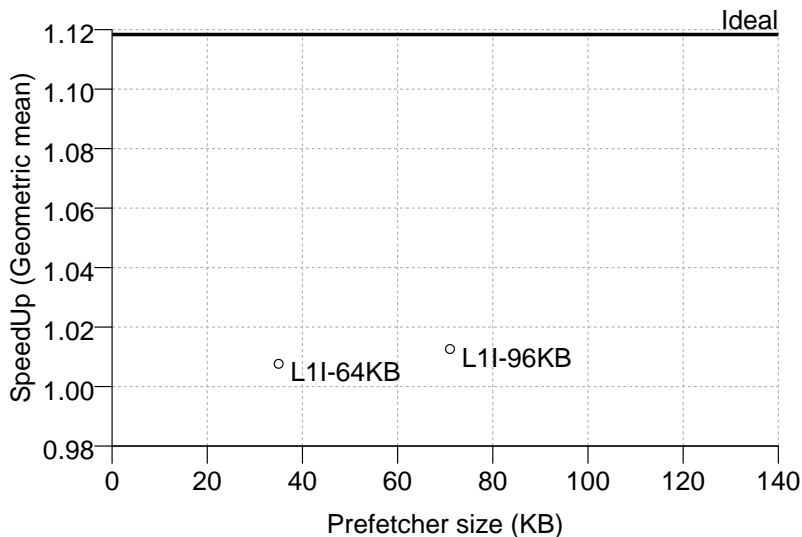
COMPRESSING DESTINATIONS



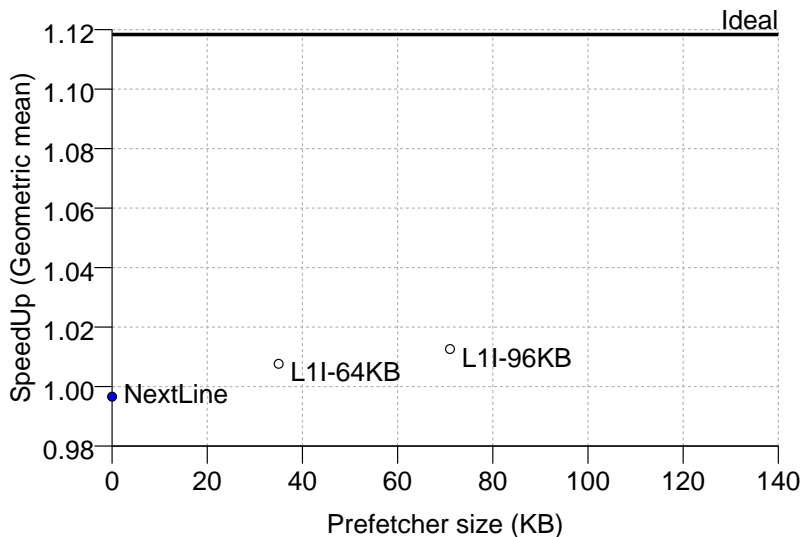
METHODOLOGY

- **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

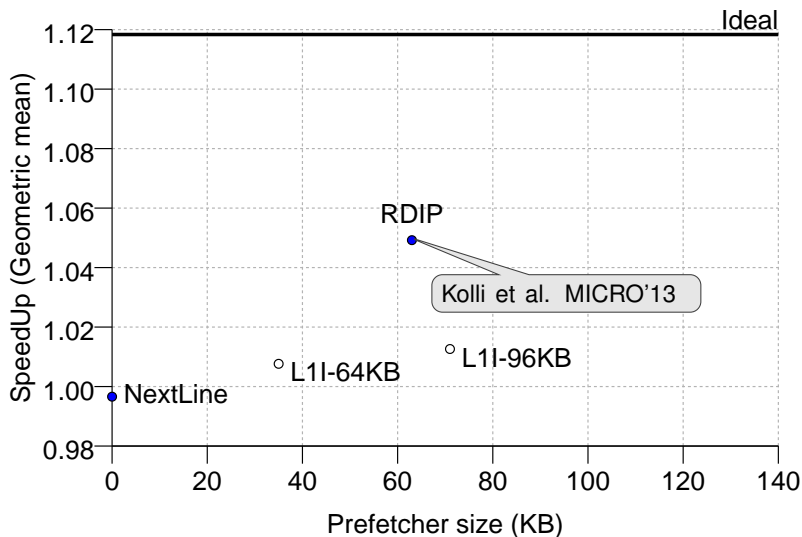
RESULTS: IPC VS MEMORY OVERHEAD



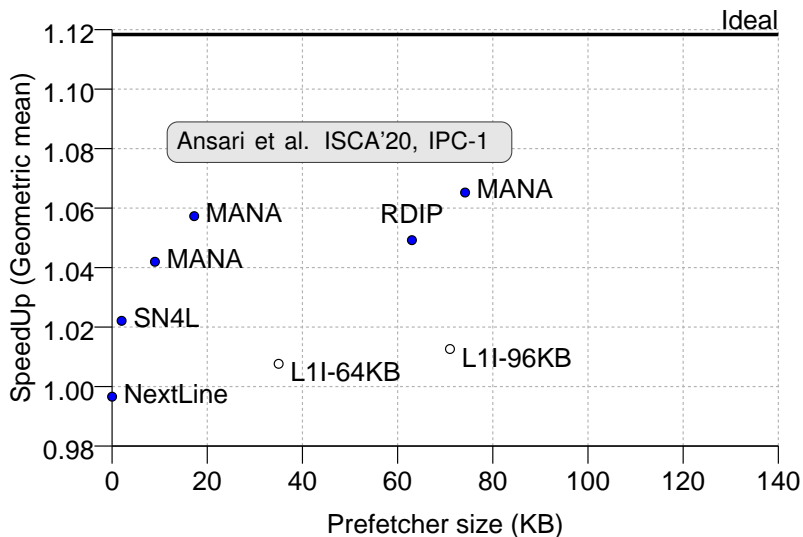
RESULTS: IPC VS MEMORY OVERHEAD



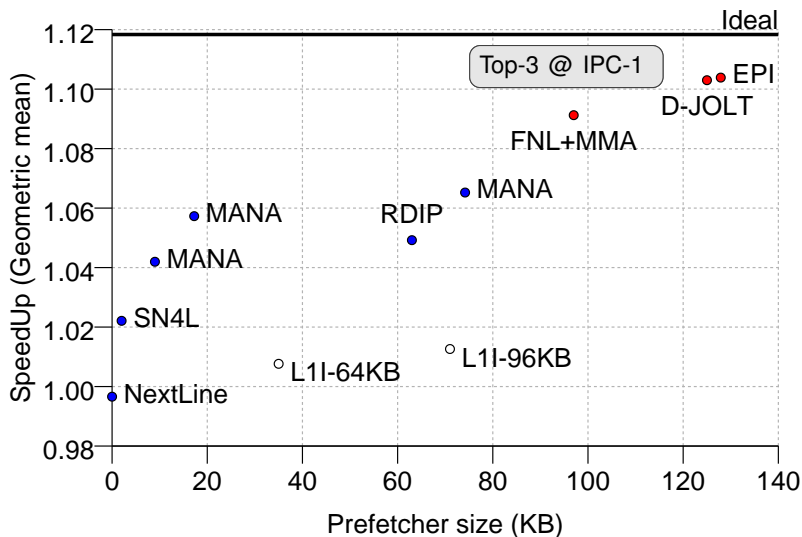
RESULTS: IPC VS MEMORY OVERHEAD



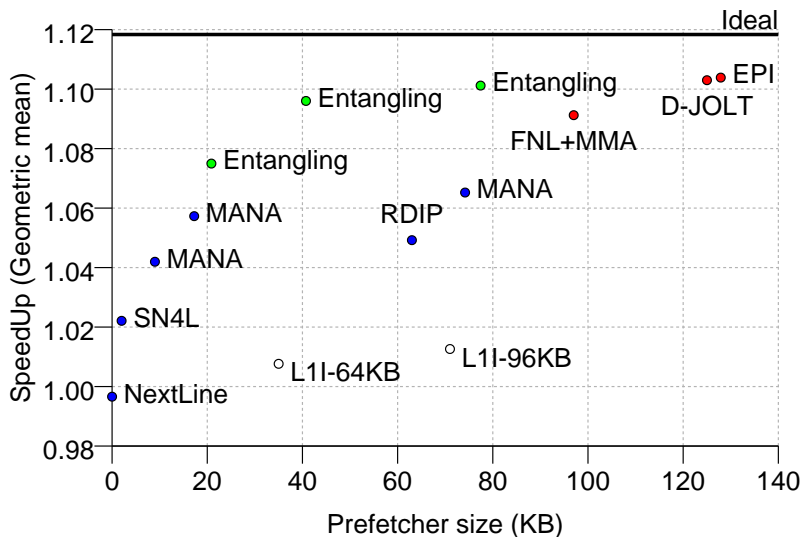
RESULTS: IPC VS MEMORY OVERHEAD



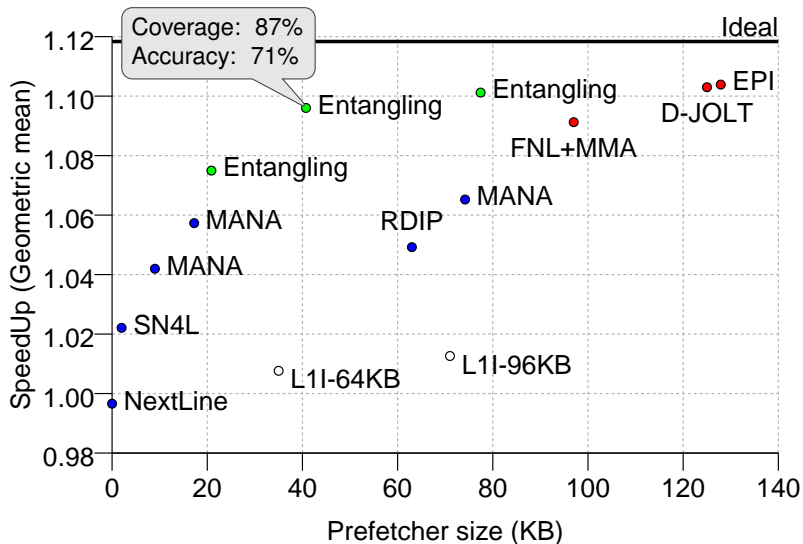
RESULTS: IPC VS MEMORY OVERHEAD



RESULTS: IPC VS MEMORY OVERHEAD



RESULTS: IPC VS MEMORY OVERHEAD



CONCLUDING REMARKS

- Data centers need good **prefetching** techniques
- **Timeliness** as a key property for a prefetcher
- **Entangle** heads of basic blocks to trigger timely prefetches
- Near **ideal** L1I performance with just 40KB

BOOSTING DATA CENTERS PERFORMANCE WITH THE ENTANGLING INSTRUCTION PREFETCHER

Alberto Ros

aros@ditec.um.es

Thank you!



ECHO, ERC Consolidator Grant (No 819134)