

BERTI: A PER-PAGE BEST-REQUEST-TIME DELTA PREFETCHER

Alberto Ros

Universidad de Murcia, Spain
aros@ditec.um.es



OUTLINE

- 1 MOTIVATION
- 2 THE BERTI PREFETCHER
- 3 EVALUATION
- 4 CONCLUSIONS AND FUTURE WORK



OUTLINE

- 1 MOTIVATION
- 2 THE BERTI PREFETCHER
- 3 EVALUATION
- 4 CONCLUSIONS AND FUTURE WORK



MOTIVATION

- **Stride** prefetcher, a simple technique but...
 - The stride can be **difficult** to find
 - The order of accesses can be altered by out-of-order cores, lower cache levels, or even prefetchers of lower cache levels
 - May lead to **late prefetches**
 - No timing considerations, although throttling can help



MOTIVATION

- **Best-offset prefetcher** (BOP) addresses the limitations of stride prefetchers
- Example: pattern found in `mcf`

Blocks accessed [0..63]	Stride
...110110110110	

MOTIVATION

- **Best-offset prefetcher** (BOP) addresses the limitations of stride prefetchers
- Example: pattern found in `mcf`

Blocks accessed [0..63]	Stride
...110110110110	-1?, -2?

Handwritten red squiggle under the binary string.



MOTIVATION

- **Best-offset prefetcher** (BOP) addresses the limitations of stride prefetchers
- Example: pattern found in `mcf`

Blocks accessed [0..63]	Stride
...110110110110	-1?, -2?, -3?



MOTIVATION

- **Best-offset prefetcher** (BOP) addresses the limitations of stride prefetchers
- Example: pattern found in `mcf`

Blocks accessed [0..63]	Stride	BO
...110110110110	-1?, -2?, -3?	-6



MOTIVATION

- **Best-offset prefetcher** (BOP) addresses the limitations of stride prefetchers
- Example: pattern found in `mcf`

Blocks accessed [0..63]	Stride	BO
...110 1 10 1 10 1 10	-1?, -2?, -3?	-6



MOTIVATION

- **Best-offset prefetcher** (BOP) addresses the limitations of stride prefetchers
- Example: pattern found in `mcf`

Blocks accessed [0..63]	Stride	BO
...110110110110	-1?, -2?, -3?	-6



MOTIVATION

- **Best-offset prefetcher** (BOP) addresses the limitations of stride prefetchers
- Example: pattern found in `mcf`

Blocks accessed [0..63]	Stride	BO
...110110110110	-1?, -2?, -3?	-6

- But BOP detects a best offset per application phase
 - And different pages have different best deltas!



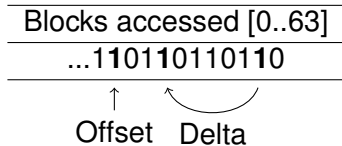
OUTLINE

- 1 MOTIVATION
- 2 THE BERTI PREFETCHER**
- 3 EVALUATION
- 4 CONCLUSIONS AND FUTURE WORK



TERMINOLOGY

- Offset vs Delta (or Stride)



KEY CONCEPT

- Per-page Berti (**B**est **r**quest **t**ime) Delta
 - Berti: the delta that would achieve more timely prefetches
 - It is calculated for each page



KEY CONCEPT

- Per-page Berti (**B**est **r**quest **t**ime) Delta
 - Berti: the delta that would achieve more timely prefetches
 - It is calculated for each page
- We track hot pages and their offsets accessed

Blocks accessed [0..63]

...110110110110



KEY CONCEPT

- Per-page Berti (**B**est **r**quest **t**ime) Delta
 - Berti: the delta that would achieve more timely prefetches
 - It is calculated for each page
- We track hot pages and their offsets accessed
- When pages become cold, Berti is calculated

Blocks accessed [0..63]	Berti
...110110110110	-6

↖ ↗
Delta

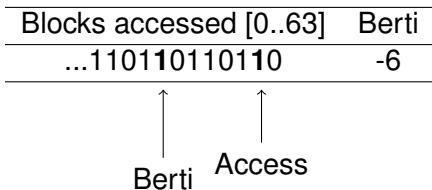
PREFETCHING MODES

- Two prefetching modes

Blocks accessed [0..63]	Berti
...110110110110	-6

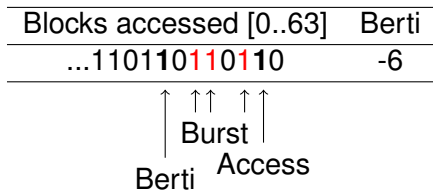
PREFETCHING MODES

- Two prefetching modes
 - Berti: One prefetch (timely)



PREFETCHING MODES

- Two prefetching modes
 - Berti: One prefetch (timely)
 - Burst: Several prefetches (first access to the page)
 - Expected to be late prefetches

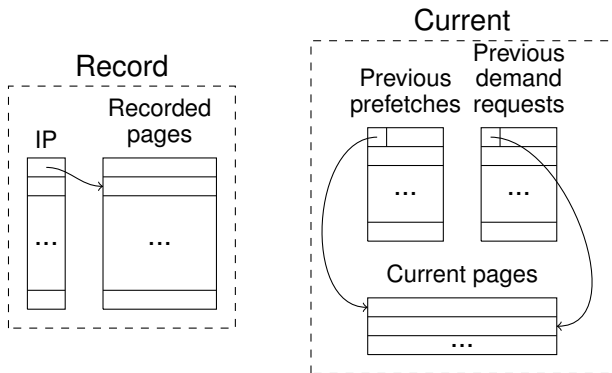


IP CLUSTERING

- Pages can be accessed by different load instructions
- We want all loads to a page to agree in the same Berti
 - Solution: Cluster the loads that access the same page

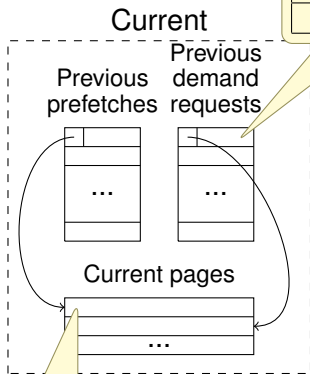
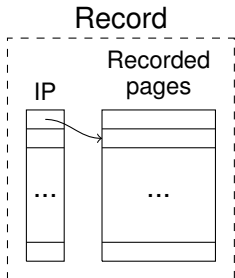


BERTI PREFETCHER OVERVIEW & EXAMPLE



BERTI PREFETCHER OVERVIEW & EXAMPLE

New page accessed and new block (offset 62)



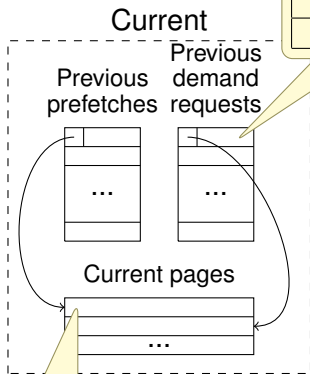
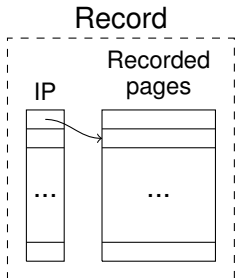
Offset	Cycle
62	10

Vector	First Offset	Berti	Counter
...0000000000010	62	0	0



BERTI PREFETCHER OVERVIEW & EXAMPLE

Access to block
with offset 61



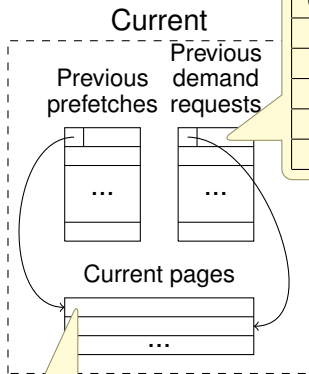
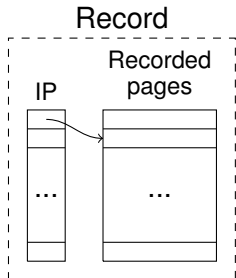
Offset	Cycle
62	10
61	20

Vector	First Offset	Berti	Counter
...000000000110	62	0	0



BERTI PREFETCHER OVERVIEW & EXAMPLE

Accesses to blocks
59, 58 and 56



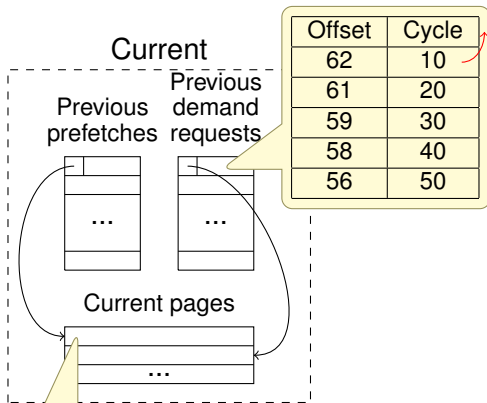
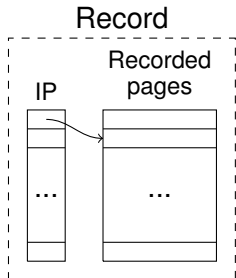
Offset	Cycle
62	10
61	20
59	30
58	40
56	50

Vector	First Offset	Berti	Counter
...000010110110	62	0	0



BERTI PREFETCHER OVERVIEW & EXAMPLE

Miss for block 62 resolved
Latency 35 cycles



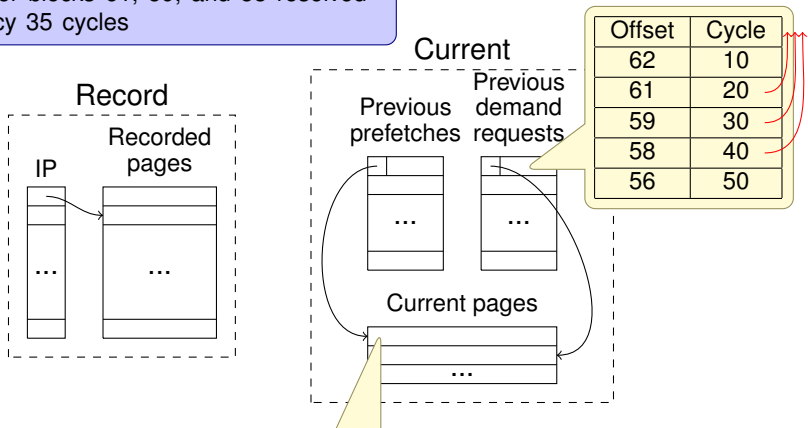
Offset	Cycle
62	10
61	20
59	30
58	40
56	50

Vector	First Offset	Berti	Counter
...000010110110	62	0	0



BERTI PREFETCHER OVERVIEW & EXAMPLE

Miss for blocks 61, 59, and 58 resolved
Latency 35 cycles

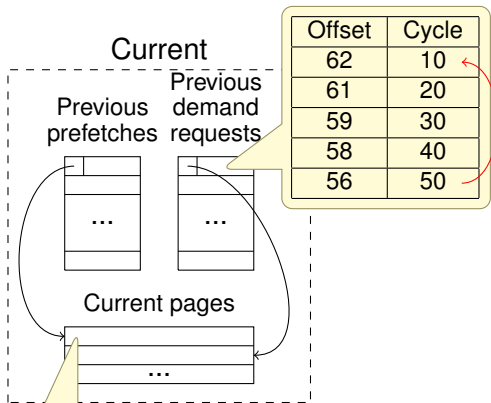
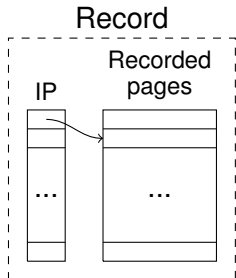


Vector	First Offset	Berti	Counter
...000010110110	62	0	0



BERTI PREFETCHER OVERVIEW & EXAMPLE

Miss for block 56 resolved
Latency 35 cycles



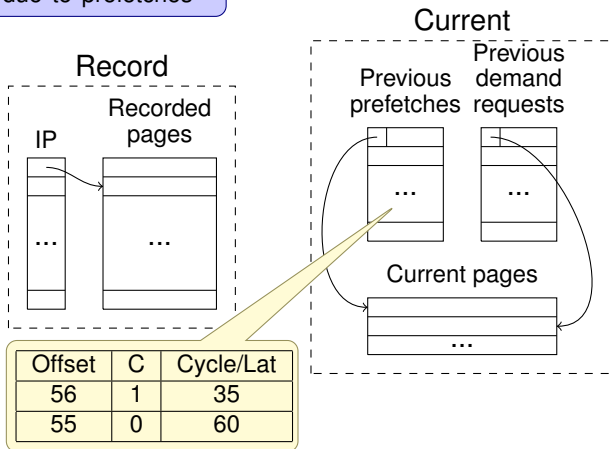
Offset	Cycle
62	10
61	20
59	30
58	40
56	50

Vector	First Offset	Berti	Counter
...000010110110	62	-6	1



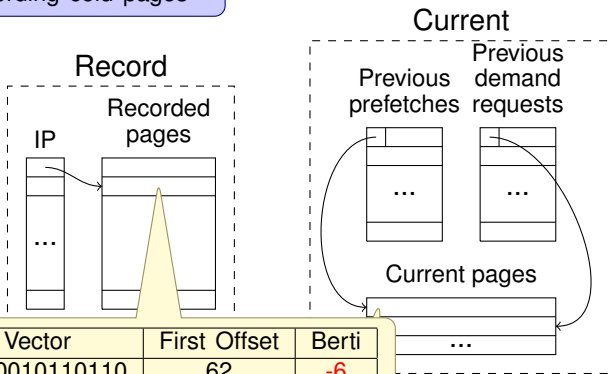
BERTI PREFETCHER OVERVIEW & EXAMPLE

Hits due to prefetches



BERTI PREFETCHER OVERVIEW & EXAMPLE

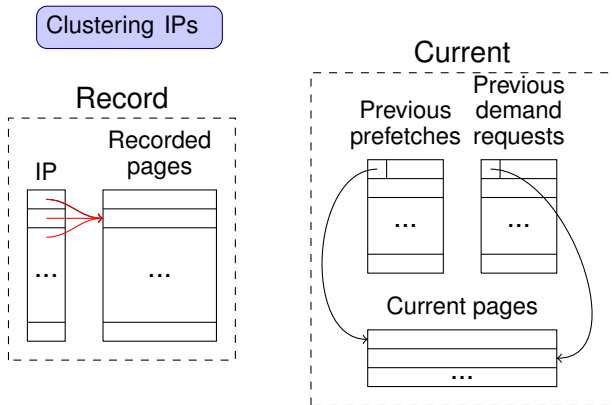
Recording cold pages



Vector	First Offset	Berti	Counter
...000010110110	62	-6	1



BERTI PREFETCHER OVERVIEW & EXAMPLE



BUILDING CONFIDENCE

- On a cache access the prefetcher checks if:
 - ① Matches page and its first access → high confidence
 - ② Matches IP and the first page access → high confidence
 - ③ Page is hot and Berti counter > 2 → medium confidence
 - ④ Matches page → low confidence
 - ⑤ Matches IP → low confidence

- If confidence is:
 - High and first access → Berti & burst
 - Medium or low → Berti

OUTLINE

- 1 MOTIVATION
- 2 THE BERTI PREFETCHER
- 3 EVALUATION**
- 4 CONCLUSIONS AND FUTURE WORK



CONFIGURATIONS AND BENCHMARKS

CONFIGURATIONS

- Cores: 1 and 4
- Prefetchers: None, NextLine, Stride, SPP, KPCP, and Berti
 - BOP not considered (porting it from DCP-2 did not give the expected results)

BENCHMARKS

- SPEC CPU 2017
- 1 core: > 1 MPKI at the LLC without prefetching
- 4 core: random mixes
- 50M instructions warm-up + 200M instructions stats



RESULTS

- Normalized with respect to no prefetch

Configuration (L1D, L2C, LLC)	1 core	4 cores
NextLine, NextLine, NextLine	1.2688	1.1045
NextLine, Stride, NextLine	1.2946	1.1147
NextLine, SPP, NextLine	1.3083	1.1043
NextLine, KPCP, NextLine	1.3142	1.1032
NextLine, Berti, NextLine	1.3211	1.1042
Berti, Berti, Berti	1.3347	1.1087
Berti+, Berti, NextLine	1.3471	1.1186
Berti+, Berti, None	1.3303	1.1268



COMMENTS/QUESTIONS

- Merged prefetches do not call `prefetcher_operate`
 - Implication: L2 and LLC prefetchers may not see the request of a block by the previous level!
 - How can this affect performance?
- Cache contention model
 - Cache port for prefetches?
 - Impact in filtering prefetches that hit in cache negligible



OUTLINE

- 1 MOTIVATION
- 2 THE BERTI PREFETCHER
- 3 EVALUATION
- 4 CONCLUSIONS AND FUTURE WORK**



CONCLUSIONS AND FUTURE WORK

CONCLUSIONS

- Per-page Berti Delta Prefetcher finds more timely prefetches

FUTURE WORK

- Adding confidence counters
- Explore Berti for LLC: late prefetches?



BERTI: A PER-PAGE BEST-REQUEST-TIME DELTA PREFETCHER

Alberto Ros

Universidad de Murcia, Spain
aros@ditec.um.es

