# Beyond Linear Read-Ahead: Logical Prefetching using Primary and Secondary Indexes in InnoDB

**preFOSDEM MySQL Belgian Days 2026**

Vitor Oliveira
Huawei Cloud Database Advanced Technology Laboratory

2026/01/30

HUAWEI

# About me

I'm a Performance Architect with background in database and in HPC performance analysis and optimization.

I work at Huawei's Cloud Database Advanced Technology Laboratory, Shannon Research Center since 2021.

My activity has been focused on database optimization for OLTP and analytics workloads, profiling tools and novel architectures involving persistent memory, low-latency networks and memory fabrics.

Prior to that I was a member of Oracle's MySQL Replication team since 2014.

There I collaborated in projects like Group Replication and MySQL Database Service (HA) on OCI, and worked closely on features like flow control and WRITESET dependency tracking, among others.

I've also been teaching sometimes, researching and enjoying music, theatre, physics and also architecture (the real one).

# Introduction

MySQL performs at peak performance when the active workset of the database is kept in main memory.

Performance degrades drastically if pages need to be swaped in and out if part of the database is offloaded to storage.
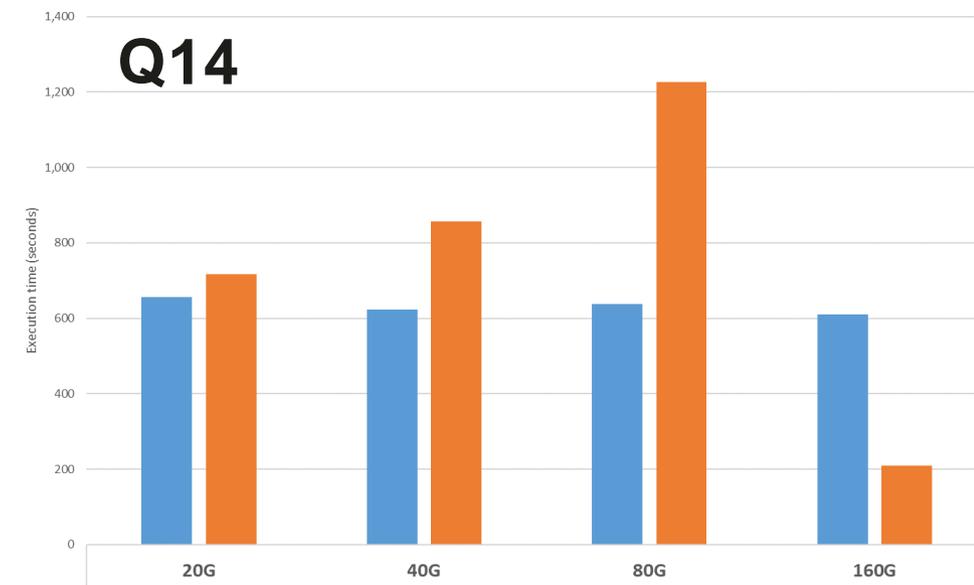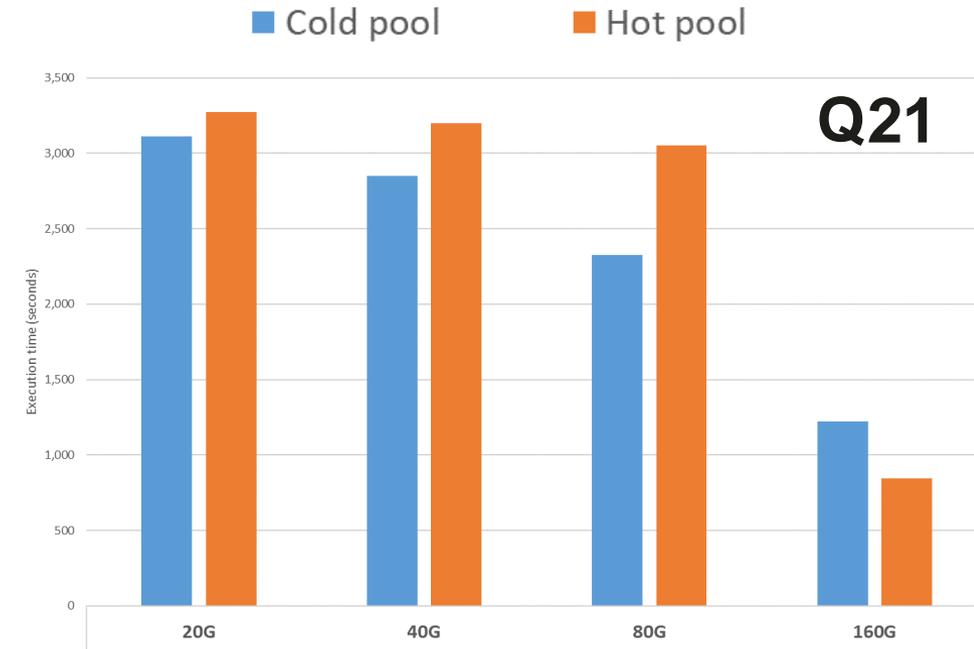
Following our interest in analytics and a strange behavior reported on network storage, we became interested on the impact of prefetching on TPC-H workloads larger than memory.

This presentation is about some of the issues we found and about an approach to address them.

# The problem

**1** TPC-H results showed that some queries would run faster when the system started (cold pool) than on subsequent runs (hot pool).

**2** That's unexpected, as having data already present in memory should improve query time, not make it worse.

**3** Even stranger, some results were actually better on smaller buffer pools than on larger buffer pools.

## TPC-H SF100 Query Execution Time

■ Cold pool    ■ Hot pool

**Q21**

Execution time (seconds): 3,500 / 3,000 / 2,500 / 2,000 / 1,500 / 1,000 / 500 / 0

20G    40G    80G    160G

**Q14**

Execution time (seconds): 1,400 / 1,200 / 1,000 / 800 / 600 / 400 / 200 / 0

20G    40G    80G    160G

# Prefetching?

Suspecting prefetcher involvement we experimented with InnoDB's two prefetcher mechanisms:

- **Linear read-ahead**
  Triggers the read of the next extent if it reads sequentially at least X pages of current extent (default X=56).

- **Random read-ahead**
  Fetches all pages if more than 13 pages in the extent are already loaded in the buffer pool (disabled by default).

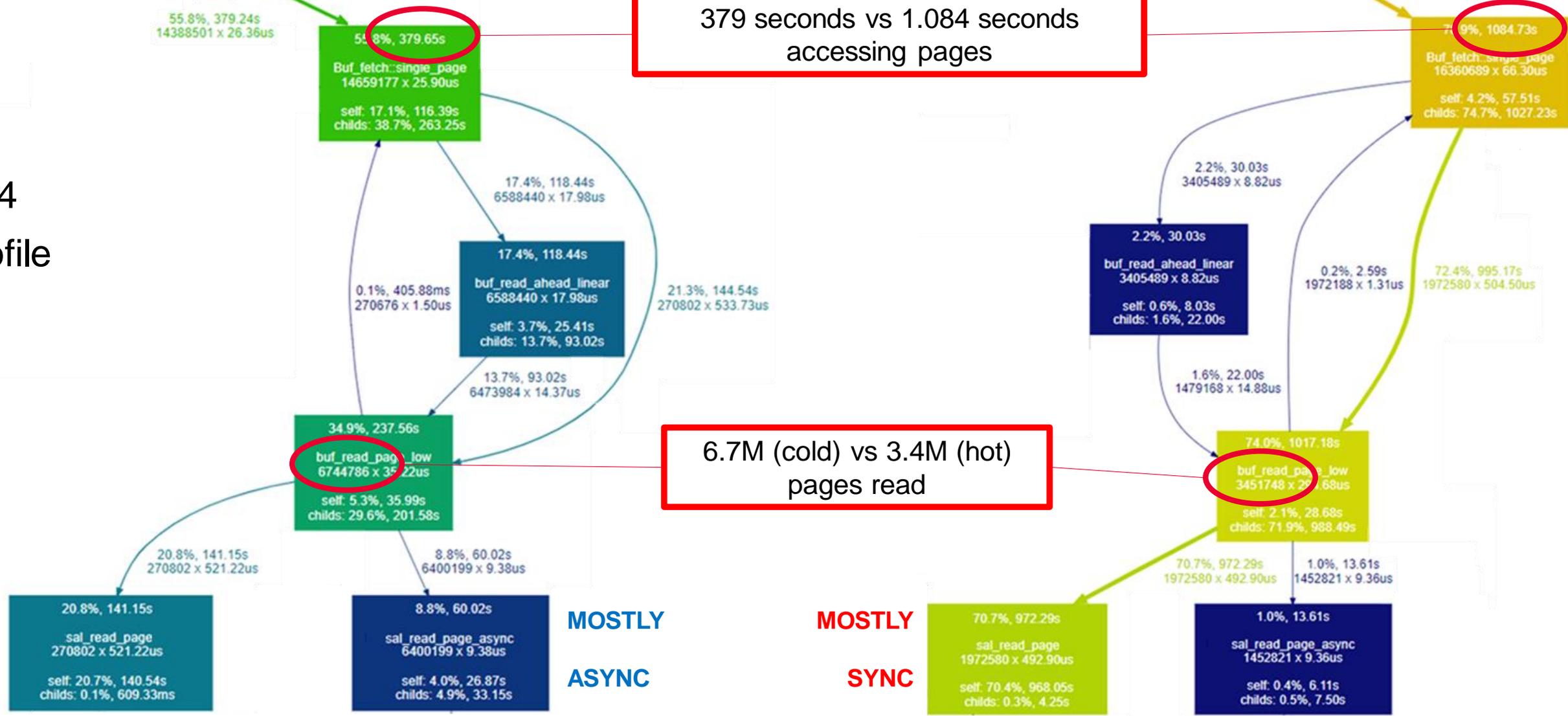| Q14 Time | | BP | |
|---|---|---|---|
| **Random RA** | **Linear RA** | **COLD** | **HOT** |
| ⊟ **OFF** | OFF (0) | 3,075 | 1,370 |
| | threshold=8 | 829 | 1,185 |
| | threshold=24 | 2,345 | 1,184 |
| | DEFAULT (56) | **529** | **1,226** |
| | threshold=64 | 2,125 | 1,267 |
| ⊟ **ON** | OFF (0) | 2,975 | 802 |
| | threshold=8 | 663 | 547 |
| | threshold=24 | 2,612 | 678 |
| | DEFAULT (56) | **2,691** | **696** |
| | threshold=64 | 2,326 | 660 |

COLD RUN

HOT RUN

Q14
Profile

100.0%, 679.92s
mysql_execute_command
1 x 679.92s
self: 43.2%, 293.61s
childs: 56.8%, 386.30s

100.0%, 1375.43s
mysql_execute_command
1 x 1375.43s
self: 20.8%, 286.12s
childs: 79.2%, 1089.31s

55.8%, 379.24s
14388501 x 26.36us

78.7%, 1082.14s
14388501 x 75.21us

55.8%, 379.65s
Buf_fetch::single_page
14659177 x 25.90us
self: 17.1%, 116.39s
childs: 38.7%, 263.25s

379 seconds vs 1.084 seconds
accessing pages

78.9%, 1084.73s
Buf_fetch::single_page
16360689 x 66.30us
self: 4.2%, 57.51s
childs: 74.7%, 1027.23s

17.4%, 118.44s
6588440 x 17.98us

2.2%, 30.03s
3405489 x 8.82us

0.1%, 405.88ms
270676 x 1.50us

17.4%, 118.44s
buf_read_ahead_linear
6588440 x 17.98us
self: 3.7%, 25.41s
childs: 13.7%, 93.02s

21.3%, 144.54s
270802 x 533.73us

2.2%, 30.03s
buf_read_ahead_linear
3405489 x 8.82us
self: 0.6%, 8.03s
childs: 1.6%, 22.00s

0.2%, 2.59s
1972188 x 1.31us

72.4%, 995.17s
1972580 x 504.50us

13.7%, 93.02s
6473984 x 14.37us

1.6%, 22.00s
1479168 x 14.88us

34.9%, 237.56s
buf_read_page_low
6744786 x 35.22us
self: 5.3%, 35.99s
childs: 29.6%, 201.58s

6.7M (cold) vs 3.4M (hot)
pages read

74.0%, 1017.18s
buf_read_page_low
3451748 x 29.68us
self: 2.1%, 28.68s
childs: 71.9%, 988.49s

20.8%, 141.15s
270802 x 521.22us

8.8%, 60.02s
6400199 x 9.38us

70.7%, 972.29s
1972580 x 492.90us

1.0%, 13.61s
1452821 x 9.36us

20.8%, 141.15s
sal_read_page
270802 x 521.22us
self: 20.7%, 140.54s
childs: 0.1%, 609.33ms

8.8%, 60.02s
sal_read_page_async
6400199 x 9.38us
self: 4.0%, 26.87s
childs: 4.9%, 33.15s

MOSTLY

ASYNC

MOSTLY

SYNC

70.7%, 972.29s
sal_read_page
1972580 x 492.90us
self: 70.4%, 968.05s
childs: 0.3%, 4.25s

1.0%, 13.61s
sal_read_page_async
1452821 x 9.36us
self: 0.4%, 6.11s
childs: 0.5%, 7.50s

# LRU insertion point

The instability shown suggests the LRU management may be playing a role, so we included also the innodb-old-block-percent (OBP).

| Q14 Time | | BP ▼ OBP ▼ | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | | ⊟COLD | | | ⊟HOT | | |
| Random RA ▼ | Linear RA ▼ | obp=5 | obp=37 | obp=95 | obp=5 | obp=37 | obp=95 |
| ⊟OFF | OFF (0) | 3,055 | 3,075 | 3,128 | 510 | 1,370 | 2,930 |
| | threshold=8 | 522 | 829 | 530 | 375 | 1,185 | 2,582 |
| | threshold=24 | 537 | 2,345 | 500 | 371 | 1,184 | 2,666 |
| | DEFAULT (56) | **1,328** | **529** | **537** | **375** | **1,226** | **2,588** |
| | threshold=64 | 546 | 2,125 | 524 | 373 | 1,267 | 2,731 |
| ⊟ON | OFF (0) | 1,119 | 2,975 | 3,042 | 413 | 802 | 2,849 |
| | threshold=8 | 1,368 | 663 | 522 | 314 | 547 | 2,584 |
| | threshold=24 | 510 | 2,612 | 506 | 302 | 678 | 2,572 |
| | DEFAULT (56) | **523** | **2,691** | **534** | **301** | **696** | **2,596** |
| | threshold=64 | 524 | 2,326 | 524 | 302 | 660 | 2,731 |

1. The LRU insertion point has indeed a big impact, with 5% showing much lower numbers than the default 37%.

2. That does not apply to all cases, and in fact it does apply to MySQL's default prefetcher configuration.

3. But it does apply to when random read ahead, which is also unexpected since the workload is mostly sequential.

# The pool is on fire

| METRIC | COLD POOL | | HOT POOL | |
|---|---|---|---|---|
| | obp5 | obp37 | obp5 | obp37 |
| page uses / storage read | 4.0 | 4.1 | 12.4 | 6.9 |
| pages loaded / buffer pool size | 127.2% | 126.4% | 34.0% | 65.3% |
| pages read synchronously | 4.3% | 3.5% | 6.3% | 30.5% |
| pages read asynchronously | 122.9% | 122.9% | 27.7% | 34.8% |
| unused read-ahead pages | 0.4% | 0.4% | 0.4% | 0.7% |

background = black -> grey -> white
number of pages in memory in that range in the **orders** table

red = read from storage

green = read from buffer pool

blue = page freed

page id

OBP=37

execution time

page id

OBP=5

# Linear read-ahead

**1** The linear read-ahead prefetcher is invoked when trying to access a page that is missing from the buffer pool, but it is only triggered if called for the boundary pages of the extent;

**2** Then it scans the pages in the current extent and counts the number of pages present in the buffer pool that have a sequential access time, either descending order if it is the first page of the extent, or ascending order if it is the last page of the extent;

**3** If that number exceeds the **innodb_read_ahead_threshold**, the successor/predecessor of the current page is fetched and if that is also a boundary page of the extent it belongs to, it will trigger the asynchronous read of all 64 pages in that extent.

# Limitations

**(1)** **It depends on the physical layout of the data**
- Fragmented databases, as long running database tend to be, will not benefit much from the prefetcher as the sequence information will be wrong.

**(2)** **Linear read-ahead can be disabled by the pages previously loaded in the BP**
- if the boundary page is present in the buffer pool, nothing is prefetched;
- non sequential access time between new and old pages counts as a sequence interruption, which counts for the maximum number in innodb_read_ahead_threshold.

**(3)** **The query may have to wait for storage even when linear read-ahead is working**
- Prefetch is only triggered on the last page of the previous extent;
- The next page will take time to load, even more when overloaded with 64 simultaneous requests.

Can we do something about this?...

# Enhanced read-ahead (ERA)

To try to improve on the previous issues, we developed an enhanced read-ahead mechanism to address, mostly as exploratory tool to evaluate the impact of different alternative approaches to prefetching.

Initial requirements:

1. prefetch pages even if theirs neighbours are already in the buffer pool;
2. Use logical prefetching to benefit the workload even on fragmented databases, not just the physical layout;
3. be friendly to storage, with patterns fit for cloud-based deployments.

Things to avoid:

- don't read too late to avoid waiting for storage;
- don't read too early to avoid prematurely discarding pages from the buffer pool;
- don't overload the system when prefetching is not needed or is ineffective.

# Logical prefetching

A main problem is knowing which page follows another without actually reading that pages from disk, which is why linear read-ahead uses physical proximity as an heuristic.

> That sequence can be read ahead from the clustered index, similar to what facebook did in https://yoshinorimatsunobu.blogspot.com/2013/10/making-full-table-scan-10x-faster-in.html.

> We also wanted to try using secondary indexes, even if those have to scan both the secondary index and the corresponding clustered index to find the pages that are next in a sequence.

> But to avoid some overhead, we also considered a page sequence cache to keep a cache of the pages sequences which have been loaded over time

Ideally the prefetcher should also work with the semantic information provided by the MySQL layer, that would be used to guide the prefetching more closely.

# Design

**1** Track all page accesses per thread and store the page and additional positioning information in a sequence history list;

**2** If at least N predecessors of the current page/position are in the history list, make sure that at least M successors of the current page/position is in the global prefetch list, except for pages that:

> are already in memory or

> pre-fetch was already been requested.

**3** Discover the predecessor/successor page sequences ahead of the requests, using page sequence information obtained from multiple sources using independent threads.

# Design

1. Control thread does history checking and selects the list of successors to pre-fetch;

2. Reader thread takes the pre-fetch requests and issues the storage reads and tracks their success;

3. Scanner thread fetches sequence information from indexes when that is found missing.



The threads communicate through messages:

- Query threads sends the page_id and the frontend context information to the control thread when accessing a page;

- When the control thread needs to read from disk it sends the list of pages to read to the read queue, and when it needs index pages to be processed it sends the index page_id to the scanner queue;

- When the scanner thread has read the sequences, it sends the list to the control thread to import.

# Evaluation

TPC-H SF100 on a 100GB buffer pool, a little less than half of the database on the buffer pool.

1. LRA saves 30% of the time, while ERA 52% of the time, compared to disabling prefetching;

2. ERA benefit more queries than LRA, or in greater amount;

3. In queries 4, 6, 8, 10, 11, 14, 17, 19 and 22 the difference is particularly significant.



Prefetch Effectiveness: TPC-H SF100, BP=100G

# Evaluation

# Evaluation



Legend: era: Buffer Pool · era: Page Cache · era: Estimate · era: Secondary Index · era: Primary Index · Rows Read · Page Requests · Page Reads · Prefetches

Q8

Q10

Primary index and more estimate on second phase

# Evaluation



Legend:
- era: Buffer Pool
- era: Page Cache
- era: Estimate
- era: Secondary Index
- era: Primary Index
- Rows Read
- Page Requests
- Page Reads
- Prefetches

Q11

Secondary index and later estimate + page cache

# Evaluation



Legend: era: Buffer Pool, era: Page Cache, era: Estimate, era: Secondary Index, era: Primary Index, Rows Read, Page Requests, Page Reads, Prefetches

Q9

Q21

Primary index on first phase, then estimate only

# Conclusion

1. Prefetching is an important tool to hide storage latency and it is well fit to the sequential nature of many large analytics queries.

2. The enhanced read-ahead mechanism proposes a mechanism that is more robust than InnoDB's linear read head, applies to more queries and in a larger scale;

3. But that comes at the cost of being complex and computationally demanding, the overhead only makes sense if large databases are used.

# Thank you.

Bring digital to every person, home and organization for a fully connected, intelligent world.

# SF10, 10G

- Overall the existing linear read-ahead (LRA) saves 8% of the time compared to running without prefetcher;

- The enhanced read-ahead (ERA) saves 52% compared to no prefetching, and 48% compared to linear read-ahead.
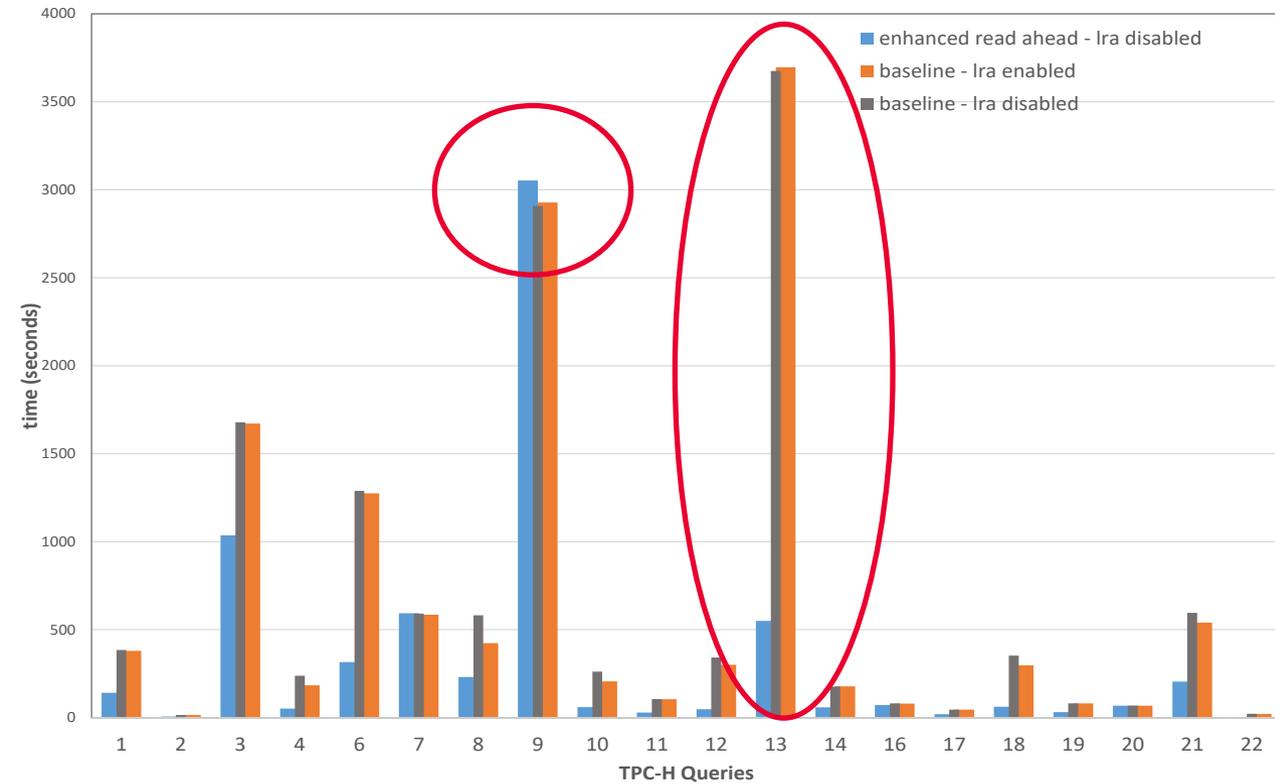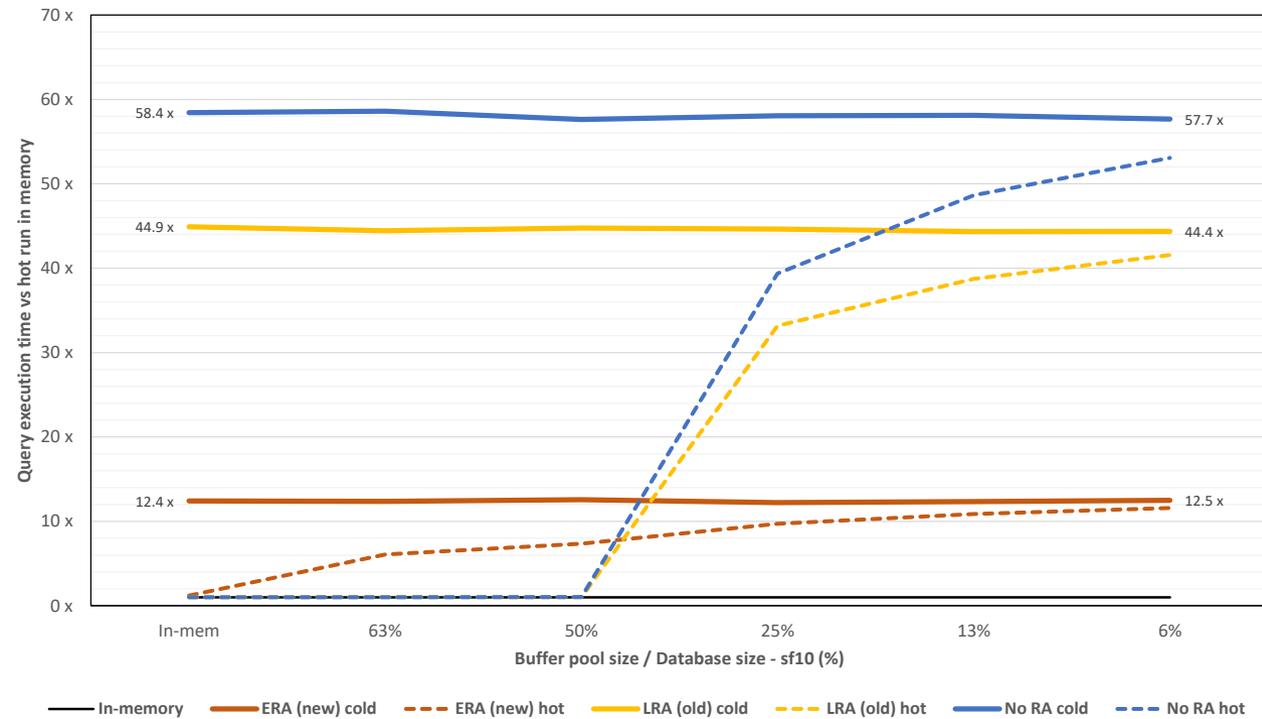
**Prefetch Effectiveness: TPC-H SF10, BP=10G**

# SF10, 2GB and 1GB



Prefetch Effectiveness: TPC-H SF10, BP=2G

Prefetch Effectiveness: TPC-H SF10, BP=1G

# SF10



**Impact of Prefetching on TPC-H Q4 Execution Time**
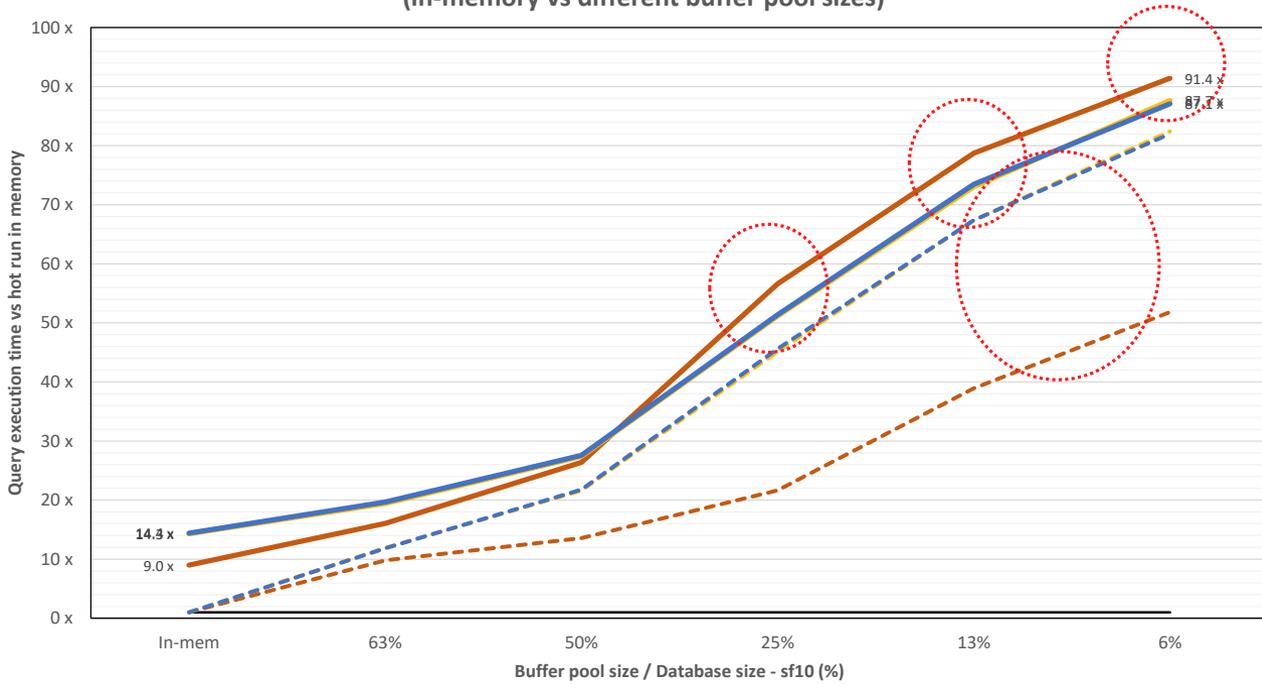**(in-memory vs different buffer pool sizes)**

**Impact of Prefetching on TPC-H Q8 Execution Time**
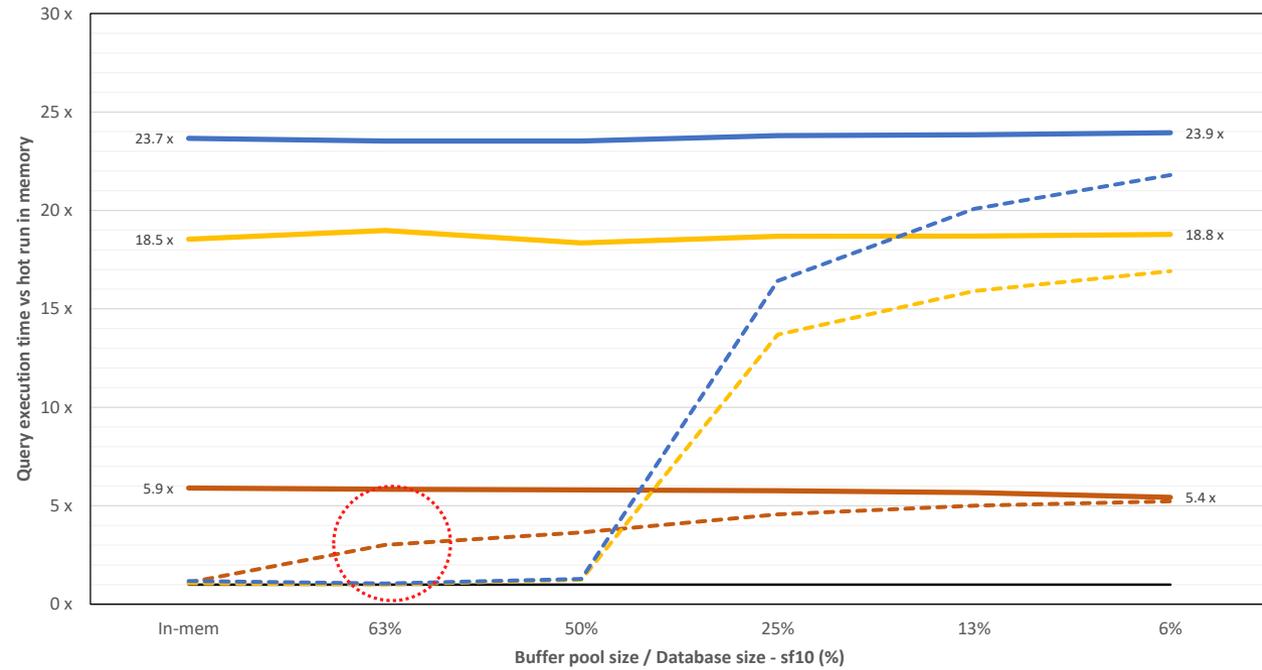**(in-memory vs different buffer pool sizes)**

# Not so good



**Impact of Prefetching on TPC-H Q9 Execution Time**
**(in-memory vs different buffer pool sizes)**

**Impact of Prefetching on TPC-H Q10 Execution Time**
**(in-memory vs different buffer pool sizes)**