Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

1 of 46

# Evaluation of Pointer Swizzling Techniques for DBMS Buffer Management

## Max Gilbert

University of Kaiserslautern

*m_gilbert13@cs.uni-kl.de*

June 8, 2017

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

2 of 46

## Table of Contents

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

3 of 46

# Section 1

# Pointer Swizzling as in "In-Memory Performance for Big Data"

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                                                    4 of 46

Subsection 1

Locate Pages in the Buffer Pool without Pointer Swizzling

# Overview of Search Strategies ([HR01])

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     Page Eviction Strategies     End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                                        5 of 46

# Overview of Search Strategies ([HR01])

Search
Strategy

## Overview of Search Strategies ([HR01])

```
                    Search
                    Strategy


    Direct Search in
    the Buffer Frames
```

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     Page Eviction Strategies     End

Locate Pages in the Buffer Pool without Pointer Swizzling     5 of 46

## Overview of Search Strategies ([HR01])

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End
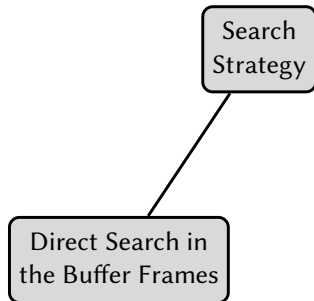
Locate Pages in the Buffer Pool without Pointer Swizzling    5 of 46

## Overview of Search Strategies ([HR01])
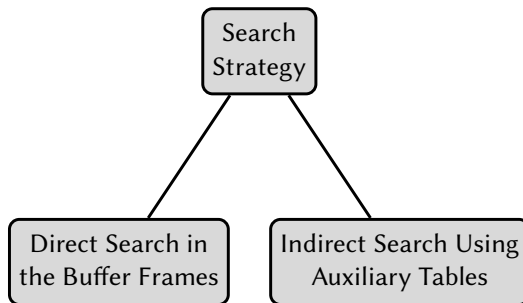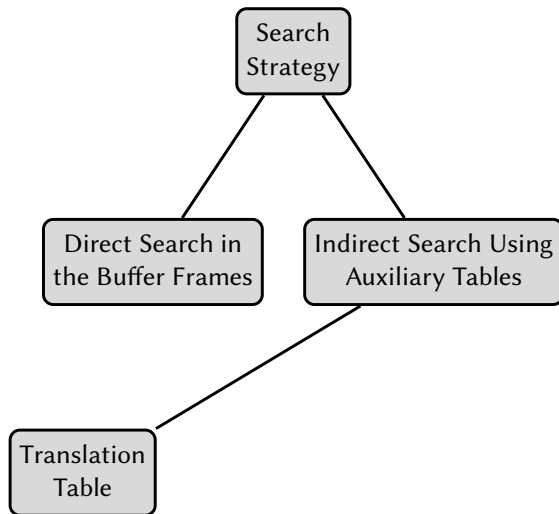
## Overview of Search Strategies ([HR01])

## Overview of Search Strategies ([HR01])

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     Page Eviction Strategies     End

Locate Pages in the Buffer Pool without Pointer Swizzling                                      5 of 46

## Overview of Search Strategies ([HR01])

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling    5 of 46

## Overview of Search Strategies ([HR01])

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                                                5 of 46

# Overview of Search Strategies ([HR01])

# Direct Search in the Buffer Frames & Unsorted Table

# Direct Search in the Buffer Frames & Unsorted Table

### Direct Search in the Buffer Frames

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                                    6 of 46

## Direct Search in the Buffer Frames & Unsorted Table

### Direct Search in the Buffer Frames

▶ Checks in each buffer frame the page ID of the contained page

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     Page Eviction Strategies     End

Locate Pages in the Buffer Pool without Pointer Swizzling        6 of 46

# Direct Search in the Buffer Frames & Unsorted Table

### Direct Search in the Buffer Frames

- ▶ Checks in each buffer frame the page ID of the contained page
- ▶ $T_{\text{avg}}^{\text{search}} \in \mathcal{O}\left(\frac{n}{2}\right)$, $T_{\text{worst}}^{\text{search}} \in \mathcal{O}\left(n\right)$

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     Page Eviction Strategies     End

Locate Pages in the Buffer Pool without Pointer Swizzling          6 of 46

# Direct Search in the Buffer Frames & Unsorted Table

### Direct Search in the Buffer Frames

- ▶ Checks in each buffer frame the page ID of the contained page
- ▶ $T_{\mathrm{avg}}^{\mathrm{search}} \in \mathcal{O}\left(\frac{n}{2}\right)$, $T_{\mathrm{worst}}^{\mathrm{search}} \in \mathcal{O}\left(n\right)$
- ▶ The usage of virtual memory management can result in extensive swapping due to read access to many pages!

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling    6 of 46

## Direct Search in the Buffer Frames & Unsorted Table

### Direct Search in the Buffer Frames

▶ Checks in each buffer frame the page ID of the contained page

▶ $T_{\text{avg}}^{\text{search}} \in \mathcal{O}\left(\frac{n}{2}\right)$, $T_{\text{worst}}^{\text{search}} \in \mathcal{O}\left(n\right)$

▶ The usage of virtual memory management can result in extensive swapping due to read access to many pages!

### Unsorted Table

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                                    6 of 46

## Direct Search in the Buffer Frames & Unsorted Table

### Direct Search in the Buffer Frames

- ▶ Checks in each buffer frame the page ID of the contained page
- ▶ $T_{\text{avg}}^{\text{search}} \in \mathcal{O}\left(\frac{n}{2}\right)$, $T_{\text{worst}}^{\text{search}} \in \mathcal{O}(n)$
- ▶ The usage of virtual memory management can result in extensive swapping due to read access to many pages!

### Unsorted Table

- ▶ Auxiliary data structure of size $S_{\text{pace}} \in \mathcal{O}(n)$

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|------|------|------|------|------|------|------|------|------|
| 7785 | 6977 | 4347 | 3380 | 5610 | 6376 | 4877 | 3332 | 3354 |

Figure: An unsorted table used to map buffer frames to page IDs.

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling    6 of 46

# Direct Search in the Buffer Frames & Unsorted Table

### Direct Search in the Buffer Frames

- ▶ Checks in each buffer frame the page ID of the contained page
- ▶ $T_{\text{avg}}^{\text{search}} \in \mathcal{O}\left(\frac{n}{2}\right)$, $T_{\text{worst}}^{\text{search}} \in \mathcal{O}(n)$
- ▶ The usage of virtual memory management can result in extensive swapping due to read access to many pages!

### Unsorted Table

- ▶ Auxiliary data structure of size $S_{\text{pace}} \in \mathcal{O}(n)$
- ▶ $T_{\text{avg}}^{\text{search}} \in \mathcal{O}\left(\frac{n}{2}\right)$, $T_{\text{worst}}^{\text{search}} \in \mathcal{O}(n)$

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| 7785 | 6977 | 4347 | 3380 | 5610 | 6376 | 4877 | 3332 | 3354 |

Figure: An unsorted table used to map buffer frames to page IDs.

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling    7 of 46

## Translation Table

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     Page Eviction Strategies     End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                                    7 of 46

## Translation Table

▶ Auxiliary data structure with one entry per page in the database
$\implies S_{\text{pace}} \in \mathcal{O}(p)$

| 0 | · |
| --- | --- |

| | |
| --- | --- |
| 3331 | · |
| 3332 | 7 |
| 3333 | · |

| 3352 | · |
| --- | --- |
| 3353 | · |
| 3354 | 8 |
| 3355 | · |
| 3356 | · |

| 3378 | · |
| --- | --- |
| 3379 | · |
| 3380 | 3 |
| 3381 | · |
| 3382 | · |

| 4345 | · |
| --- | --- |
| 4346 | · |
| 4347 | 2 |
| 4348 | · |
| 4349 | · |

| 4875 | · |
| --- | --- |
| 4876 | · |
| 4877 | 6 |
| 4878 | · |
| 4879 | · |

| 5608 | · |
| --- | --- |
| 5609 | · |
| 5610 | 4 |
| 5611 | · |
| 5612 | · |

| 6374 | · |
| --- | --- |
| 6375 | · |
| 6376 | 5 |
| 6377 | · |
| 6378 | · |

| 6975 | · |
| --- | --- |
| 6976 | · |
| 6977 | 1 |
| 6978 | · |
| 6979 | · |

| 7783 | · |
| --- | --- |
| 7784 | · |
| 7785 | 0 |
| 7786 | · |
| 7787 | · |

Figure: A translation table used to map page IDs to buffer frames.

Pointer Swizzling in the DBMS Buffer Management · · · Performance Evaluation of Pointer Swizzling · · · Page Eviction Strategies · · · End

Locate Pages in the Buffer Pool without Pointer Swizzling · · · 7 of 46

## Translation Table

- ▶ Auxiliary data structure with one entry per page in the database
  $\implies S_{\text{pace}} \in \mathcal{O}(p)$
- ▶ $T^{\text{search}} \in \mathcal{O}(1)$, $T^{\text{insert}} \in \mathcal{O}(1)$

| 0 | · |
|---|---|

| ⋮ |
|---|

| 3331 | · |
|------|---|
| 3332 | 7 |
| 3333 | · |

| 3352 | · |
|------|---|
| 3353 | · |
| 3354 | 8 |
| 3355 | · |
| 3356 | · |

| 3378 | · |
|------|---|
| 3379 | · |
| 3380 | 3 |
| 3381 | · |
| 3382 | · |

| 4345 | · |
|------|---|
| 4346 | · |
| 4347 | 2 |
| 4348 | · |
| 4349 | · |

| 4875 | · |
|------|---|
| 4876 | · |
| 4877 | 6 |
| 4878 | · |
| 4879 | · |

| 5608 | · |
|------|---|
| 5609 | · |
| 5610 | 4 |
| 5611 | · |
| 5612 | · |

| 6374 | · |
|------|---|
| 6375 | · |
| 6376 | 5 |
| 6377 | · |
| 6378 | · |

| 6975 | · |
|------|---|
| 6976 | · |
| 6977 | 1 |
| 6978 | · |
| 6979 | · |

| 7783 | · |
|------|---|
| 7784 | · |
| 7785 | 0 |
| 7786 | · |
| 7787 | · |

Figure: A translation table used to map page IDs to buffer frames.

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                    8 of 46

# Sorted & Chained Table

# Sorted & Chained Table

### Sorted Table

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     Page Eviction Strategies     End

Locate Pages in the Buffer Pool without Pointer Swizzling          8 of 46

## Sorted & Chained Table

### Sorted Table

▶ Auxiliary data structure using a table sorted by page ID only containing cached pages

| 3332 | 3354 | 3380 | 4347 | 4877 | 5610 | 6376 | 6977 | 7785 |
|------|------|------|------|------|------|------|------|------|
| → 7 | → 8 | → 3 | → 2 | → 6 | → 4 | → 5 | → 1 | → 0 |

Figure: A sorted table used to map page IDs to buffer frames.

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling    8 of 46

## Sorted & Chained Table

### Sorted Table

▶ Auxiliary data structure using a table sorted by page ID only containing cached pages

▶ $T_{\text{avg}}^{\text{search}} \in \mathcal{O}\left(\log_2 n\right)$, $T_{\text{avg}}^{\text{insert}} \in \mathcal{O}\left(n \log_2 n\right)$

| 3332 | 3354 | 3380 | 4347 | 4877 | 5610 | 6376 | 6977 | 7785 |
|------|------|------|------|------|------|------|------|------|
| → 7 | → 8 | → 3 | → 2 | → 6 | → 4 | → 5 | → 1 | → 0 |

Figure: A sorted table used to map page IDs to buffer frames.

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling      8 of 46

# Sorted & Chained Table

## Sorted Table

▶ Auxiliary data structure using a table sorted by page ID only containing cached pages

▶ $T_{\text{avg}}^{\text{search}} \in \mathcal{O}\left(\log_2 n\right)$, $T_{\text{avg}}^{\text{insert}} \in \mathcal{O}\left(n \log_2 n\right)$

| 3332 | 3354 | 3380 | 4347 | 4877 | 5610 | 6376 | 6977 | 7785 |
|------|------|------|------|------|------|------|------|------|
| → 7  | → 8  | → 3  | → 2  | → 6  | → 4  | → 5  | → 1  | → 0  |

Figure: A sorted table used to map page IDs to buffer frames.

## Chained Table

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                          8 of 46

# Sorted & Chained Table

### Sorted Table
▶ Auxiliary data structure using a table sorted by page ID only containing cached pages

▶ $T_{avg}^{search} \in \mathcal{O}\left(\log_2 n\right)$, $T_{avg}^{insert} \in \mathcal{O}\left(n\log_2 n\right)$

| 3332 | 3354 | 3380 | 4347 | 4877 | 5610 | 6376 | 6977 | 7785 |
|------|------|------|------|------|------|------|------|------|
| → 7 | → 8 | → 3 | → 2 | → 6 | → 4 | → 5 | → 1 | → 0 |

Figure: A sorted table used to map page IDs to buffer frames.

### Chained Table

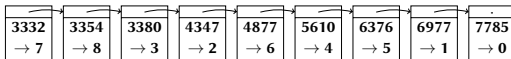▶ Auxiliary data structure using a linked list sorted by page ID only containing cached pages

| 3332 | 3354 | 3380 | 4347 | 4877 | 5610 | 6376 | 6977 | 7785 |
|------|------|------|------|------|------|------|------|------|
| → 7 | → 8 | → 3 | → 2 | → 6 | → 4 | → 5 | → 1 | → 0 |

Figure: A chained table used to map page IDs to buffer frames.

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                    8 of 46

# Sorted & Chained Table

### Sorted Table
- ▶ Auxiliary data structure using a table sorted by page ID only containing cached pages

- ▶ $T_{\text{avg}}^{\text{search}} \in \mathcal{O}\left(\log_2 n\right)$, $T_{\text{avg}}^{\text{insert}} \in \mathcal{O}\left(n \log_2 n\right)$

| 3332 | 3354 | 3380 | 4347 | 4877 | 5610 | 6376 | 6977 | 7785 |
|------|------|------|------|------|------|------|------|------|
| → 7 | → 8 | → 3 | → 2 | → 6 | → 4 | → 5 | → 1 | → 0 |

Figure: A sorted table used to map page IDs to buffer frames.

### Chained Table
- ▶ Auxiliary data structure using a linked list sorted by page ID only containing cached pages

- ▶ $T_{\text{avg}}^{\text{search}} \in \mathcal{O}\left(\log_2 n\right)$, $T_{\text{avg}}^{\text{insert}} \in \mathcal{O}\left(\log_2 n\right)$
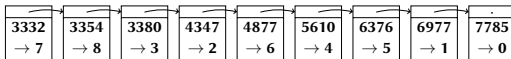
| 3332 | 3354 | 3380 | 4347 | 4877 | 5610 | 6376 | 6977 | 7785 |
|------|------|------|------|------|------|------|------|------|
| → 7 | → 8 | → 3 | → 2 | → 6 | → 4 | → 5 | → 1 | → 0 |

Figure: A chained table used to map page IDs to buffer frames.

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling    8 of 46

# Sorted & Chained Table

### Sorted Table

▶ Auxiliary data structure using a table sorted by page ID only containing cached pages

▶ $T_{\text{avg}}^{\text{search}} \in \mathcal{O}\left(\log_2 n\right)$, $T_{\text{avg}}^{\text{insert}} \in \mathcal{O}\left(n \log_2 n\right)$

| 3332 | 3354 | 3380 | 4347 | 4877 | 5610 | 6376 | 6977 | 7785 |
|------|------|------|------|------|------|------|------|------|
| → 7 | → 8 | → 3 | → 2 | → 6 | → 4 | → 5 | → 1 | → 0 |

Figure: A sorted table used to map page IDs to buffer frames.

### Chained Table

▶ Auxiliary data structure using a linked list sorted by page ID only containing cached pages

▶ $T_{\text{avg}}^{\text{search}} \in \mathcal{O}\left(\log_2 n\right)$, $T_{\text{avg}}^{\text{insert}} \in \mathcal{O}\left(\log_2 n\right)$
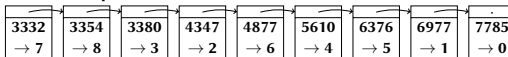
▶ Binary search requires more links!

| 3332 | 3354 | 3380 | 4347 | 4877 | 5610 | 6376 | 6977 | 7785 |
|------|------|------|------|------|------|------|------|------|
| → 7 | → 8 | → 3 | → 2 | → 6 | → 4 | → 5 | → 1 | → 0 |

Figure: A chained table used to map page IDs to buffer frames.

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling    9 of 46

# Search Trees

Pointer Swizzling in the DBMS Buffer Management   Performance Evaluation of Pointer Swizzling   Page Eviction Strategies   End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                          9 of 46
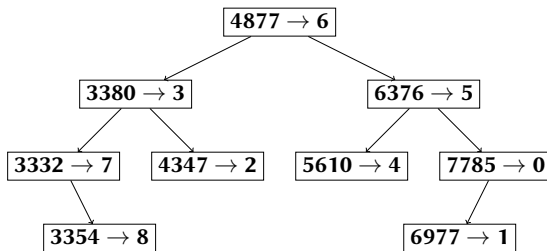
## Search Trees

▶ Auxiliary data structure is similar to the one of the chained table



Figure: A balanced search tree used to map page IDs to buffer frames.

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                    9 of 46
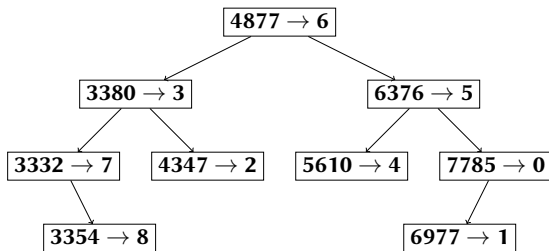
## Search Trees

▶ Auxiliary data structure is similar to the one of the chained table

▶ Many different data structures like AVL-trees, red–black trees or splay trees can be used



Figure: A balanced search tree used to map page IDs to buffer frames.

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     Page Eviction Strategies     End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                                9 of 46

## Search Trees

- ▶ Auxiliary data structure is similar to the one of the chained table
- ▶ Many different data structures like AVL-trees, red–black trees or splay trees can be used
- ▶ $T_{avg}^{search} \in \mathcal{O}(\log n)$, $T_{avg}^{insert} \in \mathcal{O}(\log n)$



Figure: A balanced search tree used to map page IDs to buffer frames.

Pointer Swizzling in the DBMS Buffer Management | Performance Evaluation of Pointer Swizzling | Page Eviction Strategies | End

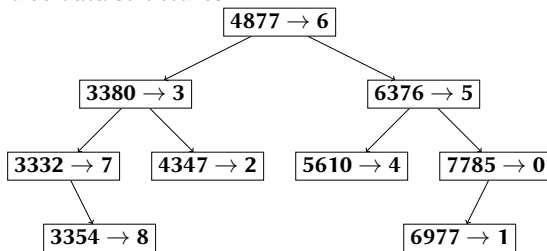Locate Pages in the Buffer Pool without Pointer Swizzling | 9 of 46

## Search Trees

▶ Auxiliary data structure is similar to the one of the chained table

▶ Many different data structures like AVL-trees, red–black trees or splay trees can be used

▶ $T_{\mathrm{avg}}^{\mathrm{search}} \in \mathcal{O}\left(\log n\right)$, $T_{\mathrm{avg}}^{\mathrm{insert}} \in \mathcal{O}\left(\log n\right)$

▶ The worst case costs and the worst cases vary between the different search tree data structures



Figure: A balanced search tree used to map page IDs to buffer frames.

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                    10 of 46

# Hash Table

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling    10 of 46

## Hash Table

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling        10 of 46

## Hash Table



▶ Each page ID is mapped to a hash bucket using a hash function

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                                      10 of 46

## Hash Table



- ► Each page ID is mapped to a hash bucket using a hash function
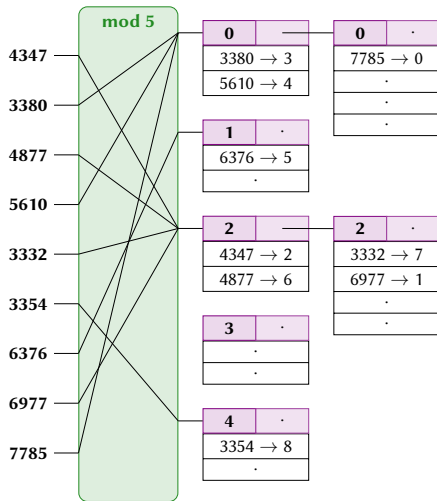- ► Only the page IDs of buffered pages are in the hash table

## Hash Table



- ▶ Each page ID is mapped to a hash bucket using a hash function
- ▶ Only the page IDs of buffered pages are in the hash table
- ▶ If a hash bucket is full, a chained bucket gets added

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling    10 of 46

## Hash Table



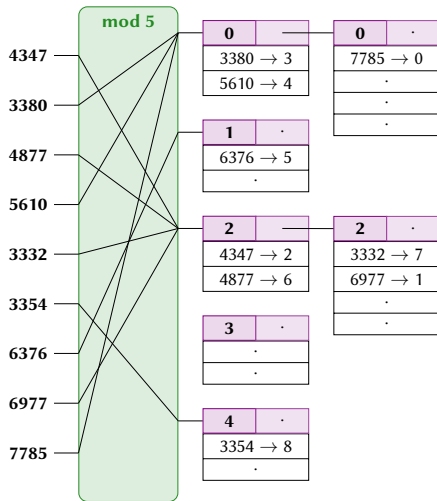- ▶ Each page ID is mapped to a hash bucket using a hash function
- ▶ Only the page IDs of buffered pages are in the hash table
- ▶ If a hash bucket is full, a chained bucket gets added
- ▶ $T_{avg}^{search} \in \mathcal{O}(1)$, $T_{avg}^{insert} \in \mathcal{O}(1)$, $T_{worst}^{search} \in \mathcal{O}(n)$

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                                11 of 46

# Locate Pages in Buffer Pool with Hash Table ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                          11 of 46

# Locate Pages in Buffer Pool with Hash Table ([Gra+14])

Buffer pool
page image

Pointer Swizzling in the DBMS Buffer Management  Performance Evaluation of Pointer Swizzling  Page Eviction Strategies  End

Locate Pages in the Buffer Pool without Pointer Swizzling  11 of 46

# Locate Pages in Buffer Pool with Hash Table ([Gra+14])

Buffer pool
page image

Search key

Pointer Swizzling in the DBMS Buffer Management        Performance Evaluation of Pointer Swizzling        Page Eviction Strategies        End

Locate Pages in the Buffer Pool without Pointer Swizzling                                                                11 of 46

# Locate Pages in Buffer Pool with Hash Table ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management | Performance Evaluation of Pointer Swizzling | Page Eviction Strategies | End

Locate Pages in the Buffer Pool without Pointer Swizzling | 11 of 46

# Locate Pages in Buffer Pool with Hash Table ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling                                    11 of 46

# Locate Pages in Buffer Pool with Hash Table ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling       11 of 46

# Locate Pages in Buffer Pool with Hash Table ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management | Performance Evaluation of Pointer Swizzling | Page Eviction Strategies | End

Locate Pages in the Buffer Pool without Pointer Swizzling | 11 of 46

# Locate Pages in Buffer Pool with Hash Table ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool without Pointer Swizzling    11 of 46

# Locate Pages in Buffer Pool with Hash Table ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management | Performance Evaluation of Pointer Swizzling | Page Eviction Strategies | End

Locate Pages in the Buffer Pool without Pointer Swizzling | 11 of 46

# Locate Pages in Buffer Pool with Hash Table ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management | Performance Evaluation of Pointer Swizzling | Page Eviction Strategies | End

Locate Pages in the Buffer Pool without Pointer Swizzling | 11 of 46

# Locate Pages in Buffer Pool with Hash Table ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management | Performance Evaluation of Pointer Swizzling | Page Eviction Strategies | End

Locate Pages in the Buffer Pool without Pointer Swizzling | 11 of 46

# Locate Pages in Buffer Pool with Hash Table ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management | Performance Evaluation of Pointer Swizzling | Page Eviction Strategies | End

Locate Pages in the Buffer Pool without Pointer Swizzling | 11 of 46

# Locate Pages in Buffer Pool with Hash Table ([Gra+14])

Subsection 2

Locate Pages in the Buffer Pool with Pointer Swizzling

## Pointer Swizzling

### Definition

To swizzle a pointer means to transform the address of the persistent object referenced there to a more direct address of the transient object in a way that this transformation could be used during multiple indirections of this pointer ([Mos92]).

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool with Pointer Swizzling                                                            14 of 46

# Classification of the Pointer Swizzling Approach following [WD95]

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool with Pointer Swizzling                                                          14 of 46

# Classification of the Pointer Swizzling Approach following [WD95]

no-swizzling

**swizzling**

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     Page Eviction Strategies     End

Locate Pages in the Buffer Pool with Pointer Swizzling                                                                14 of 46

# Classification of the Pointer Swizzling Approach following [WD95]

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool with Pointer Swizzling                                                                    14 of 46

# Classification of the Pointer Swizzling Approach following [WD95]

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool with Pointer Swizzling                                                                    14 of 46

# Classification of the Pointer Swizzling Approach following [WD95]

Pointer Swizzling in the DBMS Buffer Management | Performance Evaluation of Pointer Swizzling | Page Eviction Strategies | End

Locate Pages in the Buffer Pool with Pointer Swizzling | 14 of 46

# Classification of the Pointer Swizzling Approach following [WD95]

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool with Pointer Swizzling                                                    14 of 46

# Classification of the Pointer Swizzling Approach following [WD95]

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool with Pointer Swizzling    14 of 46

# Classification of the Pointer Swizzling Approach following [WD95]

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool with Pointer Swizzling          15 of 46

# Locate Pages in Buffer Pool w/ Pointer Swizzling ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool with Pointer Swizzling    15 of 46

# Locate Pages in Buffer Pool w/ Pointer Swizzling ([Gra+14])

Buffer pool
page image

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool with Pointer Swizzling      15 of 46

# Locate Pages in Buffer Pool w/ Pointer Swizzling ([Gra+14])

Buffer pool
page image

Search key

# Locate Pages in Buffer Pool w/ Pointer Swizzling ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Locate Pages in the Buffer Pool with Pointer Swizzling    15 of 46

# Locate Pages in Buffer Pool w/ Pointer Swizzling ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management | Performance Evaluation of Pointer Swizzling | Page Eviction Strategies | End

Locate Pages in the Buffer Pool with Pointer Swizzling | 15 of 46

# Locate Pages in Buffer Pool w/ Pointer Swizzling ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management | Performance Evaluation of Pointer Swizzling | Page Eviction Strategies | End

Locate Pages in the Buffer Pool with Pointer Swizzling | 15 of 46

# Locate Pages in Buffer Pool w/ Pointer Swizzling ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management | Performance Evaluation of Pointer Swizzling | Page Eviction Strategies | End

Locate Pages in the Buffer Pool with Pointer Swizzling | 15 of 46

# Locate Pages in Buffer Pool w/ Pointer Swizzling ([Gra+14])

# Locate Pages in Buffer Pool w/ Pointer Swizzling ([Gra+14])

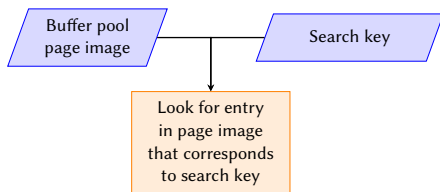Pointer Swizzling in the DBMS Buffer Management | Performance Evaluation of Pointer Swizzling | Page Eviction Strategies | End

Locate Pages in the Buffer Pool with Pointer Swizzling | 15 of 46

# Locate Pages in Buffer Pool w/ Pointer Swizzling ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management | Performance Evaluation of Pointer Swizzling | Page Eviction Strategies | End

Locate Pages in the Buffer Pool with Pointer Swizzling | 15 of 46

# Locate Pages in Buffer Pool w/ Pointer Swizzling ([Gra+14])

Pointer Swizzling in the DBMS Buffer Management · Performance Evaluation of Pointer Swizzling · Page Eviction Strategies · End

16 of 46

# Section 2

# Performance Evaluation of the Buffer Management Utilizing Pointer Swizzling

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Expected Performance                                                                                    17 of 46

Subsection 1

Expected Performance

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Expected Performance                                                                                              18 of 46

## Performance of Different Buffer Pool Sizes

Pointer Swizzling in the DBMS Buffer Management     **Performance Evaluation of Pointer Swizzling**     Page Eviction Strategies     End

Expected Performance                                                                                                          18 of 46

## Performance of Different Buffer Pool Sizes

Pointer Swizzling in the DBMS Buffer Management **Performance Evaluation of Pointer Swizzling** Page Eviction Strategies End

Expected Performance 18 of 46

# Performance of Different Buffer Pool Sizes

Pointer Swizzling in the DBMS Buffer Management  **Performance Evaluation of Pointer Swizzling**  Page Eviction Strategies  End

Expected Performance                                                                                                    18 of 46

# Performance of Different Buffer Pool Sizes



[EH84][HR01]

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Expected Performance                                                                                                    19 of 46

# Buffer Management with and without Pointer Swizzling

| Page Hits without Pointer Swizzling |
| Page Misses without Pointer Swizzling |
| Page Hits with Pointer Swizzling |
| Page Misses with Pointer Swizzling |

total execution time

$HR_{min}$                                                                                                    $HR_{CS}$

Pointer Swizzling in the DBMS Buffer Management  **Performance Evaluation of Pointer Swizzling**  Page Eviction Strategies  End

Expected Performance  19 of 46

# Buffer Management with and without Pointer Swizzling

Pointer Swizzling in the DBMS Buffer Management   **Performance Evaluation of Pointer Swizzling**   Page Eviction Strategies   End

Expected Performance                                                                                    19 of 46

# Buffer Management with and without Pointer Swizzling

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Expected Performance    19 of 46

# Buffer Management with and without Pointer Swizzling



total execution time

- Page Hits without Pointer Swizzling
- Page Misses without Pointer Swizzling
- Page Hits with Pointer Swizzling
- Page Misses with Pointer Swizzling

$HR_{min}$                                    $HR_{CS}$

Pointer Swizzling in the DBMS Buffer Management  **Performance Evaluation of Pointer Swizzling**  Page Eviction Strategies  End

Expected Performance  19 of 46

# Buffer Management with and without Pointer Swizzling

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Expected Performance                                                                                    19 of 46

# Buffer Management with and without Pointer Swizzling

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Measured Performance    20 of 46

Subsection 2

Measured Performance

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Measured Performance                                                                                                                21 of 46

# Transaction Throughput (TPC-C)



Max Gilbert                                                                                                              University of Kaiserslautern

Pointer Swizzling in the DBMS Buffer Management    **Performance Evaluation of Pointer Swizzling**    Page Eviction Strategies    End

Measured Performance      21 of 46

# Transaction Throughput (TPC-C)

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Measured Performance    21 of 46

# Transaction Throughput (TPC-C)

Pointer Swizzling in the DBMS Buffer Management   Performance Evaluation of Pointer Swizzling   Page Eviction Strategies   End

Measured Performance                                                                                     22 of 46

# Transaction Throughput (TPC-B)



Max Gilbert                                                                                University of Kaiserslautern

Evaluation of Pointer Swizzling Techniques for DBMS Buffer Management

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Measured Performance                                                                          22 of 46

# Transaction Throughput (TPC-B)

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Measured Performance                                                                                                22 of 46

# Transaction Throughput (TPC-B)

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Measured Performance                                                                                      23 of 46

# Buffer Pool Performance Acquiring Shared Latches



| Traditional Buffer Pool | Pointer Swizzling Buffer Pool |

500 MiB (w/ eviction)                                    20 GiB (w/o eviction)

Pointer Swizzling in the DBMS Buffer Management    **Performance Evaluation of Pointer Swizzling**    Page Eviction Strategies    End

Measured Performance                                                                                                              23 of 46

# Buffer Pool Performance Acquiring Shared Latches



500 MiB (w/ eviction)

20 GiB (w/o eviction)

Pointer Swizzling in the DBMS Buffer Management    **Performance Evaluation of Pointer Swizzling**    Page Eviction Strategies    End

Measured Performance                                                                                      23 of 46

# Buffer Pool Performance Acquiring Shared Latches



500 MiB (w/ eviction)

20 GiB (w/o eviction)

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Measured Performance                                                                                      23 of 46

# Buffer Pool Performance Acquiring Shared Latches



500 MiB (w/ eviction)

20 GiB (w/o eviction)

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Measured Performance                                                                                                23 of 46

# Buffer Pool Performance Acquiring Shared Latches



500 MiB (w/ eviction)                    20 GiB (w/o eviction)

Pointer Swizzling in the DBMS Buffer Management  **Performance Evaluation of Pointer Swizzling**  Page Eviction Strategies  End

Conclusion  24 of 46

Subsection 3

Conclusion

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Conclusion    25 of 46

# Conclusion

Pointer Swizzling in the DBMS Buffer Management    **Performance Evaluation of Pointer Swizzling**    Page Eviction Strategies    End

Conclusion      25 of 46

# Conclusion

## Overall Performance

## Conclusion

### Overall Performance

▶ Pointer swizzling couldn't improve the performance on TPC-C benchmark runs with a duration of 10 min.

## Conclusion

### Overall Performance

▶ Pointer swizzling couldn't improve the performance on TPC-C
  benchmark runs with a duration of 10 min.

▶ The page hits after the cold start couldn't compensate the
  overhead of pointer swizzling during the cold start.

## Conclusion

### Overall Performance

- ▶ Pointer swizzling couldn't improve the performance on TPC-C benchmark runs with a duration of 10 min.
- ▶ The page hits after the cold start couldn't compensate the overhead of pointer swizzling during the cold start.
- ▶ A continuously running DB with large buffer pool could profit from pointer swizzling.

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Conclusion      25 of 46

## Conclusion

### Overall Performance

- ▶ Pointer swizzling couldn't improve the performance on TPC-C benchmark runs with a duration of 10 min.
- ▶ The page hits after the cold start couldn't compensate the overhead of pointer swizzling during the cold start.
- ▶ A continuously running DB with large buffer pool could profit from pointer swizzling.

### Buffer Pool Performance

Pointer Swizzling in the DBMS Buffer Management   Performance Evaluation of Pointer Swizzling   Page Eviction Strategies   End

Conclusion                                                                                                25 of 46

## Conclusion

### Overall Performance

▶ Pointer swizzling couldn't improve the performance on TPC-C benchmark runs with a duration of 10 min.

▶ The page hits after the cold start couldn't compensate the overhead of pointer swizzling during the cold start.

▶ A continuously running DB with large buffer pool could profit from pointer swizzling.

### Buffer Pool Performance

▶ A page hit is faster when pointer swizzling is activated.

Pointer Swizzling in the DBMS Buffer Management   Performance Evaluation of Pointer Swizzling   Page Eviction Strategies   End

Conclusion          25 of 46

# Conclusion

### Overall Performance

- ▶ Pointer swizzling couldn't improve the performance on TPC-C benchmark runs with a duration of 10 min.
- ▶ The page hits after the cold start couldn't compensate the overhead of pointer swizzling during the cold start.
- ▶ A continuously running DB with large buffer pool could profit from pointer swizzling.

### Buffer Pool Performance

- ▶ A page hit is faster when pointer swizzling is activated.
- ▶ A page miss is slower when pointer swizzling is activated.

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Conclusion                                                                                    25 of 46

# Conclusion

## Overall Performance

- ▶ Pointer swizzling couldn't improve the performance on TPC-C benchmark runs with a duration of 10 min.

- ▶ The page hits after the cold start couldn't compensate the overhead of pointer swizzling during the cold start.

- ▶ A continuously running DB with large buffer pool could profit from pointer swizzling.

## Buffer Pool Performance

- ▶ A page hit is faster when pointer swizzling is activated.

- ▶ A page miss is slower when pointer swizzling is activated.

- ▶ After the cold start phase, activated pointer swizzling will improve the buffer pool performance for large buffer pools.

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

26 of 46

# Section 3

# Page Eviction Strategies in the Context of Pointer Swizzling

Pointer Swizzling in the DBMS Buffer Management  Performance Evaluation of Pointer Swizzling  Page Eviction Strategies  End

27 of 46

# Motivation not to Analyze Different Page Eviction Strategies

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

27 of 46

# Motivation <u>not</u> to Analyze Different Page Eviction Strategies

- ▶ Even LRU results in decent hit rates

TPC-C with Warehouses: 100, Threads: 25

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     **Page Eviction Strategies**     End

28 of 46

# But ...

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

28 of 46

## But ...

- ▶ Page reference pattern containing a loop slightly to long to fit in the buffer pool:

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

28 of 46

## But ...

- ▶ Page reference pattern containing a loop slightly to long to fit in the buffer pool:
    - ▶ **OPT:** Hit rate close to 1

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

28 of 46

# But ...

- ▶ Page reference pattern containing a loop slightly to long to fit in the buffer pool:
    - ▶ **OPT:** Hit rate close to 1
    - ▶ **LRU:** Hit rate of 0

Pointer Swizzling in the DBMS Buffer Management   Performance Evaluation of Pointer Swizzling   **Page Eviction Strategies**   End

28 of 46

## But ...

- ▶ Page reference pattern containing a loop slightly to long to fit in the buffer pool:
  - ▶ **OPT:** Hit rate close to 1
  - ▶ **LRU:** Hit rate of 0
- ▶ Some pages gets referenced very frequent for a limited time:

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

28 of 46

## But ...

- ▶ Page reference pattern containing a loop slightly to long to fit in the buffer pool:
    - ▶ **OPT:** Hit rate close to 1
    - ▶ **LRU:** Hit rate of 0
- ▶ Some pages gets referenced very frequent for a limited time:
    - ▶ **OPT:** Pages would be evicted after their last reference

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     **Page Eviction Strategies**     End

28 of 46

# But ...

- ▶ Page reference pattern containing a loop slightly to long to fit in the buffer pool:
    - ▶ **OPT:** Hit rate close to 1
    - ▶ **LRU:** Hit rate of 0
- ▶ Some pages gets referenced very frequent for a limited time:
    - ▶ **OPT:** Pages would be evicted after their last reference
    - ▶ **LFU:** Pages waste buffer frames probably during the whole running time of the DB

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

28 of 46

## But ...

- ▶ Page reference pattern containing a loop slightly to long to fit in the buffer pool:
  - ▶ **OPT:** Hit rate close to 1
  - ▶ **LRU:** Hit rate of 0
- ▶ Some pages gets referenced very frequent for a limited time:
  - ▶ **OPT:** Pages would be evicted after their last reference
  - ▶ **LFU:** Pages waste buffer frames probably during the whole running time of the DB
- ▶ Huge access time gap $\implies$ Every saved page miss significantly improves the performance

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     **Page Eviction Strategies**     End

28 of 46

## But ...

- ▶ Page reference pattern containing a loop slightly to long to fit in the buffer pool:
    - ▶ **OPT:** Hit rate close to 1
    - ▶ **LRU:** Hit rate of 0
- ▶ Some pages gets referenced very frequent for a limited time:
    - ▶ **OPT:** Pages would be evicted after their last reference
    - ▶ **LFU:** Pages waste buffer frames probably during the whole running time of the DB
- ▶ Huge access time gap $\implies$ Every saved page miss significantly improves the performance
- ▶ Pointer swizzling even amplifies that effect

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Probable pitfalls when Implementing a Page Eviction Strategy for a DBMS Buffer Manager                                    29 of 46

Subsection 1

Probable pitfalls when Implementing a Page Eviction
Strategy for a DBMS Buffer Manager

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     **Page Eviction Strategies**     End

Probable pitfalls when Implementing a Page Eviction Strategy for a DBMS Buffer Manager     30 of 46

# General Problems Concerning DBMS Buffer Managers

# General Problems Concerning DBMS Buffer Managers

▶ Fixed pages cannot be evicted but a long timespan between a fix and an unfix of a page could make it a candidate for eviction.

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Probable pitfalls when Implementing a Page Eviction Strategy for a DBMS Buffer Manager    30 of 46

# General Problems Concerning DBMS Buffer Managers

▶ Fixed pages cannot be evicted but a long timespan between a fix and an unfix of a page could make it a candidate for eviction.

▶ A page pinned for refix cannot be evicted but a long timespan in which a page is pinned could make it a candidate for eviction.

## General Problems Concerning DBMS Buffer Managers

▶ Fixed pages cannot be evicted but a long timespan between a fix and an unfix of a page could make it a candidate for eviction.

▶ A page pinned for refix cannot be evicted but a long timespan in which a page is pinned could make it a candidate for eviction.

▶ Dirty pages cannot be evicted but a page being dirty for a long timespan due to the update propagation using write-back policy could make it a candidate for eviction.

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     Page Eviction Strategies     End

Probable pitfalls when Implementing a Page Eviction Strategy for a DBMS Buffer Manager     31 of 46

# Additional Problem When Using Pointer Swizzling

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Probable pitfalls when Implementing a Page Eviction Strategy for a DBMS Buffer Manager    31 of 46

# Additional Problem When Using Pointer Swizzling

▶ A page containing swizzled pointer cannot be evicted but a page unfixed before the last unfix of one of its child pages could make it a candidate for eviction before its child pages got evicted.

# Solutions

## Solutions

▶ Check each of the restrictions before the eviction of a page.

## Solutions

- ▶ Check each of the restrictions before the eviction of a page.
- ▶ Update the statistics of the eviction strategy during an unfix, too.

## Solutions

- ▶ Check each of the restrictions before the eviction of a page.
- ▶ Update the statistics of the eviction strategy during an unfix, too.
- ▶ Update the statistics of the eviction strategy during an pin and unpin, too.

## Solutions

- ▶ Check each of the restrictions before the eviction of a page.
- ▶ Update the statistics of the eviction strategy during an unfix, too.
- ▶ Update the statistics of the eviction strategy during an pin and unpin, too.
- ▶ Use write-thru for update propagation or a page cleaner decoupled from the buffer pool as proposed in [SHG16].

## Solutions

- ▶ Check each of the restrictions before the eviction of a page.
- ▶ Update the statistics of the eviction strategy during an unfix, too.
- ▶ Update the statistics of the eviction strategy during an pin and unpin, too.
- ▶ Use write-thru for update propagation or a page cleaner decoupled from the buffer pool as proposed in [SHG16].
- ▶ Use a page eviction strategy that takes into account the content of pages (like the structure of an B tree).

Pointer Swizzling in the DBMS Buffer Management | Performance Evaluation of Pointer Swizzling | **Page Eviction Strategies** | End

Evaluated Page Replacement Strategies | 33 of 46

Subsection 2

Evaluated Page Replacement Strategies

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Evaluated Page Replacement Strategies                                                                34 of 46

# RANDOM

Overview

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Evaluated Page Replacement Strategies                                                                        34 of 46

# RANDOM

### Overview

▶ Simplest page eviction strategy

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     **Page Eviction Strategies**     End

Evaluated Page Replacement Strategies          34 of 46

# RANDOM

### Overview

▶ Simplest page eviction strategy

▶ Evicts a random page that can be evicted

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Evaluated Page Replacement Strategies                                                                34 of 46

# RANDOM

### Overview

▶ Simplest page eviction strategy

▶ Evicts a random page that can be evicted

▶ Won't evict frequently used pages as they're latched all the time

# GCLOCK

Overview

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     **Page Eviction Strategies**    End

Evaluated Page Replacement Strategies     35 of 46

# GCLOCK

### Overview

▶ Slight enhancement of the CLOCK algorithm: *generalized CLOCK*

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     **Page Eviction Strategies**     End

Evaluated Page Replacement Strategies                                 35 of 46

# GCLOCK

### Overview

- ▶ Slight enhancement of the CLOCK algorithm: *generalized CLOCK*
- ▶ Uses finer-grained statistics about the recency of page references

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Evaluated Page Replacement Strategies                                                                                                    35 of 46

# GCLOCK

### Overview

- ▶ Slight enhancement of the CLOCK algorithm: *generalized CLOCK*
- ▶ Uses finer-grained statistics about the recency of page references
- ▶ Parameter $k$ defines granulation of statistics

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Evaluated Page Replacement Strategies                                                            35 of 46

# GCLOCK

### Overview

▶ Slight enhancement of the CLOCK algorithm: *generalized CLOCK*

▶ Uses finer-grained statistics about the recency of page references

▶ Parameter $k$ defines granulation of statistics
  ▶ $k = 1$: CLOCK

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     **Page Eviction Strategies**     End

Evaluated Page Replacement Strategies          35 of 46

# GCLOCK

### Overview

- ▶ Slight enhancement of the CLOCK algorithm: *generalized CLOCK*
- ▶ Uses finer-grained statistics about the recency of page references
- ▶ Parameter $k$ defines granulation of statistics
  - ▶ $k = 1$: CLOCK
  - ▶ $k = \#$**frames**: Similar to LRU

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Evaluated Page Replacement Strategies                                                                                36 of 46

## GCLOCK

### Example

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Evaluated Page Replacement Strategies                                                                    37 of 46

# GCLOCK

Advantage of Higher $k$-Values

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Evaluated Page Replacement Strategies                                                                37 of 46

# GCLOCK

### Advantage of Higher $k$-Values

### Advantages of Lower $k$-Values

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Evaluated Page Replacement Strategies                                                                                    37 of 46

# GCLOCK

### Advantage of Higher *k*-Values

- ▶ More detailed statistics about page references
  - $\implies$ Higher hit rate
  - $\implies$ Higher performance

### Advantages of Lower *k*-Values

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Evaluated Page Replacement Strategies                                                                    37 of 46

# GCLOCK

### Advantage of Higher $k$-Values

- ▶ More detailed statistics about page references
  - $\implies$ Higher hit rate
  - $\implies$ Higher performance

### Advantages of Lower $k$-Values

- ▶ Lower processing time required to find an eviction victim
  - $\implies$ Higher performance

# GCLOCK

### Advantage of Higher *k*-Values

- ▶ More detailed statistics about page references
  - $\implies$ Higher hit rate
  - $\implies$ Higher performance

### Advantages of Lower *k*-Values

- ▶ Lower processing time required to find an eviction victim
  - $\implies$ Higher performance
- ▶ Lower memory overhead due to shorter referenced-numbers

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Evaluated Page Replacement Strategies    37 of 46

# GCLOCK

### Advantage of Higher $k$-Values

- More detailed statistics about page references
  $\implies$ Higher hit rate
  $\implies$ Higher performance

### Advantages of Lower $k$-Values

- Lower processing time required to find an eviction victim
  $\implies$ Higher performance
- Lower memory overhead due to shorter referenced-numbers

$\implies$ Trade-off between CPU- and I/O-optimization

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     **Page Eviction Strategies**     End

Evaluated Page Replacement Strategies        38 of 46

# CAR

Overview

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Evaluated Page Replacement Strategies        38 of 46

# CAR

### Overview

▶ Extensive enhancement of the CLOCK algorithm: *Clock with Adaptive Replacement* [BM04]

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     **Page Eviction Strategies**     End

Evaluated Page Replacement Strategies        38 of 46

# CAR

### Overview

▶ Extensive enhancement of the CLOCK algorithm: *Clock with Adaptive Replacement* [BM04]

▶ Approximation of the ARC page eviction strategy

Pointer Swizzling in the DBMS Buffer Management | Performance Evaluation of Pointer Swizzling | **Page Eviction Strategies** | End

Evaluated Page Replacement Strategies | 38 of 46

# CAR

### Overview

- ▶ Extensive enhancement of the CLOCK algorithm: *Clock with Adaptive Replacement* [BM04]
- ▶ Approximation of the ARC page eviction strategy
- ▶ Uses two clocks and two LRU-lists

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Evaluated Page Replacement Strategies                                                                                        38 of 46

# CAR

### Overview

- ▶ Extensive enhancement of the CLOCK algorithm: *Clock with Adaptive Replacement* [BM04]
- ▶ Approximation of the ARC page eviction strategy
- ▶ Uses two clocks and two LRU-lists
- ▶ Advantages of CAR compared to CLOCK:

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Evaluated Page Replacement Strategies    38 of 46

# CAR

### Overview

- ▶ Extensive enhancement of the CLOCK algorithm: *Clock with Adaptive Replacement* [BM04]
- ▶ Approximation of the ARC page eviction strategy
- ▶ Uses two clocks and two LRU-lists
- ▶ Advantages of CAR compared to CLOCK:
  - ▶ Weighted consideration of reference recency and frequency

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Evaluated Page Replacement Strategies                                                                      38 of 46

# CAR

### Overview

- ▶ Extensive enhancement of the CLOCK algorithm: *Clock with Adaptive Replacement* [BM04]
- ▶ Approximation of the ARC page eviction strategy
- ▶ Uses two clocks and two LRU-lists
- ▶ Advantages of CAR compared to CLOCK:
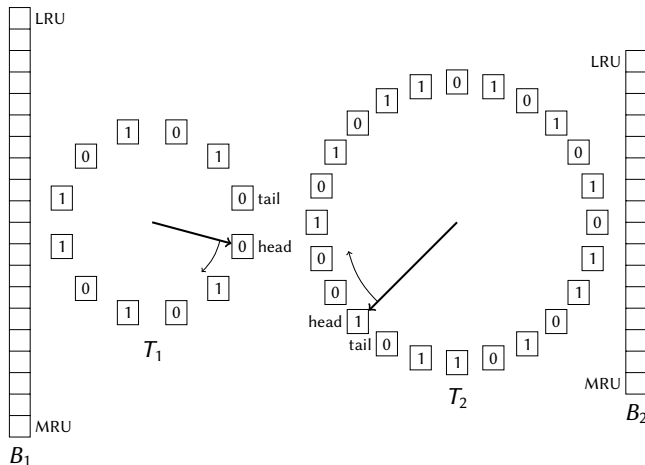  - ▶ Weighted consideration of reference recency and frequency
  - ▶ Scan-resistence

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Evaluated Page Replacement Strategies                                                                                              39 of 46

# CAR

### Example

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Performance Evaluation                                                                                              40 of 46

Subsection 3

Performance Evaluation

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End
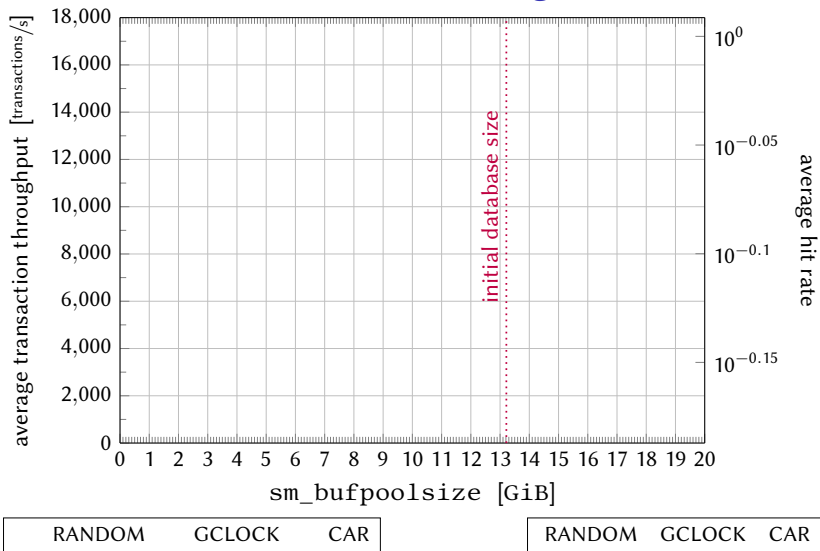
Performance Evaluation                                                                    41 of 46

# Buffer Pool Without Pointer Swizzling (TPC-C)



RANDOM    GCLOCK    CAR          RANDOM    GCLOCK    CAR

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Performance Evaluation                                                                                            41 of 46

# Buffer Pool Without Pointer Swizzling (TPC-C)



| | RANDOM | GCLOCK | CAR | | RANDOM | GCLOCK | CAR |

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Performance Evaluation                                                                41 of 46

# Buffer Pool Without Pointer Swizzling (TPC-C)

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Performance Evaluation    41 of 46

# Buffer Pool Without Pointer Swizzling (TPC-C)

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Performance Evaluation                                                                41 of 46

# Buffer Pool Without Pointer Swizzling (TPC-C)

Pointer Swizzling in the DBMS Buffer Management       Performance Evaluation of Pointer Swizzling       **Page Eviction Strategies**       End

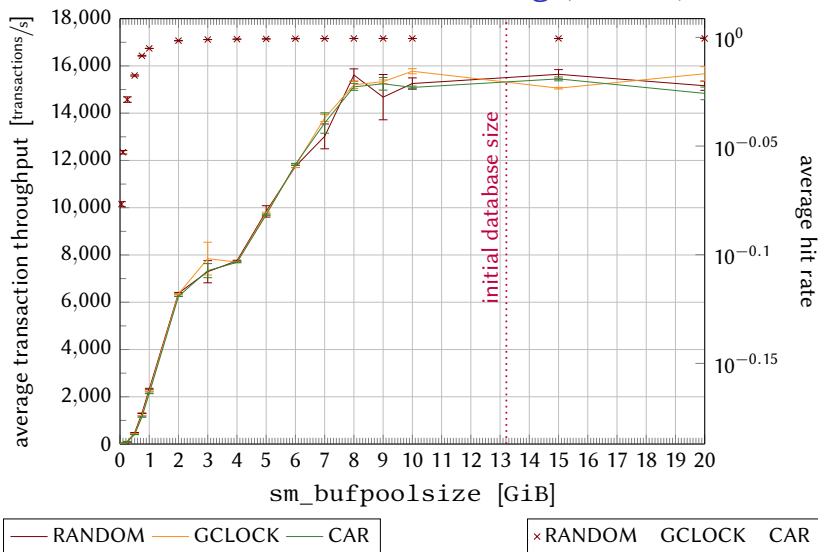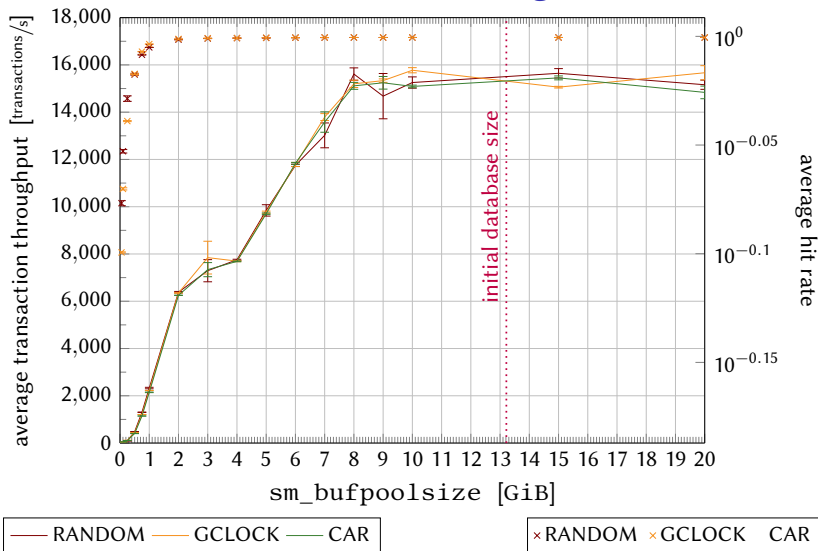Performance Evaluation                                                                                                41 of 46

# Buffer Pool Without Pointer Swizzling (TPC-C)

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     **Page Eviction Strategies**     End

Performance Evaluation                                                                                      41 of 46

# Buffer Pool Without Pointer Swizzling (TPC-C)

Pointer Swizzling in the DBMS Buffer Management   Performance Evaluation of Pointer Swizzling   **Page Eviction Strategies**   End

Performance Evaluation                                                                                    42 of 46

# Buffer Pool With Pointer Swizzling (TPC-C)

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Performance Evaluation                                                                                          42 of 46

# Buffer Pool With Pointer Swizzling (TPC-C)

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Performance Evaluation                                                                                          42 of 46

# Buffer Pool With Pointer Swizzling (TPC-C)

Pointer Swizzling in the DBMS Buffer Management　　　Performance Evaluation of Pointer Swizzling　　**Page Eviction Strategies**　　End

Performance Evaluation　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　42 of 46

# Buffer Pool With Pointer Swizzling (TPC-C)

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     **Page Eviction Strategies**     End

Performance Evaluation     42 of 46

# Buffer Pool With Pointer Swizzling (TPC-C)

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Performance Evaluation                                                                                42 of 46

# Buffer Pool With Pointer Swizzling (TPC-C)

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

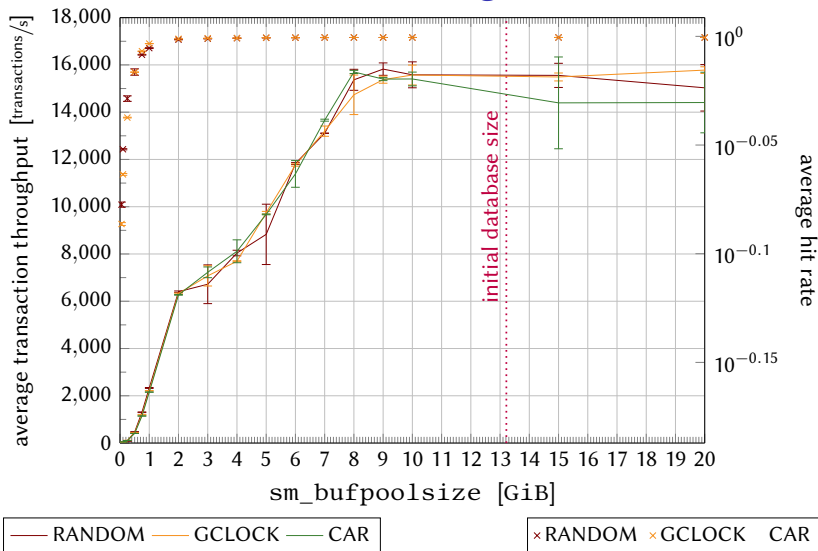Performance Evaluation                                                                                          42 of 46

# Buffer Pool With Pointer Swizzling (TPC-C)

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Performance Evaluation    43 of 46

## Operation Performance

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Performance Evaluation                                                                                                      43 of 46

## Operation Performance

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Performance Evaluation                                                                                      43 of 46

## Operation Performance

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     **Page Eviction Strategies**     End

Performance Evaluation     43 of 46

## Operation Performance

Pointer Swizzling in the DBMS Buffer Management · · · Performance Evaluation of Pointer Swizzling · · · **Page Eviction Strategies** · · · End

Performance Evaluation · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · 43 of 46

# Operation Performance

Max Gilbert · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · · University of Kaiserslautern

Evaluation of Pointer Swizzling Techniques for DBMS Buffer Management

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Conclusion                                                                                                    44 of 46

Subsection 4

Conclusion

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Conclusion    45 of 46

# Conclusion

## Performance

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Conclusion                                                                                                    45 of 46

## Conclusion

### Performance

▶ CAR has a significantly higher hit rate than RANDOM or GCLOCK

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Conclusion                                                                                                              45 of 46

# Conclusion

### Performance

▶ CAR has a significantly higher hit rate than RANDOM or GCLOCK

▶ The hit rate of GCLOCK isn't significantly higher than the one of RANDOM

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Conclusion                                                                                                          45 of 46

## Conclusion

### Performance

▶ CAR has a significantly higher hit rate than RANDOM or GCLOCK

▶ The hit rate of GCLOCK isn't significantly higher than the one of RANDOM

▶ Major differences in hit rate are only for buffer pool sizes of $\leq \frac{1}{10}$ of the database size

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    End

Conclusion                                                                                          45 of 46

# Conclusion

### Performance

- ▶ CAR has a significantly higher hit rate than RANDOM or GCLOCK
- ▶ The hit rate of GCLOCK isn't significantly higher than the one of RANDOM
- ▶ Major differences in hit rate are only for buffer pool sizes of $\leq \frac{1}{10}$ of the database size
- ▶ The computational effort spent to do CAR eviction is 27–58 times higher

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    **Page Eviction Strategies**    End

Conclusion                                                                                      45 of 46

# Conclusion

### Performance

- ▶ CAR has a significantly higher hit rate than RANDOM or GCLOCK
- ▶ The hit rate of GCLOCK isn't significantly higher than the one of RANDOM
- ▶ Major differences in hit rate are only for buffer pool sizes of $\leq \frac{1}{10}$ of the database size
- ▶ The computational effort spent to do CAR eviction is 27–58 times higher
- ▶ The overall performance of CAR isn't better than the one of RANDOM or GCLOCK

Pointer Swizzling in the DBMS Buffer Management    Performance Evaluation of Pointer Swizzling    Page Eviction Strategies    **End**

Conclusion                                                                                                  45 of 46

## References I

📄  Sorav Bansal and Dharmendra S. Modha. "CAR: Clock with Adaptive Replacement". Mar. 31, 2004.

📄  Wolfgang Effelsberg and Theo Härder. "Principles of Database Buffer Management". Dec. 1984.

📄  Goetz Graefe et al. "In-Memory Performance for Big Data". Sept. 2014.

📕  Theo Härder and Erhard Rahm. *Datenbanksysteme - Konzepte und Techniken der Implementierung.* 2001. ISBN: 978-3-642-62659-3.

📄  J. Eliot B. Moss. "Working with Persistent Objects: To Swizzle or Not to Swizzle". Aug. 1992.

Pointer Swizzling in the DBMS Buffer Management      Performance Evaluation of Pointer Swizzling      Page Eviction Strategies      End

Conclusion                                                                                                    46 of 46

## References II

📄    Lucas Sauer Caetano Lersch, Theo Härder, and Goetz Graefe.
      "Update propagation strategies for high-performance OLTP".
      Aug. 14, 2016.

📄    Seth J. White and David J. DeWitt. "QuickStore: A High
      Performance Mapped Object Store".   Oct. 1995.

Pointer Swizzling in the DBMS Buffer Management     Performance Evaluation of Pointer Swizzling     Page Eviction Strategies     End

47 of 46

# Your Turn to Ask ...