

# Database Buffer Management

Yanlei Diao  
UMass Amherst  
Feb 20, 2007

Slides Courtesy of R. Ramakrishnan and J. Gehrke

1

## DBMS vs. OS File System

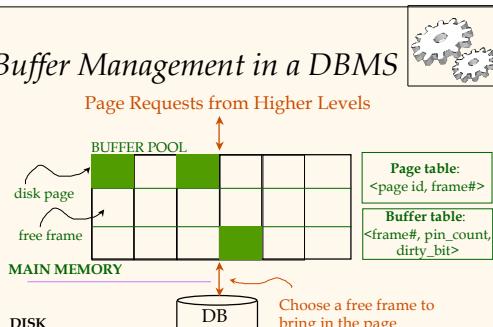
OS does buffer management: why not let OS manage it?

- ❖ Differences in OS support: portability issues
- ❖ Buffer mgmt in DBMS needs to work with transaction management.
  - pin a page in buffer pool, force a page to disk (important for implementing CC & recovery)
- ❖ Buffer mgmt in DBMS needs to understand query-specific data access patterns.
  - adjust replacement policy, and pre-fetch pages based on access patterns in typical DB operations.

2

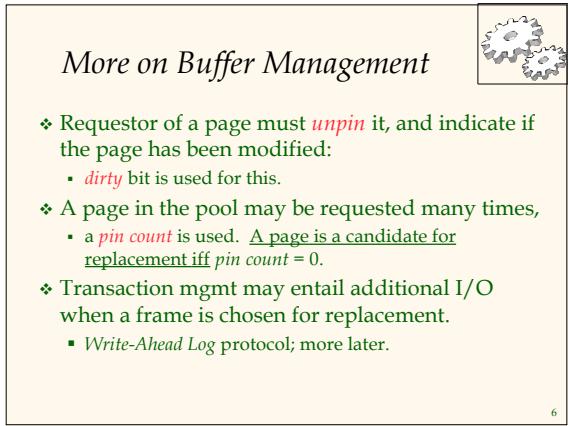
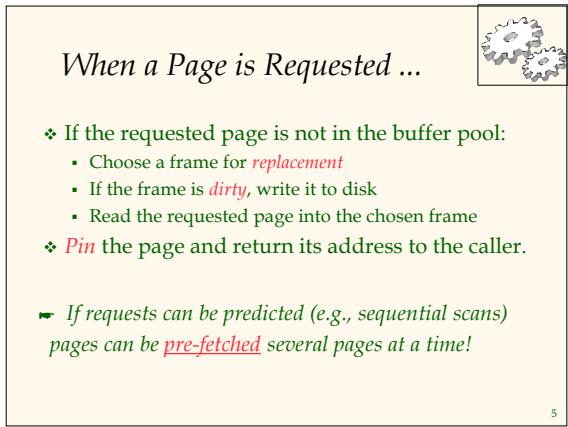
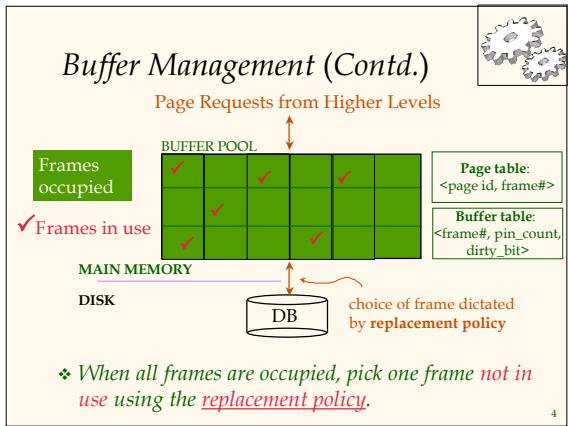
## Buffer Management in a DBMS

Page Requests from Higher Levels



- ❖ Data must be in RAM for DBMS to operate on it!

3



## Buffer Replacement Policies

- ❖ A frame is chosen for replacement by a *replacement policy*:
  - Least-recently-used (LRU)
  - Clock, a variant of LRU
  - Most-recently-used (MRU)
  - First-in-first-out (FIFO)
  - Last-in-first-out (LIFO)
  - LRU2
  - ...

7



---

---

---

---

---

## Impact of the Replacement Policy

- ❖ The policy can have a big impact on # of I/O's; depends on the *access pattern*.
- ❖ None of the stochastic policies is uniformly appropriate for typical DBMS access patterns.
- ❖ Sequential flooding: Nasty situation caused by LRU + repeated sequential scans.
  - # buffer frames < # pages in file means each page request causes an I/O.
  - MRU much better in this situation (but not in all situations, of course).

8



---

---

---

---

---

## Advanced Buffer Management

- ❖ Just a few basic *data access patterns* in relational DBs, with noticeable *locality behaviors*.
- ❖ A DB operation can be decomposed into a subset of these access patterns.
- ❖ To reduce I/Os, expose these patterns to the buffer manager for correct estimation of:
  - *buffer size*: many queries share the buffer pool; need to know how to allocate frames to each query.
  - *replacement policy*: evict unused pages to make room for newly requested pages.

9



---

---

---

---

---

**DBMIN: QLSM**



❖ Query Locality Set Model (QLSM)

<i>Sequential data access</i>	<i>Straight Sequential (SS)</i>	<i>Clustered Sequential (CS)</i>	<i>Looping Sequential (LS)</i>
<b>Behavior</b>	One access w/o repetition	Local rescan	Repetitive sequential scan of a file
<b>Examples</b>	File scan	Merge join with point of plurality	Nested loops join: inner
<b>Locality set</b>	1	Largest cluster (# pages with the same key value)	# pages of the file
<b>Replacement policy</b>	Any; pre-fetching is good	If enough memory, FIFO/LRU (evict old data)	If file doesn't fit, MRU

10

---



---



---



---



---



---



---



---



---



---

**DBMIN: QLSM (Contd.)**



<i>Random data access</i>	<i>Independent Random</i>	<i>Clustered Random</i>
<b>Behavior</b>	A series of independent accesses	A series of random accesses
<b>Examples</b>	Indexed scan: access to data pages through a non-clustered index	Indexed join: outer is a clustered file with duplicate join values, inner is through a non-clustered index.
<b>Locality set</b>	1 or b (# of pages that will be re-visited)	Largest cluster (maximum # records in the cluster)
<b>Replacement policy</b>	Any	Any

11

---



---



---



---



---



---



---



---



---



---

**DBMIN: QLSM (Contd.)**



<i>Hierarchical Index Access</i>	<i>Straight Hierarchical</i>	<i>Hierarchical w. SS</i>	<i>Hierarchical w. CS</i>	<i>Looping Hierarchical</i>
<b>Behavior</b>	Index is traversed once	Tree traversal+ scan of leaves	Local rescan in page access	Repeated access to the index structure
<b>Examples</b>	Equality search using an index	Range search using the index	Inner of indexed nested loops join, w.o. duplicate values from the outer	Inner of indexed nested loops, w.o. duplicate values from the outer
<b>Locality set</b>	1	1	= CS	Closer to root, often only sure about the root for high F, 3~4.
<b>Replacement policy</b>	Any	Any	If enough memory, FIFO or LRU	LIFO

12

---



---



---



---



---



---



---



---



---



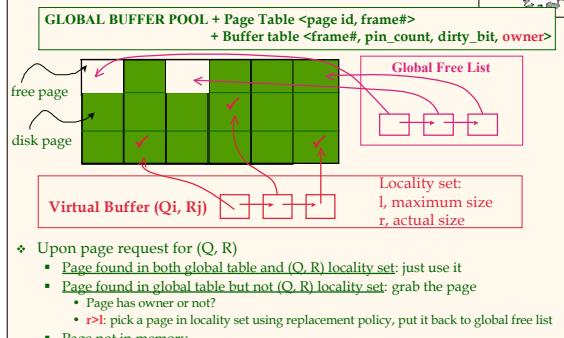
---

DBMIN Buffer Management

- ❖ A buffer manager that uses QLSM
    - e.g. all those access patterns, sizes of locality sets, suitable replacement policies
  - ❖ To reduce I/O: a “virtual buffer” per query-file (Q, R)
    - Can allocate # frames as the locality set requires
    - Use the right replacement policy
  - ❖ Data sharing via the use of a global buffer pool
    - Two queries accessing the same page: one copy of the page in the global buffer, two pointers.
    - Fairness?
  - ❖ Load control
    - *Overload:* can not take more queries
    - *Underload:* a query can use the entire global pool!

13

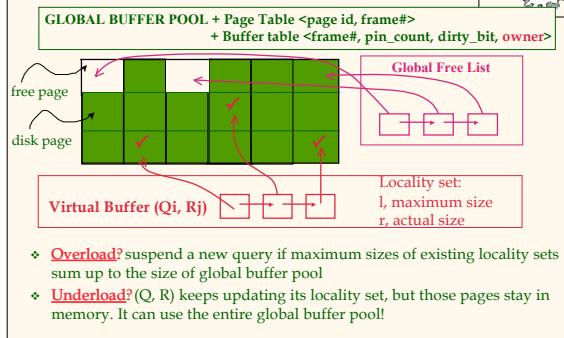
## *DBMIN Algorithm*



- Upon page request for (Q, R)
    - Page found in both global table and (Q, R) locality set: just use it
    - Page found in global table but not (Q, R) locality set: grab the page
      - Page has owner or not?
      - pick a page in locality set using replacement policy, put it back to global free list
    - Page not in memory

14

### *DBMIN Algorithm (Contd.)*



- ❖ **Overload**? suspend a new query if maximum sizes of existing locality sets sum up to the size of global buffer pool
  - ❖ **Underload**? (Q, R) keeps updating its locality set, but those pages stay in memory. It can use the entire global buffer pool!

15

## Summary



- ❖ Buffer manager brings pages into RAM.
  - Page stays in RAM until released by requestor.
  - Written to disk when frame chosen for replacement (which is sometime after requestor releases the page).
  - Choice of frame to replace based on *replacement policy*.
  - Tries to *pre-fetch* several pages at a time.
- ❖ Advanced buffer management
  - Understand access patterns to allocate enough buffer and choose the right replacement policy
  - Load control for overload and underload

16

---

---

---

---

---

---

## Summary (Contd.)



- ❖ DBMS vs. OS File Support: DBMS needs features not found in many OS's, e.g.,
  - forcing a page to disk (transaction mgmt)
  - controlling the order of page writes to disk (WAL for recovery)
  - ability to control pre-fetching
  - page replacement policy based on predictable access patterns, etc.

17

---

---

---

---

---

---

## Questions



18

---

---

---

---

---

---