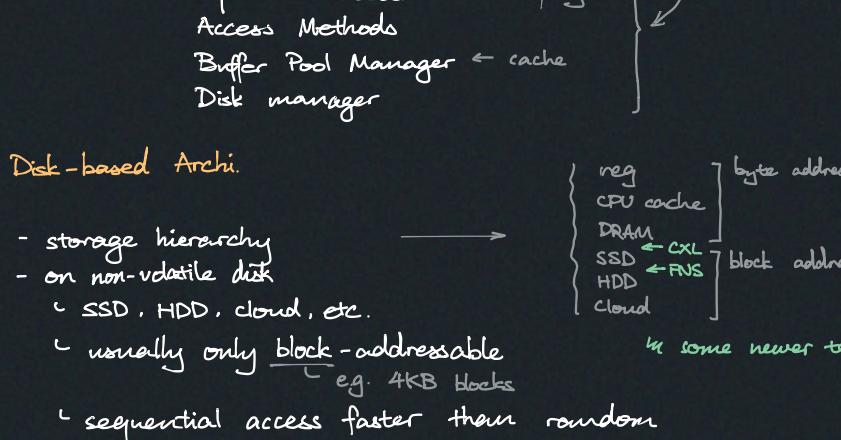


Lec 3 DB Storage I

DB Management



Disk-based Archi.

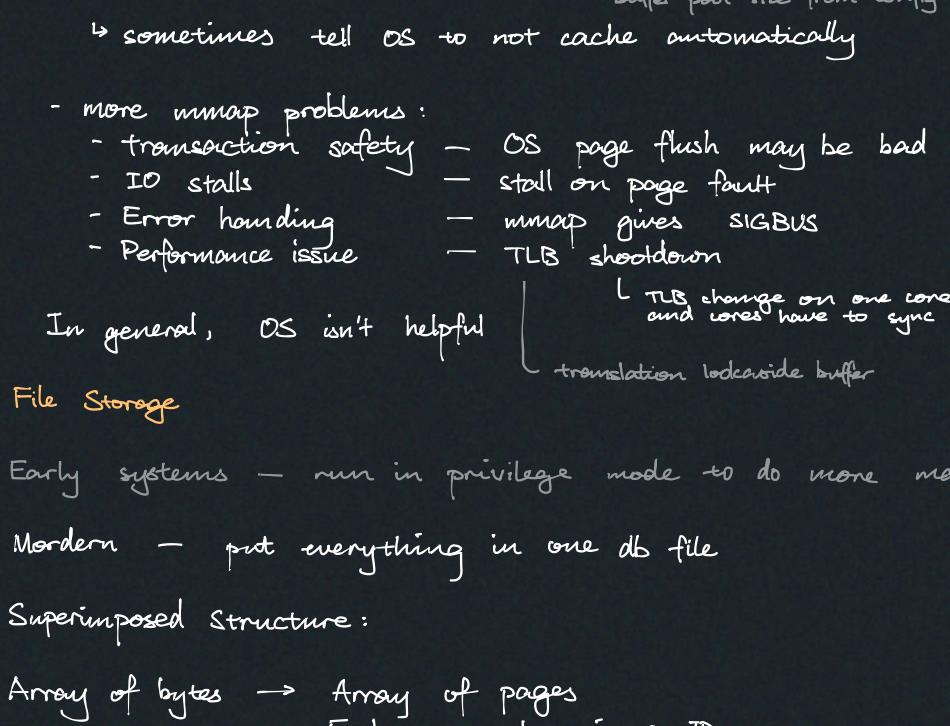
- storage hierarchy →
 - on non-volatile disk
 - ↳ SSD, HDD, cloud, etc.
 - ↳ usually only block-addressable
 - e.g. 4KB blocks
 - ↳ sequential access faster than random
 - ↳ need good locality for performance
 - ↳ multithreading to use multiple cores
 - ↳ write conflict handling



↳ some newer tech

DB Files

e.g. one file per table / all tables in one file



↳ bring disk file to virtual mem (at least appear so)

? OS provides mmap, why not use it?

- OS mmap doesn't handle file eviction the way we want
 - we want to be in control using semantic context in the DB
 - ↳ usually DB process just malloc its entire cache
- ↳ sometimes tell OS to not cache automatically

- more mmap problems:

- transaction safety
- IO stalls
- Error handling
- Performance issue
- OS page flush may be bad
- stall on page fault
- mmap gives SIGBUS
- TLB shootdown

In general, OS isn't helpful

↳ TLB change on one core and cores have to sync

↳ translation lookaside buffer

File Storage

Early systems — run in privilege mode to do more memory

Modern — put everything in one db file

Superimposed structure:

Array of bytes → Array of pages
Each page gets unique ID
Set some page size

Page layers [Hardware page
OS page
DB page]

↳ in the "dump it there" sense

Heap File Structure

↳ in the logical sense, not that they lack index

Unordered collection of pages

- create / delete / write / get page

- supports iterating

BB [page0 | page1 | page2 | ...]

Directory : - use one of the pages as directory that refer to other pages

- keep track of metadata

- keep track of free pages

Page header : - checksum

- page size

- DB version

- transaction data

- metadata / aggregation

- e.g. min/max date for quicker query

Self-contained page — all info in the page
some systems do this

Tuple oriented page

If fixed len — just array of records

In practice for variable len: slotted pages

Diagram of slotted pages:

tuple 1 | tuple 2 | tuple 3 | ...

tuple 4 | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header | data | header | data | ...

header |