



**Linux IO internals
for PostgreSQL administrators
in 2020**



- Linux is a most common OS for databases
- DBAs often run into IO problems
- Most of the information on topic is written by kernel developers (for kernel developers) or is checklist-style
- Checklists are useful, but up to certain workload



- **How to maximize page throughput between memory and disks**
- Things involved:
 - ▶ Disks
 - ▶ Memory
 - ▶ CPU
 - ▶ IO Schedulers
 - ▶ Filesystems
 - ▶ Database itself
- IO problems for databases are not always only about disks



- **How to maximize page throughput between memory and disks**
- Things involved:
 - ▶ **Disks - because latency of this part was very significant**
 - ▶ Memory
 - ▶ CPU
 - ▶ IO Schedulers
 - ▶ Filesystems
 - ▶ Database itself
- IO problems for databases are not always only about disks

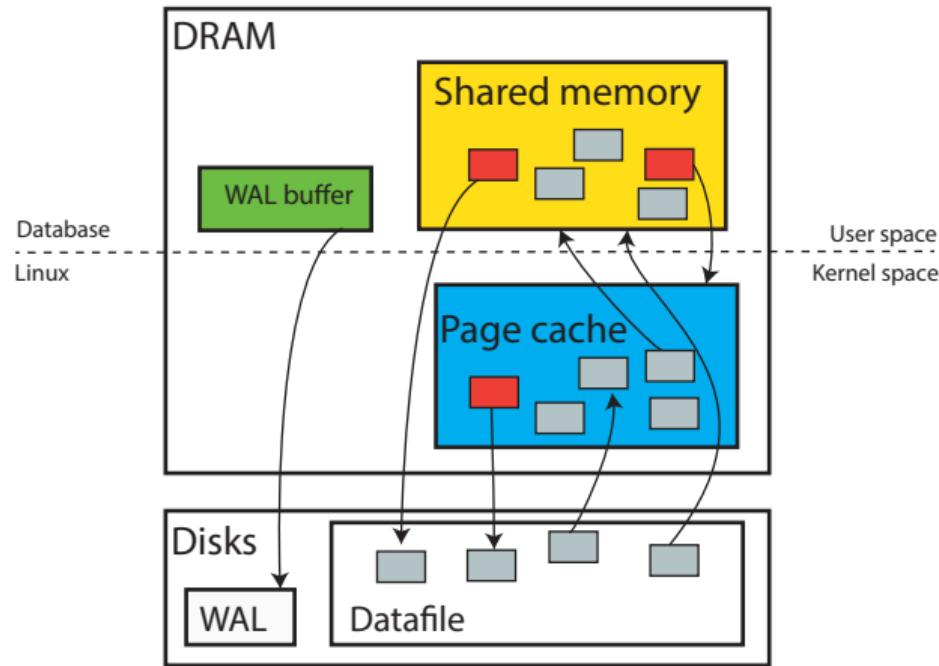


- Maximizing IO performance through maximizing throughput is easy up to certain moment
- Minimizing latency of IO usually is tricky
- With large adoption of proper SSDs, hardware latency dropped dramatically



- Database development was concentrated around maximization of throughput
- So did Linux kernel development
- Many rotating disks era IO optimization techniques are not that good for SSDs



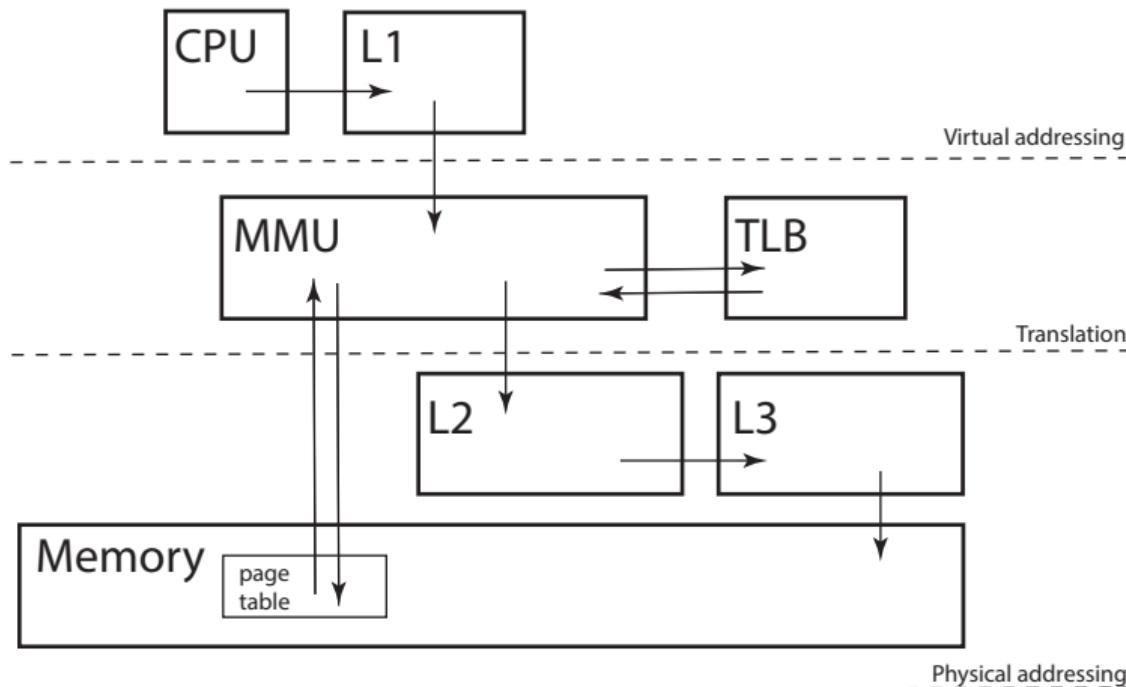


- Shared memory segment can be very large
- Keeping in-memory pages synchronized with disk generates huge IO
- WAL should be written fast and safe
- One and every layer of OS IO stack involved



Memory allocation and mapping

9

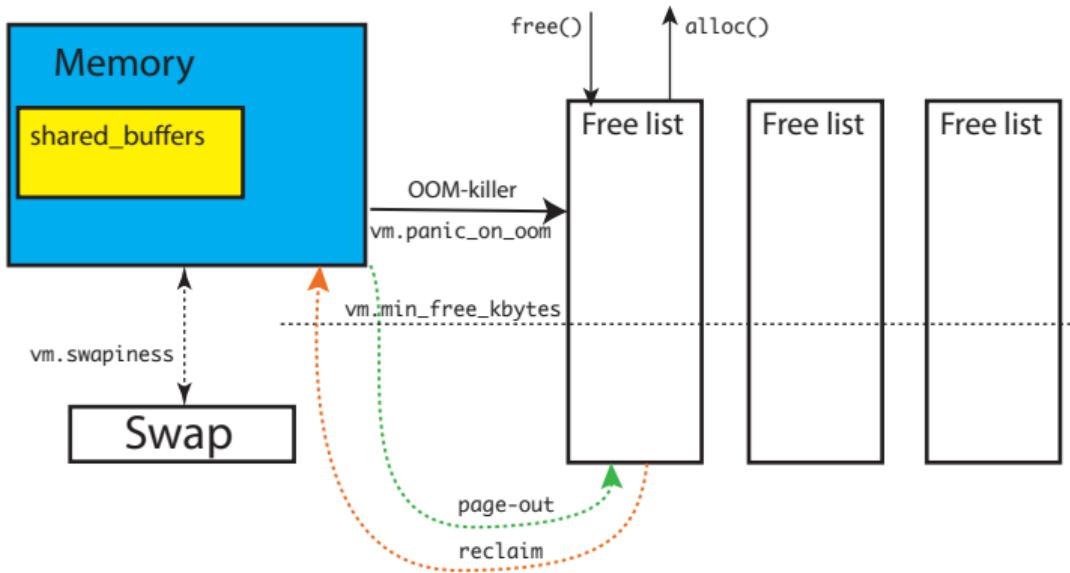


- Database with huge shared memory segment benefits from huge pages
- But not from transparent huge pages
 - ▶ Databases operate large continuous shared memory segments
 - ▶ THP defragmentation can lead to severe performance degradation in such cases



Freeing memory

11



- Page-out happens if:
 - ▶ someone calls fsync
 - ▶ 30 sec timeout exceeded (*vm.dirty_expire_centisecs*)
 - ▶ Too many dirty pages (*vm.dirty_background_ratio* and *vm.dirty_ratio*)
- It is reasonable to tune *vm.dirty_** on rotating disks with RAID controller, but helps a little on server class SSDs

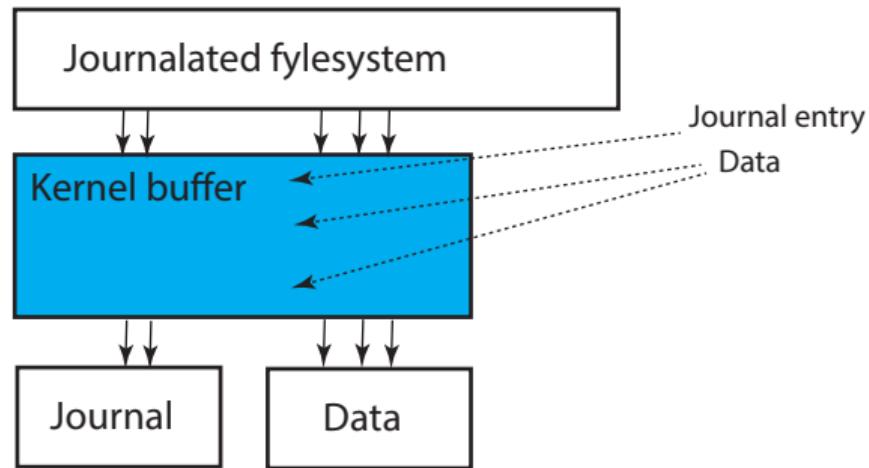


- *vm.overcommit_memory*
 - ▶ 0 - heuristic overcommit, reduces swap usage
 - ▶ 1 - always overcommit
 - ▶ **2 - do not overcommit** (*vm.overcommit_ratio = 50* by default)
- *vm.min_free_kbytes* - reasonably high (can be as low as 1000000 on a server with enough memory)
- **vm.swappiness = 1**
 - ▶ 0 - swap disabled
 - ▶ 60 - default
 - ▶ 100 - swap preferred instead of other reaping mechanisms
- Your database would not like OOM-killer



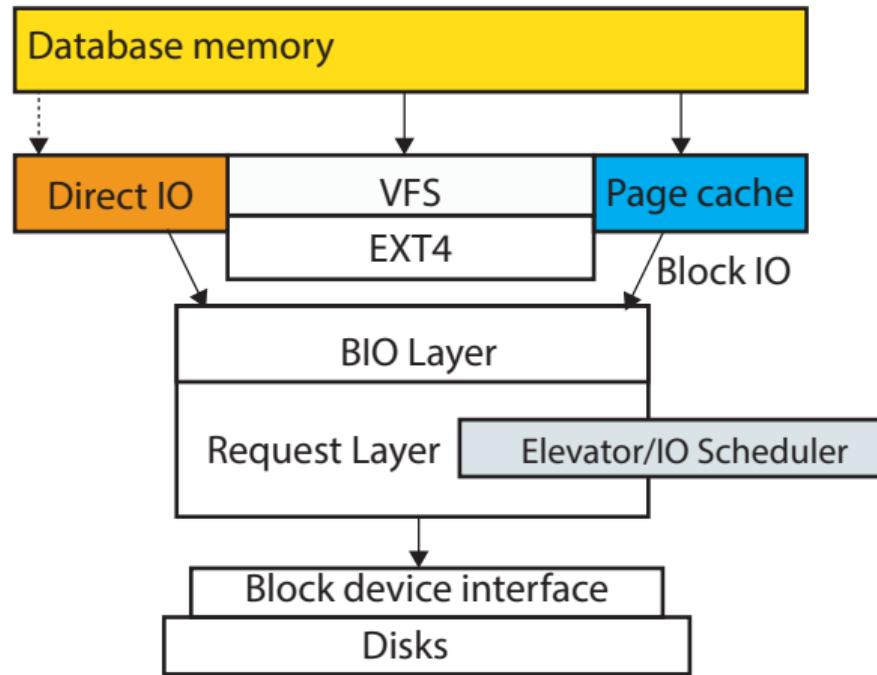
- *vm.panic_on_oom* effectively disables OOM-killer, but that is probably not the result you desire
- Or for a certain process: echo -17 > /proc/12465/oom_adj but again



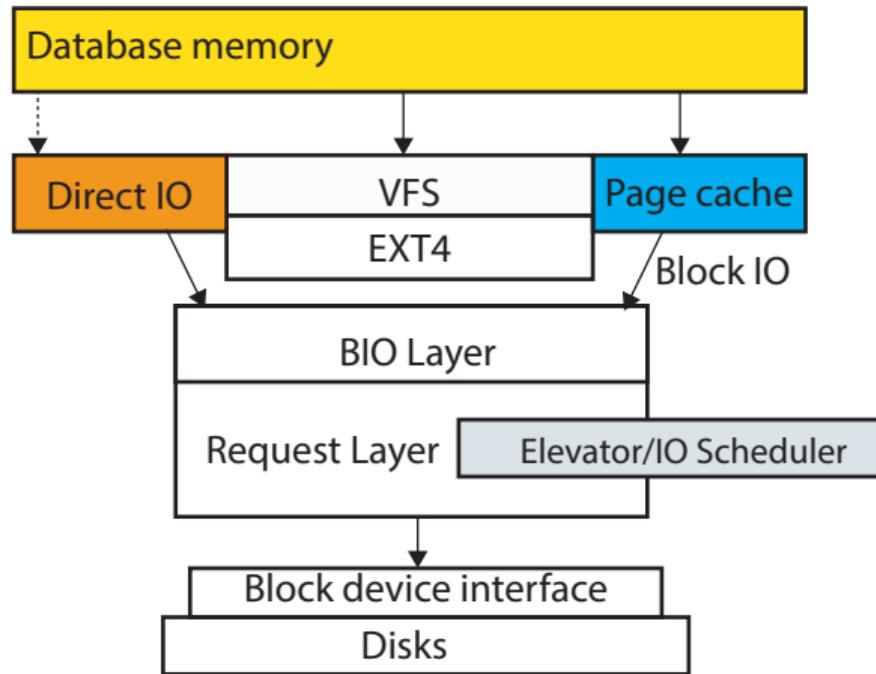


- ext4 or xfs
- Disable write barrier (only if SSD/controller cache is protected by capacitor/battery)



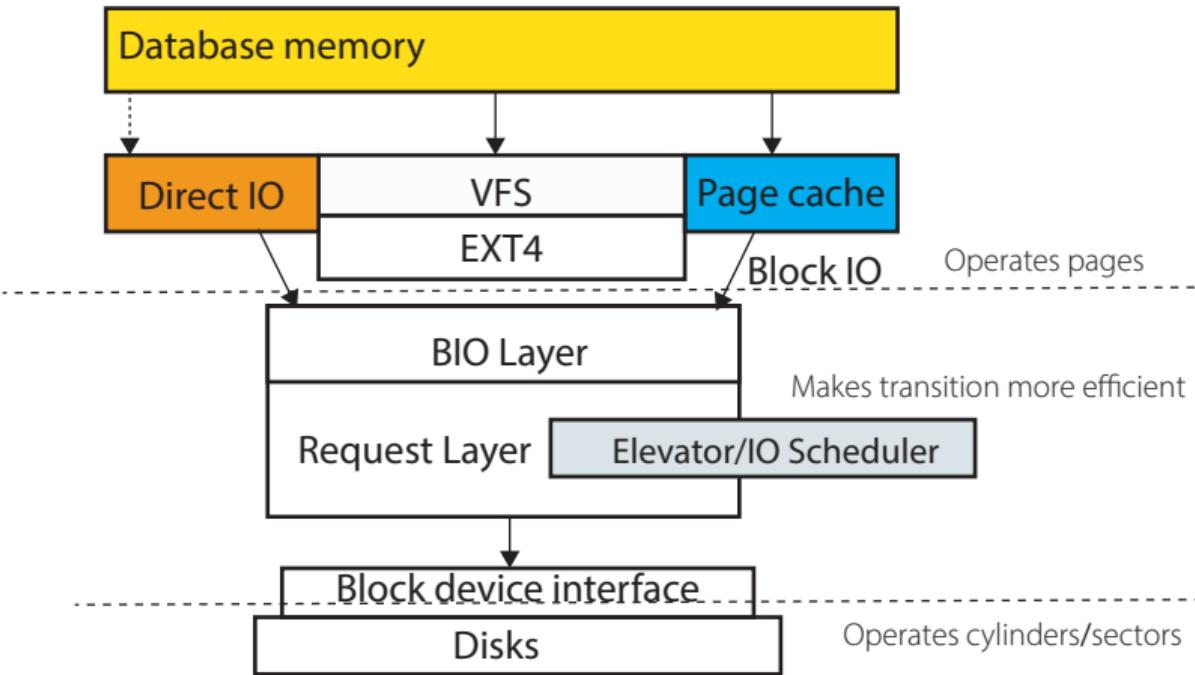


IO stack (as it used to look like)



IO stack (as it used to look like)

19



- **Linus Elevator** - the only one in times of 2.4
- *merging and sorting* request queues
- Had **lots** of problems

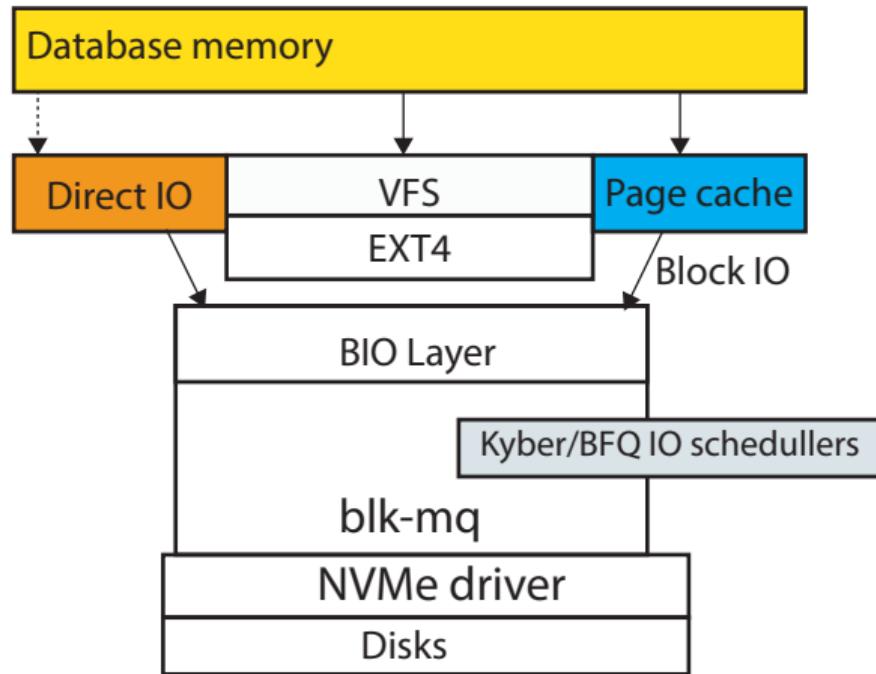


- **CFQ** - universal, default one
- **deadline** - rotating disks
- **noop** or **none** - then disks throughput is so high, that it can not benefit from keen scheduling
 - ▶ PCIe SSDs
 - ▶ SAN disk arrays



- Effectiveness of **noop** clearly shows ineffectiveness of others, or ineffectiveness of smart sorting as an approach
- **blk-mq** scheduler was merged into 3.13 kernel
- Much better deals with parallelism of modern SSD - basically separate IO queue for each CPU
- The best option for good SSDs right now
- **blk-mq and NVMe driver is actually more than scheduler, but a system aimed to substitute whole request layer**





- **https://www.thomas-krenn.com/en/wiki/Linux_Storage_Stack_Diagram**
- Regular updates
- Some things are difficult to draw, but it is a complex topic



- Sets of standards, which helps to use modern SSDs more effectively
- For Linux it is first of all NVMe driver (or subsystem)
- Most common example of NVMe SSDs are PCIe NAND drives
- With NVMe v.5 (currently 3 is ready for production) can work up to 32GB/sec
- Are databases NVMe ready?



- IO polling
- New IO schedulers Kyber and BFQ (Kernel 4.12)
- IO tagging
- Direct IO improvements
- io_uring (Kernel 5.1)



- Currently PostgreSQL supports DirectIO only for WAL, but it is unusable on practice
- Requires a lots of development
- Very OS specific
- Allows to use specific things, like O_ATOMIC
- PostgreSQL is the only database, which is not using Direct IO
- Many people consider Direct IO as a dirty hack
- Looks like io_uring has better perspectives



ik@dataegret.com

