WRITING LINUX FS 4FUN

Part 2 Short story of Pages and Virtual Memory

Outline

Why

Main Concepts and bit of history

- VM as a start of the journey
- Paging, Swapping
- Caching vs Latencies
- VM caching problem statement
 - Solaris Page Cache
 - Linux Page Cache
- IO paths
- Code Fragments:
 - Read/Write Page
 - mmap

Why this talk?!

Cons

- Writing FS is quite time consuming (approx. 10 years...)
- Just few production ready FS, many abandoned or not truly maintained
- File Systems considered as passé: now is Object Stores era!

Pros

- Address specific Education gap
- Solving other complicated problems
 - Memory Management is considered as difficult topic
 - Storage stack is complicated and usually became a bottleneck
 - Data is foundation of most todays application

VM in UNIX

- Virtual Memory: "high level" abstraction
- Multiple Processes with different address space
- Concept of Primary Storage: Main Memory (Fast - RAM) and Secondary Storage (Slow -Backing Device)

UNIX, like all time-sharing systems, and some multiprogramming systems uses "program swapping" (also called "rollin/roll-out") to share the limited resource of the main physical memory among several processes.

Processes which are suspended may be selectively "swapped out" by writing their data segments (including the "per process data") into a "swap area" on disk

From J.Lions[77] c14. Program Swapping

Paging

- Mechanism of Memory Management: Divide whole physical memory to small chunks called pages
- Pages referred by their number i.e. PFN
- Three important paging policies [1973]:
 - Fetch policy
 - Placement policy
 - Replacement Policy



Demand Paging

- Introduced in : Demand paging made its appearance in UNIX with the introduction of the VAX-11/780 in 1978
- All versions of UNIX used demand paging as the primary memory management technique, with swapping relegated to a secondary role in 80s
- When referencing memory which is not present Page Fault: Minor, Major



Computing is all about caching

- Intel Laptop CPU can execute 50 Bilions instruction per second
 - Unless it takes a cache miss
- DRAM can deliver 500 Milions cache lines per second
- SSD can deliver 800 thousand 4KiB pages per second
- PDP 11/70 could execute
 400 Thousand IPS
- Its Memory could handle
 300 Thousand cache misses p/s
- The RK05 drive could deliver
 4 thousand 4KiB pages p/s



Latencies numbers that you should know

```
L1 cache reference ..... 0.5 ns
Branch mispredict ..... 5 ns
L2 cache reference ..... 7 ns
```

Main memory reference 100 ns

```
SSD random read ..... 150,000 ns = 150 \mus
Read 1 MB sequentially from memory .... 250,000 ns = 250 \mus
```

```
Read 1 MB sequentially from SSD* .... 1,000,000 ns = 1 ms
Disk seek ..... 10,000,000 ns = 10 ms
Read 1 MB sequentially from disk .... 20,000,000 ns = 20 ms
```

Lets multiply all these durations by a billion:

L1 cache reference	0.5 s	One heart beat (0.5 s)
Branch mispredict	5 s	Yawn
L2 cache reference machine	7 s	Dispatching expresso in coffee
Main memory reference	100 s	Doing a toasted sandwich
SSD random read	1.7 days	A normal weekend
Read 1 MB sequentially from memory	2.9 days	A long weekend
Read 1 MB sequentially from SSD	11.6 days	Waiting for 2 weeks for a delivery
Disk seek	16.5 weeks	A semester in university
Read 1 MB sequentially from disk	7.8 months	Almost producing a new human being
The above 2 together	1 year	



File memory in pages



SOLARIS

- Page Cache implemented as a Hash Table
- Hashing of Pages done by key [vnode, offset]
- Each vnode maintain list of pages associated to this vnode



Linux

- Page Cache: inode connected to the address space
- Address space mapping of the pages
- Radix tree as Page Cache
- One Page Cache per inode (file)



IO Paths:

- Data can be sent to the Storage device in many ways
- Multiple components inside kernel space
- Userspace can advise the kernel



Buffered IO

- Application do frequent R/W
- Is fine for data to sync with disk after some time
- Read/Write directly from/to the Cache
- Cache do synchronization with storage by its own
- If cache miss during the Read need to wait



Synchronous Buffered IO

- Application want Write the data to reach the disk and get confirmation.
- Cache do store this data as well as it needs to sync data between it and Storage medium



Direct IO

- Application wants to manage caching the data by its own
- Whenever IO performed it should go directly to the disk, no point for caching the same data inside kernel



Memory mapped IO

- Application wants to use kernel Page Cache, but not keep a copy of the data on their site
- Same pages can be shared between Processes



Page Cache and Storage Device

- Page Cache is a set of pages per address space
- Possible to GET/PUT page to the page cache



What are we trying to achieve?

- Make use of the Page Cache for:
 - Buffered IO
 - Sync IO
- Make sure that Direct IO works and omit Page Cache
- Support mmap/munmap

```
Fragments:
           mmap
    "support mmap(2)"
// file.c fops for dm file
file operations dummy file ops = {
•••
  .mmap = dm file mmap,
•••
// Handler of mmap
static int
dm file mmap(struct file *file,
           struct vm area struct *vma)
 return generic file mmap(file, vma);
```

```
//mm/filemap.c
/* This is used for a general mmap of a disk file
int generic file mmap(struct file * file, struct
vm area struct * vma)
{
...//
       file accessed(file);
        vma->vm ops = &generic file vm ops;
        return 0;
}
/* These are filemap generic methods */
vm operations struct generic file vm ops = {
  .fault
                  = filemap fault,
  .map pages = filemap map pages,
  .page mkwrite = filemap page mkwrite,
};
```

Fragments: R/W through PC

```
// mm/filemap.c
generic file read iter(struct kiocb *iocb, struct iov iter
*iter)
   struct file *file = iocb->ki filp;
   struct address space *mapping = file->f mapping;
   struct inode *inode = mapping->host;
   loff t size;
   size = i size read(inode);
   if (iocb->ki flags & IOCB NOWAIT) {
      if (filemap range has page(mapping, iocb->ki pos,
               iocb->ki pos + count - 1))
       return -EAGAIN;
   } else {
      return filemap write and wait range(mapping,
            iocb->ki pos, iocb->ki pos + count - 1);
    }
   file_accessed(file);
   retval = mapping->a ops->direct IO(iocb, iter);
 }
```

return generic_file_buffered_read(iocb, iter, retval);

Fragments: pages "support r/w pages"

// inode as ops for dm inode struct address space operations dm as ops = $\{$.readpage = dm readpage, .readpages = dm readpages, .writepage = dm writepage, .writepages = dm writepages, .write begin = dm write begin, .write end = dm write end, .bmap = dm bmap,

```
//inode.c
static int dm readpage(struct file *file,
                       struct page *page)
{
 return mpage readpage(page, dm get block);
}
static int
dm readpages(struct file *file, struct address space
*mapping, struct list head *pages, unsigned nr pages)
{
 return mpage readpages (mapping, pages, nr pages,
                         dm get block);
}
static int dm writepage(struct page *page,
                        struct writeback control *wbc)
{
 return block write full page(page, dm get block, wbc);
}
static int
dm writepages(struct address space *mapping,
              struct writeback control *wbc)
{
 return mpage writepages(mapping, wbc, dm get block);
```

Fragments: get block

{

}

```
/**
 * Get Block with block index
 * iblock
 * map it to the buffer head
 * This is FS specific routine
 */
int dm get block(
    struct inode *inode,
    sector t iblock,
    struct buffer head *bh_result,
    int create)
```

```
int dm get block(struct inode *inode, sector t iblock,
                      struct buffer head *bh result, int create)
           unsigned max blocks =
                         bh result->b size >> inode->i blkbits;
           bool new = false, boundary = false;
           u32 bno;
           int ret;
/ * *
* return > 0, 'N' of blocks mapped or allocated.
* return = 0, if plain lookup failed.
* return < 0, error case.</pre>
 */
           ret = dm get blocks(inode, iblock, max blocks, &bno,
                               &new, &boundary, create);
           if (ret <= 0)
                      return ret;
           map bh(bh result, inode->i sb, bno);
           bh result->b size = (ret << inode->i blkbits);
           if (new)
                      set buffer new(bh result);
           if (boundary)
                      set buffer boundary(bh result);
           return 0;
```

Fragments: Get Blocks

/**

```
* Return the blocks
* routine for Filesystem
* iblock is start block
* maxblocks max number of blocks
*/
static int dm_get_blocks(
    struct inode *inode,
    sector_t iblock,
    unsigned long maxblocks,
    u32 *bno,
    bool *new,
    bool *new,
    int create)
```

```
/* Iterate all extends see if we can find a block range inside existing range \ast/
```

```
while (count < maxblocks && count <= blocks_to_boundary)</pre>
```

/* Simple case block found in existing blocks and do
not overflow return it */

```
if (err != -EAGAIN)
    goto got_it;
```

/* Next simple case - plain lookup or failed read of
indirect block so we need to chain them. */

```
if (!create || err == -EIO)
    goto cleanup;
```

```
/* Truncate if size bigger than requested */
```

```
if (count > blocks_to_boundary)
    *boundary = true;
```

```
err = count;
```

Fragments: DirectIO

"support for DirectIO"

// inode as ops for dm_inode

struct address_space_operations
dm_as_ops = {

• • •

.direct_IO = dm_direct_IO,

// We can again use magic get_blok function
// and plug to FS framework

{

}

```
static ssize_t
dm_direct_IO(struct kiocb *iocb, struct iov_iter *iter)
```

```
struct file *file = iocb->ki_filp;
struct address_space *mapping = file->f_mapping;
struct inode *inode = mapping->host;
size_t count = iov_iter_count(iter);
loff_t offset = iocb->ki_pos;
```

Other resources:

- J.Lions: "A commentary on the sixth edition UNIX Operating System" [1977]
- R McDougall, J. Mauro: "Solaris Internals Second edition" [2006]
- Uresh Vahalia: "UNIX Internals The new Frontiers" [1996]
- Steve D. Pate: "UNIX Filesystems: Evolution, Design and Implementation"
- Ext2, Ufs: Linux source code
- github.com/gotoco/dummyfs

