Linux Kernel Programming
**Memory Management Pt.2**

Sang-Hoon Kim

Systems Software Research Group @ Virginia Tech

April 13, 2017

## Everyday questions

- ▶ I found my smartphone/laptop is low on free memory
- ▶ This app greatly increases the free memory from my phone!! Five stars!!
- ▶ I hate Google Chrome because it hogs memory
- ▶ Linux installer suggests a swap partition which is twice of the RAM size
  - ▶ By the way, I have 32GB RAM and 128GB SSD... what?
- ▶ I don't understand the performance evaluation results
  - ▶ 300 MB/sec from hard disk drive????
- ▶ I could `malloc()` 2 GB from 1 GB Raspberry Pi...?

Virginia
Tech

# Myth

- High memory utilization makes me feel bad ☹
  - Swapping will destroy the system performance
  - Will increase the memory allocation time
  - Will trigger *some* background operations

- Are they all true?
- Will be happy if 20% is used and the rest is free?
- What do you want to do with the *"free memory"*?

# Outline

Virginia Tech

# Outline

# System memory utilization
Metrics for memory utilization

- ▶ The portion of memory that is currently used in the system
- ▶ Memory in use + free memory = total memory
- ▶ Allocated memory / total memory or
  (total memory - free memory) / total memory
- ▶ E.g., using 3 GB out of 4 GB: 3 GB / 4 GB = 75%
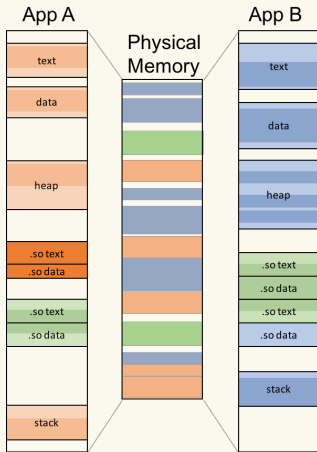
Virginia Tech

# Per-process memory utilization
Metrics for memory utilization

- ▶ Processes request the kernel for virtual memory areas (VMAs)
- ▶ A read causes the corresponding page to be mapped to the zero page.
- ▶ A write causes the copy-on-write to the zero page
    - ▶ Allocate the actual page
    - ▶ Update the mapping to the new page
- ▶ VMAs can be shared between processes
- ▶ The sharing can be private

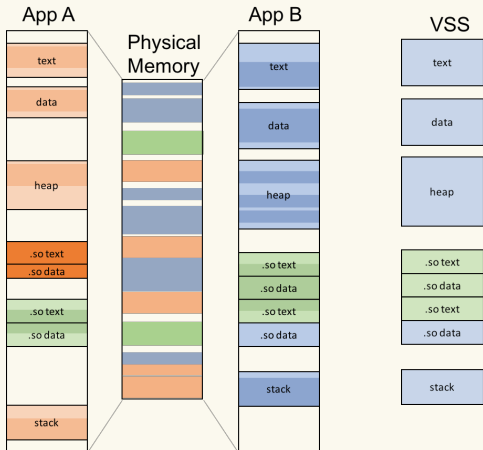Virginia Tech

# Per-process memory utilization

Metrics for memory utilization

# Virtual set size (VSS)

Metrics for memory utilization

▶ The size of virtually allocated memory

# Virtual set size (VSS)
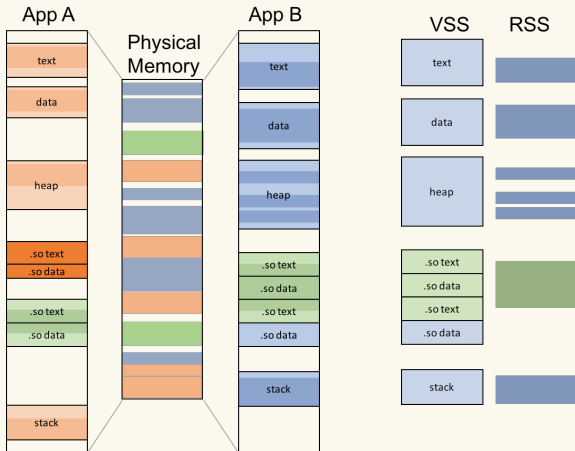Metrics for memory utilization

- ▶ The size of virtually allocated memory
- ▶ Usually, way to large than the actual memory footprint
- ▶ Can obtain by summing the size of each VMA
  - ▶ $\sum_{\text{Each VMA}} |\text{VMA}|$

```c
/* in include/linux/mm_types.h */

struct mm_struct {
  /* ... */
  unsigned long total_vm; /* Total pages mapped */
  /* ... */
};
```

Virginia Tech

# Resident set size (RSS)
Metrics for memory utilization

▶ The memory footprint that the process actually occupies

# Resident set size (RSS)

Metrics for memory utilization

- ▶ The memory footprint that the process actually occupies
- ▶ Can get by counting the valid page table entries
    - ▶ Page size $\times$ # of mapped pages

```
1  /* In include/linux/mm.h */
2
3  static inline unsigned long get_mm_rss(struct mm_struct *mm);
```
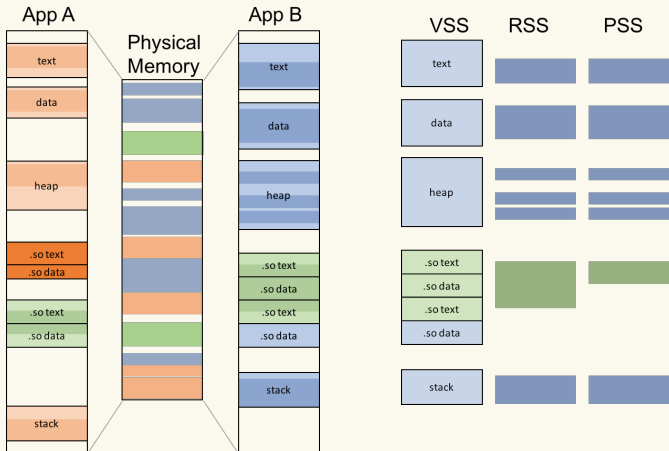
# Proportional set size (PSS) / Unique set size (USS)
Metrics for memory utilization

- ▶ VSS and RSS do not consider the memory sharing between processes
- ▶ Note: VMA is the unit of memory sharing
- ▶ Proportional set size (PSS)
  - ▶ Account for the number of sharing processes
  - ▶ Indicate the memory contribution of the process
  - ▶ Page size $\times \sum_{\text{Each VMA}}$ # of mapped pages / # of share
  - ▶ E.g., VMA contains 10 pages and shared by 4 processes
    = 4 KB $\times$ 10 / 4 = 10 KB
- ▶ Unique set size (USS)
  - ▶ The memory footprint that is not shared with any other processes
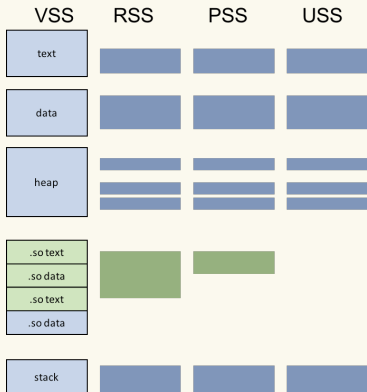  - ▶ The amount of memory that can be reclaimed by killing the process

Virginia Tech
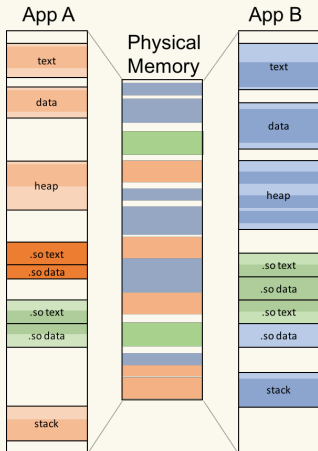
# Proportional set size (PSS) / Unique set size (USS)
## Metrics for memory utilization

# Proportional set size (PSS) / Unique set size (USS)
## Metrics for memory utilization

# Accessing the metrics

Metrics for memory utilization

▶ top

# Accessing the metrics

Metrics for memory utilization

- ▶ /proc/[pid]/smap

```
1  00400000-004f2000 r-xp 00000000 08:01 21233800          /bin/bash
2  Size:               968 kB
3  Rss:                888 kB
4  Pss:                177 kB   /* 888 / 177 = 5 */
5  Shared_Clean:       888 kB   /* Rss == Shared */
6  Shared_Dirty:         0 kB
7
8  ...
9
10 006f1000-006f2000 r--p 000f1000 08:01 21233800          /bin/bash
11 Size:                 4 kB
12 Rss:                  4 kB
13 Pss:                  4 kB
14
15 ...
16
17 7f92acb43000-7f92acd42000 ---p 0000b000 08:01 14682101    /lib/x86_64-linux
       -gnu/libnss_files-2.19.so
18 Size:              2044 kB
19 Rss:                  0 kB
20 Pss:                  0 kB
```

# Outline

Virginia Tech

# Is high memory utilization bad?

## Page cache

- ▶ In fact, you are wasting your memory if you don't utilize it

- ▶ Keep dissipating 5–30% of system energy [1, 2, 3]
    - ▶ The cells in DRAM are effectively capacitors
    - ▶ The value is determined by the amount of electron charges (i.e., voltage) in the cell
    - ▶ Should be refreshed periodically (∼64ms)

- ▶ Free pages and used pages are all the same from the memory modules' perspective

- ▶ Miss the opportunity to utilize the *fastest* storage in the system that you/OS can fully control its contents
    - ▶ Explicitly utilize as RAM disks
    - ▶ Else?

Wordline

Bitline

Virginia Tech

# Page cache
Page cache

- ► Utilize free pages to cache the data from slow storage devices
- ► Implemented at the virtual file system (VFS) layer
    - ► Will be discussed in the following lecture
- ► Maintain file-backed pages
    - ► Pages are indexed by ({device, file}, offset)

- ► When `read()`/`fread()` from a file,
    - ► Look up the page cache
    - ► If exist, copy the data from the page cache page. Fast!!
    - ► If not,
        - ► Allocate page(s)
        - ► Put the page into the page cache and lock properly
        - ► Generate I/O request(s) to corresponding block device/file system
        - ► Copy the read content to the user-space buffer

Virginia
Tech

# Page cache
Page cache

- ▶ When `write()`/`fwrite()` to a file,
    - ▶ Look up the page cache
    - ▶ If exist, apply the write to the page and flag it dirty. Fast!!
    - ▶ If not,
        - ▶ Allocate page(s)
        - ▶ Generate I/O request(s) to corresponding block/file system
        - ▶ Apply the write to those pages, make them dirty
    - ▶ Dirty pages will be *magically* written back to the storage
    - ▶ The write-back makes the page clean

Virginia
Tech

# Page cache
## Page cache

- ► Page cache significantly improves the system performance
  - ► Can read without accessing the slow storage devices
  - ► Can buffer writes
  - ► Merge overwrites before writing to the storage
  - ► Can shape the write pattern to the storage device
- ► In enterprise and HPC environments, performance-critical servers are equipped with a huge amount of memory for the page cache

# How large is the page cache?
Page cache

```
1  sanghoon@echo_debian:~$ free -m
2              total      used      free    shared   buffers    cached
3  Mem:        32133     26789      5344        24      1592     23720
4  -/+ buffers/cache:     1475     30658
5  Swap:           0         0         0
```

- Memory utilizatoin = 26,789 / 32,133 = 83%
- Page cache = 23,720 MB (88% / 74%)
- Available = 5,344 + (1,592 + 23,720) $\pm\epsilon$ = 30,656 MB

```
1  sanghoon@cerberus_ubuntu:~$ free -m
2              total      used      free    shared  buff/cache  available
3  Mem:        12012       156      3158        88        8697      11154
4  Swap:       15250        48      1520
5  sanghoon@echo:~$
```

- Memory utilizatoin = (156 + 8,697) / 12,012 = 74%
- Page cache = +-8,697 MB (98% / 72%)
- Available = 3,158 + 8,697 $\pm\epsilon$ = 11,154 MB
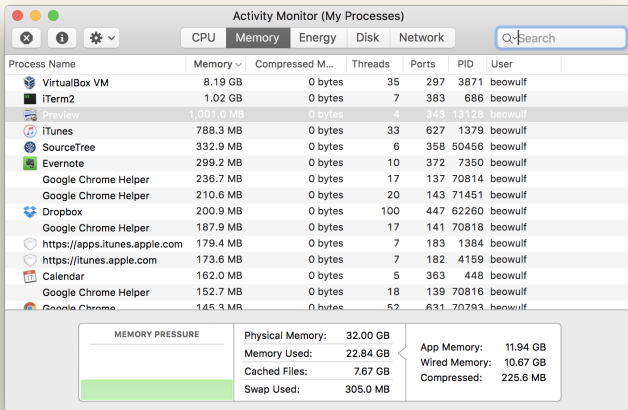
Virginia Tech

# How large is the page cache?
Page cache

```
1  # cat /proc/meminfo
2  MemTotal:       32904980 kB
3  MemFree:         4971852 kB
4  MemAvailable:   32399416 kB
5  Buffers:         1634068 kB
6  Cached:         24768688 kB <--
7    ...
```

# How large is the page cache?

Page cache

# QnA
## Page cache

**Q:** I don't understand the performance evaluation results
- ▶ 300 MB/sec from hard disk drive????

**A:** Check whether the data is accessed from the page cache
- ▶ Unfair between the 1st run and subsequent runs
- ▶ Drop or warm up the page cache
- ▶ You can explicitly drop page caches
- ▶ Bypass the page cache by `open()` with `O_DIRECT`

```
1  /* To free pagecache: */
2  $ echo 1 > /proc/sys/vm/drop_caches
3
4  /* To free reclaimable slab objects (includes dentries and inodes): */
5  $ echo 2 > /proc/sys/vm/drop_caches
6
7  /* To free slab objects and pagecache: */
8  $ echo 3 > /proc/sys/vm/drop_caches
```

Virginia Tech

# QnA
Page cache

**Q:** I found my smartphone/laptop is low on free memory

**A:** Don't panic

- ▶ Your operating system is doing its best to maximize the system performance
- ▶ By the way, check the swap status for sure

**Q:** This app increases the free memory from my phone!! Five stars!!

**A:** ...

# Outline

Virginia Tech

Sang-Hoon Kim (SSRG@VT)  Memory Management Pt.2  April 13, 2017  28 / 45

# Memory is limited
Page reclamation

- ▶ Eventually, memory will be filled up with page cache pages and process-allocated pages
- ▶ The kernel cannot keep allocating pages
- ▶ Should reclaim allocated memory somehow

- ▶ When the kernel should start the reclamation?
- ▶ Who will reclaim the memory?
- ▶ Which memory should be reclaimed?
- ▶ How reclaim?

# When to start the memory reclamation?
Page reclamation

- ▶ When fail to allocate a page from `alloc_pages()`
    - ▶ Mostly due to zone imbalancing (e.g., alloc_page(GFP_DMA))
    - ▶ `ZONE_DMA`, `ZONE_DMA32`, `ZONE_NORMAL`, `ZONE_HIGHMEM`
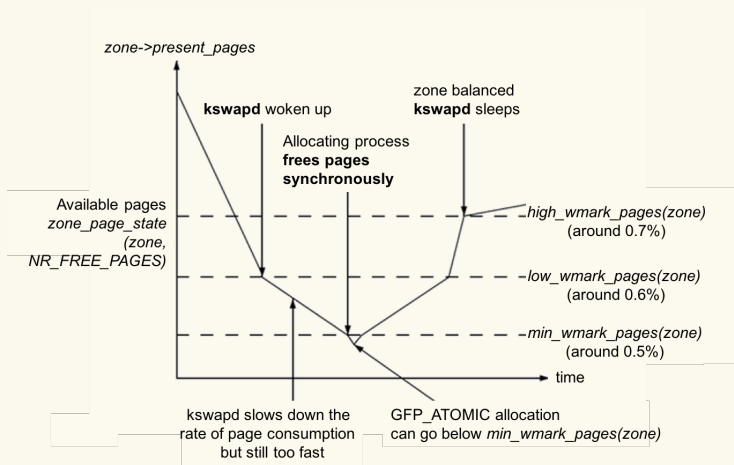- ▶ Based on the amount of free memory

Virginia Tech

# Who will reclaim?
## Page reclamation

- Direct reclamation
    - When `alloc_pages()` fails
    - Synchronously reclaim pages for critical page allocation

- Background reclamation
    - Triggered when `alloc_pages()` fails
    - A kernel thread a.k.a `kswapd` reclaims the page frames
    - Keep the free memory above thresholds

# Background reclamation
## Page reclamation

# What to reclaim?
Page reclamation

- ▶ Should evict file pages
    - ▶ Memory is for anonymous pages originally
    - ▶ Reclaiming clean file pages is fast
        - ▶ Drop the page from the page cache
    - ▶ Require slow I/Os to swap-out and -in anonymous pages
    - ▶ Many files are only accessed once
        - ▶ E.g., multimedia files

- ▶ Should keep file pages
    - ▶ Many clean pages are performance critical (e.g., code)
    - ▶ The write-back might block the reclamation
    - ▶ Cold anonymous pages are really cold
    - ▶ Storage devices are (were) extremely slow

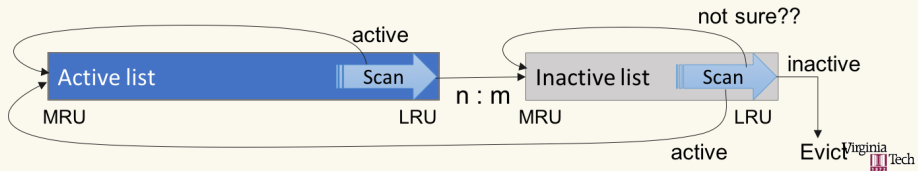Virginia Tech

# What to reclaim?
Page reclamation

- ▶ Keep *important* pages
    - ▶ Keep important pages whatever they are for
    - ▶ Can apply the traditional cache management policies
        - ▶ LRU, LFU, LRFU, CLOCK, ARC, CAR, MQ, ...
    - ▶ Anonymous pages and file-back pages have different characteristics
        - ▶ Many files are only accessed once
        - ▶ Code pages are performance critical
        - ▶ Cold anonymous pages are really cold

- ▶ OK. let manage them separately
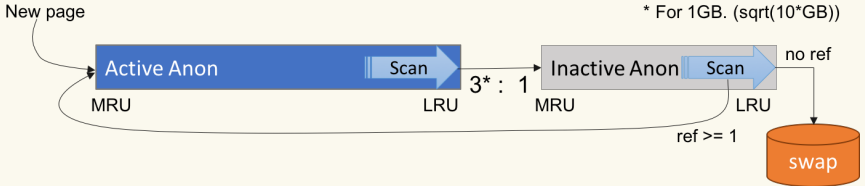
# 2Q
Page reclamation

- ▶ None of OS/cache uses the LRU/LFU in practice
  - ▶ Updating the position upon every access requires a significant overhead

- ▶ Linux employs a modified 2Q policy
  - ▶ Known good to separate inactive pages from active ones
  - ▶ Maintain an LRU list for inactive pages and an LRU list for active pages
  - ▶ Scan a part of the lists
  - ▶ Migrate pages between the lists according to some policies

# Managing the page lists
## Page reclamation

- ▶ Anonymous pages



New page

* For 1GB. (sqrt(10*GB))

Active Anon    Scan    3* : 1    Inactive Anon    Scan    no ref

MRU    LRU    MRU    LRU

ref >= 1

swap

- ▶ File pages



ref >= 1 && VM_EXEC

New page

ref == 1

Active File    Scan    1 : 1    Inactive File    Scan    no ref

MRU    LRU    MRU    LRU

ref >= 2 ||
ref==1 && VM_EXEC

file
system

# Swappiness
Page reclamation

- ▶ Tendency to reclaim anonymous pages with swap-out
- ▶ Integer between 0 to 200 (60 by default)
    - ▶ 60 means anon:file = 60:140
    - ▶ 100 means anon:file = 100:100
- ▶ Effectively 0 when no swap device exists
- ▶ Can control via `/proc/sys/vm/swappiness`

Virginia Tech

# Have a look at the source code
## Page reclamation

- ▶ `mm/vmstat.c` and `mm/vmscan.c`
- ▶ Isolate pages from the lists
    - ▶ `shrink_inactive_list()`
    - ▶ `shrink_active_list()`
- ▶ Process the isolated page
    - ▶ `shrink_page_list()`

Virginia Tech

# Check the size of page lists
Page reclamation

```
1  # cat /proc/meminfo
2  MemTotal:       32904980 kB
3  MemFree:         4971852 kB
4  MemAvailable:   32399416 kB
5    ...
6  Active:         14175700 kB
7  Inactive:       12260264 kB
8  Active(anon):      34020 kB
9  Inactive(anon):    24728 kB
10 Active(file):   14141680 kB
11 Inactive(file): 12235536 kB
12   ...
```

▶ Comment: We cannot tell whether the pages in the active/inactive page list are actually active/inactive or not

   ▶ The activeness can be determined by examining the page
   ▶ The pages evicted from the inactive list are inactive.

Virginia Tech
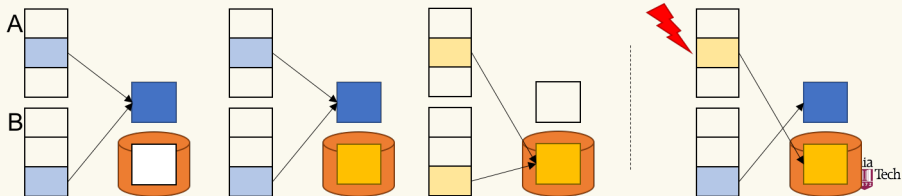
## How to reclaim?
Page reclamation

- ▶ File pages
  - ▶ Clean: Drop from the page cache
    - ▶ The data is in the filesystem
  - ▶ Dirty: May drop after writing them back, or keep it
- ▶ Anonymous pages
  - ▶ Perform the swap-out

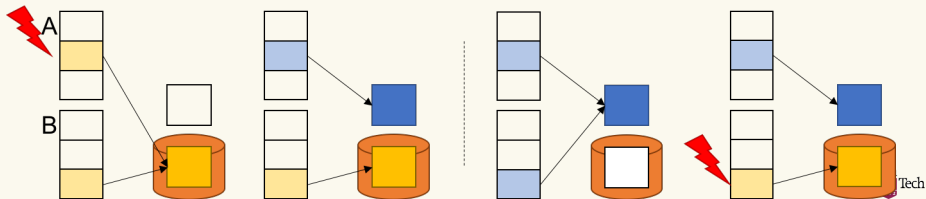Virginia Tech

# Swap-out v0.1
Page reclamation

- ► Note: A page can be shared between processes

- ► Allocate a swap entry from the swap device
- ► Write the page to the swap entry
- ► For each PTE points to the page (Q: How?)
    - ► Mark the PTE as invalid and update to the swap entry

- ► What if a process accesses the swapped-out page during the swap-out?
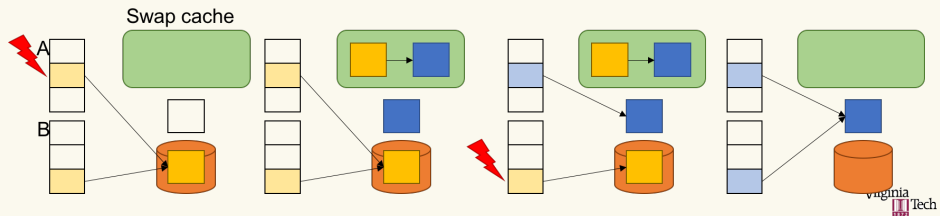
# Swap-in v0.1
Page reclamation

- ▶ When a process accesses the swapped out PTE
- ▶ Allocate a page
- ▶ Read in the page from the swap device
- ▶ Update the PTE accordingly

- ▶ Should the kernel fix the PTE from other processes and release swap entry?
- ▶ Which one is the current one?

# Swap cache
Page reclamation

- ▶ Prevent the race conditions during the swap-in and out
- ▶ Keep (swap entry ID, page) mappings for transient pages
- ▶ Consult to the swap cache first to swap-in
- ▶ If a corresponding mapping exists, update the PTE to the page in the memory instead of loading the swap entry

## Conclusion

- ▶ Do not under-utilize your memory
- ▶ Trust the techniques that are reviewed, evaluated, and proved

- ▶ Yes, memory management is dirty and complicated
- ▶ However, it is connected to everything
- ▶ Non-volatile memory and flash memory provide new opportunities

Virginia Tech

# Bibliography I

[1]  CARROLL, A., AND HEISER, G.
     An analysis of power consumption in a smartphone.
     In *Proceedings of the 2010 USENIX Annual Technical Conference (USENIX ATC'10)* (2010), pp. 21–35.

[2]  CARROLL, A., AND HEISER, G.
     The systems hacker's guide to the galaxy: Energy usage in a modern smartphone.
     In *Proceedings of the 4th ACM SIGOPS Asia-Pacific Workshop on Systems (APSys'13)* (2013).

[3]  LIU, S., PATTABIRAMAN, K., MOSCIBRODA, T., AND ZORN, B. G.
     Flikker: Saving DRAM refresh-power through critical data partitioning.
     In *Proceedings of the 16th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS XVI)* (2011), pp. 213–224.

Virginia Tech