# Linux Character Device Driver

*Sunbeam Infotech*

# Linux mem mgmt

① Paging
    ↳ MMU - TLB hw
    ↳ page table - multi-level paging
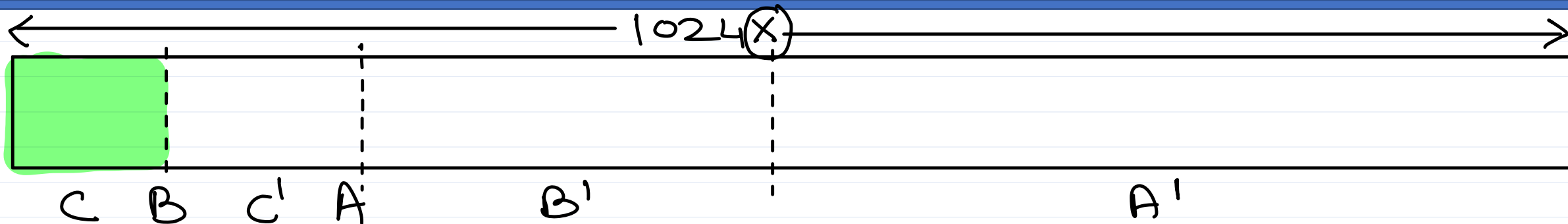        - 3 level pasing

② Page fault handling
    ↳ CPU req logical page - but page not in RAM.

                  OS page fault handler ← page fault ex.
           ↳ ✓ check if addr is valid & perms are valid.
             ✓ allocate empty frame → ?
              - if not empty frame, page replacement.
            ✓ if page on disk/swap, load it in RAM (in the frame)
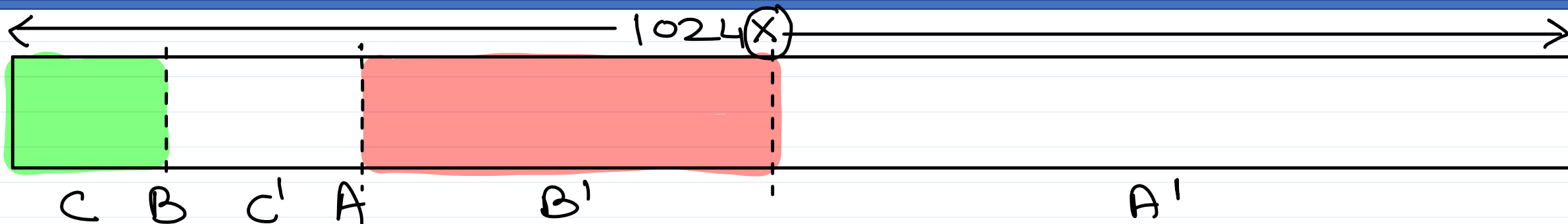            ✓ update page table & restart instn.

# Buddy allocator



$$1024\ \textcircled{X}$$

C   B   C'   A          B'                              A'

alloc req

① 128 → C

$1024 \rightarrow \cancel{X}$
$512 \rightarrow A', \cancel{A}$
$256 \rightarrow B', \cancel{B}$
$128 \rightarrow C', \cancel{C}$
$64 \rightarrow$

# Buddy allocator

$$1024\;\text{X}$$

C  B  C'  A  B'  A'

$\underline{alloc\;req}$

① $128 \rightarrow$ C

② $256 \rightarrow$ B'

$1024 \rightarrow \cancel{X}$

$512 \rightarrow A', \cancel{A}$

$258 \rightarrow \cancel{B}, \cancel{B}$

$128 \rightarrow C', \cancel{C}$

$64 \rightarrow$

# Buddy allocator



$1024\ \text{(X)}$

C   B   C'   A'   B'   A'

alloc req

1. $128 \rightarrow$ C
2. $256 \rightarrow$ B'
3. $64 \rightarrow$ D

$$1024 \rightarrow X$$
$$512 \rightarrow A', A$$
$$256 \rightarrow B, B$$
$$128 \rightarrow C, C$$
$$64 \rightarrow D', D$$

# Buddy allocator



alloc req

① 128 → C

② 256 → B'

③ 64 → D

④ 128 → F

$1024 \to \cancel{X}$

$512 \to \cancel{A'}, \cancel{A}$

$258 \to \cancel{B}, \cancel{D} \; E', \cancel{E}$

$128 \to \cancel{C}, \cancel{C'} \; F', \cancel{F}$

$64 \to D', \cancel{D}$

# Buddy allocator



alloc req
① 128 → C
② 256 → B'
③ 64 → D
④ 128 → F

dealloc req
① C ✗

$$1024 \to \cancel{X}$$
$$512 \to \cancel{A'}, \cancel{A}$$
$$256 \to \cancel{B'}, \cancel{B} \; E', \cancel{F}$$
$$128 \to \cancel{D'}, \cancel{D} \; F', \cancel{F}, \underline{C}$$
$$64 \to \cancel{D'}, \cancel{D}$$

# Buddy allocator



alloc reqr

① 128 → C

② 256 → B'

③ 64 → D

④ 128 → F

dealloc reqr

① C ✗

② B' ✗

$1024 \to \cancel{X}$

$512 \to \cancel{D'}, \cancel{A}$

$256 \to \cancel{B}, \cancel{D} \quad E', \cancel{F}, \underline{B'}$

$128 \to \cancel{C'}, \cancel{F} \quad F', \cancel{E} \quad \underline{C}$

$64 \to D', \cancel{D}$

# Buddy allocator



alloc reqr

① 128 → C
② 256 → B'
③ 64 → D
④ 128 → F

dealloc reqr

① C ✗
② B' ✗
③ F ✗

$1024 \rightarrow \cancel{X}$
$512 \rightarrow \cancel{A'}, \cancel{A} \quad A' \rightarrow A'$
$256 \rightarrow \cancel{B'}, \cancel{B} \quad \boxed{E'}, \cancel{F}, B' \quad \boxed{E}$
$128 \rightarrow \cancel{D}, \cancel{E} \quad \boxed{F'}, \cancel{F}, C \quad \boxed{F}$
$64 \rightarrow D', \cancel{D}$

E

# Buddy allocator



alloc reg

① 128 → C

② 256 → B'

③ 64 → D

④ 128 → F

dealloc reg

① C ✗

② B' ✗

③ F ✗

④ D ✗

↳ home work

1024 → ✗

512 → A', A → A'

256 → B', B → E', F, B', E

128 → C', C → F', F, C, F

64 → D', D

# Buddy allocator in Linux

RAM (4GB/8GB)

on x86-32 arch → Three Zones

① DMA → 0 - 16 MB → Buddy 1

② NORMAL → 16MB - 896MB → Buddy2

③ HIGHMEM → 896MB and above → Buddy3

cmd) cat /proc/buddyinfo

```
0 ──────────────
DMA
16MB ────────────
NORMAL
896MB ───────────
HIGH MEM
4GB/8GB
```

$2^0$ →

$2^1$ →

$2^2$ →

$2^3$ →

.
.

$2^{10}$ →

base addr.

to allocate mem from buddy allocator :

① get_free_pages (mask, order);

allocator behaviour

num of Consecutive pages - power 2

② free_pages (addr);

# Linux Buddy allocator

x86_32

$addr = \_\_get\_free\_pages(mask, order);$

$free\_pages(addr, order);$

Buddy Alloc
free list

$2^{10} \rightarrow$

$2^9 \rightarrow$

$2^8 \rightarrow$

$2^7 \rightarrow$

$\cdots$

$2^3 \rightarrow$

$2^2 \rightarrow$

$2^1 \rightarrow$

$2^0 \rightarrow$

Combn of
action mod
+ zone mod
= Type flag

→ e.g. GFP_KERNEL

$= \_\_GFP\_WAIT$
$+ \_\_GFP\_IO$
$+ \_\_GFP\_FS$

→ e.g. GFP_ATOMIC

$= \_\_GFP\_HIGH$

$0 \rightarrow 1$ Page
$1 \rightarrow 2$ Pages
$2 \rightarrow 4$ "
$3 \rightarrow 8$ "
$4 \rightarrow 16$ "
$\cdots$
$10 \rightarrow 1024$ "

Buddy Alloc
free list

$2^{10} \rightarrow$
$2^9 \rightarrow$
$2^8 \rightarrow$
$2^7 \rightarrow$
$\cdots$
$2^3 \rightarrow$
$2^2 \rightarrow$
$2^1 \rightarrow$
$2^0 \rightarrow$

Buddy Alloc
free list

$2^{10} \rightarrow$
$2^9 \rightarrow$
$2^8 \rightarrow$
$2^7 \rightarrow$
$\cdots$
$2^3 \rightarrow$
$2^2 \rightarrow$
$2^1 \rightarrow$
$2^0 \rightarrow$

Buddy Alloc
free list

$2^{10} \rightarrow$
$2^9 \rightarrow$
$2^8 \rightarrow$
$2^7 \rightarrow$
$\cdots$
$2^3 \rightarrow$
$2^2 \rightarrow$
$2^1 \rightarrow$
$2^0 \rightarrow$

RAM – Zones

0

ZONE_DMA

16MB

ZONE_NORMAL

896MB

ZONE_HIGHMEM

4GB/8GB

# Slab allocator

3KB

```
ptr = kmalloc(size, mask);
    :
kfree(ptr);
```

↳ GFP_XYZ

- allocates smaller contiguous block from slab cache.
- lesser internal fragmentation.
- faster allocation for commonly required (well-known) objects. e.g. task_struct, mm_struct, inode, ...

Example:
- object size = 3KB
- pages per slab = 4
- objects per slab = 5
- slabs per cache = 3
- num of objs = 15
- active obj = 8

← full

← partial

← empty

Slab cache

- A slab is set of contiguous pages.
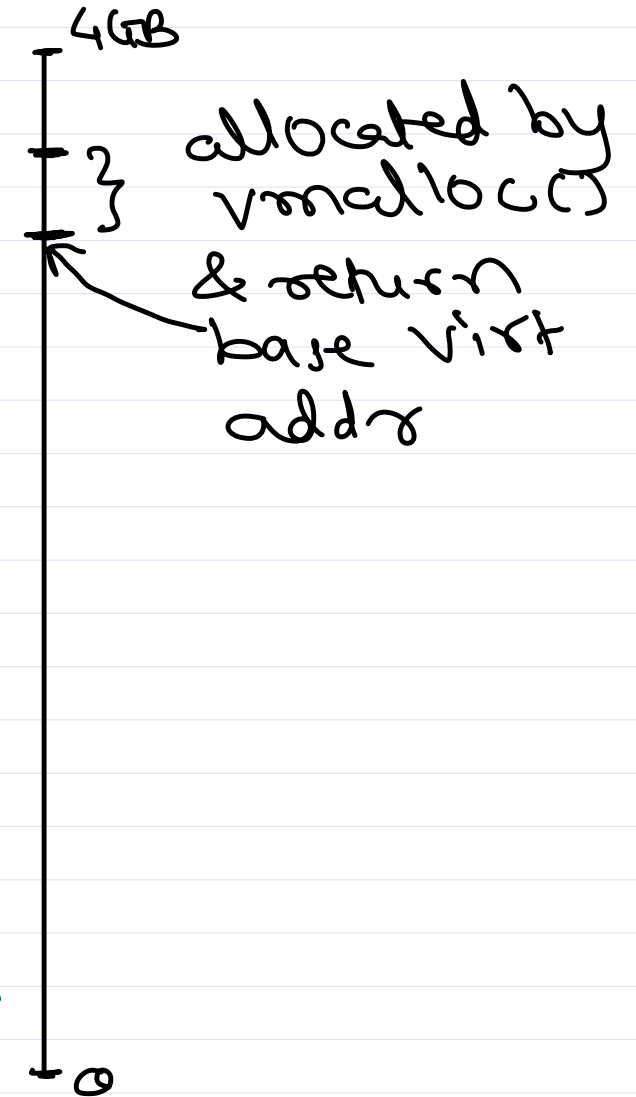- Slab allocator allocates slab - using buddy allocator.

# vmalloc()

In user space, contiguous virtual addresses are allocated by mmap(). It internally calls vmalloc().

To allocate contiguous virtual memory in kernel modules/code, use vmalloc().

```
ptr = vmalloc (size);
    :
    :
vfree (ptr);
```

✓ Allocates contiguous range of virtual address.
　　✓ allocates a new VAD (vm_area_struct).
　　✓ allocates corresponding page table entries.
✓ Actual physical pages will be allocated, when pages are accessed by CPU (read/write).

4GB

} allocated by vmalloc() & return base virt addr

0

# Page replacement

If RAM is full, existing page(s) needs to be swapped out - to make space available for new pages.

Local Replacement:
- To make space avail for process swap out a few pages of that process itself.

Global Replacement:
- To make space avail for process swap out a few pages of any process.

Kernel runs a daemon process which ensures that avail memory is not going below a threshold level. If needed, it will swap out pages from RAM (as per page replacement algo). This process is called as "Page stealing process". In Linux kernel, this is done by "kswapd" daemon.
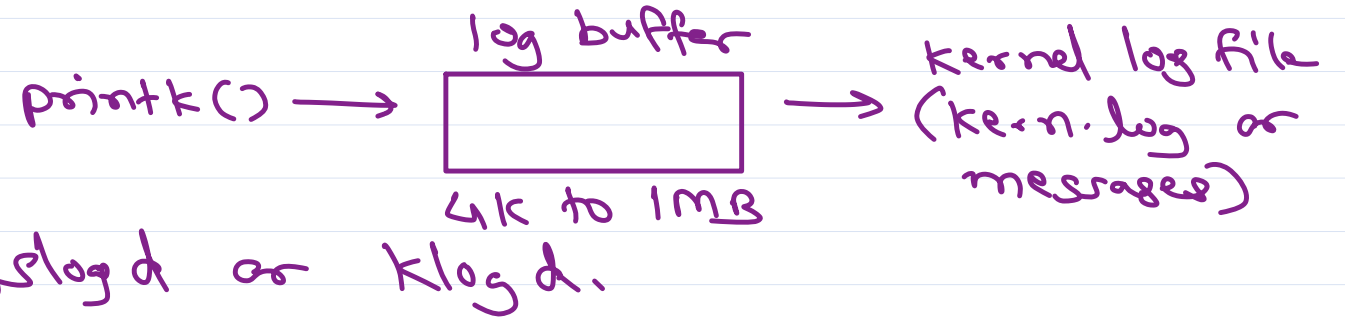
# Kernel debugging techniques

① **debugging by printing**
printk() → log levels
@ .... 7
EMERG...DEBUG
kernel log daemon → syslogd or klogd.

printk() ⟶ | log buffer | ⟶ kernel log file
4k to 1MB (kern.log or messages)

② **debugging by querying**
① ioctl() operation      ② procfs entry.

③ **debugging by watching**
strace command → shows which sys calls are called

④ **System faults/hangs**
① kernel OOPs message → register values + stack trace.
② kernel panic - crash - scheduler stopped

⑤ **kernel debuggers**
① gdb      ② kdb

# *Thank you!*

Nilesh Ghule <nilesh@sunbeaminfo.com>