

Kernel_4

노트북: SW
만든 날짜: 2018-04-10 오후 2:17
작성자: fstopdg
태그: :cs find 0
URL: <http://complicated0idea.tistory.com/14>

수정한 날짜: 2018-04-11 오전 12:39

2018. 4. 10 화 - 34회차

과정 : TI, DSP, Xilinx Zynq FPGA, MCU 기반의 프로그래밍 전문가 과정

Prof. 이상훈
gcccompil3r@gmail.com
Stu. 정상용
fstopdg@gmail.com

Chapter 4 : 메모리관리

-> The information about memory : mm_struct, slack(when we practiced a lot of programme, we put buf_size like buf[32], buf[64]. slack is about this size), buddy)

-> Virtual Memory, Physical Memory, Paging

1. 메모리 관리 기법과 가상메모리

-> 물리메모리의 한계를 극복하기 위하여 '가상메모리' 개발(i.e Cortex-A series 는 MMU를 달고 나온다. MMU Memory Management Unit)

-> A size of virtual memory : 4GB(32bit)

-> A size of physical memory : 28byte(4byte X 7)

-> The pros : 메모리배치정책이 불필요(kernel이 알아서 해결), 태스크의 빠른 생성이 가능(physical memory를 직접 복사가 아닌 task_struct만 복사하면 되므로), 태스크간 메모리 공유/보호가 쉬움

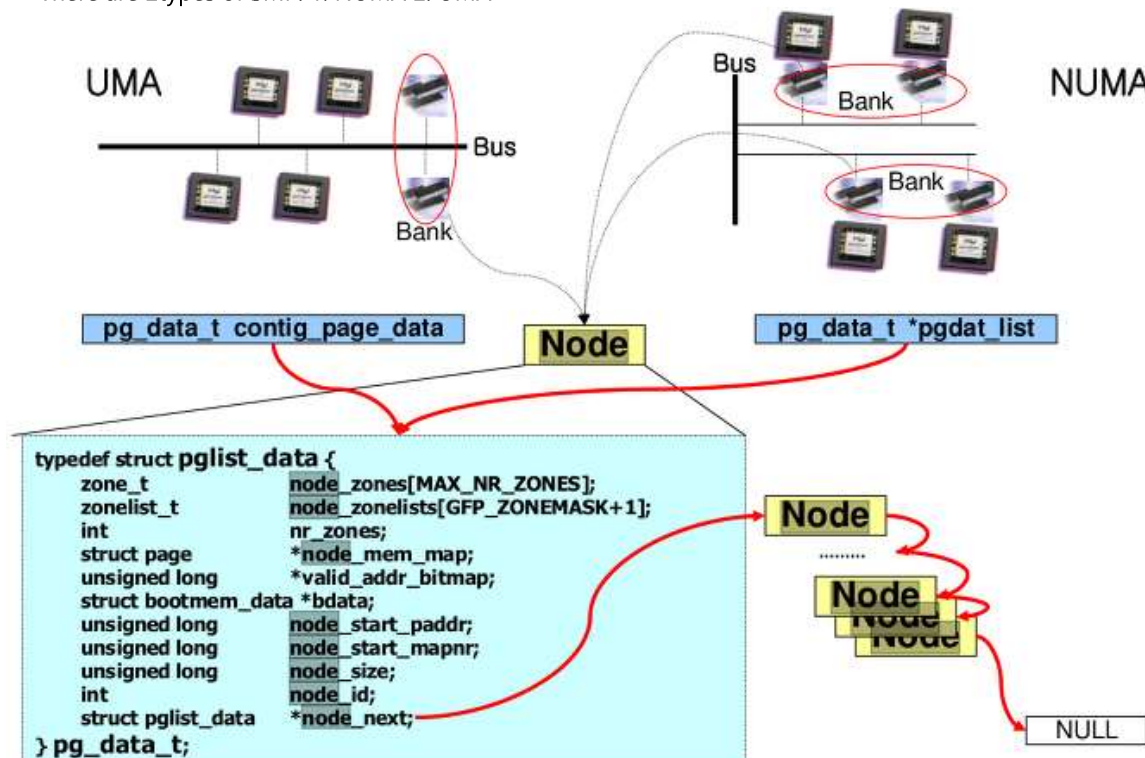
2. 물리메모리 관리 자료구조(How to work physical memory in linux : Node -> Zone -> Page frame)

-> Bootloader : 부팅시 물리메모리의 최소공간을 할당해놓음(i.e wait queue를 받을 공간, run queue를 받을 공간등 나중에 사용될 공간을 미리 할당)

-> 2 Purposes : First, How to express a physical memory in linux. Second, Which policies do people use to handle efficiently the limited physical memory.(The policies : paging, buddy, UMA, NUMA...)

-> SMP(Symmetric Multiprocessing) : All CPU and the memories share I/O bus. It causes a bottleneck phenomenon.

-> There are 2types of SMP. 1. NUMA 2. UMA



2.1 Node

-> bank : A group of memories have a same access velocity.

-> if there is a bank : UMA
-> if there is banks : NUMA
-> **node** is expressed by pglist_data (= pg_data_t) : **bank** struct
i.e A number of node = A number of bank

```

634 * per-zone basis
635 */
636 struct bootmem_data;
637 typedef struct pglist_data {
638     struct zone node_zones[MAX_NR_ZONES];
639     struct zonelist node_zones_lists[MAX_ZONELISTS];
640     int nr_zones;
641 #ifdef CONFIG_FLAT_NODE_MEM_MAP /* means !SPARSEMEM */
642     struct page *node_mem_map;
643 #endif
644 #ifdef CONFIG_PAGE_EXTENSION
645     struct page_ext *node_page_ext;
646 #endif
647 #ifdef CONFIG_NO_BOOTMEM
648     struct bootmem_data *bdata;
649 #endif
650 #ifdef CONFIG_MEMORY_HOTPLUG
651     /*
652      * Must be held any time you expect node_start_pfn, node_present_pages
653      * or node_spanned_pages stay constant. Holding this will also
654      * guarantee that any pfn_valid() stays that way.
655      */
656     pgdat_resize_lock() and pgdat_resize_unlock() are provided to
657     manipulate node_size_lock without checking for CONFIG_MEMORY_HOTPLUG.
658     *
659     * Nests above zone->lock and zone->span_seqlock
660     */
661     spinlock_t node_size_lock;
662 #endif
663     unsigned long node_start_pfn;
664     unsigned long node_present_pages; /* total number of physical pages */
665     unsigned long node_spanned_pages; /* total size of physical page
666     range, including holes */
667     int node_id;
668     wait_queue_head_t kswapd_wait;
669     wait_queue_head_t pfmemalloc_wait;
670     struct task_struct *kswapd; /* Protected by
671     mem_hotplug_begin/end() */
672     int kswapd_max_order;
673     enum zone_type classzone_idx;
674 #ifdef CONFIG_NUMA_BALANCING
675     /* Lock serializing the migrate rate limiting window */
676     spinlock_t numabalancing_migrate_lock;
677     /* Rate limiting time interval */
678     unsigned long numabalancing_migrate_next_window;
679     /* Number of pages migrated during the rate limiting time interval */
680     unsigned long numabalancing_migrate_nr_pages;
681 #endif
682 #ifdef CONFIG_DEFERRED_STRUCT_PAGE_INIT
683     /*
684     * If memory initialisation on large machines is deferred then this
685     * is the first PFN that needs to be initialised.
686     */
687     unsigned long first_deferred_pfn;
688 #endif
689 #endif /* CONFIG_DEFERRED_STRUCT_PAGE_INIT */
690 } pg_data_t;
691 #define node_present_pages(nid) (NODE_DATA(nid)->node_present_pages)

```

-> a node of UMA : contig_page_data = pglist_data
-> Conclusion : pglist_data cares all node. & pglist_data = bank struct
-> p.98 2nd paragraph : Cache

2.2 Zone

-> It's located in struct pglist_data
-> The information of a zone.
zone : cat /proc/zoneinfo - What we can see : DMA, DM32, Normal, Movable(There is no ISA.)
-> zone : 메모리의 특정한 영역(메모리를 역할에 따라서 나눠주는 역할)

Name	role	Physical memory
ZONE_DMA	DMA stands for Direct Memory Access -> It means DMA is a feature of computer systems that allows certain hardware subsystems to access main system memory independently of the CPU. When a size of data is kind of big, it is used. e.g network, video data, sound data -> V4L2, ALSA	x < 16MB
ZONE_NORMAL	1:1 match	16MB < x < 895MB
ZONE_HIGHMEM		896MB < x

-> 16MB이상 : ZONE_NORMAL
-> ZONE_HIGHMEM 896MB기준으로 그 이하 : 1:1로 대응 / 그 이상 : 간접참조(Paging과는 다른 :10bit로 인덱스찾고, 10bit로 인덱스 찾고, 12bit로 실제 물리메모리 찾아가는 방식)

2.3 Page frame

-> The smallest unit of memory
-> Page frame is under the struct page
-> This is, each page frame has struct page(therefore, a page's size < a page frame's)

3.Buddy와 Slab

->Q. When we allocate certain data has very small size(30Byte, 60Byte) or too big one compared with the size of page frame(4KB), what happen to page frame?

A. Slab allocator and Buddy allocator solve the problem(Internal Fragmentation and External Fragmentation).

3.1 버디 할당자(Buddy Allocator : node -> zone -> free_area)

-> it is located on array free_area which is put in struct zone.

it means each free_area has a buddy.

```
/ include / linux / mmzone.h
94
95
96 struct free_area {
97     struct list_head    free_list[MIGRATE_TYPES];
98     unsigned long        nr_free;
99 };
100
```

-> Array free_area's index : 0~9(0~9는 해당 free_area가 관리하는 할당의 크기를 나타낸다)

최대할당크기 : $2^{10} \times 4KB$ (512 X 4KB)

cf

1. How to use ':cs find 0'

when it is impossible use ':cs find 0 dl_rq', type ctag -R on terminal. And, type 'vi -t task_struct' in kernel/linux-4.4

2. <http://atsequence.tistory.com/30>

3. <http://complicated0idea.tistory.com/14>