

# Per-cpu -1- (Basic)

📅 2016-02-11 (<http://jake.dothome.co.kr/per-cpu/>) 👤 Moon Young-il

(<http://jake.dothome.co.kr/author/admin/>) 📁 Linux Kernel (<http://jake.dothome.co.kr/category/linux/>)

<kernel v5.0>

## Per-cpu -1- (Basic)

Providing independent memory space for each CPU for the same variable is called a per-cpu variable. In SMP (Symmetry Multi Processing), synchronization mechanisms such as cache consistency and spinlock are required for resources that the cpu accesses at the same time, but per-cpu variables allow access to individual spaces, so there is no need for a synchronization mechanism, which is faster.

Per-CPU is a memory management technique that allocates variables separately by the number of CPUs. This separation of space for each CPU eliminates the need for a locking mechanism between CPUs, which provides fast performance and increases the cache hit rate of the processor.

## Key Features of Per-CPU

The per-CPU area was added in kernel 2.6. By creating a per-cpu zone of data so that each cpu has a copy, there is less to consider synchronization if you only have access to data that is relevant to your own cpu in an SMP environment. It is used in designs that do not use locks and want to repeat fast writes. Protection against task preemption and interrupts is necessary.

It is often used for counting network packets, disk, and kernel objects, and it is fast because it is performed without the use of a locking mechanism when thousands of updates per second are required. It caches effectively on each CPU and is fast because it doesn't have cache line bouncing issues with shared memory.

In the initial design, each per-CPU data was used as byte aligned as the number of cache lines, but now the kernel declares per-CPU data separately according to its type, so it is used in conjunction with other per-CPU data without sorting (read-most, etc.).

Per-CPU data is defined in a separate per-CPU address space rather than the regular code or data address space used by the kernel, and can be used to filter out zone invasion errors when using a sparse static code analysis tool. The per-CPU data statically defined in the kernel image is organized in the first chunk. Dynamic per-CPU allocation was added as a patch in 2009, introducing chunk. The asynchronous chunk population was added to kernel 2014.9-rc3 in September 18.

- Consultation:

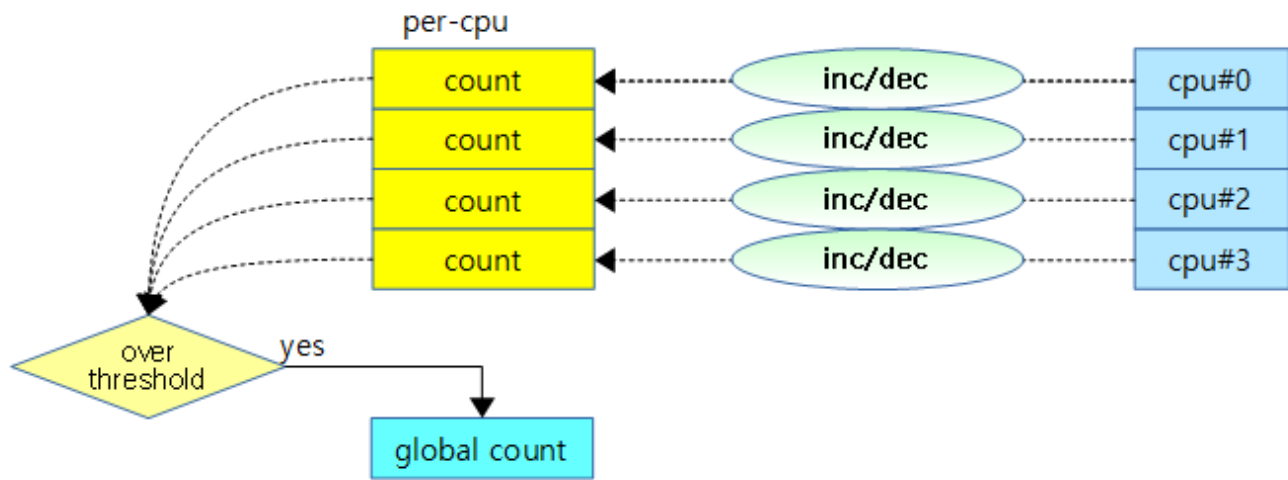
- percpu: implement asynchronous chunk population  
(<https://github.com/torvalds/linux/commit/1a4d76076cda69b0abf15463a8cebc172406da25>)
- percpu: implmeent pcpu\_nr\_empty\_pop\_pages and chunk->nr\_populated  
(<https://github.com/torvalds/linux/commit/b539b87fed37ffc16c89a6bc3beca2d7aed82e1c>)

## Per-CPU Uses

Let's take a look at some of the most common uses for per-CPUs.

- Cache Objects
  - Lists and trees used throughout the world are allocated to each CPU and used as caches. While the performance is fast, the disadvantage is that it uses as much memory as the number of CPUs, which increases the memory requirements. To compensate for this, it acts as a cache that does not manage the entire list or tree, but only provides a limited part of it ( PCP cache on the buddy system, Per-CPU cache applied to slab objects, LRU cache).
- Stats Counter
  - Various statistical counters are very frequent with simple increases and decreases. If you use the stat counter variable as a per-CPU, you can use it without a lock mechanism to remember the value of each CPU increase or decrease, and if the threshold is exceeded, it will be added to the global stat counter variable. If it is okay to use an approximate value of the counter value, only the value of the global statistic counter variable should be read, and if a precise value is required, the value of the per-cpu statistic counter that has been increased or decreased on each CPU should be applied to the value of the global stat counter. In order to guarantee the accuracy of the global counter to a certain extent, if the size of the increase or decrease on each CPU exceeds a certain threshold, the global counter is updated (the update of the global counter uses atomic operations).
- Others
  - RCU infrastructure
  - Profiling, Ftrace
  - VFS components

The following figure shows a case where it is used as a per-cpu counter.



(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-12b.png>)

## Unit: How much space is used to store per-CPU data per CPU?

A unit is a space where per-CPU data corresponding to a single CPU is stored. The unit is divided into three types of storage areas: The size of the unit is the value of the following three areas added and sorted upwards on a page-by-page basis.

- Static Area
  - Static per-CPU data is stored in the kernel using the `DEFINE_PER_CPU( )` macro.
  - RPI2: Yes) 0x3ec0
- reserved area
  - Static per-CPU data using the `DEFINE_PER_CPU( )` macro in the module is stored.
  - rpi2: e.g. default 0x2000 (8K)
- Dynamic Zones
  - **Dynamic** per-CPU data is stored via the `alloc_percpu( )` function in the kernel or module.
  - rpi2 e.g. default 0x5000 (20K)
    - On 64-bit systems, the default 28K

There are two main variables involved

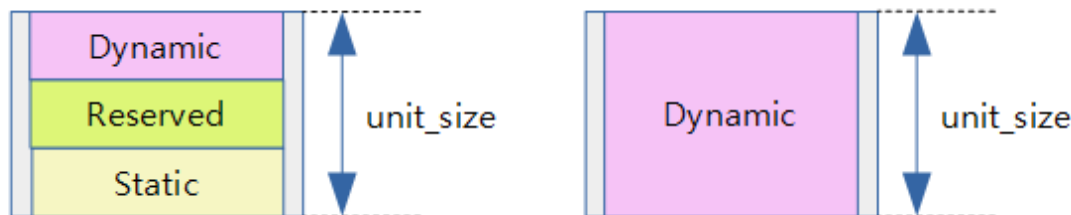
- `pcpu_unit_pages`
  - It contains the number of pages used by one unit.
  - RPI2: Yes) 11
- `pcpu_unit_size`
  - It contains the byte size used by one unit.
  - rpi2: Yes) 0xb000 (44 KB = 11 \* 4 KB (PAGE\_SIZE))

For systems with 4K small pages, such as ARM, multiple pages are used in a single unit when allocating per-cpu area. And use that area in a page-by-page arrangement. Add all the static\_size + reserved\_size + dyn\_size and align the 1K and add the remaining space to the dynamic area.

- rpi2: e.g.  $0x3ec0 + 0x2000 + 0x5000 + 0x140$  (4K align to add extra space to the dynamic area)

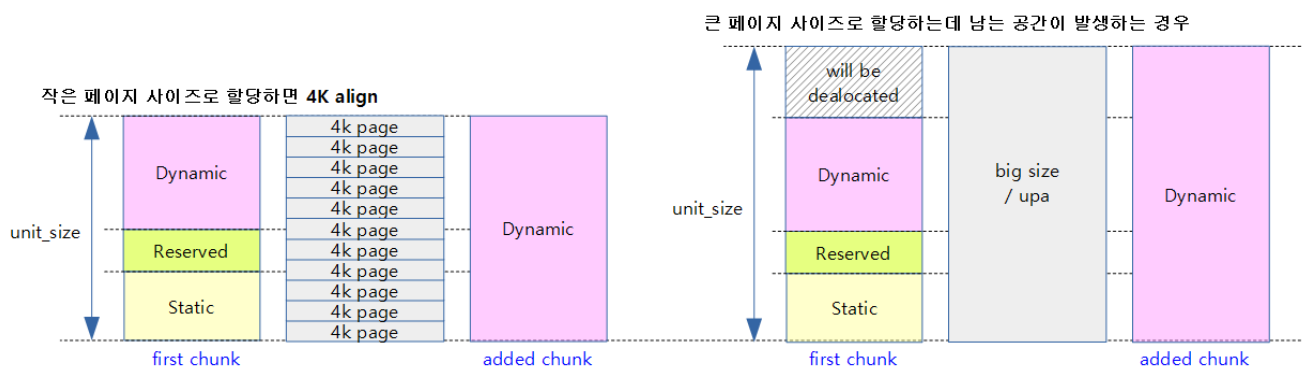
In some architectures, such as x86\_64 and ARM64, large pages, such as 2M pages, can be allocated using them. Due to the large size of this space, multiple units can be divided into one allotment size. In the process of placing a unit, there is a surplus space at the end of the unit, which is released from memory allocation.

The following figure shows that when the per-CPU area was first created, it was divided into three zones. Also, when creating an additional per-cpu region, it shows a unit containing only the dynamic region with the same unit size.



The following illustration shows that the assignment page is divided into two types: small and large.

- Compare the dynamic areas between the first chunk and the added chunk.
- You can see that only the first chunk has static and reserved areas.



([http://jake.dothome.co.kr/wp-content/uploads/2016/02/setup\\_per\\_cpu\\_areas-17a.png](http://jake.dothome.co.kr/wp-content/uploads/2016/02/setup_per_cpu_areas-17a.png))

## CPU -> Unit Mapping

Per-CPU data uses NR\_CPUS arrays, so using a larger number than the actual number of CPUs in existence can result in wasted memory. This is because a system that is configured to support thousands of CPUs (NR\_CPUS) sometimes has only a few CPUs that are capable of running at actual boot.

- Note: cpu alloc v1: Optimize by removing arrays of pointers to per cpu objects  
(<https://lwn.net/Articles/258248/>) | LWN.net

In order to reduce wasted per-CPU arrays, additional patches have been applied in some architectures to make efficient use of space. In a NUMA system, each node has a different number of CPUs, and if you use a large page, the number of units is the same for each allocated memory. Instead of always mapping the cpu and units the same, we mapped them to a non-linear/sparse >unit for NUMA.

This can be easily implemented by adding a CPU->Unit mapping array and can be used to find matched units. When configured with an asymmetric NUMA system, the number of units and the number of CPUs are not the same because some units are not mapped.

The arrangement used here is as follows:

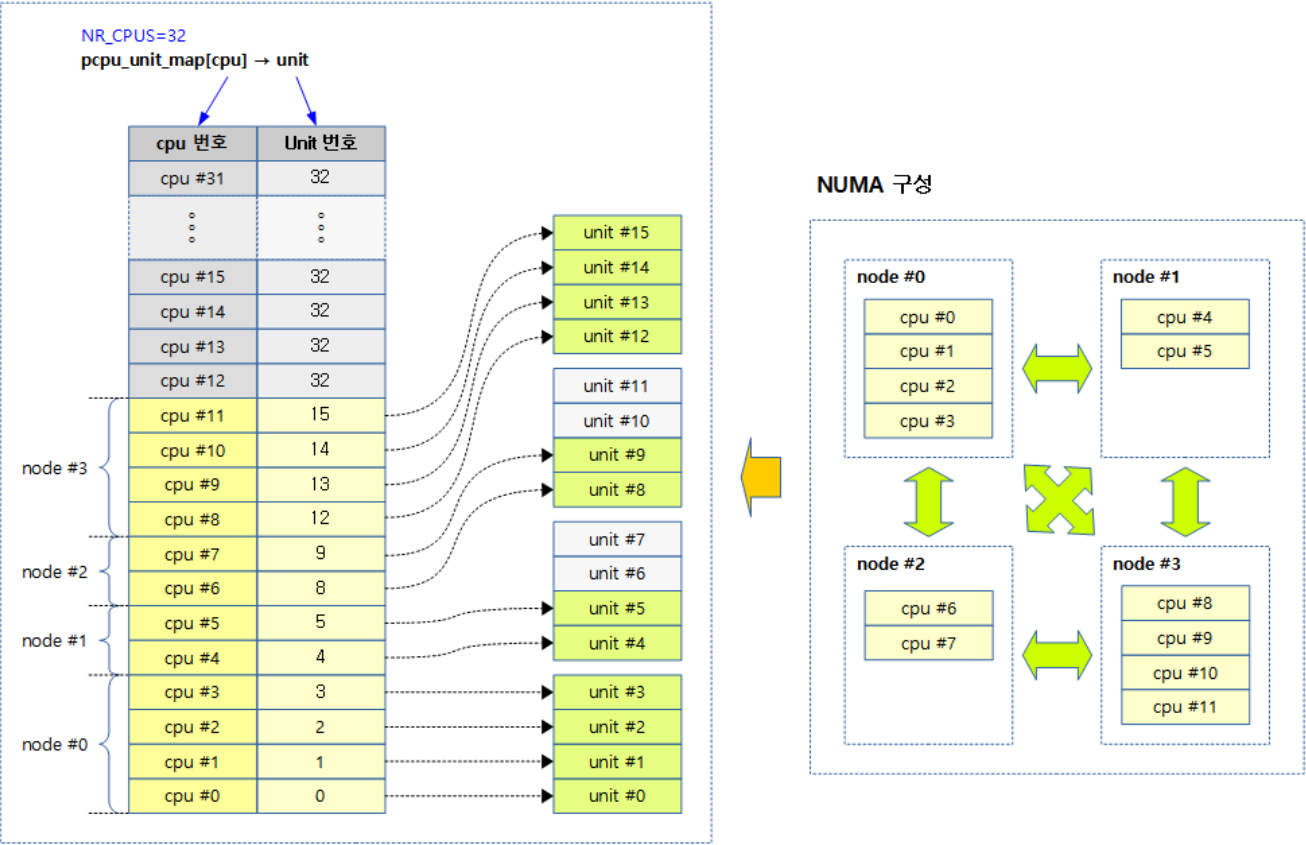
- `pcpu_unit_map[]`
  - The CPU ID is mapped to the Unit Index.
  - The mapping is valid for all chunks.
  - e.g. 0,1,2,3
- `pcpu_unit_offsets[]`
  - The baseaddr of that chunk contains the offset to each unit.
  - The unit offset is the same for all chunks.
  - e.g. 0, 0x8000, 0x10000, 0x18000

## NUMA support for per-CPU's

To support NUMA, per-cpu adds non-linear/sparse `cpu->unit` mapping to the same mapping of the `cpu` and `units`. This can be easily implemented by adding a CPU->Unit mapping array and can be used to find matched units. The base address of `unit0` is not fixed. If the mapping changes, the base address of `unit0` changes.

In ARM64, the group mapping structure information below is not used after kernel initialization, but rather global variables are used for the actual `cpu->unit` mapping. In the x86 and SPARC64 architectures, the CPU->unit mapping is done using the following two structure information even after kernel initialization.

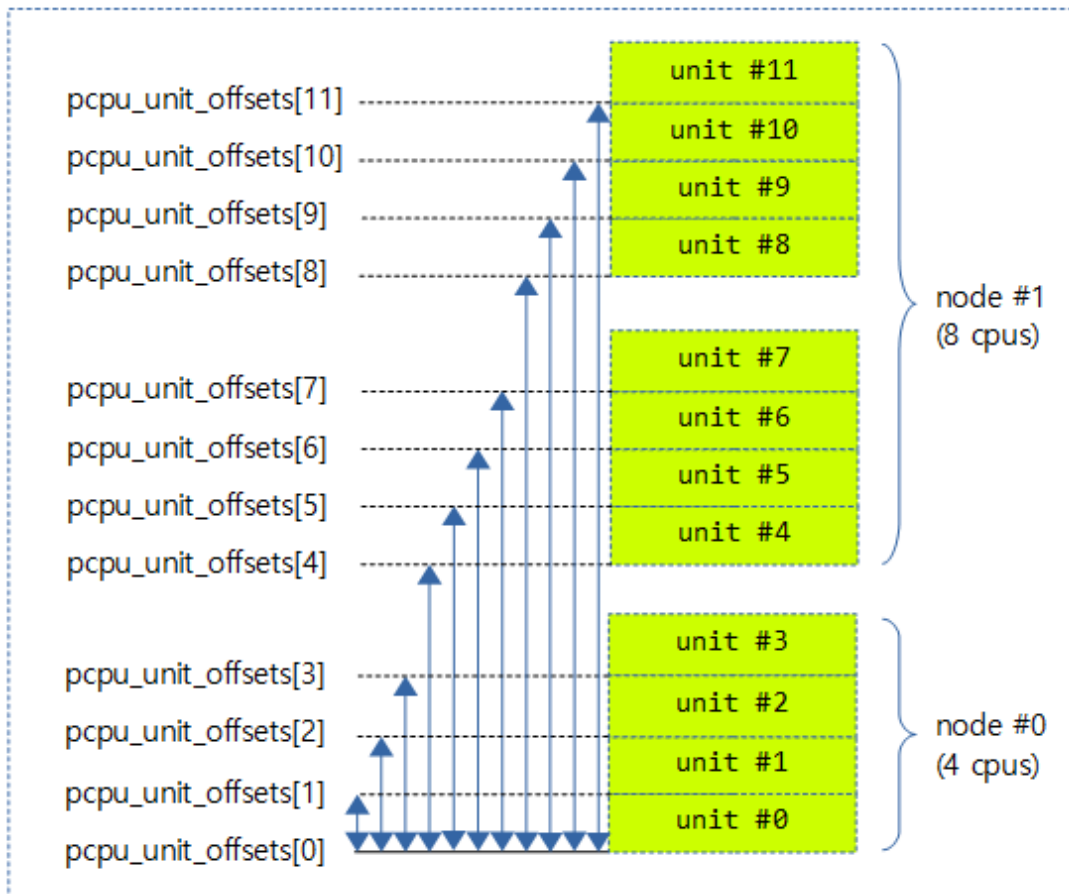
The following figure shows the unit number mapped to the CPU number available in the NUMA system. (Of the 32 `NR_CPUS`, there are 12 units that correspond to the actual 16 possible CPUs.)



(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-7a.png>)

The following diagram shows that the starting offset of each unit is specified from the unit at the lowest address.

- The initialization of the `pcpu_unit_offset[]` array is done by reference to the `pcpu_setup_first_chunk()` function.



(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-3b.png>)

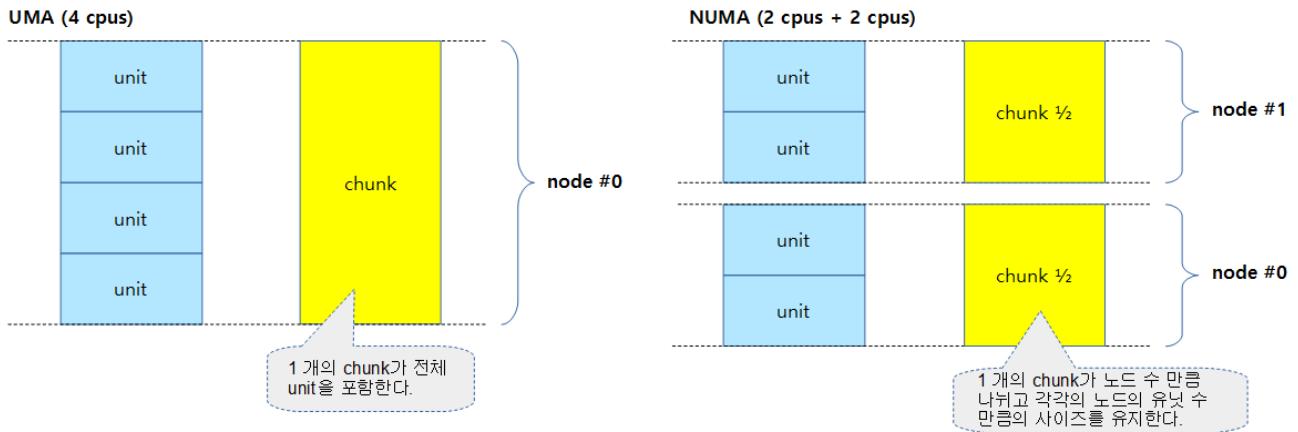
## Chunk

A chunk is a unit that collects all units corresponding to all CPUs, and the characteristics are as follows:

- One chunk is divided by `nr_units` number
- If you run out of chunks, create and use additional chunks
- Each unit in every added chunk that is not the first chunk only has a dynamic area
- All chunks are managed as `pcpu_chunk` structs
- In a NUMA system, one chunk is allocated using its own node (group) memory

The first chunk is the first chunk created during the kernel initialization process. Each unit in the First Chunk is divided into three areas: Static, Reserved, and Dynamic.

The following figure shows that a chunk is divided and allocated by the number of nodes.

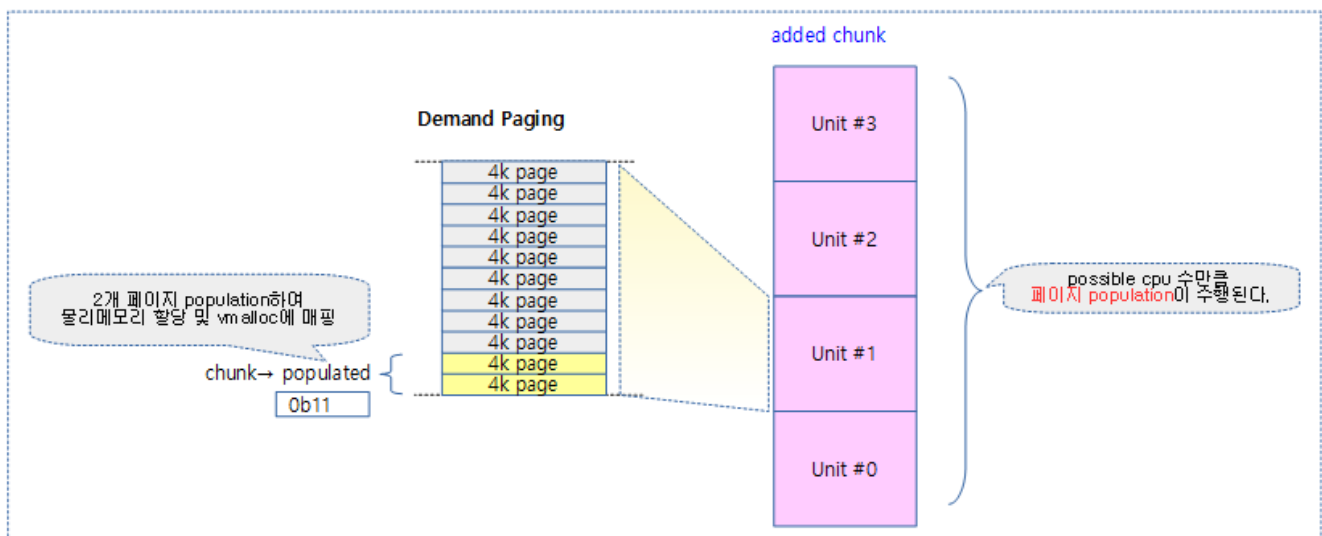


(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-13.png>)

## Demand Paging (Population by Page)

- All additional chunks use the Demand Paging method to save memory.
  - Exclude the first chunk because it is already mapped
- When a new chunk is allocated, it does not allocate and map pages contained in every unit, but only when there is an actual per-cpu allocation request, which allocates physical memory for those pages and maps them to vmalloc space. This is called the page population.
- Whether or not it is populated for the entire page configured within a chunk is managed by a bitmap, a populated member variable.
  - When per-CPU data is allocated and used, if the pages corresponding to the specified area of the chunk are not set to a populated bitmap, the allocation of physical memory pages and mapping by vmap is performed.

The following illustration shows the populated pages.



(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-30.png>)



# Configuring First Chunk

The first chunk is a chunk that plays a key role in determining the layout of each group and unit, and is specially managed. At compile time, `.data..percpu` section and copied to each static area of the first chunk in the boot init routine `setup_per_cpu_areas ( )`.

reserved areas are initialized to `PERCPU_MODULE_RESERVE` size. If you use a module, it will be 8K, and if you don't use a module, it will be 0K. The dynamic area is initialized to the `PERCPU_DYNAMIC_RESERVE` size. On a 32-bit system, it will be 20K, and on a 64-bit system, it will be 28K. Until the v3.17 kernel, they were smaller than they are now, at 12K and 20K, respectively.

There are two ways to do this, and you can choose one of them as a `percpu_alloc` option at boot time.

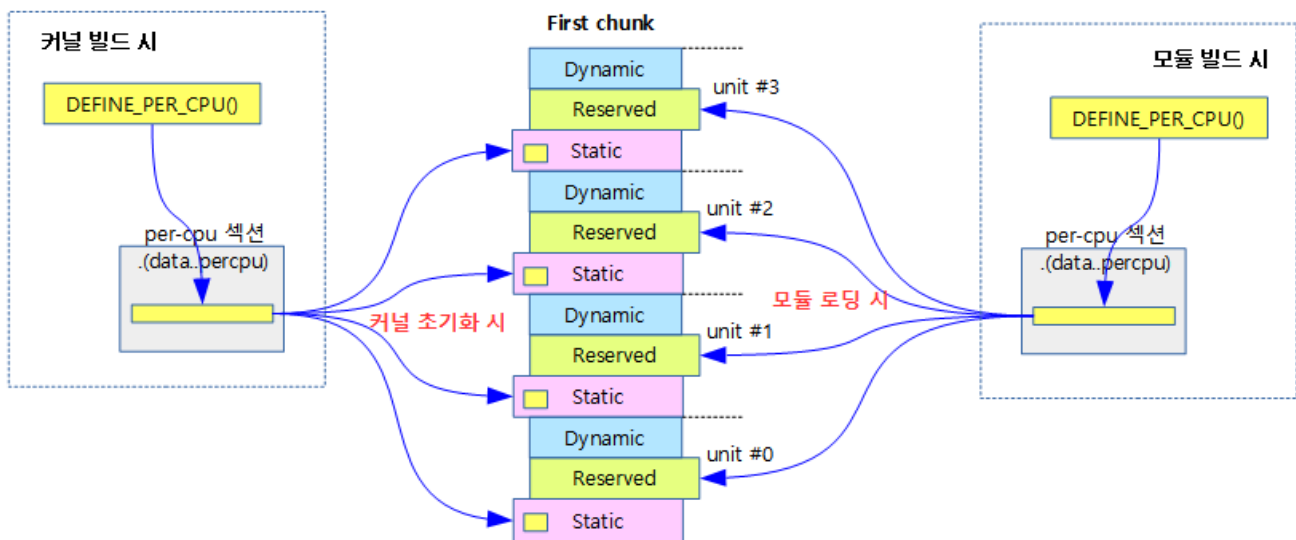
- Embed method
  - If your architecture supports it, use huge page sizes for unit allocation to increase the efficiency of the TLB cache.
    - ARM64 uses the embed method, but does not yet support the use of huge pages on per-cpu.
  - When constructing a first chunk, space is left behind the dynamic of each unit. This space can be deallocated and used for other purposes
  - In the case of UMA systems, the first chunk is allocated and used at a time.
  - In the case of a NUMA system, the first chunk is divided and allocated to each node's memory.
  - If you add chunks, you should allocate them from top to bottom in the `vmalloc` space so that they are spaced at the same spacing as the node-by-node configuration when you created the first chunk.
    - In the case of 32-bit NUMA systems, the `vmalloc` space is very small (`arm=240M`, `x86=120M`), so if the embed method cannot be used, it is switched to the page method.
      - In order to account for the placement of additional chunks of `vmalloc` space, the embed method is restricted if the farthest `max_distance` between the base addresses of the chunks allocated to each node exceeds 75% of the `vmalloc` space.
        - Note: `percpu: make embedding first chunk allocator check vmalloc space size`  
(<https://github.com/torvalds/linux/commit/6ea529a2037ce662fc6bfa572b46d47407d08805>)
      - On x86, if it can't be generated in an embed way, it will automatically retry using the page method.
  - Search for the empty `vmalloc` space from top to bottom.
    - When allocating `vmalloc` space by the `vmalloc()` and `vmap()` functions, the `vmalloc` space is searched from bottom to top, and the chunk of the per-cpu does the opposite, preventing chunks from overlapping each other with each additional chunk as much as possible.
    - When allocating the `vmalloc` space in the NUMA system, the per-cpu should be generated in the opposite direction to prevent it from mixing with the normal

vmalloc area, so that it can be generated quickly and the probability of failure is reduced.

◦ Paged Method

- When constructing the first chunk, it is allocated and mapped to the VMALLOC space for a minimum number of pages in the vmalloc space.
  - Since the slab memory allocator is not yet working at kernel initialization, mapping to the vmalloc region is not possible, so the early registration method is used so that the slab memory allocator can be registered after it has been activated.
  - After the slab memory allocator is running, call the `percpu_init_late()` function to read, allocate, and map all the chunk maps that you have already created.
- Instead of dividing chunks into nodes in vmalloc space, the number of units is configured in succession to the number of CPUs possible.
- Since it doesn't use large pages, it may be slower than the embed method, but it can still be used on 32-bit NUMA systems with less vmalloc space.
- Search for the empty vmalloc space from the bottom to the top.
  - Since you don't need to divide chunks by the number of nodes in the vmalloc space, you can use the same way you allocate vmalloc space by the `vmalloc()` and `vmap()` functions.

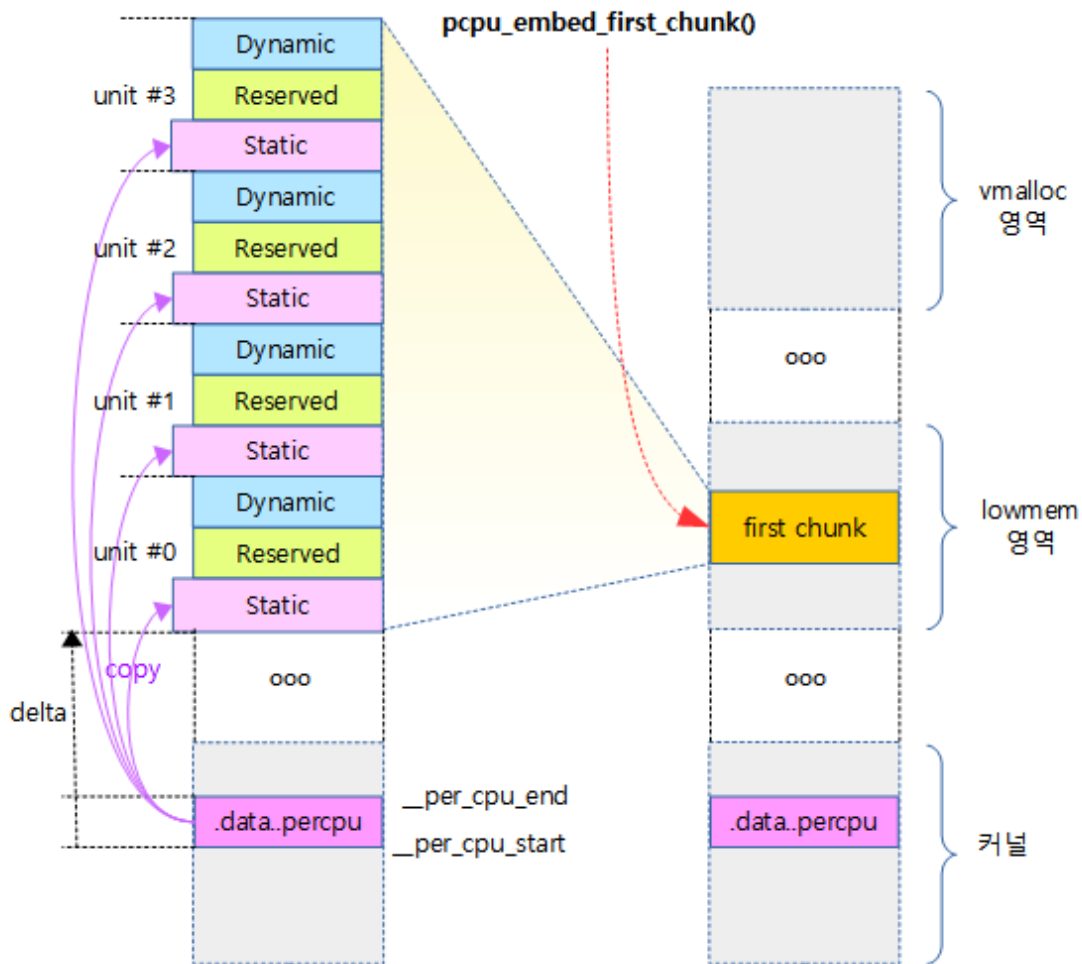
The following figure shows that static per-cpu data is declared and copied to the static area of each unit of the first chunk when the kernel is initialized after the build, and also copied to the reserved area of each unit of the first chunk when the module is loaded.



(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-14.png>)

The following figure shows the location of the space where the first chunk is created using the embed method in the UMA system.

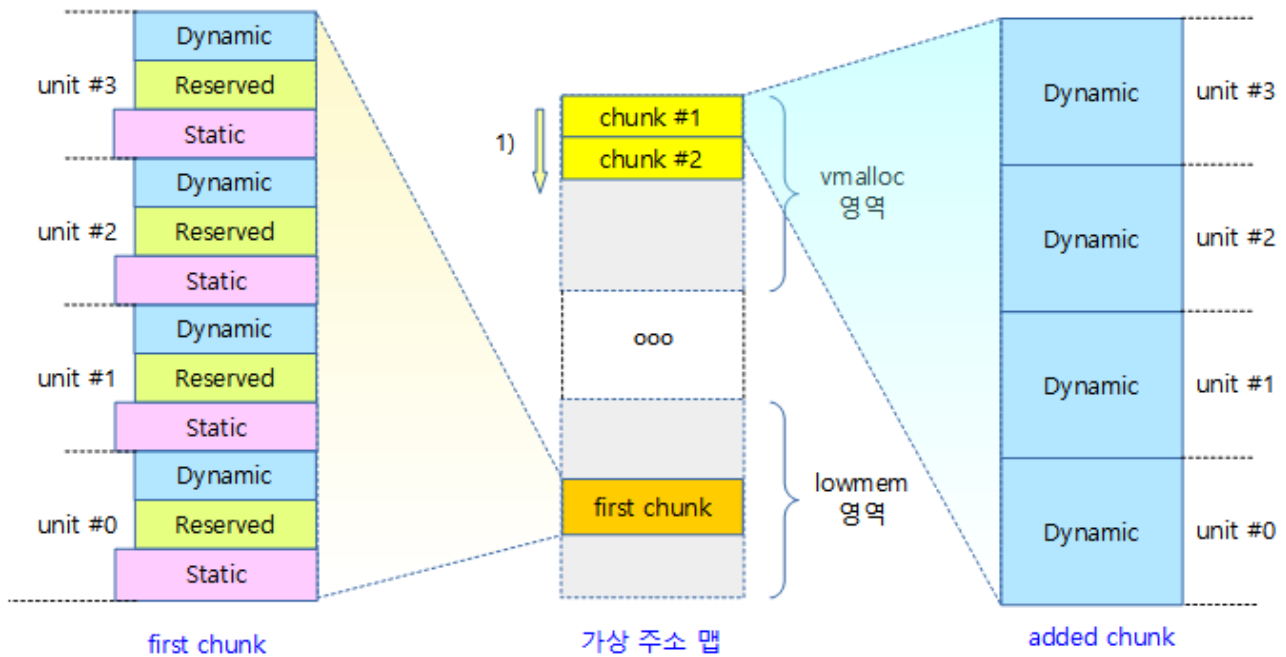
## UMA 시스템 - Embed 방식에서 first chunk 빌드



(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-15.png>)

The following figure shows the location of the space where additional chunks are created using the embed method in the UMA system.

## UMA 시스템 - Embed 방식에서 chunk 추가

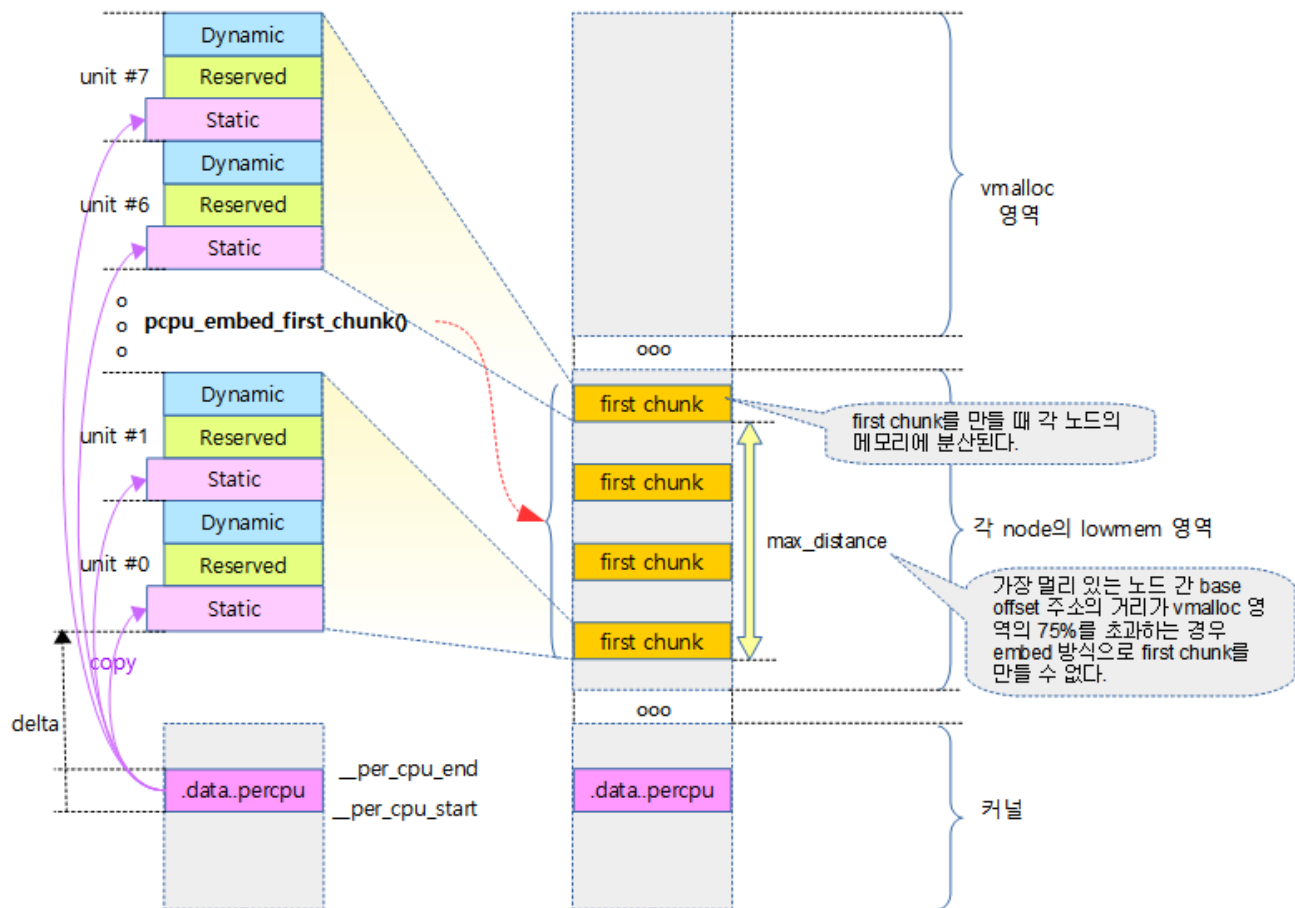


1) vmalloc space space는 일반적으로 아래에서 위로 할당해 올라가는데 embed 방식 first chunk를 만든 경우 per-cpu에 사용하는 chunk 공간은 거꾸로 top→down 방향으로 할당하여 사용한다.

(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-16.png>)

The following figure shows the location of the space where the first chunk is created using the embed method in the NUMA system.

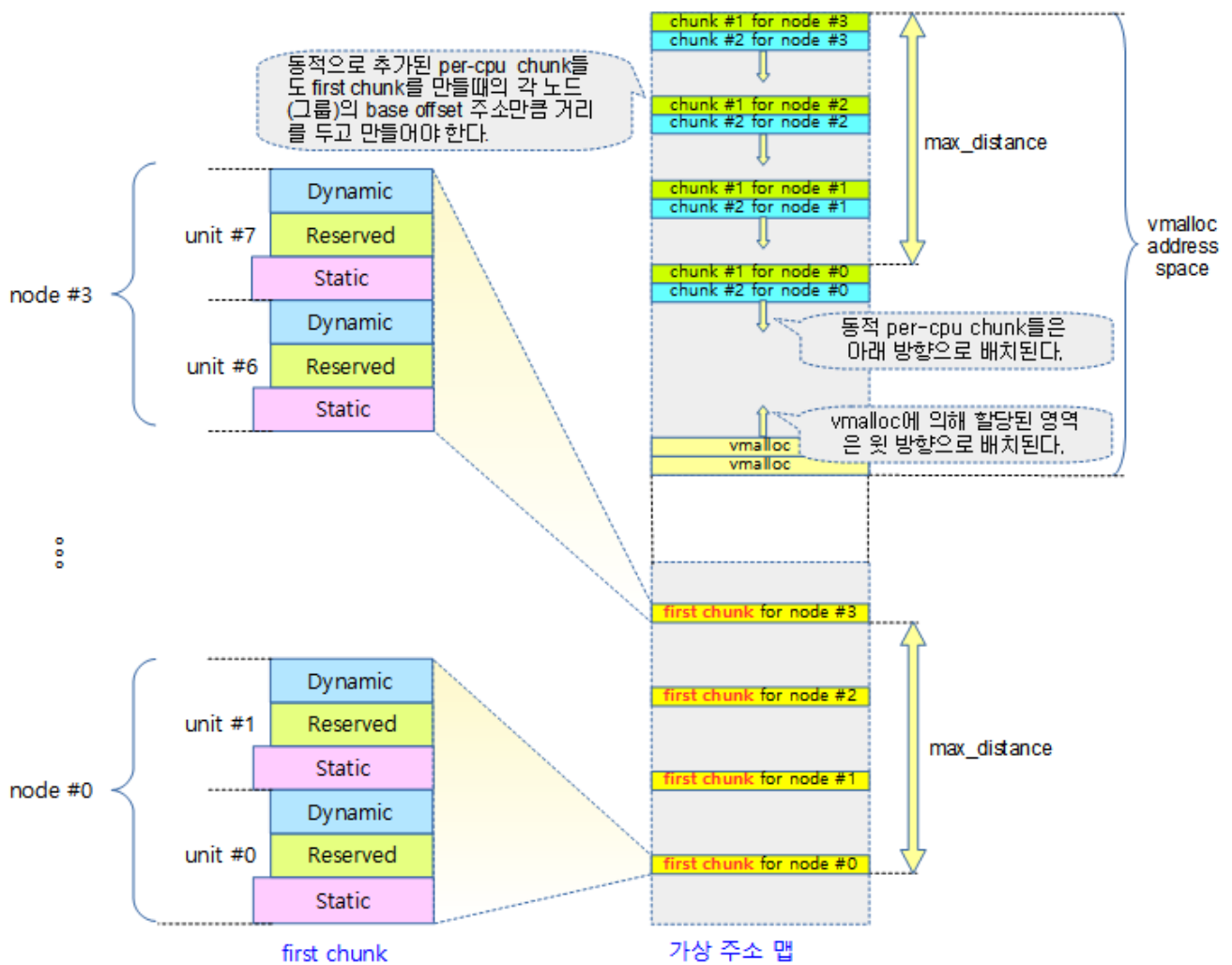
## NUMA 시스템 - Embed 방식에서 first chunk 빌드



(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-17.png>)

The following figure shows the location of the space in which additional chunks are created using the embed method in the NUMA system.

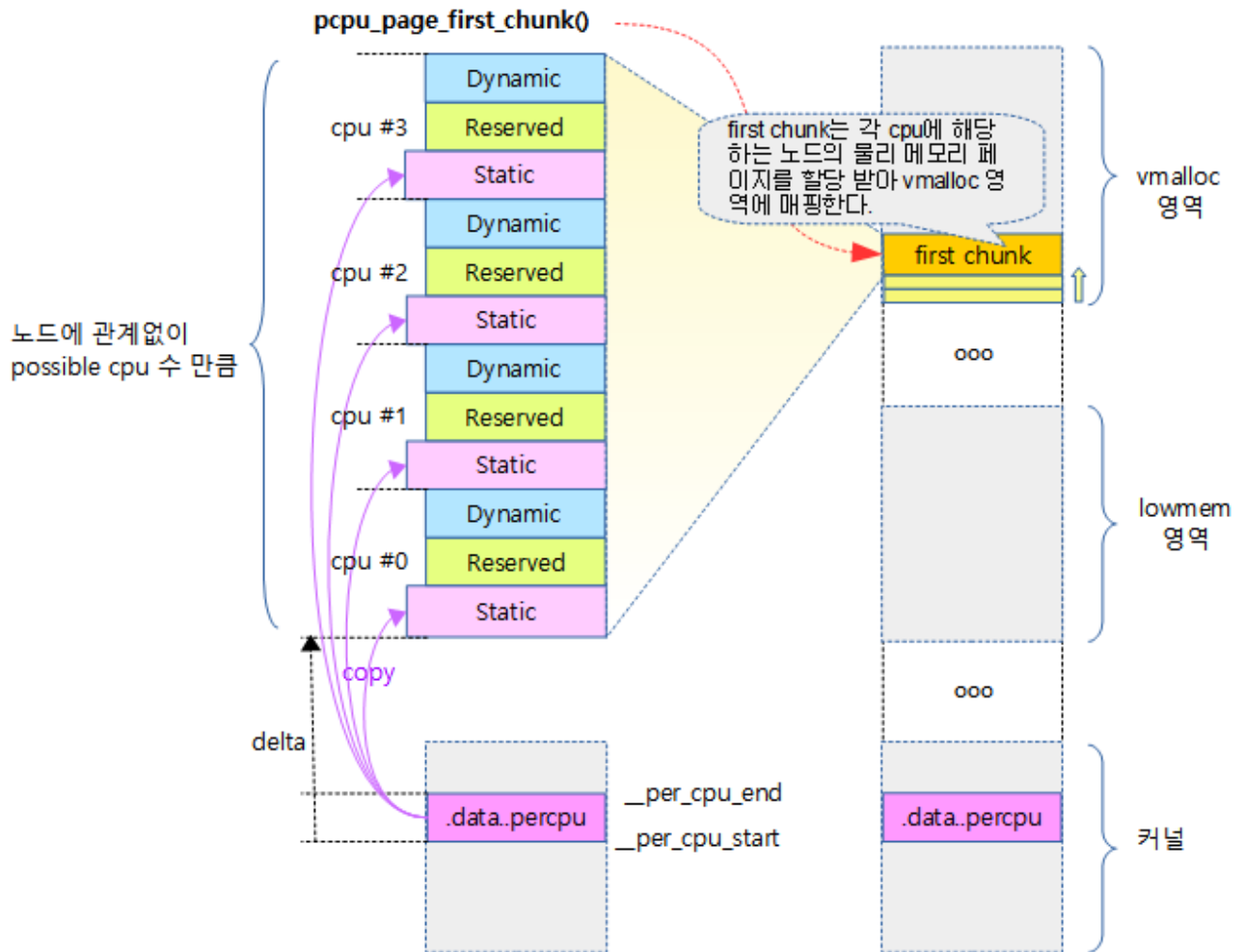
## NUMA 시스템 - Embed 방식에서 Chunk 추가



(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-18.png>)

The following figure shows the location of the space in which the first chunk is made using the Paged method.

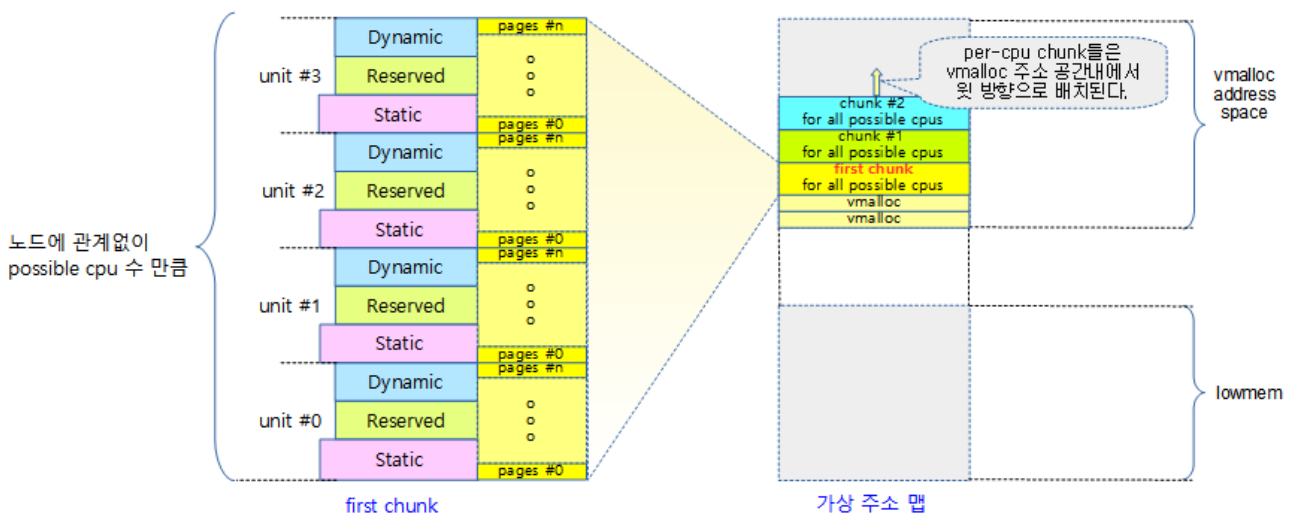
## UMA/NUMA 시스템 - Page 방식에서 chunk 빌드



(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-19.png>)

The following figure shows the location of the space where additional chunks are made using the Paged method.

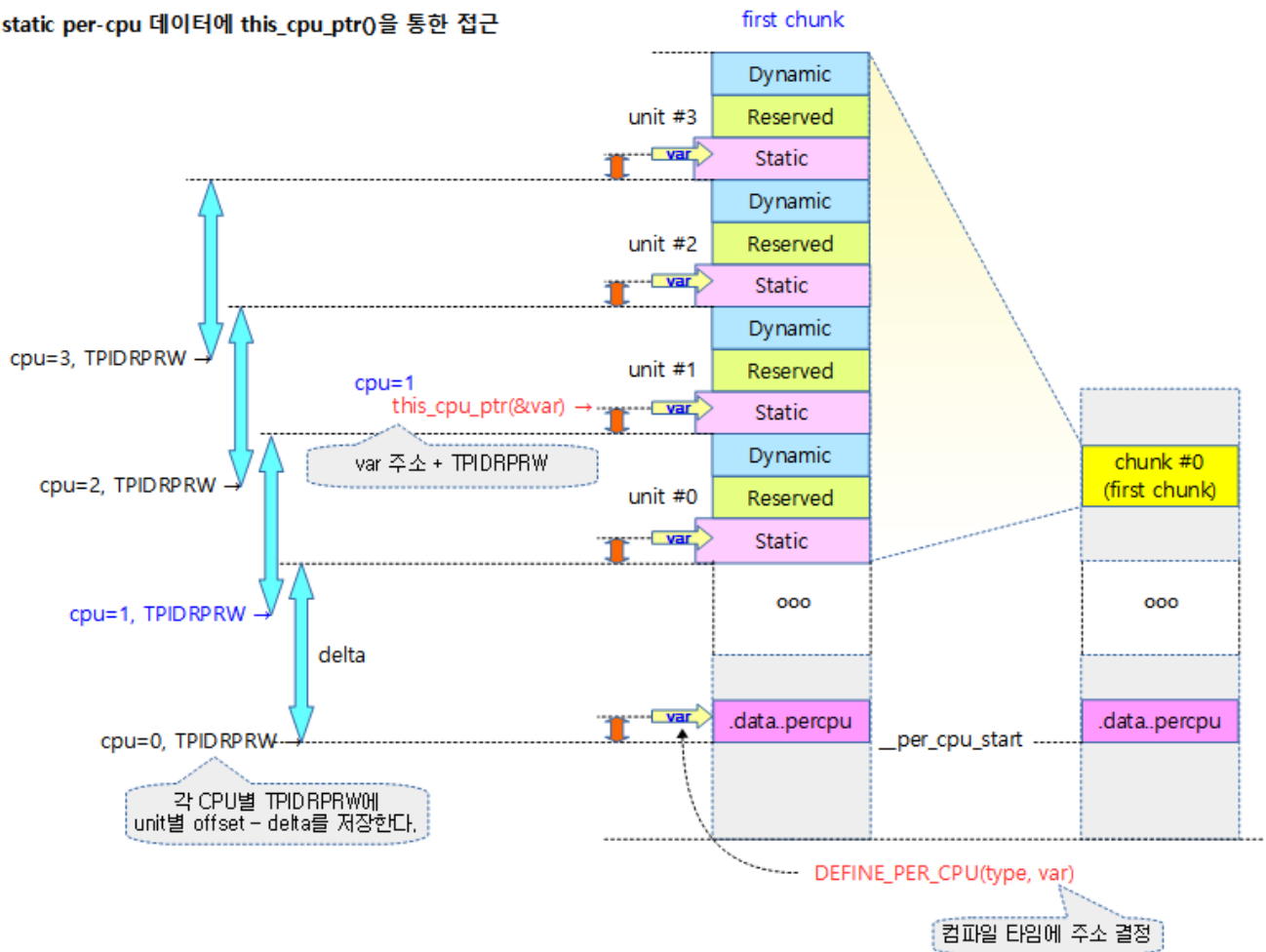
## NUMA 시스템 - Page 방식에서 chunk 추가



(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-20.png>)

The following figure shows the static per-cpu data accessed via the `this_cpu_ptr()` function.

static per-cpu 데이터에 `this_cpu_ptr()`을 통한 접근

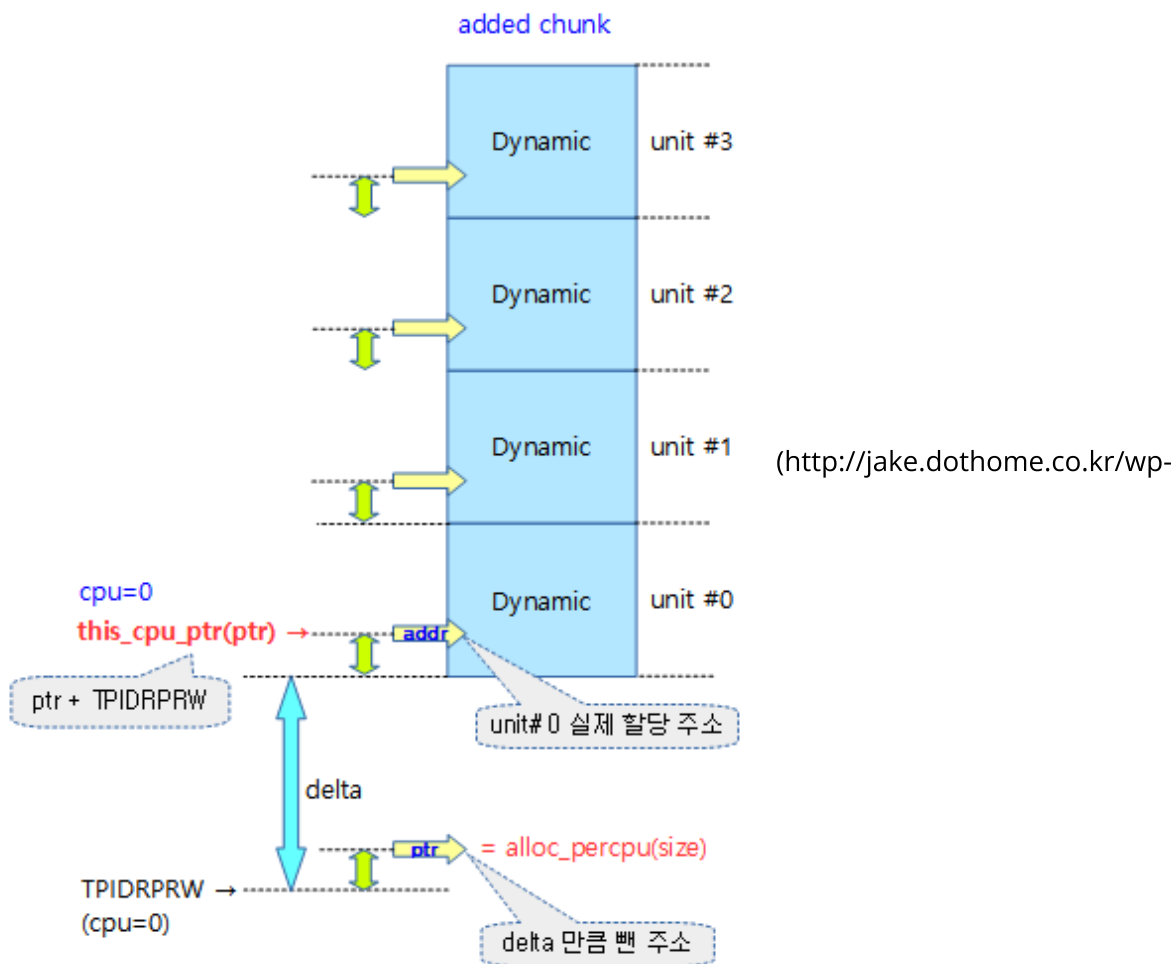


(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-21.png>)

The following figure shows accessing the dynamic per-cpu data of a chunk added via the `this_cpu_ptr()` function.



dynamic per-cpu 데이터에 `this_cpu_ptr()`을 통한 접근

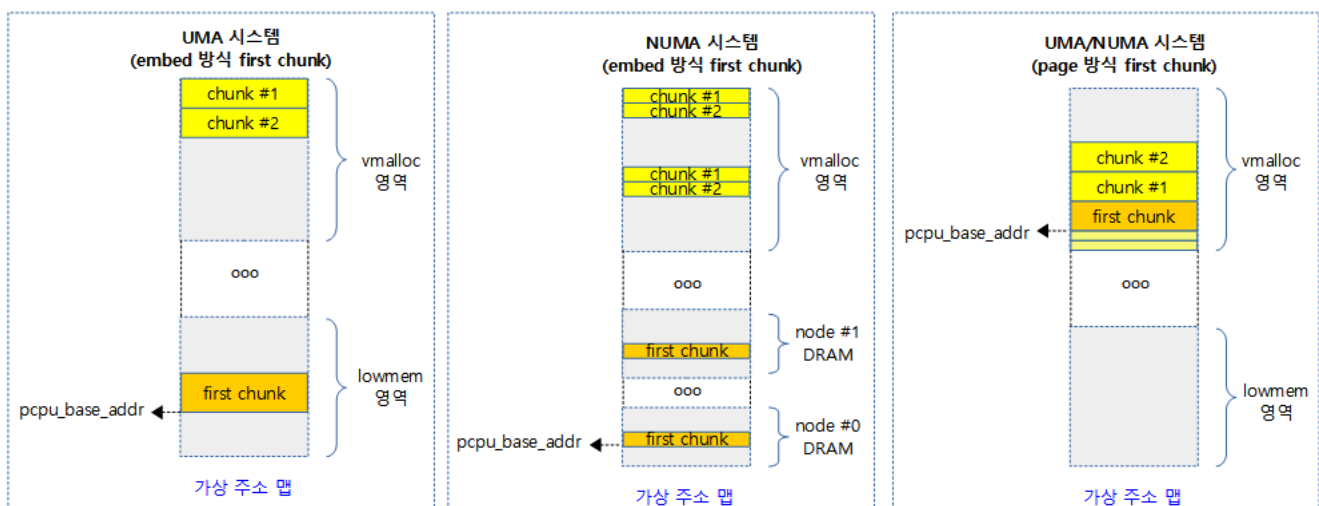


content/uploads/2016/02/percpu-22.png)

**pcpu\_base\_addr**

As shown in the following figure, the `pcpu_base_addr` contains the base address of the first chunk.

- Depending on the architecture, the vmalloc space may be located below the lowmem area. Note that ARM32 has a vmalloc space above lowmem, and ARM64 has a vmalloc space below lowmem.



(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-4a.png>)

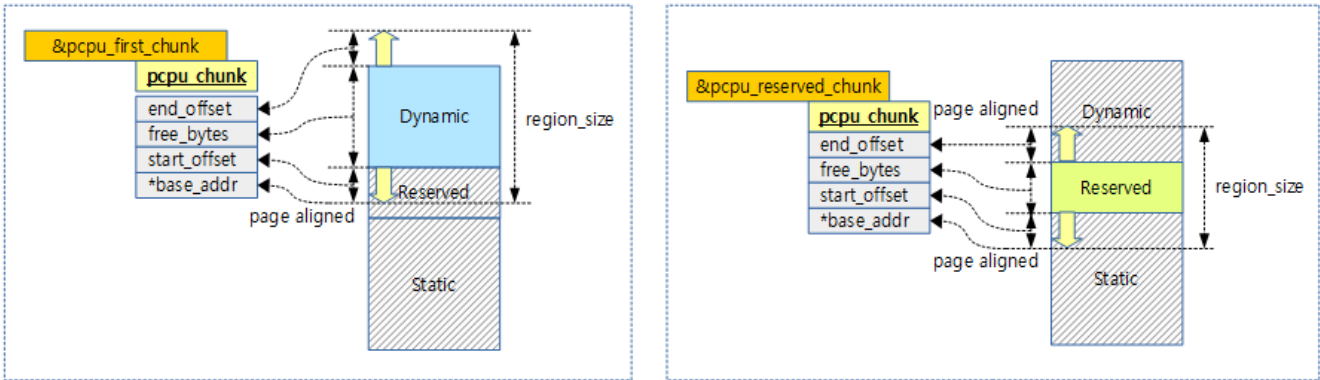
# Managing First Chunk's Map Entries

## Mapping Management in First chunk

These chunks are specially managed, playing an important role in determining the placement of each group and unit. For one allocation area of the first chunk, the management chunk map can consist of one or two. `pcpu_first_chunk` manages the dynamic space in the same way as other chunks, and is used in addition to `pcpu_slot[ ]`. `pcpu_reserved_chunk` is not added to the `pcpu_slot[ ]`, but is used by the module management mechanism.

If you're using modules in the kernel, `pcpu_reserved_chunk` have global variables manage the reserved areamap entries for the first chunk. Each time a kernel module is loaded or unloaded, the per-cpu data of the `DEFINE_PER_CPU( )` macro used by the kernel module is added to the reserved area.

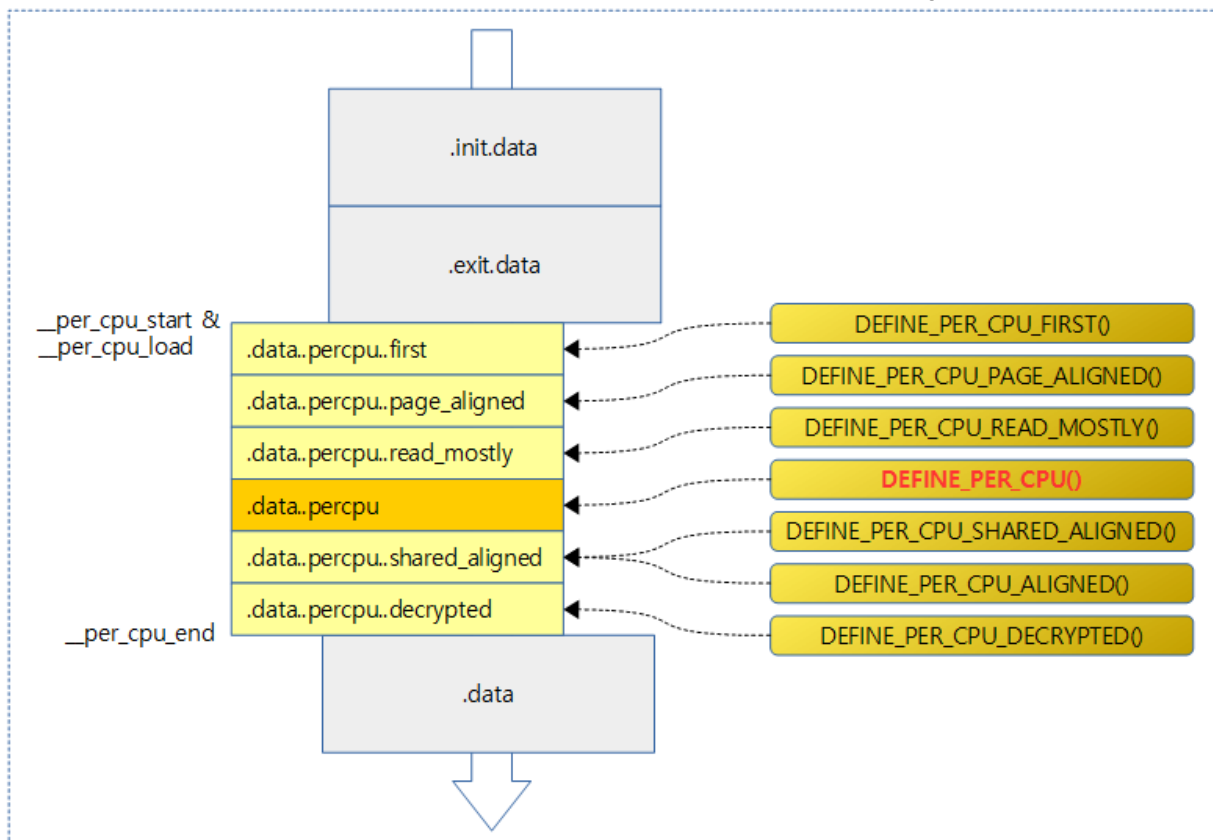
The following figure shows how the management structure of the first chunk is managed in either one or two parts, depending on whether the kernel module is used or not.



([http://jake.dothome.co.kr/wp-content/uploads/2016/02/setup\\_per\\_cpu\\_areas-11d-1.png](http://jake.dothome.co.kr/wp-content/uploads/2016/02/setup_per_cpu_areas-11d-1.png))

## Per-CPU Area Storage Section

`_per_cpu_load` is the virtual address at which the per-CPU area is loaded into memory. To improve performance, various macros are available to store per-CPU data in different areas, as shown in the following figure.



(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-8a.png>)

## Per-CPU Section Location

The following vmlinux.lds.S linker script will show that the per-cpu section is located under the .init.data section.

### for ARM32

arch/arm/kernel/vmlinux.lds.S – ARM32

```

01 | .init.data : {
02 |     INIT_DATA
03 |     INIT_SETUP(16)
04 |     INIT_CALLS
05 |     CON_INITCALL
06 |     INIT_RAM_FS
07 |     *(.init.rodata.* .init.bss)    /* from the EFI stub */
08 | }
09 | .exit.data : {
10 |     ARM_EXIT_KEEP(EXIT_DATA)
11 | }
12 |
13 | PERCPU_SECTION(L1_CACHE_BYTES)

```

### for ARM64

arch/arm64/kernel/vmlinux.lds.S – ARM64

```

1 |     INIT_DATA_SECTION(16)
2 |
3 |     .exit.data : {
4 |         ARM_EXIT_KEEP(EXIT_DATA)
5 |     }

```

```

6
7 #ifdef CONFIG_SMP
8     PERCPU_SECTION(L1_CACHE_BYTES)
9 #endif

```

## PERCPU\_SECTION()

include/asm-generic/vmlinux.lds.h

```

01 /**
02  * PERCPU_SECTION - define output section for percpu area, simple versio
n
03  * @cacheline: cacheline size
04  *
05  * Align to PAGE_SIZE and outputs output section for percpu area. This
06  * macro doesn't manipulate @vaddr or @phdr and __per_cpu_load and
07  * __per_cpu_start will be identical.
08  *
09  * This macro is equivalent to ALIGN(PAGE_SIZE); PERCPU_VADDR(@cachelin
e,,)
10  * except that __per_cpu_load is defined as a relative symbol against
11  * .data..percpu which is required for relocatable x86_32 configuration.
12  */
13 #define PERCPU_SECTION(cacheline)
14 \
15     . = ALIGN(PAGE_SIZE);
16 \
17     .data..percpu : AT(ADDR(.data..percpu) - LOAD_OFFSET) {
18 \
19         __per_cpu_load = .;
20         PERCPU_INPUT(cacheline)
21 \
22     }

```

Define an output section for the per-CPU area.

- **LOAD\_OFFSET**
  - The value defined in asm-generic is 0 and is used unchanged in the ARM architecture.

## PERCPU\_INPUT()

include/asm-generic/vmlinux.lds.h

```

01 /**
02  * PERCPU_INPUT - the percpu input sections
03  * @cacheline: cacheline size
04  *
05  * The core percpu section names and core symbols which do not rely
06  * directly upon load addresses.
07  *
08  * @cacheline is used to align subsections to avoid false cacheline
09  * sharing between subsections for different purposes.
10  */
11 #define PERCPU_INPUT(cacheline)
12 \
13     __per_cpu_start = .;
14 \
15     *(.data..percpu..first)
16 \
17     . = ALIGN(PAGE_SIZE);
18 \
19     *(.data..percpu..page_aligned)
20 \

```

```

16 |         . = ALIGN(cacheline);
17 |         \
18 |         . = ALIGN(cacheline);
19 |         \
20 |         *(.data..percpu..read_mostly)
21 |         \
22 |         PERCPU_DECRYPTED_SECTION
    |         \
    |         __per_cpu_end = .;

```

Define the input section for the per-cpu area.

## per-cpu data allocation API

### Static per-cpu data allocation

Learn which APIs to declare to allocate and use per-cpu data at compile time.

#### DEFINE\_PER\_CPU()

include/linux/percpu-defs.h

```

1 | #define DEFINE_PER_CPU(type, name) \
2 |     DEFINE_PER_CPU_SECTION(type, name, "")

```

per\_cpu declare the data to be used. The per-CPU data declared in the kernel core uses the static area of the first chunk.

- You can also use arrays for types.
  - e.g. DEFINE\_PER\_CPU(int[3], my\_percpu\_array);

#### DEFINE\_PER\_CPU\_SECTION()

include/linux/percpu-defs.h

```

1 | #define DEFINE_PER_CPU_SECTION(type, name, sec) \
2 |     __PCPU_ATTRS(sec) __typeof__(type) name

```

#### \_\_PCPU\_ATTRS()

include/linux/percpu-defs.h

```

01 | /*
02 |  * Base implementations of per-CPU variable declarations and definition
03 |  * s, where
04 |  * the section in which the variable is to be placed is provided by the
05 |  * 'sec' argument. This may be used to affect the parameters governing
    | the
05 |  * variable's storage.

```

```

06  *
07  * NOTE! The sections for the DECLARE and for the DEFINE must match, le
st
08  * linkage errors occur due the compiler generating the wrong code to ac
cess
09  * that section.
10  */
11 #define __PCPU_ATTRS(sec)
12 \
13 \
    __percpu __attribute__((section(PER_CPU_BASE_SECTION sec)))
    PER_CPU_ATTRIBUTES

```

Specifies the section location where the per-cpu data will be stored.

```

1 #ifndef PER_CPU_BASE_SECTION
2 #ifdef CONFIG_SMP
3 #define PER_CPU_BASE_SECTION ".data..percpu"
4 #else
5 #define PER_CPU_BASE_SECTION ".data"
6 #endif
7 #endif

```

This is the section that is stored in the per-cpu data declaration.

- On SMP machines, .data.. percpu section.
- On UP machines, per-cpu data is stored in the .data section because it is no different from normal data.

## Declaring a static per-cpu data extern

To reference static per-cpu data from other files, use the following API:

include/linux/percpu-defs.h”

```

1 #define DECLARE_PER_CPU(type, name)
2 \
    DECLARE_PER_CPU_SECTION(type, name, "")

```

include/linux/percpu-defs.h”

```

1 #define DECLARE_PER_CPU_SECTION(type, name, sec)
2 \
    extern __PCPU_ATTRS(sec) __typeof__(type) name

```

## Create/Disable Dynamic

Learn about the API to dynamically allocate and deallocate per-CPU data at runtime.

- alloc\_percpu()
- free\_percpu()

## alloc\_percpu()

include/linux/percpu.h

```

1 | #define alloc_percpu(type)
2 | \
   (typeof(type) __percpu *)__alloc_percpu(sizeof(type),
3 | \
                                   __alignof__(type))

```

Allocates a per-CPU data area equal to the specified type size. At first, it uses the dynamic of the first chunk to allocate it, and if the dynamic area of the first chunk is insufficient, it creates and uses additional chunks. The added chunk is entirely a dynamic area.

## \_\_alloc\_percpu()

mm/percpu.c

```

1 | /**
2 |  * __alloc_percpu - allocate dynamic percpu area
3 |  * @size: size of area to allocate in bytes
4 |  * @align: alignment of area (max PAGE_SIZE)
5 |  *
6 |  * Equivalent to __alloc_percpu_gfp(size, align, %GFP_KERNEL).
7 |  */

1 | void __percpu *__alloc_percpu(size_t size, size_t align)
2 | {
3 |     return pcpu_alloc(size, align, false, GFP_KERNEL);
4 | }
5 | EXPORT_SYMBOL_GPL(__alloc_percpu);

```

@align to the requested @size and allocate a per-cpu data area.

## Using Per-CPU Data in Modules

The way a module uses per-CPU is handled slightly differently in kernel internals.

- Dynamic Assignment
  - The method of allocation and management in the existing dynamic area is the same.
- static assignment
  - It is not assigned to the static area of the first chunk, but dynamically manages the allocation to the reserve area.

## per-cpu data access API

It is designed with preemption in mind, and since per-cpu data does not point to the address of the actual per-cpu data, it must be accessed through the appropriate API without directly accessing and using it.

## 1) L-value operation

Use the `get_cpu_var()` function to increment the per-cpu data, as shown in the following example. After use, you must call the `put_cpu_var()` function.

- `get_cpu_var(sockets_in_use)++;`
- `put_cpu_var(sockets_in_use);`

### `get_cpu_var()`

`include/linux/percpu-defs.h`

```

1  /*
2   * Must be an lvalue. Since @var must be a simple identifier,
3   * we force a syntax error here if it isn't.
4   */
5  #define get_cpu_var(var)
6  ({
7      preempt_disable();
8      this_cpu_ptr(&var);
9  })

```

- In SMP machines, when you want to get a value, you `preempt_disable` it.
- The argument must be l-value.

### `put_cpu_var()`

`include/linux/percpu-defs.h`

```

1  /*
2   * The weird & is necessary because sparse considers (void)(var) to be
3   * a direct dereference of percpu variable (var).
4   */
5  #define put_cpu_var(var)
6  do {
7      (void)&(var);
8      preempt_enable();
9  } while (0)

```

- Write the value on the SMP machine and then `preempt_enable()`.

## 2) Pointer Manipulation

If you need to access it with a pointer, as shown in the following example, you first get the current CPU ID, then use the `per_cpu_ptr()` function to get the pointer and then use the pointer to perform the desired user action. When the per-cpu data is done, `put_cpu()` must be used.



```

1 | int cpu;
2 |
3 | cpu = get_cpu( )
4 | ptr = per_cpu_ptr(per_cpu_var, cpu);
5 | /* work with ptr */
6 | put_cpu( );

```

## per\_cpu\_ptr()

include/linux/percpu-defs.h

```

1 | #define per_cpu_ptr(ptr, cpu)
2 | \
3 | \
4 | \      __verify_pcpu_ptr(ptr);
5 | \      SHIFT_PERCPU_PTR((ptr), per_cpu_offset((cpu)));
6 | \
7 | \
8 | \
9 | \
10 | \
11 | \
12 | \
13 | \
14 | \
15 | \
16 | \
17 | \
18 | \
19 | \
20 | \
21 | \
22 | \
23 | \
24 | \
25 | \
26 | \
27 | \
28 | \
29 | \
30 | \
31 | \
32 | \
33 | \
34 | \
35 | \
36 | \
37 | \
38 | \
39 | \
40 | \
41 | \
42 | \
43 | \
44 | \
45 | \
46 | \
47 | \
48 | \
49 | \
50 | \
51 | \
52 | \
53 | \
54 | \
55 | \
56 | \
57 | \
58 | \
59 | \
60 | \
61 | \
62 | \
63 | \
64 | \
65 | \
66 | \
67 | \
68 | \
69 | \
70 | \
71 | \
72 | \
73 | \
74 | \
75 | \
76 | \
77 | \
78 | \
79 | \
80 | \
81 | \
82 | \
83 | \
84 | \
85 | \
86 | \
87 | \
88 | \
89 | \
90 | \
91 | \
92 | \
93 | \
94 | \
95 | \
96 | \
97 | \
98 | \
99 | \
100 | \

```

When using the Sparse static code analysis tool, check if there is any problem with the ptr (pointer) through \_\_verify\_pcpu\_ptr(). It then returns an address that combines the ptr with the offset of the CPU.

## per\_cpu\_offset()

include/asm-generic/percpu.h

```

01 | /*
02 |  * per_cpu_offset() is the offset that has to be added to a
03 |  * percpu variable to get to the instance for a certain processor.
04 |  *
05 |  * Most arches use the __per_cpu_offset array for those offsets but
06 |  * some arches have their own ways of determining the offset (x86_64, s3
07 |  * 90).
08 |  */
09 | #ifndef __per_cpu_offset
10 | extern unsigned long __per_cpu_offset[NR_CPUS];
11 | #define per_cpu_offset(x) (__per_cpu_offset[x])
12 | #endif

```

Returns the per-cpu offset value for the requested CPU @x number.

- \_\_per\_cpu\_offset[]
  - To access per-CPU data, you need to add a specified offset value for each CPU number, which is stored in an array of global variables.

## this\_cpu\_ptr()

include/linux/percpu-defs.h

```

1 | #ifdef CONFIG_DEBUG_PREEMPT
2 | #define this_cpu_ptr(ptr)
3 | \
4 | \
5 | \
6 | \
7 | \
8 | \
9 | \
10 | \
11 | \
12 | \
13 | \
14 | \
15 | \
16 | \
17 | \
18 | \
19 | \
20 | \
21 | \
22 | \
23 | \
24 | \
25 | \
26 | \
27 | \
28 | \
29 | \
30 | \
31 | \
32 | \
33 | \
34 | \
35 | \
36 | \
37 | \
38 | \
39 | \
40 | \
41 | \
42 | \
43 | \
44 | \
45 | \
46 | \
47 | \
48 | \
49 | \
50 | \
51 | \
52 | \
53 | \
54 | \
55 | \
56 | \
57 | \
58 | \
59 | \
60 | \
61 | \
62 | \
63 | \
64 | \
65 | \
66 | \
67 | \
68 | \
69 | \
70 | \
71 | \
72 | \
73 | \
74 | \
75 | \
76 | \
77 | \
78 | \
79 | \
80 | \
81 | \
82 | \
83 | \
84 | \
85 | \
86 | \
87 | \
88 | \
89 | \
90 | \
91 | \
92 | \
93 | \
94 | \
95 | \
96 | \
97 | \
98 | \
99 | \
100 | \

```



**RELOC\_HIDE()**

include/linux/compiler-gcc.h

```

01  /*
02  * This macro obfuscates arithmetic on a variable address so that gcc
03  * shouldn't recognize the original var, and make assumptions about it.
04  *
05  * This is needed because the C standard makes it undefined to do
06  * pointer arithmetic on "objects" outside their boundaries and the
07  * gcc optimizers assume this is the case. In particular they
08  * assume such arithmetic does not wrap.
09  *
10  * A miscompilation has been observed because of this on PPC.
11  * To work around it we hide the relationship of the pointer and the obj
    ect
12  * using this macro.
13  *
14  * Versions of the ppc64 compiler before 4.1 had a bug where use of
15  * RELOC_HIDE could trash r30. The bug can be worked around by changing
16  * the inline assembly constraint from =g to =r, in this particular
17  * case either is valid.
18  */
19  #define RELOC_HIDE(ptr, off)                                \
20  ({ unsigned long __ptr;                                     \
21     __asm__ (" : "=r"(__ptr) : "0"(ptr));                  \
22     (typeof(ptr)) (__ptr + (off)); })

```

- Returns the address of the \_\_ptr plus off.
- A literal translation of the above comment is that this macro obscures the arithmetic operation of a variable, so that gcc cannot recognize the original value and therefore cannot make assumptions about it, in other words, it takes only the address of the variable, regardless of the type of the variable, adds an offset to it, and casts it to the type of the original variable.
- Compilers prior to GCC 4.1 had a bug with the PPC64 architecture, which was resolved by applying this walkaround.

**raw\_cpu\_ptr()**

include/linux/percpu-defs.h

```

1  #define raw_cpu_ptr(ptr)
2  \
3  ({
4  \
5  \
6  \
7  \
8  \
9  \
10 \
11 \
12 \
13 \
14 \
15 \
16 \
17 \
18 \
19 \
20 \
21 \
22 \
23 \
24 \
25 \
26 \
27 \
28 \
29 \
30 \
31 \
32 \
33 \
34 \
35 \
36 \
37 \
38 \
39 \
40 \
41 \
42 \
43 \
44 \
45 \
46 \
47 \
48 \
49 \
50 \
51 \
52 \
53 \
54 \
55 \
56 \
57 \
58 \
59 \
60 \
61 \
62 \
63 \
64 \
65 \
66 \
67 \
68 \
69 \
70 \
71 \
72 \
73 \
74 \
75 \
76 \
77 \
78 \
79 \
80 \
81 \
82 \
83 \
84 \
85 \
86 \
87 \
88 \
89 \
90 \
91 \
92 \
93 \
94 \
95 \
96 \
97 \
98 \
99 \
100 \
101 \
102 \
103 \
104 \
105 \
106 \
107 \
108 \
109 \
110 \
111 \
112 \
113 \
114 \
115 \
116 \
117 \
118 \
119 \
120 \
121 \
122 \
123 \
124 \
125 \
126 \
127 \
128 \
129 \
130 \
131 \
132 \
133 \
134 \
135 \
136 \
137 \
138 \
139 \
140 \
141 \
142 \
143 \
144 \
145 \
146 \
147 \
148 \
149 \
150 \
151 \
152 \
153 \
154 \
155 \
156 \
157 \
158 \
159 \
160 \
161 \
162 \
163 \
164 \
165 \
166 \
167 \
168 \
169 \
170 \
171 \
172 \
173 \
174 \
175 \
176 \
177 \
178 \
179 \
180 \
181 \
182 \
183 \
184 \
185 \
186 \
187 \
188 \
189 \
190 \
191 \
192 \
193 \
194 \
195 \
196 \
197 \
198 \
199 \
200 \
201 \
202 \
203 \
204 \
205 \
206 \
207 \
208 \
209 \
210 \
211 \
212 \
213 \
214 \
215 \
216 \
217 \
218 \
219 \
220 \
221 \
222 \
223 \
224 \
225 \
226 \
227 \
228 \
229 \
230 \
231 \
232 \
233 \
234 \
235 \
236 \
237 \
238 \
239 \
240 \
241 \
242 \
243 \
244 \
245 \
246 \
247 \
248 \
249 \
250 \
251 \
252 \
253 \
254 \
255 \
256 \
257 \
258 \
259 \
260 \
261 \
262 \
263 \
264 \
265 \
266 \
267 \
268 \
269 \
270 \
271 \
272 \
273 \
274 \
275 \
276 \
277 \
278 \
279 \
280 \
281 \
282 \
283 \
284 \
285 \
286 \
287 \
288 \
289 \
290 \
291 \
292 \
293 \
294 \
295 \
296 \
297 \
298 \
299 \
300 \
301 \
302 \
303 \
304 \
305 \
306 \
307 \
308 \
309 \
310 \
311 \
312 \
313 \
314 \
315 \
316 \
317 \
318 \
319 \
320 \
321 \
322 \
323 \
324 \
325 \
326 \
327 \
328 \
329 \
330 \
331 \
332 \
333 \
334 \
335 \
336 \
337 \
338 \
339 \
340 \
341 \
342 \
343 \
344 \
345 \
346 \
347 \
348 \
349 \
350 \
351 \
352 \
353 \
354 \
355 \
356 \
357 \
358 \
359 \
360 \
361 \
362 \
363 \
364 \
365 \
366 \
367 \
368 \
369 \
370 \
371 \
372 \
373 \
374 \
375 \
376 \
377 \
378 \
379 \
380 \
381 \
382 \
383 \
384 \
385 \
386 \
387 \
388 \
389 \
390 \
391 \
392 \
393 \
394 \
395 \
396 \
397 \
398 \
399 \
400 \
401 \
402 \
403 \
404 \
405 \
406 \
407 \
408 \
409 \
410 \
411 \
412 \
413 \
414 \
415 \
416 \
417 \
418 \
419 \
420 \
421 \
422 \
423 \
424 \
425 \
426 \
427 \
428 \
429 \
430 \
431 \
432 \
433 \
434 \
435 \
436 \
437 \
438 \
439 \
440 \
441 \
442 \
443 \
444 \
445 \
446 \
447 \
448 \
449 \
450 \
451 \
452 \
453 \
454 \
455 \
456 \
457 \
458 \
459 \
460 \
461 \
462 \
463 \
464 \
465 \
466 \
467 \
468 \
469 \
470 \
471 \
472 \
473 \
474 \
475 \
476 \
477 \
478 \
479 \
480 \
481 \
482 \
483 \
484 \
485 \
486 \
487 \
488 \
489 \
490 \
491 \
492 \
493 \
494 \
495 \
496 \
497 \
498 \
499 \
500 \
501 \
502 \
503 \
504 \
505 \
506 \
507 \
508 \
509 \
510 \
511 \
512 \
513 \
514 \
515 \
516 \
517 \
518 \
519 \
520 \
521 \
522 \
523 \
524 \
525 \
526 \
527 \
528 \
529 \
530 \
531 \
532 \
533 \
534 \
535 \
536 \
537 \
538 \
539 \
540 \
541 \
542 \
543 \
544 \
545 \
546 \
547 \
548 \
549 \
550 \
551 \
552 \
553 \
554 \
555 \
556 \
557 \
558 \
559 \
560 \
561 \
562 \
563 \
564 \
565 \
566 \
567 \
568 \
569 \
570 \
571 \
572 \
573 \
574 \
575 \
576 \
577 \
578 \
579 \
580 \
581 \
582 \
583 \
584 \
585 \
586 \
587 \
588 \
589 \
590 \
591 \
592 \
593 \
594 \
595 \
596 \
597 \
598 \
599 \
600 \
601 \
602 \
603 \
604 \
605 \
606 \
607 \
608 \
609 \
610 \
611 \
612 \
613 \
614 \
615 \
616 \
617 \
618 \
619 \
620 \
621 \
622 \
623 \
624 \
625 \
626 \
627 \
628 \
629 \
630 \
631 \
632 \
633 \
634 \
635 \
636 \
637 \
638 \
639 \
640 \
641 \
642 \
643 \
644 \
645 \
646 \
647 \
648 \
649 \
650 \
651 \
652 \
653 \
654 \
655 \
656 \
657 \
658 \
659 \
660 \
661 \
662 \
663 \
664 \
665 \
666 \
667 \
668 \
669 \
670 \
671 \
672 \
673 \
674 \
675 \
676 \
677 \
678 \
679 \
680 \
681 \
682 \
683 \
684 \
685 \
686 \
687 \
688 \
689 \
690 \
691 \
692 \
693 \
694 \
695 \
696 \
697 \
698 \
699 \
700 \
701 \
702 \
703 \
704 \
705 \
706 \
707 \
708 \
709 \
710 \
711 \
712 \
713 \
714 \
715 \
716 \
717 \
718 \
719 \
720 \
721 \
722 \
723 \
724 \
725 \
726 \
727 \
728 \
729 \
730 \
731 \
732 \
733 \
734 \
735 \
736 \
737 \
738 \
739 \
740 \
741 \
742 \
743 \
744 \
745 \
746 \
747 \
748 \
749 \
750 \
751 \
752 \
753 \
754 \
755 \
756 \
757 \
758 \
759 \
760 \
761 \
762 \
763 \
764 \
765 \
766 \
767 \
768 \
769 \
770 \
771 \
772 \
773 \
774 \
775 \
776 \
777 \
778 \
779 \
780 \
781 \
782 \
783 \
784 \
785 \
786 \
787 \
788 \
789 \
790 \
791 \
792 \
793 \
794 \
795 \
796 \
797 \
798 \
799 \
800 \
801 \
802 \
803 \
804 \
805 \
806 \
807 \
808 \
809 \
810 \
811 \
812 \
813 \
814 \
815 \
816 \
817 \
818 \
819 \
820 \
821 \
822 \
823 \
824 \
825 \
826 \
827 \
828 \
829 \
830 \
831 \
832 \
833 \
834 \
835 \
836 \
837 \
838 \
839 \
840 \
841 \
842 \
843 \
844 \
845 \
846 \
847 \
848 \
849 \
850 \
851 \
852 \
853 \
854 \
855 \
856 \
857 \
858 \
859 \
860 \
861 \
862 \
863 \
864 \
865 \
866 \
867 \
868 \
869 \
870 \
871 \
872 \
873 \
874 \
875 \
876 \
877 \
878 \
879 \
880 \
881 \
882 \
883 \
884 \
885 \
886 \
887 \
888 \
889 \
890 \
891 \
892 \
893 \
894 \
895 \
896 \
897 \
898 \
899 \
900 \
901 \
902 \
903 \
904 \
905 \
906 \
907 \
908 \
909 \
910 \
911 \
912 \
913 \
914 \
915 \
916 \
917 \
918 \
919 \
920 \
921 \
922 \
923 \
924 \
925 \
926 \
927 \
928 \
929 \
930 \
931 \
932 \
933 \
934 \
935 \
936 \
937 \
938 \
939 \
940 \
941 \
942 \
943 \
944 \
945 \
946 \
947 \
948 \
949 \
950 \
951 \
952 \
953 \
954 \
955 \
956 \
957 \
958 \
959 \
960 \
961 \
962 \
963 \
964 \
965 \
966 \
967 \
968 \
969 \
970 \
971 \
972 \
973 \
974 \
975 \
976 \
977 \
978 \
979 \
980 \
981 \
982 \
983 \
984 \
985 \
986 \
987 \
988 \
989 \
990 \
991 \
992 \
993 \
994 \
995 \
996 \
997 \
998 \
999 \
1000 \
1001 \
1002 \
1003 \
1004 \
1005 \
1006 \
1007 \
1008 \
1009 \
1010 \
1011 \
1012 \
1013 \
1014 \
1015 \
1016 \
1017 \
1018 \
1019 \
1020 \
1021 \
1022 \
1023 \
1024 \
1025 \
1026 \
1027 \
1028 \
1029 \
1030 \
1031 \
1032 \
1033 \
1034 \
1035 \
1036 \
1037 \
1038 \
1039 \
1040 \
1041 \
1042 \
1043 \
1044 \
1045 \
1046 \
1047 \
1048 \
1049 \
1050 \
1051 \
1052 \
1053 \
1054 \
1055 \
1056 \
1057 \
1058 \
1059 \
1060 \
1061 \
1062 \
1063 \
1064 \
1065 \
1066 \
1067 \
1068 \
1069 \
1070 \
1071 \
1072 \
1073 \
1074 \
1075 \
1076 \
1077 \
1078 \
1079 \
1080 \
1081 \
1082 \
1083 \
1084 \
1085 \
1086 \
1087 \
1088 \
1089 \
1090 \
1091 \
1092 \
1093 \
1094 \
1095 \
1096 \
1097 \
1098 \
1099 \
1100 \
1101 \
1102 \
1103 \
1104 \
1105 \
1106 \
1107 \
1108 \
1109 \
1110 \
1111 \
1112 \
1113 \
1114 \
1115 \
1116 \
1117 \
1118 \
1119 \
1120 \
1121 \
1122 \
1123 \
1124 \
1125 \
1126 \
1127 \
1128 \
1129 \
1130 \
1131 \
1132 \
1133 \
1134 \
1135 \
1136 \
1137 \
1138 \
1139 \
1140 \
1141 \
1142 \
1143 \
1144 \
1145 \
1146 \
1147 \
1148 \
1149 \
1150 \
1151 \
1152 \
1153 \
1154 \
1155 \
1156 \
1157 \
1158 \
1159 \
1160 \
1161 \
1162 \
1163 \
1164 \
1165 \
1166 \
1167 \
1168 \
1169 \
1170 \
1171 \
1172 \
1173 \
1174 \
1175 \
1176 \
1177 \
1178 \
1179 \
1180 \
1181 \
1182 \
1183 \
1184 \
1185 \
1186 \
1187 \
1188 \
1189 \
1190 \
1191 \
1192 \
1193 \
1194 \
1195 \
1196 \
1197 \
1198 \
1199 \
1200 \
1201 \
1202 \
1203 \
1204 \
1205 \
1206 \
1207 \
1208 \
1209 \
1210 \
1211 \
1212 \
1213 \
1214 \
1215 \
1216 \
1217 \
1218 \
1219 \
1220 \
1221 \
1222 \
1223 \
1224 \
1225 \
1226 \
1227 \
1228 \
1229 \
1230 \
1231 \
1232 \
1233 \
1234 \
1235 \
1236 \
1237 \
1238 \
1239 \
1240 \
1241 \
1242 \
1243 \
1244 \
1245 \
1246 \
1247 \
1248 \
1249 \
1250 \
1251 \
1252 \
1253 \
1254 \
1255 \
1256 \
1257 \
1258 \
1259 \
1260 \
1261 \
1262 \
1263 \
1264 \
1265 \
1266 \
1267 \
1268 \
1269 \
1270 \
1271 \
1272 \
1273 \
1274 \
1275 \
1276 \
1277 \
1278 \
1279 \
1280 \
1281 \
1282 \
1283 \
1284 \
1285 \
1286 \
1287 \
1288 \
1289 \
1290 \
1291 \
1292 \
1293 \
1294 \
1295 \
1296 \
1297 \
1298 \
1299 \
1300 \
1301 \
1302 \
1303 \
1304 \
1305 \
1306 \
1307 \
1308 \
1309 \
1310 \
1311 \
1312 \
1313 \
1314 \
1315 \
1316 \
1317 \
1318 \
1319 \
1320 \
1321 \
1322 \
1323 \
1324 \
1325 \
1326 \
1327 \
1328 \
1329 \
1330 \
1331 \
1332 \
1333 \
1334 \
1335 \
1336 \
1337 \
1338 \
1339 \
1340 \
1341 \
1342 \
1343 \
1344 \
1345 \
1346 \
1347 \
1348 \
1349 \
1350 \
1351 \
1352 \
1353 \
1354 \
1355 \
1356 \
1357 \
1358 \
1359 \
1360 \
1361 \
1362 \
1363 \
1364 \
1365 \
1366 \
1367 \
1368 \
1369 \
1370 \
1371 \
1372 \
1373 \
1374 \
1375 \
1376 \
1377 \
1378 \
1379 \
1380 \
1381 \
1382 \
1383 \
1384 \
1385 \
1386 \
1387 \
1388 \
1389 \
1390 \
1391 \
1392 \
1393 \
1394 \
1395 \
1396 \
1397 \
1398 \
1399 \
1400 \
1401 \
1402 \
1403 \
1404 \
1405 \
1406 \
1407 \
1408 \
1409 \
1410 \
1411 \
1412 \
1413 \
1414 \
1415 \
1416 \
1417 \
1418 \
1419 \
1420 \
1421 \
1422 \
1423 \
1424 \
1425 \
1426 \
1427 \
1428 \
1429 \
1430 \
1431 \
1432 \
1433 \
1434 \
1435 \
1436 \
1437 \
1438 \
1439 \
1440 \
1441 \
1442 \
1443 \
1444 \
1445 \
1446 \
1447 \
1448 \
1449 \
1450 \
1451 \
1452 \
1453 \
1454 \
1455 \
1456 \
1457 \
1458 \
1459 \
1460 \
1461 \
1462 \
1463 \
1464 \
1465 \
1466 \
1467 \
1468 \
1469 \
1470 \
1471 \
1472 \
1473 \
1474 \
1475 \
1476 \
1477 \
1478 \
1479 \
1480 \
1481 \
1482 \
1483 \
1484 \
1485 \
1486 \
1487 \
1488 \
1489 \
1490 \
1491 \
1492 \
1493 \
1494 \
1495 \
1496 \
1497 \
1498 \
1499 \
1500 \
1501 \
1502 \
1503 \
1504 \
1505 \
1506 \
1507 \
1508 \
1509 \
1510 \
1511 \
1512 \
1513 \
1514 \
1515 \
1516 \
1517 \
1518 \
1519 \
1520 \
1521 \
1522 \
1523 \
1524 \
1525 \
1526 \
1527 \
1528 \
1529 \
1530 \
1531 \
1532 \
1533 \
1534 \
1535 \
1536 \
1537 \
1538 \
1539 \
1540 \
1541 \
1542 \
1543 \
1544 \
1545 \
1546 \
1547 \
1548 \
1549 \
1550 \
1551 \
1552 \
1553 \
1554 \
1555 \
1556 \
1557 \
1558 \
1559 \
1560 \
1561 \
1562 \
1563 \
1564 \
1565 \
1566 \
1567 \
1568 \
1569 \
1570 \
1571 \
1572 \
1573 \
1574 \
1575 \
1576 \
1577 \
1578 \
1579 \
1580 \
1581 \
1582 \
1583 \
1584 \
1585 \
1586 \
1587 \
1588 \
1589 \
1590 \
1591 \
1592 \
1593 \
1594 \
1595 \
1596 \
1597 \
1598 \
1599 \
1600 \
1601 \
1602 \
1603 \
1604 \
1605 \
1606 \
1607 \
1608 \
1609 \
1610 \
1611 \
1612 \
1613 \
1614 \
1615 \
1616 \
1617 \
1618 \
1619 \
1620 \
1621 \
1622 \
1623 \
1624 \
1625 \
1626 \
1627 \
1628 \
1629 \
1630 \
1631 \
1632 \
1633 \
1634 \
1635 \
1636 \
1637 \
1638 \
1639 \
1640 \
1641 \
1642 \
1643 \
1644 \
1645 \
1646 \
1647 \
1648 \
1649 \
1650 \
1651 \
1652 \
1653 \
1654 \
1655 \
1656 \
1657 \
1658 \
1659 \
1660 \
1661 \
1662 \
1663 \
1664 \
1665 \
1666 \
1667 \
1668 \
1669 \
1670 \
1671 \
1672 \
1673 \
1674 \
1675 \
1676 \
1677 \
1678 \
1679 \
1680 \
1681 \
1682 \
1683 \
1684 \
1685 \
1686 \
1687 \
1688 \
1689 \
1690 \
1691 \
1692 \
1693 \
1694 \
1695 \
1696 \
1697 \
1698 \
1699 \
1700 \
1701 \
1702 \
1703 \
1704 \
1705 \
1706 \
1707 \
1708 \
1709 \
1710 \
1711 \
1712 \
1713 \
1714 \
1715 \
1716 \
1717 \
1718 \
1719 \
1720 \
1721 \
1722 \
1723 \
1724 \
1725 \
1726 \
1727 \
1728 \
1729 \
1730 \
1731 \
1732 \
1733 \
1734 \
1735 \
1736 \
1737 \
1738 \
1739 \
1740 \
1741 \
1742 \
1743 \
1744 \
1745 \
1746 \
1747 \
1748 \
1749 \
1750 \
1751 \
1752 \
1753 \
1754 \
1755 \
1756 \
1757 \
1758 \
1759 \
1760 \
1761 \
1762 \
1763 \
1764 \
1765 \
1766 \
1767 \
1768 \
1769 \
1770 \
1771 \
1772 \
1773 \
1774 \
1775 \
1776 \
1777 \
1778 \
1779 \
1780 \
1781 \
1782 \
1783 \
1784 \
1785 \
1786 \
1787 \
1788 \
1789 \
1790 \
1791 \
1792 \
1793 \
1794 \
1795 \
1796 \
1797 \
1798 \
1799 \
1800 \
1801 \
1802 \
1803 \
1804 \
1805 \
1806 \
1807 \
1808 \
1809 \
1810 \
1811 \
1812 \
1813 \
1814 \
1815 \
1816 \
1817 \
1818 \
1819 \
1820 \
1821 \
1822 \
1823 \
1824 \
1825 \
1826 \
1827 \
1828 \
1829 \
1830 \
1831 \
1832 \
1833 \
1834 \
1835 \
1836 \
1837 \
1838 \
1839 \
1840 \
1841 \
1842 \
1843 \
1844 \
1845 \
1846 \
1847 \
1848 \
1849 \
1850 \
1851 \
1852 \
1853 \
1854 \
1855 \
1856 \
1857 \
1858 \
1859 \
1860 \
1861 \
1862 \
1863 \
1864 \
1865 \
1866 \
1867 \
1868 \
1869 \
1870 \
1871 \
1872 \
1873 \
1874 \
1875 \
1876 \
1877 \
1878 \
1879 \
1880 \
1881 \
1882 \
1883 \
1884 \
1885 \
1886 \
1887 \
1888 \
1889 \
1890 \
1891 \
1892 \
1893 \
1894 \
1895 \
1896 \
1897 \
1898 \
1899 \
1900 \
1901 \
1902 \
1903 \
1904 \
1905 \
1906 \
1907 \
1908 \
1909 \
1910 \
1911 \
1912 \
1913 \
1914 \
1915 \
1916 \
1917 \
1918 \
1919 \
1920 \
1921 \
1922 \
1923 \
1924 \
1925 \
1926 \
1927 \
1928 \
1929 \
1930 \
1931 \
1932 \
1933 \
1934 \
1935 \
1936 \
1937 \
1938 \
1939 \
1940 \
1941 \
1942 \
1943 \
1944 \
1945 \
1946 \
1947 \
1948 \
1949 \
1950 \
1951 \
1952 \
1953 \
1954 \
1955 \
1956 \
1957 \
1958 \
1959 \
1960 \
1961 \
1962 \
1963 \
1964 \
1965 \
1966 \
1967 \
1968 \
1969 \
1970 \
1971 \
1972 \
1973 \
1974 \
1975 \
1976 \
1977 \
1978 \
1979 \
1980 \
1981 \
1982 \
1983 \
1984 \
1985 \
1986 \
1987 \
1988 \
1989 \
1990 \
1991 \
1992 \
1993 \
1994 \
1995 \
1996 \
1997 \
1998 \
1999 \
2000 \
2001 \
2002 \
2003 \
2004 \
2005 \
2006 \
2007 \
2008 \
2009 \
2010 \
2011 \
2012 \
2013 \
2014 \
2015 \
2016 \
2017 \
2018 \
2019 \
2020 \
2021 \
2022 \
2023 \
2024 \
2025 \
2026 \
2027 \
2028 \
2029 \
2030 \
2031 \
2032 \
2033 \
2034 \
2035 \
2036 \
2037 \
2038 \
2039 \
2040 \
2041 \
2042 \
2043 \
2044 \
2045 \
2046 \
2047 \
2048 \
2049 \
2050 \
2051 \
2052 \
2053 \
2054 \
2055 \
2056 \
2057 \
2058 \
2059 \
2060 \
2061 \
2062 \
2063 \
2064 \
2065 \
2066 \
2067 \
2068 \
2069 \
2070 \
2071 \
2072 \
2073 \
2074 \
2075 \
2076 \
2077 \
2078 \
2079 \
2080 \
2081 \
2082 \
2083 \
2084 \
2085 \
2086 \
2087 \
2088 \
2089 \
2090 \
2091 \
2092 \
2093 \
2094 \
2095 \
2096 \
2097 \
2098 \
2099 \
2100 \
2101 \
2102 \
2103 \
2104 \
2105 \
2106 \
2107 \
2108 \
2109 \
2110 \
2111 \
2112 \
2113 \
2114 \
2115 \
2116 \
211
```

```

1  /*
2   * Arch may define arch_raw_cpu_ptr() to provide more efficient address
3   * translations for raw_cpu_ptr().
4   */
5  #ifndef arch_raw_cpu_ptr
6  #define arch_raw_cpu_ptr(ptr) SHIFT_PERCPU_PTR(ptr, __my_cpu_offset)
7  #endif

```

Returns the pointer address of the @ptr plus the per-cpu offset for the current CPU.

### \_\_my\_cpu\_offset() – ARM32

arch/arm/include/asm/percpu.h

```

01  #define __my_cpu_offset __my_cpu_offset()
02
03  static inline unsigned long __my_cpu_offset(void)
04  {
05      unsigned long off;
06
07      /*
08       * Read TPIDRPRW.
09       * We want to allow caching the value, so avoid using volatile a
10  nd
11       * instead use a fake stack read to hazard against barrier().
12       */
13      asm("mrc p15, 0, %0, c13, c0, 4" : "=r" (off)
14          : "Q" (*(const unsigned long *)current_stack_pointer));
15      return off;
16  }

```

Read the value of the offset corresponding to the cpu stored in the TPIDRPRW register.

- This value is used to access the per-cpu data in addition to the per-cpu ptr address.
- In the ARM architecture, the TPIDRPRW register, which is not used among the CP15 registers, is used for CP<> offset archiving per CPU, which increases per-cpu computational performance.

### \_\_my\_cpu\_offset() – ARM64

arch/arm64/include/asm/percpu.h

```

01  #define __my_cpu_offset __my_cpu_offset()
02
03  static inline unsigned long __my_cpu_offset(void)
04  {
05      unsigned long off;
06
07      /*
08       * We want to allow caching the value, so avoid using volatile a
09  nd
10       * instead use a fake stack read to hazard against barrier().
11       */
12      asm(ALTERNATIVE("mrs %0, tpidr_el1",
13                      "mrs %0, tpidr_el2",
14                      ARM64_HAS_VIRT_HOST_EXTN)
15          : "=r" (off) :
16          "Q" (*(const unsigned long *)current_stack_pointer));
17      return off;
18  }

```

Reads the per-cpu offset value corresponding to the CPU stored in the TPIDR register.

## Other key functions

### \_\_percpu Attribute Macros

include/linux/compiler\_types.h

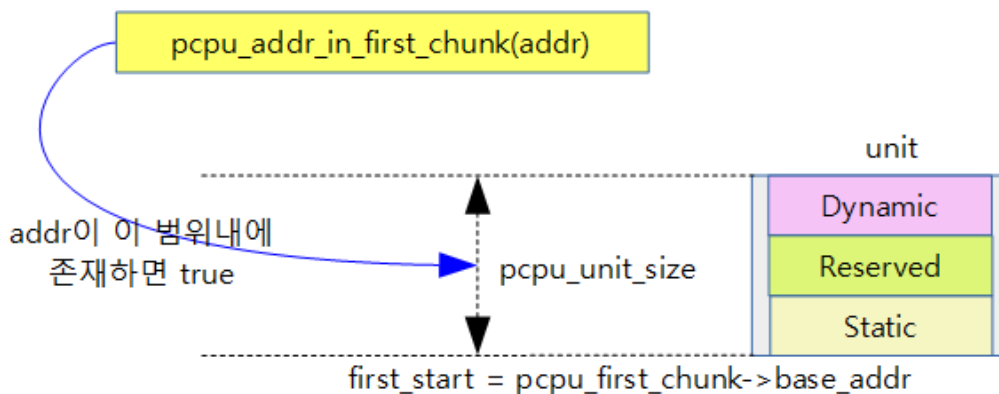
```
1 | # define __percpu      __attribute__((noderef, address_space(3)))
```

Sparse is an attribute specified by the static code analysis tool.

- address\_space(3):
  - Specify the per-CPU address\_space to be used.
- noderef
  - Stipulate that dereferencing pointers cannot be used.
    - You can't use a generic pointer variable, \*ptr, etc., but you can use the &val form that directly provides the address of the variable.

### pcpu\_addr\_in\_first\_chunk()

mm/percpu.c – Removed as of kernel v4.14



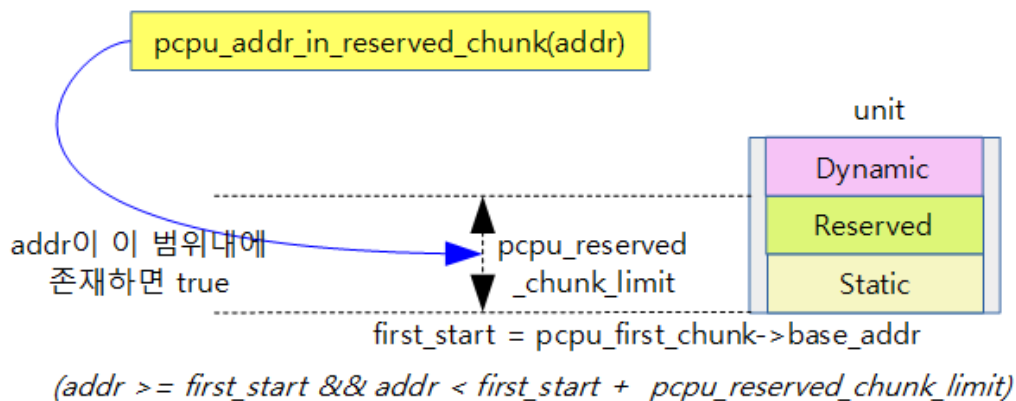
*(addr >= first\_start && addr < first\_start + pcpu\_unit\_size)*

(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-9.png>)

- Returns whether the given address is in the first chunk range.

### pcpu\_addr\_in\_reserved\_chunk()

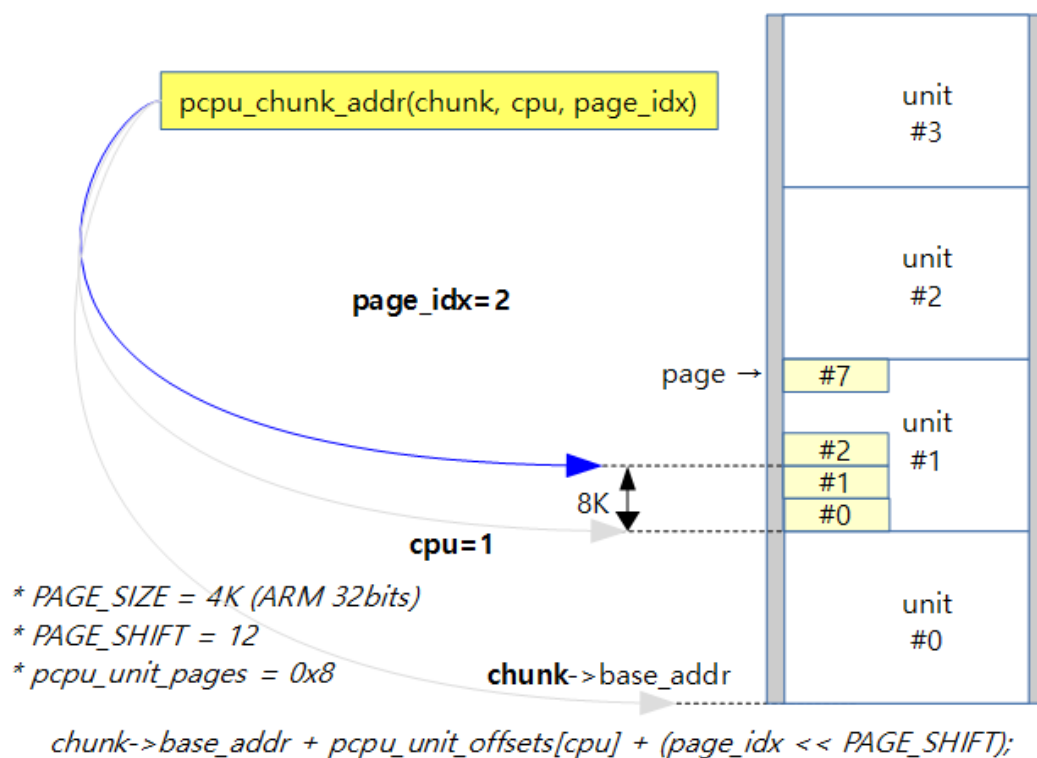
mm/percpu.c – Removed as of kernel v4.14



(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-10.png>)

- Returns whether a given address is included in the Static and Reserved zones.
- If no reserved chunk is given, i.e. there is no reserved zone, the `pcpu_reserved_chunk_limit` value is 0.

### pcpu\_chunk\_addr()



(<http://jake.dothome.co.kr/wp-content/uploads/2016/02/percpu-11.png>)

- The given chunk number, CPU number, and page number are used to get the start address of the page.

## Example of using per-cpu data used to count slab objects

`include/linux/slub_def.h`

```
1 | /*
2 |  * Slab cache management.
```

```

3  */
4  struct kmem_cache {
5      struct kmem_cache_cpu __percpu *cpu_slab;

1 struct kmem_cache_cpu {
2     void **freelist;           /* Pointer to next available object */
3     unsigned long tid;        /* Globally unique transaction id */
4     struct page *page;        /* The slab from which we are allocating */
5     struct page *partial;     /* Partially allocated frozen slabs */
6 #ifdef CONFIG_SLUB_STATS
7     unsigned stat[NR_SLUB_STAT_ITEMS];
8 #endif
9 };

1 static inline void flush_slab(struct kmem_cache *s, struct kmem_cache_cp
  u *c)
2 {
3     stat(s, CPUSLAB_FLUSH);
4     deactivate_slab(s, c->page, c->freelist);
5
6     c->tid = next_tid(c->tid);
7     c->page = NULL;
8     c->freelist = NULL;
9 }

01 static inline void stat(const struct kmem_cache *s, enum stat_item si)
02 {
03 #ifdef CONFIG_SLUB_STATS
04     /*
05      * The rmw is racy on a preemptible kernel but this is acceptabl
06 e, so
07      * avoid this_cpu_add()'s irq-disable overhead.
08      */
09     raw_cpu_inc(s->cpu_slab->stat[si]);
10 #endif
11 }

```

include/linux/slub\_def.h

```

01 enum stat_item {
02     ALLOC_FASTPATH,           /* Allocation from cpu slab */
03     ALLOC_SLOWPATH,          /* Allocation by getting a new cpu slab */
04     /*
05     FREE_FASTPATH,           /* Free to cpu slab */
06     FREE_SLOWPATH,          /* Freeing not to cpu slab */
07     FREE_FROZEN,            /* Freeing to frozen slab */
08     FREE_ADD_PARTIAL,        /* Freeing moves slab to partial list */
09     FREE_REMOVE_PARTIAL,     /* Freeing removes last object */
10     ALLOC_FROM_PARTIAL,      /* Cpu slab acquired from node partial l
11     ist */
12     ALLOC_SLAB,              /* Cpu slab acquired from page allocator */
13     /*
14     ALLOC_REFILL,            /* Refill cpu slab from slab freelist */
15     ALLOC_NODE_MISMATCH,     /* Switching cpu slab */
16     FREE_SLAB,               /* Slab freed to the page allocator */
17     CPUSLAB_FLUSH,           /* Abandoning of the cpu slab */
18     (...생략...)
19 }

```

include/linux/percpu-defs.h

```

1 #define raw_cpu_inc(pcp)          raw_cpu_add(pcp, 1)

```

```

1 | #define raw_cpu_add(pcp, val)      __pcpu_size_call(raw_cpu_add_, pcp,
   | val)
01 | #define __pcpu_size_call(stem, variable, ...)
   | \
02 | do {
   | \
03 |     __verify_pcpu_ptr(&(variable));
   | \
04 |     switch(sizeof(variable)) {
   | \
05 |         case 1: stem##1(variable, __VA_ARGS__);break;
   | \
06 |         case 2: stem##2(variable, __VA_ARGS__);break;
   | \
07 |         case 4: stem##4(variable, __VA_ARGS__);break;
   | \
08 |         case 8: stem##8(variable, __VA_ARGS__);break;
   | \
09 |         default:
   | \
10 |             __bad_size_call_parameter();break;
   | \
11 |     }
   | \
12 | } while (0)

```

include/asm-generic/percpu.h

```

1 | #ifndef raw_cpu_add_4
2 | #define raw_cpu_add_4(pcp, val)      raw_cpu_generic_to_op(pcp, val,
   | +=)
3 | #endif
1 | #define raw_cpu_generic_to_op(pcp, val, op)
   | \
2 | do {
   | \
3 |     *raw_cpu_ptr(&(pcp)) op val;
   | \
4 | } while (0)

```

## per-cpu-related structs

### pcpu\_chunk Structure

mm/percpu-internal.h

```

01 | struct pcpu_chunk {
02 | #ifdef CONFIG_PERCPU_STATS
03 |     int nr_alloc; /* # of allocations */
04 |     size_t max_alloc_size; /* largest allocation size */
   | ze */
05 | #endif
06 |
07 |     struct list_head list; /* linked to pcpu_slot list */
   | ists */
08 |     int free_bytes; /* free bytes in the chunk */
   | nk */
09 |     int contig_bits; /* max contiguous size in bits */
   | int */

```



```

10         int
11         ng
12         void
13         chunk */
14         unsigned long
15         unsigned long
16         struct pcpu_block_md
17
18         void
19         int
20         bool
21         owed */
22         int
23         previous
24         aligned
25         int
26         red to
27         age
28         int
29         this chunk */
30         int
31         ted pages */
32         unsigned long
    
```

```

contig_bits_start; /* contig_bits start
                                offset */
*base_addr; /* base address of this
chunk */
*alloc_map; /* allocation map */
*bound_map; /* boundary map */
*md_blocks; /* metadata blocks */
*data; /* chunk data */
first_bit; /* no free below this */
immutable; /* no [de]population all
owed */
start_offset; /* the overlap with the
previous
region to have a page
aligned
base_addr */
end_offset; /* additional area requi
red to
have the region end p
age
aligned */
nr_pages; /* # of pages served by
this chunk */
nr_populated; /* # of populated pages
*/
nr_empty_pop_pages; /* # of empty popula
ted pages */
populated[]; /* populated bitmap */
    
```

- nr\_alloc
  - Number of allocations requested in this chunk
- max\_alloc\_size
  - Each time an allocation request is made to this chunk, the maximum allocation size of the requested allocation size is updated.
- list
  - This is a list of pcpu\_slot maps used to manage chunks.
  - All chunks are linked to a list of pcpu\_slot maps and used to manage dynamic spaces, sorted by the size of the free space in each chunk, with the allocation starting with the chunk with the fewest free space.
- free\_bytes
  - It is the sum of the free space in the chunk. A unit is a byte.
  - Each allocation and release changes the free\_size and moves from the chunk slot to the other appropriate slot according to the free\_size.
- contig\_bits (formerly called contig\_hints, and the size was in bytes.)
  - It is the largest contiguous free space size within a chunk. (4 bytes per bit)
- contig\_bits\_start
  - This is the offset bit from which the contig\_bits begins.
- \*alloc\_map
  - A bitmap that expresses an assignment, with each bit representing whether or not it is allocated in units of 4 bytes.
- \*bound\_map

- A bitmap that represents a boundary, with each bit managed in 4-byte increments, representing the starting bit of the allocated data as 1.
- The last bit of the bound\_map must also be represented by 1, so it is always 1 bit larger than the size of the alloc\_map.
- \*md\_blocks
  - It refers to metadata that is managed in per-CPU blocks.
  - Within the chunk, this meter data is used to quickly scan each per-CPU block for allocation.
- data
  - The address of the page structure pointer array is stored to which the chunk is assigned.
- immutable
  - Indicates whether the chunk in question can be changed.
- base\_addr
  - It is the lowest virtual starting address in the area where the chunk is assigned, and uses an address that is aligned on a page-by-page basis.
  - Since it is aligned on a page-by-page basis, the number of offset addresses that are aligned is contained in the start\_offset below.
    - Example: If you are using 4K page units, the start address is 0x1\_2340, and the size is 0x1000.
      - base\_addr=0x1\_2000, start\_offset=0x340, end\_offset=0xcb0
- start\_offset
  - This is the gap between the start address of the area to be managed by the chunk and the start address sorted down on a page-by-page basis.
- end\_offset
  - This is the gap between the end address of the area to be managed by the chunk and the end address sorted up on a page-by-page basis.
- nr\_pages
  - A chunk is the number of pages that can be used.
- nr\_populated
  - Within the chunk, the actual physics pages have been allocated and mapped, and the number of pages currently in use.
  - If you allocated a new chunk, this value is 0 because the actual page is not mapped.
- nr\_empty\_pop\_pages
  - The number of pages within the chunk where the actual physical pages have been allocated and mapped, but are currently unused.
  - If you allocated a new chunk, this value is 0 because the actual page is not mapped.
- populated[]
  - This arrangement is for managing the areas within the chunk where page assignments are made and which are not.
  - The first chunk uses memory that the actual page has been allocated and mapped to allocated, so all pages are set to 1.
  - The added chunk will have all pages set to 0 until the actual area is allocated and used.
    - The vmalloc space is managed as in use, but no physical pages are allocated and no maps are made.

### Members used in the legacy Areamap method

- `contig_hints` (X)
  - It is the largest contiguous space size within a chunk.
  - In the case of the first chunk, the `contig_hint` is the same size as the `free_size`.
  - When an allocation is required, sizes larger than this value are not possible in this chunk.
- `map_used` (X)
  - Indicates the number of items in the `map[]` array within the chunk.
  - The entry of the number +1 contains the last address offset+1 (in-use flag).
    - For example, if `map_used=4`, use 4+1 entries.
- `map_alloc` (X)
  - The maximum number of items in the map array allocated within a chunk. (The initial default value starts at 128 and expands if insufficient.)
- `map[]` (X)
  - The `map[]` array is a mutable array that expands dynamically at runtime.
    - As the map array size expands, the number of `map_alloc` increases accordingly.
    - While the Early Memory Allocator is running, `PERCPU_DYNAMIC_EARLY_SLOTS` (128) arrays are initially declared in the map array.
    - After the Slub Memory Allocator is running, the map can grow larger if it needs more.
  - Distinguish between map array values
    - How to distinguish old array values (legacy method)
      - A positive number is the size of the empty space within a unit.
      - A negative number is the size of the space being used.
    - The new array value separation method was introduced in kernel v2014.3-rc3 in March 15.
      - Note: `percpu`: store offsets instead of lengths in `->map[]`  
(<https://github.com/torvalds/linux/commit/723ad1d90b5663ab623bb3bfba3e4ee7101795d7>)
      - Each byte represents a size, and a value aligned with `sizeof(int)` is used. However, it is used for the last other purpose of the size, so it must be cut and used.
        - The last bit separates the allocation status:
          - 1=in-use status
          - 0=free status
      - e.g. 4 bytes free, 8 in use, 4 in use, 4 free, 12 in use, 100 free, total unit size=132, `map_used=3`
        - Old method: `map[] = {4, -8, -4, 4, -12, 100, 0, }`
        - new method: `map[] = {0, 5, 13, 16, 21, 32, 133, 0, }`
          - `<0,0>`, `<4,1>`, `<12,1>`, `<16,0>`, `<20,1>`, `<32,0>`, `<132,1>` and `<offset, in-use flag>` and 1 byte of offset is used as the in-use flag in actual storage.
          - 0 bytes free from 3~4, 4 in use from 11~8, 12 in use from 15~4, 16 in use from 19~4, 20 in use from 31~32, 32 free from 131~100

## **percpu\_block\_md Struct**

## mm/percpu-internal.h

```

1  /*
2  * pcpu_block_md is the metadata block struct.
3  * Each chunk's bitmap is split into a number of full blocks.
4  * All units are in terms of bits.
5  */

01 struct pcpu_block_md {
02     int                contig_hint;    /* contig hint for block
03     */
04     int                contig_hint_start; /* block relative sta
05     rting                position of the co
06     ntig hint */
07     int                left_free;      /* size of free space al
08     ong                the left side of the
09     block */
10     int                right_free;     /* size of free space al
11     ong                the right side of the
12     block */
13     int                first_free;     /* block position of fir
14     st free */
15 };

```

This is the metadata that is managed by the bitmap in the chunk in per-CPU blocks (the bitmap representation of 1 page in 4-byte increments).

- contig\_hint
  - It contains the number of bits for the maximum free area within the block.
    - e.g. contig\_hint=16 means that there is a maximum free area of 16 x 4 bytes = 64 bytes in this block.
- contig\_hint\_start
  - The bit position where the maximum free area starts within the block
- left\_free
  - If the beginning of this block is the free area, it contains the number of bits in the free area.
- right\_free
  - If the end of this block is in the free region, it contains the number of bits in the free region.
- first\_free
  - It contains an offset bit that points to the first free area of this block.
  - Note: Unlike left\_free and right\_free, it contains positions, not counts.

**pcpu\_group\_info Structure**

include/linux/percpu.h

```

1  struct pcpu_group_info {
2      int                nr_units;      /* aligned # of units */
3      unsigned long      base_offset;   /* base address offset
4      */
5      unsigned int       *cpu_map;      /* unit->cpu map, empty
6      US */
7      /* entries contain NR_CP
8  */
9  };

```

The nodes in the NUMA system are managed as groups, and each group manages the CPUs that belong to that node. The UMA system is a single node, so it uses one group.

- nr\_units
  - The number of units used by that group.
- base\_offset
  - The base\_offset address used by the group
- \*cpu\_map[]
  - Mapping the unit->cpu used by that group

## pcpu\_alloc\_info Struct

include/linux/percpu.h

```

01 struct pcpu_alloc_info {
02     size_t      static_size;
03     size_t      reserved_size;
04     size_t      dyn_size;
05     size_t      unit_size;
06     size_t      atom_size;
07     size_t      alloc_size;
08     size_t      __ai_size;      /* internal, don't use */
09     /* */
10     int          nr_groups;      /* 0 if grouping unneces
11     sary */
12     struct pcpu_group_info groups[];
13 };

```

This struct is stored and used as a local variable within the function during the initial initialization of the per-cpu data in the setup\_per\_cpu\_areas() function. Since the information disappears after initialization, you need to use global variables. (e.g. when doing unit->cpu mapping, etc...)

- static\_size
  - Static area size within 1 unit
  - RPI2: Yes) 0x3ec0
- reserved\_size
  - The reserved area size within one unit can be used when CONFIG\_MODULE kernel options are used.
  - rpi2: 8K.
- dyn\_size
  - The size of the dyn\_size area within one unit varies from one architecture to another.
  - ARM also has a slightly larger value than before, so it basically asks for 32K for 20bit and 64K for 28bit, but if you align 4K with static\_size + reserved\_size + dyn\_size, there will be extra space, which is added to the dyn\_size.
  - rpi2: 0x5140 (initial request value 32x0 to 5000K)
- atom\_size
  - Minimum Assigned Size
  - ARM 4K
- alloc\_size
  - Size to assign

- rpi2: 0xb000 x 4 = 0x3\_c000
- nr\_groups
  - Number of groups (nodes)
  - rpi2: 1

## Key Global Variables

- pcpu\_unit\_pages
  - Number of pages in 1 unit
    - RPI2: e.g. 0xb (1 unit page=11 pages)
- pcpu\_unit\_size
  - 1 Unit Size
    - rpi2: e.g. 0xb000 (1 unit size=44K)
- pcpu\_nr\_units
  - Total number of units
    - RPI2: 4 pcs
- pcpu\_atom\_size
  - As the minimum page size to use for allocation, ARM and ARM64 use PAGE\_SIZE.
  - x86-64-bit systems and a few architectures use huge pages such as 2M pages.
- pcpu\_nr\_slots
  - Number of slots for management
  - Change it to the slot number that corresponds to the unit size and add 2
    - rpi2: If the unit size is 44K, the slot number is 13, so adding 2 is pcpu\_nr\_slots=15
- pcpu\_chunk\_struct\_size
  - The value set when the first chunk is created is the chunk structure size + the number of bytes of the populated bitmap (determined by the unit size).
- pcpu\_low\_unit\_cpu
  - The CPU number where base\_addr is located at the bottom
  - rpi2: 0
- pcpu\_high\_unit\_cpu
  - The CPU number where the base\_addr is located at the top
  - rpi2: 3
- pcpu\_base\_addr
  - The virtual start address to which the first chunk is assigned
- pcpu\_unit\_map[]
  - CPU -> Unit Mapping
  - rpi2: { 0, 1, 2, 3 }
- pcpu\_unit\_offsets[]
  - An array that stores the start offset of each unit
    - rpi2: { 0, 0xb000, 0x1\_6000, 0x2\_1000 }
- pcpu\_nr\_groups
  - Number of groups (number of NUMA nodes)
    - rpi2: 1

- `pcpu_group_offsets[]`
  - An array that stores the starting offset of each group.
    - `rpi2: 0`
- `pcpu_group_sizes[]`
  - An array that stores the size of each group.
    - `rpi2: { 0x2_c000 }`
- `pcpu_first_chunk`
  - `pcpu_chunk` struct pointer to manage dynamic area maps in the First Chunk
- `pcpu_reserved_chunk`
  - `pcpu_chunk` struct pointer to manage reserved area maps in the First Chunk
  - null when not using a module
- `pcpu_reserved_chunk_limit`
  - If there is a reserved zone, the static size plus the reserved\_size
- `pcpu_slot[]`
  - `pcpu_nr_slots` controls the array size, and each array contains a `list_head` that points to a chunk.
  - The last slot is for empty chunks.
  - This arrangement of slots manages chunks that use dynamic areas.
  - `rpi2: static_size=0x3ec0, reserved_size=0x2000, dyn_size=0x5140` example)
    - The `free_size` of the first chunk = `0x5140`, and the slot is slot 12.
    - The reserved chunk has `free_size=0x5ec0`, but the reserved chunk is allocated when the module is loaded, so it is not managed by the `pcpu_slot[]` that manages the dynamic region.
    - Thus, when the per-CPU data is initialized for the first time, `pcpu_slot[12]` is configured with one first chunk as a list, and a spare empty chunk is added to the last slot.
- `pcpu_nr_empty_pop_pages`
  - Indicates the number of empty populated pages in the added chunks that manage only the dynamic area.
  - In this case, the reserved chunk, which manages the reserve area of the first chunk, of course, does not have an effect.
- `pcpu_async_enabled`
  - If you allow async for free in the per-cpu area, the `pcpu_balance_workfn()` function will be scheduled and late.
  - It is enabled when the kernel is initialized and `initcall` is called.
- `pcpu_atomic_alloc_failed`
  - If atomic allocation fails, it will be changed to true.

## consultation

- Per-cpu -1- (Basic) | Question C – Current Article
- Per-cpu -2-(initialize) ([http://jake.dothome.co.kr/setup\\_per\\_cpu\\_areas](http://jake.dothome.co.kr/setup_per_cpu_areas)) | Qc
- per-cpu -3- (dynamic allocation (<http://jake.dothome.co.kr/per-cpu-dynamic>)) | Qc

- Per-cpu -4- (atomic operations) (<http://jake.dothome.co.kr/per-cpu-atomic>) | Qc
- Per-CPU Variables (<http://www.makelinux.net/ldd3/chp-8-sect-5>)
- robust per\_cpu allocation for modules (<https://lwn.net/Articles/180101/>) | LWN.net
- Per-CPU Memory Management (1) (<http://egloos.zum.com/studyfoss/v/5375570>) | F/OSS
- Per-CPU Memory Management (2) (<http://studyfoss.egloos.com/5377666>) | F/OSS
- Documentation/preempt-locking.txt (<https://www.kernel.org/doc/Documentation/preempt-locking.txt>) | kernel.org
- Better per-CPU variables (<https://lwn.net/Articles/258238/>) | LWN.net
- What every programmer should know about memory, Part 1 (<http://lwn.net/Articles/250967/>) | LWN.net
- ARM: implement optimized percpu variable access (<http://lwn.net/Articles/527927/>) | LWN.net

## 9 thoughts to “Per-cpu -1- (Basic)”



**IPARAN (HTTPS://PARANLEE.GITHUB.IO/)**

2021-04-13 11:48 (<http://jake.dothome.co.kr/per-cpu/#comment-305107>)

Hello. Thank you for always keeping an eye on me~

I was watching it again today during the review.

# CPU -> Unit Mapping

Per-CPU data uses NR\_CPUS arrays, which  
can waste memory on <>.

This is because a system that is configured to support thousands of CPUs (NR\_CPUS)  
sometimes has only a few CPUs that are capable of running at actual boot.

I don't know if I'm misunderstanding the part I <> in, but I think the meaning is  
different from the line below.

"If the number of CPUs that the system will use after booting is less than NR\_CPUS"

Is this what I understand correct?

RESPONSE (/PER-CPU/?REPLYTOCOM=305107#RESPOND)



**IPARAN (HTTPS://PARANLEE.GITHUB.IO/)**

2021-04-13 11:50 (<http://jake.dothome.co.kr/per-cpu/#comment-305108>)

[[If you are using a number greater than the number of CPUs that actually exist]]

[[ ... ]] I don't know if I'm misunderstanding the part I put in, or if it's different in  
meaning from the bottom line.



"If the number of CPUs that the system will use after booting is less than NR\_CPUS"

RESPONSE (/PER-CPU/?REPLYTOCOM=305108#RESPOND)



**MOON YOUNG-IL ([HTTP://JAKE.DOTHOME.CO.KR](http://jake.dothome.co.kr))**

2021-04-13 13:40 (<http://jake.dothome.co.kr/per-cpu/#comment-305109>)

Hello?

To explain it a little bit, they mean the same thing:

Since

per-CPU data uses an array of NR\_CPUS numbers, 1) a number greater than the number of CPUs that actually exist > if the number of NR\_CPUS is greater than the number of CPUs that actually exist

2) A number that is greater than the number of CPUs that actually exist > the NR\_CPUS used by the built kernel is a large number, but the booted system has fewer

CPUs 3) A number that is greater than the number of CPUs that actually exist > The number of CPUs that the system will use after booting is less than NR\_CPUS (Iparan)

Note that CONFIG\_NR\_CPUS kernel options have changed a few times in the case of arm64 defaults, and have been increased from 5 to 1 since kernel v64.256. Since then, it hasn't changed until the latest v5.12.

I appreciate it.

RESPONSE (/PER-CPU/?REPLYTOCOM=305109#RESPOND)



**이파란 ([HTTPS://PARANLEE.GITHUB.IO/](https://paranlee.github.io/))**

2021-04-13 18:00 (<http://jake.dothome.co.kr/per-cpu/#comment-305110>)

질문으로 더 깊은 지식을 얻어가네요! 항상 감사드립니다~

응답 (/PER-CPU/?REPLYTOCOM=305110#RESPOND)



**이파란 ([HTTPS://PARANLEE.GITHUB.IO/](https://paranlee.github.io/))**

2021-04-13 18:13 (<http://jake.dothome.co.kr/per-cpu/#comment-305111>)

댓글에 추가적으로 말씀해주신 내용에서

CONFIG\_NR\_CPUS 관련 호기심이 생겨서 질문하나 추가로 드려도 될까요?

제가 만약에 CPU 코어가 4000개 짜리 프로세서가 있으면,

직접 수정해서 NR\_CPUS 값을 4000개로 만들어주어야 per-cpu 데이터를 효과적으로 쓸 수 있나요?

응답 (/PER-CPU/?REPLYTOCOM=305111#RESPOND)



**문영일 (HTTP://JAKE.DOTHOME.CO.KR)**

2021-04-15 08:42 (<http://jake.dothome.co.kr/per-cpu/#comment-305126>)

CONFIG\_NR\_CPUS 값은 빌드할 커널이 지원하는 최대 수의 cpu를 의미합니다.

디폴트 값이 4096일 때 4000개의 cpu를 가진 경우 수정할 필요 없습니다.

만일 8000개의 cpu를 가진 경우 8000개 이상으로 수정해야 시스템이 동작할 수 있습니다.

물론 시스템의 cpu 수와 동일한 CONFIG\_NR\_CPUS 값을 사용하면 메모리 낭비가 없어 좋긴 하지만,

하나의 커널 이미지로 여러 시스템에서 돌릴 수 있게 하기 위해 조금 큰 수를 사용합니다.

감사합니다.

응답 (/PER-CPU/?REPLYTOCOM=305126#RESPOND)



**이파란 (HTTPS://PARANLEE.GITHUB.IO/)**

2021-04-15 09:11 (<http://jake.dothome.co.kr/per-cpu/#comment-305128>)

네 초보적인 질문인데도 잘 설명해주셔서 감사합니다!

좋은 하루되세요~

응답 (/PER-CPU/?REPLYTOCOM=305128#RESPOND)



**이파란 (HTTPS://PARANLEE.GITHUB.IO/)**

2021-04-14 17:30 (<http://jake.dothome.co.kr/per-cpu/#comment-305123>)

항상 감사드립니다.

다음 그림은 가장 낮은 주소의 unit으로 부터 각각의 unit의 시작 offset이 지정되는 것을 보여준다.  
pcpu\_unit\_offset[] 배열의 초기화는 pcpu\_setup\_first\_chunk() 함수를 참고한다.

percpu-3a.png 그림에서

예) node #0 -> 2 cpu

[[ 2 cpu ]] -> [[ 4 cpu ]] 인가요?

응답 (/PER-CPU/?REPLYTOCOM=305123#RESPOND)

**문영일 (HTTP://JAKE.DOTHOME.CO.KR)**2021-04-15 08:46 (<http://jake.dothome.co.kr/per-cpu/#comment-305127>)

지적해주신 것 처럼 첫 번째 노드의 cpu 수는 2가 아니라 4입니다. 오타입니다. ^^;  
node #0 -> 4 cpus로 수정합니다.

감사합니다.

[응답 \(/PER-CPU/?REPLYTOCOM=305127#RESPOND\)](/PER-CPU/?REPLYTOCOM=305127#RESPOND)

## 댓글 남기기

이메일은 공개되지 않습니다. 필수 입력창은 \* 로 표시되어 있습니다

댓글

이름 \*

이메일 \*

웹사이트

댓글 작성

◀ [kernel/head.S – \\_\\_enable\\_mmu: \(http://jake.dothome.co.kr/\\_enable\\_mmu/\)](http://jake.dothome.co.kr/_enable_mmu/)

[Per-cpu -2- \(초기화\) ▶ \(http://jake.dothome.co.kr/setup\\_per\\_cpu\\_areas/\)](http://jake.dothome.co.kr/setup_per_cpu_areas/)

문c 블로그 (2015 ~ 2024)