

# Dynamic DMA mapping Guide

**Author:** David S. Miller <davem@redhat.com>  
**Author:** Richard Henderson <rth@cygnus.com>  
**Author:** Jakub Jelinek <jakub@redhat.com>

This is a guide to device driver writers on how to use the DMA API with example pseudo-code. For a concise description of the API, see DMA-API.txt.

## CPU and DMA addresses

There are several kinds of addresses involved in the DMA API, and it's important to understand the differences.

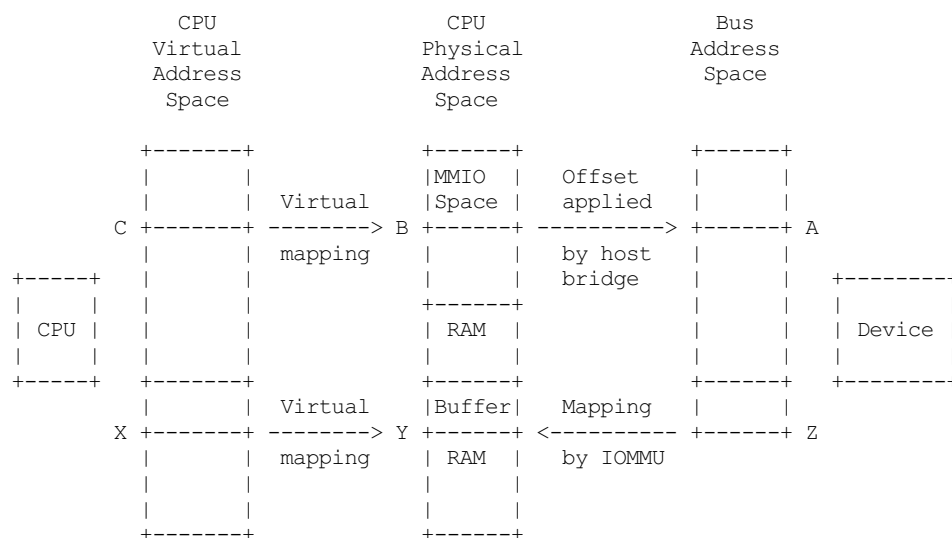
The kernel normally uses virtual addresses. Any address returned by `kmalloc()`, `vmalloc()`, and similar interfaces is a virtual address and can be stored in a `void *`.

The virtual memory system (TLB, page tables, etc.) translates virtual addresses to CPU physical addresses, which are stored as "phys\_addr\_t" or "resource\_size\_t". The kernel manages device resources like registers as physical addresses. These are the addresses in `/proc/iomem`. The physical address is not directly useful to a driver; it must use `ioremap()` to map the space and produce a virtual address.

I/O devices use a third kind of address: a "bus address". If a device has registers at an MMIO address, or if it performs DMA to read or write system memory, the addresses used by the device are bus addresses. In some systems, bus addresses are identical to CPU physical addresses, but in general they are not. IOMMUs and host bridges can produce arbitrary mappings between physical and bus addresses.

From a device's point of view, DMA uses the bus address space, but it may be restricted to a subset of that space. For example, even if a system supports 64-bit addresses for main memory and PCI BARs, it may use an IOMMU so devices only need to use 32-bit DMA addresses.

Here's a picture and some examples:



During the enumeration process, the kernel learns about I/O devices and their MMIO space and the host bridges that connect them to the system. For example, if a PCI device has a BAR, the kernel reads the bus address (A) from the BAR and converts it to a CPU physical address (B). The address B is stored in a struct resource and usually exposed via `/proc/iomem`. When a driver claims a device, it typically uses `ioremap()` to map physical address B at a virtual address (C). It can then use, e.g., `ioread32(C)`, to access the device registers at bus address A.

If the device supports DMA, the driver sets up a buffer using `kmalloc()` or a similar interface, which returns a virtual address (X). The virtual memory system maps X to a physical address (Y) in system RAM. The driver can use virtual address X to access the buffer, but the device itself cannot because DMA doesn't go through the CPU virtual memory system.

In some simple systems, the device can do DMA directly to physical address Y. But in many others, there is IOMMU hardware that translates DMA addresses to physical addresses, e.g., it translates Z to Y. This is part of the reason for the DMA API: the driver can give a virtual address X to an interface like `dma_map_single()`, which sets up any required IOMMU mapping and returns the DMA address Z. The driver then tells the device to do DMA to Z, and the IOMMU maps it to the buffer at address Y in system RAM.

So that Linux can use the dynamic DMA mapping, it needs some help from the drivers, namely it has to take into account that DMA addresses should be mapped only for the time they are actually used and unmapped after the DMA transfer.

The following API will work of course even on platforms where no such hardware exists.

Note that the DMA API works with any bus independent of the underlying microprocessor architecture. You should use the DMA API rather than the bus-specific DMA API, i.e., use the `dma_map_*`() interfaces rather than the `pci_map_*`() interfaces.

First of all, you should make sure:

```
#include <linux/dma-mapping.h>
```

is in your driver, which provides the definition of `dma_addr_t`. This type can hold any valid DMA address for the platform and should be used everywhere you hold a DMA address returned from the DMA mapping functions.

## What memory is DMA'able?

The first piece of information you must know is what kernel memory can be used with the DMA mapping facilities. There has been an unwritten set of rules regarding this, and this text is an attempt to finally write them down.

If you acquired your memory via the page allocator (i.e. `__get_free_page*()`) or the generic memory allocators (i.e. `kmalloc()` or `kmem_cache_alloc()`) then you may DMA to/from that memory using the addresses returned from those routines.

This means specifically that you may not use the memory/addresses returned from `vmalloc()` for DMA. It is possible to DMA to the underlying memory mapped into a `vmalloc()` area, but this requires walking page tables to get the physical addresses, and then translating each of those pages back to a kernel address using something like `__va()`. [ EDIT: Update this when we integrate Gerd Knorr's generic code which does this. ]

This rule also means that you may use neither kernel image addresses (items in data/text/bss segments), nor module image addresses, nor stack addresses for DMA. These could all be mapped somewhere entirely different than the rest of physical memory. Even if those classes of memory could physically work with DMA, you'd need to ensure the I/O buffers were cacheline-aligned. Without that, you'd see cacheline sharing problems (data corruption) on CPUs with DMA-incoherent caches. (The CPU could write to one word, DMA would write to a different one in the same cache line, and one of them could be overwritten.)

Also, this means that you cannot take the return of a `kmap()` call and DMA to/from that. This is similar to `vmalloc()`.

What about block I/O and networking buffers? The block I/O and networking subsystems make sure that the buffers they use are valid for you to DMA from/to.

## DMA addressing capabilities

By default, the kernel assumes that your device can address 32-bits of DMA addressing. For a 64-bit capable device, this needs to be increased, and for a device with limitations, it needs to be decreased.

Special note about PCI: PCI-X specification requires PCI-X devices to support 64-bit addressing (DAC) for all transactions. And at least one platform (SGI SN2) requires 64-bit consistent allocations to operate correctly when the IO bus is in PCI-X mode.

For correct operation, you must set the DMA mask to inform the kernel about your devices DMA addressing capabilities.

This is performed via a call to `dma_set_mask_and_coherent()`:

```
int dma_set_mask_and_coherent(struct device *dev, u64 mask);
```

which will set the mask for both streaming and coherent APIs together. If you have some special requirements, then the following two separate calls can be used instead:

The setup for streaming mappings is performed via a call to `dma_set_mask()`:

```
int dma_set_mask(struct device *dev, u64 mask);
```

The setup for consistent allocations is performed via a call to `dma_set_coherent_mask()`:

```
int dma_set_coherent_mask(struct device *dev, u64 mask);
```

Here, `dev` is a pointer to the device struct of your device, and `mask` is a bit mask describing which bits of an address your device supports. Often the device struct of your device is embedded in the bus-specific device struct of your device. For example, `&pdev->dev` is a pointer to the device struct of a PCI device (`pdev` is a pointer to the PCI device struct of your device).

These calls usually return zero to indicate your device can perform DMA properly on the machine given the address mask you provided, but they might return an error if the mask is too small to be supportable on the given system. If it returns non-zero, your device cannot perform DMA properly on this platform, and attempting to do so will result in undefined behavior. You must not use DMA on this device unless the `dma_set_mask` family of functions has returned success.

This means that in the failure case, you have two options:

1. Use some non-DMA mode for data transfer, if possible.
2. Ignore this device and do not initialize it.

It is recommended that your driver print a kernel `KERN_WARNING` message when setting the DMA mask fails. In this manner, if a user of your driver reports that performance is bad or that the device is not even detected, you can ask them for the kernel messages to find out exactly why.

The standard 64-bit addressing device would do something like this:

```
if (dma_set_mask_and_coherent(dev, DMA_BIT_MASK(64))) {
    dev_warn(dev, "mydev: No suitable DMA available\n");
    goto ignore_this_device;
}
```

If the device only supports 32-bit addressing for descriptors in the coherent allocations, but supports full 64-bits for streaming mappings it would look like this:

```
if (dma_set_mask(dev, DMA_BIT_MASK(64))) {
    dev_warn(dev, "mydev: No suitable DMA available\n");
    goto ignore_this_device;
}
```

The coherent mask will always be able to set the same or a smaller mask as the streaming mask. However for the rare case that a device driver only uses consistent allocations, one would have to check the return value from `dma_set_coherent_mask()`.

Finally, if your device can only drive the low 24-bits of address you might do something like:

```
if (dma_set_mask(dev, DMA_BIT_MASK(24))) {
    dev_warn(dev, "mydev: 24-bit DMA addressing not available\n");
    goto ignore_this_device;
}
```

When `dma_set_mask()` or `dma_set_mask_and_coherent()` is successful, and returns zero, the kernel saves away this mask you have provided. The kernel will use this information later when you make DMA mappings.

There is a case which we are aware of at this time, which is worth mentioning in this documentation. If your device supports multiple functions (for example a sound card provides playback and record functions) and the various different functions have different DMA addressing limitations, you may wish to probe each mask and only provide the functionality which the machine can handle. It is important that the last call to `dma_set_mask()` be for the most specific mask.

Here is pseudo-code showing how this might be done:

```
#define PLAYBACK_ADDRESS_BITS    DMA_BIT_MASK(32)
#define RECORD_ADDRESS_BITS     DMA_BIT_MASK(24)

struct my_sound_card *card;
struct device *dev;

...
if (!dma_set_mask(dev, PLAYBACK_ADDRESS_BITS)) {
    card->playback_enabled = 1;
} else {
    card->playback_enabled = 0;
    dev_warn(dev, "%s: Playback disabled due to DMA limitations\n",
             card->name);
}
if (!dma_set_mask(dev, RECORD_ADDRESS_BITS)) {
    card->record_enabled = 1;
} else {
    card->record_enabled = 0;
    dev_warn(dev, "%s: Record disabled due to DMA limitations\n",
             card->name);
}
```

A sound card was used as an example here because this genre of PCI devices seems to be littered with ISA chips given a PCI front end, and thus retaining the 16MB DMA addressing limitations of ISA.

## Types of DMA mappings

There are two types of DMA mappings:

- Consistent DMA mappings which are usually mapped at driver initialization, unmapped at the end and for which the hardware should guarantee that the device and the CPU can access the data in parallel and will see updates made by each other without any explicit software flushing.

Think of "consistent" as "synchronous" or "coherent".

The current default is to return consistent memory in the low 32 bits of the DMA space. However, for future compatibility you should set the consistent mask even if this default is fine for your driver.

Good examples of what to use consistent mappings for are:

- Network card DMA ring descriptors.
- SCSI adapter mailbox command data structures.
- Device firmware microcode executed out of main memory.

The invariant these examples all require is that any CPU store to memory is immediately visible to the device, and vice versa. Consistent mappings guarantee this.

### Important

Consistent DMA memory does not preclude the usage of proper memory barriers. The CPU may reorder stores to consistent memory just as it may normal memory. Example: if it is important for the device to see the first word of a descriptor updated before the second, you must do something like:

```
desc->word0 = address;
wmb();
desc->word1 = DESC_VALID;
```

in order to get correct behavior on all platforms.

Also, on some platforms your driver may need to flush CPU write buffers in much the same way as it needs to flush write buffers found in PCI bridges (such as by reading a register's value after writing it).

- Streaming DMA mappings which are usually mapped for one DMA transfer, unmapped right after it (unless you use `dma_sync_*` below) and for which hardware can optimize for sequential accesses.

Think of "streaming" as "asynchronous" or "outside the coherency domain".

Good examples of what to use streaming mappings for are:

- Networking buffers transmitted/received by a device.
- Filesystem buffers written/read by a SCSI device.

The interfaces for using this type of mapping were designed in such a way that an implementation can make whatever performance optimizations the hardware allows. To this end, when using such mappings you must be explicit about what you want to happen.

Neither type of DMA mapping has alignment restrictions that come from the underlying bus, although some devices may have such restrictions. Also, systems with caches that aren't DMA-coherent will work better when the underlying buffers don't share cache lines with other data.

## Using Consistent DMA mappings

To allocate and map large (`PAGE_SIZE` or so) consistent DMA regions, you should do:

```
dma_addr_t dma_handle;

cpu_addr = dma_alloc_coherent(dev, size, &dma_handle, gfp);
```

where device is a `struct device *`. This may be called in interrupt context with the `GFP_ATOMIC` flag.

Size is the length of the region you want to allocate, in bytes.

This routine will allocate RAM for that region, so it acts similarly to `__get_free_pages()` (but takes size instead of a page order). If your driver needs regions sized smaller than a page, you may prefer using the `dma_pool` interface, described below.

The consistent DMA mapping interfaces, will by default return a DMA address which is 32-bit addressable. Even if the device indicates (via the DMA mask) that it may address the upper 32-bits, consistent allocation will only return > 32-bit addresses for DMA if the consistent DMA mask has been explicitly changed via `dma_set_coherent_mask()`. This is true of the `dma_pool` interface as well.

`dma_alloc_coherent()` returns two values: the virtual address which you can use to access it from the CPU and `dma_handle` which you pass to the card.

The CPU virtual address and the DMA address are both guaranteed to be aligned to the smallest `PAGE_SIZE` order which is greater than or equal to the requested size. This invariant exists (for example) to guarantee that if you allocate a chunk which is smaller than or equal to 64 kilobytes, the extent of the buffer you receive will not cross a 64K boundary.

To unmap and free such a DMA region, you call:

```
dma_free_coherent(dev, size, cpu_addr, dma_handle);
```

where `dev`, `size` are the same as in the above call and `cpu_addr` and `dma_handle` are the values `dma_alloc_coherent()` returned to you. This function may not be called in interrupt context.

If your driver needs lots of smaller memory regions, you can write custom code to subdivide pages returned by `dma_alloc_coherent()`, or you can use the `dma_pool` API to do that. A `dma_pool` is like a `kmem_cache`, but it uses `dma_alloc_coherent()`, not `__get_free_pages()`. Also, it understands common hardware constraints for alignment, like queue heads needing to be aligned on N byte boundaries.

Create a `dma_pool` like this:

```
struct dma_pool *pool;

pool = dma_pool_create(name, dev, size, align, boundary);
```

The "name" is for diagnostics (like a `knmem_cache` name); `dev` and `size` are as above. The device's hardware alignment requirement for this type of data is "align" (which is expressed in bytes, and must be a power of two). If your device has no boundary crossing restrictions, pass 0 for boundary; passing 4096 says memory allocated from this pool must not cross 4KByte boundaries (but at that time it may be better to use `dma_alloc_coherent()` directly instead).

Allocate memory from a DMA pool like this:

```
cpu_addr = dma_pool_alloc(pool, flags, &dma_handle);
```

flags are GFP\_KERNEL if blocking is permitted (not in interrupt nor holding SMP locks), GFP\_ATOMIC otherwise. Like `dma_alloc_coherent()`, this returns two values, `cpu_addr` and `dma_handle`.

Free memory that was allocated from a `dma_pool` like this:

```
dma_pool_free(pool, cpu_addr, dma_handle);
```

where `pool` is what you passed to `dma_pool_alloc()`, and `cpu_addr` and `dma_handle` are the values `dma_pool_alloc()` returned. This function may be called in interrupt context.

Destroy a `dma_pool` by calling:

```
dma_pool_destroy(pool);
```

Make sure you've called `dma_pool_free()` for all memory allocated from a pool before you destroy the pool. This function may not be called in interrupt context.

## DMA Direction

The interfaces described in subsequent portions of this document take a DMA direction argument, which is an integer and takes on one of the following values:

```
DMA_BIDIRECTIONAL
DMA_TO_DEVICE
DMA_FROM_DEVICE
DMA_NONE
```

You should provide the exact DMA direction if you know it.

DMA\_TO\_DEVICE means "from main memory to the device" DMA\_FROM\_DEVICE means "from the device to main memory" It is the direction in which the data moves during the DMA transfer.

You are strongly encouraged to specify this as precisely as you possibly can.

If you absolutely cannot know the direction of the DMA transfer, specify DMA\_BIDIRECTIONAL. It means that the DMA can go in either direction. The platform guarantees that you may legally specify this, and that it will work, but this may be at the cost of performance for example.

The value DMA\_NONE is to be used for debugging. One can hold this in a data structure before you come to know the precise direction, and this will help catch cases where your direction tracking logic has failed to set things up properly.

Another advantage of specifying this value precisely (outside of potential platform-specific optimizations of such) is for debugging. Some platforms actually have a write permission boolean which DMA mappings can be marked with, much like page protections in the user program address space. Such platforms can and do report errors in the kernel logs when the DMA controller hardware detects violation of the permission setting.

Only streaming mappings specify a direction, consistent mappings implicitly have a direction attribute setting of DMA\_BIDIRECTIONAL.

The SCSI subsystem tells you the direction to use in the 'sc\_data\_direction' member of the SCSI command your driver is working on.

For Networking drivers, it's a rather simple affair. For transmit packets, map/unmap them with the DMA\_TO\_DEVICE direction specifier. For receive packets, just the opposite, map/unmap them with the DMA\_FROM\_DEVICE direction specifier.

## Using Streaming DMA mappings

The streaming DMA mapping routines can be called from interrupt context. There are two versions of each map/unmap, one which will map/unmap a single memory region, and one which will map/unmap a scatterlist.

To map a single region, you do:

```
struct device *dev = &my_dev->dev;
dma_addr_t dma_handle;
void *addr = buffer->ptr;
```

```

size_t size = buffer->len;

dma_handle = dma_map_single(dev, addr, size, direction);
if (dma_mapping_error(dev, dma_handle)) {
    /*
     * reduce current DMA mapping usage,
     * delay and try again later or
     * reset driver.
     */
    goto map_error_handling;
}

```

and to unmap it:

```

dma_unmap_single(dev, dma_handle, size, direction);

```

You should call `dma_mapping_error()` as `dma_map_single()` could fail and return error. Doing so will ensure that the mapping code will work correctly on all DMA implementations without any dependency on the specifics of the underlying implementation. Using the returned address without checking for errors could result in failures ranging from panics to silent data corruption. The same applies to `dma_map_page()` as well.

You should call `dma_unmap_single()` when the DMA activity is finished, e.g., from the interrupt which told you that the DMA transfer is done.

Using CPU pointers like this for single mappings has a disadvantage: you cannot reference HIGHMEM memory in this way. Thus, there is a map/unmap interface pair akin to `dma_{map,unmap}_single()`. These interfaces deal with page/offset pairs instead of CPU pointers. Specifically:

```

struct device *dev = &my_dev->dev;
dma_addr_t dma_handle;
struct page *page = buffer->page;
unsigned long offset = buffer->offset;
size_t size = buffer->len;

dma_handle = dma_map_page(dev, page, offset, size, direction);
if (dma_mapping_error(dev, dma_handle)) {
    /*
     * reduce current DMA mapping usage,
     * delay and try again later or
     * reset driver.
     */
    goto map_error_handling;
}

...

dma_unmap_page(dev, dma_handle, size, direction);

```

Here, "offset" means byte offset within the given page.

You should call `dma_mapping_error()` as `dma_map_page()` could fail and return error as outlined under the `dma_map_single()` discussion.

You should call `dma_unmap_page()` when the DMA activity is finished, e.g., from the interrupt which told you that the DMA transfer is done.

With scatterlists, you map a region gathered from several regions by:

```

int i, count = dma_map_sg(dev, sglist, nents, direction);
struct scatterlist *sg;

for_each_sg(sglist, sg, count, i) {
    hw_address[i] = sg_dma_address(sg);
    hw_len[i] = sg_dma_len(sg);
}

```

where `nents` is the number of entries in the `sglist`.

The implementation is free to merge several consecutive `sglist` entries into one (e.g. if DMA mapping is done with `PAGE_SIZE` granularity, any consecutive `sglist` entries can be merged into one provided the first one ends and the second one starts on a page boundary - in fact this is a huge advantage for cards which either cannot do scatter-gather or have very limited number of scatter-gather entries) and returns the actual number of `sg` entries it mapped them to. On failure 0 is returned.

Then you should loop `count` times (note: this can be less than `nents` times) and use `sg_dma_address()` and `sg_dma_len()` macros where you previously accessed `sg->address` and `sg->length` as shown above.

To unmap a scatterlist, just call:

```

dma_unmap_sg(dev, sglist, nents, direction);

```

Again, make sure DMA activity has already finished.

**Note**

The 'nents' argument to the `dma_unmap_sg` call must be the `_same_` one you passed into the `dma_map_sg` call, it should `_NOT_` be the 'count' value `_returned_` from the `dma_map_sg` call.

Every `dma_map_{single,sg}()` call should have its `dma_unmap_{single,sg}()` counterpart, because the DMA address space is a shared resource and you could render the machine unusable by consuming all DMA addresses.

If you need to use the same streaming DMA region multiple times and touch the data in between the DMA transfers, the buffer needs to be synced properly in order for the CPU and device to see the most up-to-date and correct copy of the DMA buffer.

So, firstly, just map it with `dma_map_{single,sg}()`, and after each DMA transfer call either:

```
dma_sync_single_for_cpu(dev, dma_handle, size, direction);
```

or:

```
dma_sync_sg_for_cpu(dev, sglist, nents, direction);
```

as appropriate.

Then, if you wish to let the device get at the DMA area again, finish accessing the data with the CPU, and then before actually giving the buffer to the hardware call either:

```
dma_sync_single_for_device(dev, dma_handle, size, direction);
```

or:

```
dma_sync_sg_for_device(dev, sglist, nents, direction);
```

as appropriate.

**Note**

The 'nents' argument to `dma_sync_sg_for_cpu()` and `dma_sync_sg_for_device()` must be the same passed to `dma_map_sg()`. It is `_NOT_` the count returned by `dma_map_sg()`.

After the last DMA transfer call one of the DMA unmap routines `dma_unmap_{single,sg}()`. If you don't touch the data from the first `dma_map_*`() call till `dma_unmap_*`(), then you don't have to call the `dma_sync_*`() routines at all.

Here is pseudo code which shows a situation in which you would need to use the `dma_sync_*`() interfaces:

```
my_card_setup_receive_buffer(struct my_card *cp, char *buffer, int len)
{
    dma_addr_t mapping;

    mapping = dma_map_single(cp->dev, buffer, len, DMA_FROM_DEVICE);
    if (dma_mapping_error(cp->dev, mapping)) {
        /*
         * reduce current DMA mapping usage,
         * delay and try again later or
         * reset driver.
         */
        goto map_error_handling;
    }

    cp->rx_buf = buffer;
    cp->rx_len = len;
    cp->rx_dma = mapping;

    give_rx_buf_to_card(cp);
}

...

my_card_interrupt_handler(int irq, void *devid, struct pt_regs *regs)
{
    struct my_card *cp = devid;

    ...
    if (read_card_status(cp) == RX_BUF_TRANSFERRED) {
        struct my_card_header *hp;

        /* Examine the header to see if we wish
         * to accept the data. But synchronize
         * the DMA transfer with the CPU first
         * so that we see updated contents.
         */
        dma_sync_single_for_cpu(&cp->dev, cp->rx_dma,
```

```

        cp->rx_len,
        DMA_FROM_DEVICE);

/* Now it is safe to examine the buffer. */
hp = (struct my_card_header *) cp->rx_buf;
if (header_is_ok(hp)) {
    dma_unmap_single(&cp->dev, cp->rx_dma, cp->rx_len,
                    DMA_FROM_DEVICE);
    pass_to_upper_layers(cp->rx_buf);
    make_and_setup_new_rx_buf(cp);
} else {
    /* CPU should not write to
     * DMA_FROM_DEVICE-mapped area,
     * so dma_sync_single_for_device() is
     * not needed here. It would be required
     * for DMA_BIDIRECTIONAL mapping if
     * the memory was modified.
     */
    give_rx_buf_to_card(cp);
}
}
}

```

Drivers converted fully to this interface should not use `virt_to_bus()` any longer, nor should they use `bus_to_virt()`. Some drivers have to be changed a little bit, because there is no longer an equivalent to `bus_to_virt()` in the dynamic DMA mapping scheme - you have to always store the DMA addresses returned by the `dma_alloc_coherent()`, `dma_pool_alloc()`, and `dma_map_single()` calls (`dma_map_sg()` stores them in the scatterlist itself if the platform supports dynamic DMA mapping in hardware) in your driver structures and/or in the card registers.

All drivers should be using these interfaces with no exceptions. It is planned to completely remove `virt_to_bus()` and `bus_to_virt()` as they are entirely deprecated. Some ports already do not provide these as it is impossible to correctly support them.

## Handling Errors

DMA address space is limited on some architectures and an allocation failure can be determined by:

- checking if `dma_alloc_coherent()` returns NULL or `dma_map_sg` returns 0
- checking the `dma_addr_t` returned from `dma_map_single()` and `dma_map_page()` by using `dma_mapping_error()`:

```

dma_addr_t dma_handle;

dma_handle = dma_map_single(dev, addr, size, direction);
if (dma_mapping_error(dev, dma_handle)) {
    /*
     * reduce current DMA mapping usage,
     * delay and try again later or
     * reset driver.
     */
    goto map_error_handling;
}

```

- unmap pages that are already mapped, when mapping error occurs in the middle of a multiple page mapping attempt. These example are applicable to `dma_map_page()` as well.

Example 1:

```

dma_addr_t dma_handle1;
dma_addr_t dma_handle2;

dma_handle1 = dma_map_single(dev, addr, size, direction);
if (dma_mapping_error(dev, dma_handle1)) {
    /*
     * reduce current DMA mapping usage,
     * delay and try again later or
     * reset driver.
     */
    goto map_error_handling1;
}
dma_handle2 = dma_map_single(dev, addr, size, direction);
if (dma_mapping_error(dev, dma_handle2)) {
    /*
     * reduce current DMA mapping usage,
     * delay and try again later or
     * reset driver.
     */
    goto map_error_handling2;
}

...

```



```
map_error_handling2:
    dma_unmap_single(dma_handle1);
map_error_handling1:
```

### Example 2:

```
/*
 * if buffers are allocated in a loop, unmap all mapped buffers when
 * mapping error is detected in the middle
 */

dma_addr_t dma_addr;
dma_addr_t array[DMA_BUFFERS];
int save_index = 0;

for (i = 0; i < DMA_BUFFERS; i++) {

    ...

    dma_addr = dma_map_single(dev, addr, size, direction);
    if (dma_mapping_error(dev, dma_addr)) {
        /*
         * reduce current DMA mapping usage,
         * delay and try again later or
         * reset driver.
         */
        goto map_error_handling;
    }
    array[i].dma_addr = dma_addr;
    save_index++;
}

...

map_error_handling:

for (i = 0; i < save_index; i++) {

    ...

    dma_unmap_single(array[i].dma_addr);
}

}
```

Networking drivers must call `dev_kfree_skb()` to free the socket buffer and return `NETDEV_TX_OK` if the DMA mapping fails on the transmit hook (`ndo_start_xmit`). This means that the socket buffer is just dropped in the failure case.

SCSI drivers must return `SCSI_MLQUEUE_HOST_BUSY` if the DMA mapping fails in the `queuecommand` hook. This means that the SCSI subsystem passes the command to the driver again later.

## Optimizing Unmap State Space Consumption

On many platforms, `dma_unmap_{single,page}()` is simply a nop. Therefore, keeping track of the mapping address and length is a waste of space. Instead of filling your drivers up with `ifdefs` and the like to "work around" this (which would defeat the whole purpose of a portable API) the following facilities are provided.

Actually, instead of describing the macros one by one, we'll transform some example code.

1. Use `DEFINE_DMA_UNMAP_{ADDR,LEN}` in state saving structures. Example, before:

```
struct ring_state {
    struct sk_buff *skb;
    dma_addr_t mapping;
    __u32 len;
};
```

after:

```
struct ring_state {
    struct sk_buff *skb;
    DEFINE_DMA_UNMAP_ADDR(mapping);
    DEFINE_DMA_UNMAP_LEN(len);
};
```

2. Use `dma_unmap_{addr,len}_set()` to set these values. Example, before:

```
ringp->mapping = FOO;
ringp->len = BAR;
```

after:

```
dma_unmap_addr_set(ringp, mapping, FOO);
dma_unmap_len_set(ringp, len, BAR);
```

3. Use `dma_unmap_{addr,len}()` to access these values. Example, before:

```
dma_unmap_single(dev, ringp->mapping, ringp->len,
DMA_FROM_DEVICE);
```

after:

```
dma_unmap_single(dev,
dma_unmap_addr(ringp, mapping),
dma_unmap_len(ringp, len),
DMA_FROM_DEVICE);
```

It really should be self-explanatory. We treat the ADDR and LEN separately, because it is possible for an implementation to only need the address in order to perform the unmap operation.

## Platform Issues

If you are just writing drivers for Linux and do not maintain an architecture port for the kernel, you can safely skip down to "Closing".

1. Struct scatterlist requirements.

You need to enable `CONFIG_NEED_SG_DMA_LENGTH` if the architecture supports IOMMUs (including software IOMMU).

2. `ARCH_DMA_MINALIGN`

Architectures must ensure that `kmallo`'ed buffer is DMA-safe. Drivers and subsystems depend on it. If an architecture isn't fully DMA-coherent (i.e. hardware doesn't ensure that data in the CPU cache is identical to data in main memory), `ARCH_DMA_MINALIGN` must be set so that the memory allocator makes sure that `kmallo`'ed buffer doesn't share a cache line with the others. See `arch/arm/include/asm/cache.h` as an example.

Note that `ARCH_DMA_MINALIGN` is about DMA memory alignment constraints. You don't need to worry about the architecture data alignment constraints (e.g. the alignment constraints about 64-bit objects).

## Closing

This document, and the API itself, would not be in its current form without the feedback and suggestions from numerous individuals. We would like to specifically mention, in no particular order, the following people:

```
Russell King <rmk@arm.linux.org.uk>
Leo Dagum <dagum@barrel.engr.sgi.com>
Ralf Baechle <ralf@oss.sgi.com>
Grant Grundler <grundler@cup.hp.com>
Jay Estabrook <Jay.Estabrook@compaq.com>
Thomas Sailer <sailer@ife.ee.ethz.ch>
Andrea Arcangeli <andrea@suse.de>
Jens Axboe <jens.axboe@oracle.com>
David Mosberger-Tang <davidm@hpl.hp.com>
```