



Officially these flags can also be used in combination with READ-ONLY and WRITE-ONLY:

```
CL_MEM_WRITE_ONLY | CL_MEM_COPY_HOST_READ_ONLY  
CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_WRITE_ONLY
```

Please let me know, if you have encountered cases where one of these two combinations resulted in faster transfers.

## 3: Host Pointer Flags

Where it gets interesting (and messy) is when is defined how host and device memory interact. This we can describe with host-pointer flags.

First let have a look on the normal way to send over a buffer of floats to the device:

```
[raw]cl_mem cl_input = clCreateBuffer(context, CL_MEM_READ_ONLY,  
    sizeof(float) * count, NULL, NULL);  
int err = clEnqueueWriteBuffer(commands, cl_input, CL_TRUE,  
    0, sizeof(float) * dataCount, p_input, 0, NULL, NULL);  
err = clSetKernelArg(kernel, 0, sizeof(cl_mem), &cl_input);[/raw]
```

Here p\_input is a pointer to an array of floats, cl\_input an OpenCL memory object.

Let us see the definitions of the flags, which influences how data is transferred between host and device.

<empty>. Normal behaviour, with need for explicit writing and reading buffers, as in the above example.

**CL\_MEM\_ALLOC\_HOST\_PTR**. Allocate memory at device, accessible from host. Used for shared memory devices (only ARM MALI has support, afaik).

**CL\_MEM\_COPY\_HOST\_PTR**. Allocate memory at device and initialise with data at host\_ptr.

**CL\_MEM\_USE\_HOST\_PTR**. Allocate memory at device and pin it to host\_ptr.

NB: Even if not explicitly defined, the host access flag **CL\_MEM\_COPY\_HOST\_NO\_ACCESS** cannot be used in combination with a host pointer flag.

## Let OpenCL handle the copy

The flag **CL\_MEM\_COPY\_HOST\_PTR** results in only one big difference with the default. It does the allocation of device-memory and copying to the device in one step. The code looks like:

```
[raw]cl_mem cl_input = clCreateBuffer(context, CL_MEM_READ_ONLY |  
    CL_MEM_COPY_HOST_PTR, sizeof(float) * count, p_input, NULL);  
err = clSetKernelArg(kernel, 0, sizeof(cl_mem), &cl_input);[/raw]
```

It is a matter of taste, what you use (and I personally don't use this version). It has no advantages for faster transfer, known to me. This flag is not meant for getting data back.

There are cases known that buffer-object 'p\_input' in this example cannot be used by clEnqueueWriteBuffer or clEnqueueReadBuffer; in that case **CL\_MEM\_ALLOC\_HOST\_PTR** needs also be specified too. Be careful with this, as this is a contradictive contract.

## Pinned memory

The flag **CL\_MEM\_USE\_HOST\_PTR** enables the so called "pinned memory". Pinned memory makes it possible to use DMA-transfers over PCIe, which is much faster.

The idea of pinned memory is that there is an exact copy of an object at both the device and the host. Exclusive access rights is given to or the host or the device.

This is the above example with pinned memory:

```
[raw]cl_mem cl_input = clCreateBuffer(context, CL_MEM_READ_ONLY |
    CL_MEM_USE_HOST_PTR, sizeof(float) * count, p_input, NULL);
void* p_map_input = clEnqueueMapBuffer(queue, cl_input, CL_TRUE,
    CL_MAP_READ, 0, sizeof(float) * count, 0, NULL, NULL, &err);
// here we can write data to the buffer from the host
clEnqueueUnmapMemObject(queue, cl_input, p_map_input, 0, NULL, NULL);
// here all changes have been sent back to the device
err = clSetKernelArg(kernel, 0, sizeof(cl_mem), &cl_input);
clEnqueueNDRangeKernel(queue, dimension, NULL, global_size,
    local_size, 0, NULL, NULL);
err = clReleaseMemObject(input);[/raw]
```

With `clEnqueueMapBuffer` the access rights are given to the host. With `clEnqueueUnmapMemObject` the exclusive access rights to the memory-objects are given to the device.

For reading from the device, most noticeable difference is `CL_MAP_WRITE`. For read&write, use "`CL_MAP_READ | CL_MAP_WRITE`". This flag triggers memory-transfers:

**<empty>**. undefined

**CL\_MAP\_WRITE**. Will not do any transfers during mapping, will copy data back to the host during unmapping.

**CL\_MAP\_READ**. Will copy data to the device during mapping, will not do any transfers during unmapping.

**CL\_MAP\_READ | CL\_MAP\_WRITE**. Will do transfers during both mapping and unmapping.

Note that the function [clEnqueueMapBuffer](#) has many options for blocking and non-blocking execution - so above example might not be optimal for your case.

## SoCs: CPUs with embedded GPU

If host and device are the same, what to do then? Best is if the pointer to one memory-area could be shared. Problem is that this changes the whole idea of the current computer-model and this takes time. Say you have a long program and you want the GPU to compute on a part of a buffer - this implies that the GPU should have temporary full access to CPU-memory. This is potentially very dangerous (GPU-viruses) and will not become available over night.

Best is to use pinned memory. If the new models have been worked out, the code can be faster without any code-change. Now it copies memory over at the speed of main memory (25GB/s).

Update: ARM MALI supports exactly this. You need to set the flag `CL_MEM_ALLOC_HOST_PTR`.

## More to learn

AMD has introduced a new flag to be added with `CL_MEM_USE_HOST_PTR`: `CL_MEM_USE_PERSISTENT_MEM_AMD`. It claims faster transfer-speeds under Windows 7 only. I expect the usage will be merged with `CL_MEM_USE_HOST_PTR`.

Next step is figuring out asynchronous and non-blocking transfers. Have fun!