

# LINMA2710 - Scientific Computing

## Graphics processing unit (GPU)

*P.-A. Absil and B. Legat*

☐ Full Width Mode    ☐ Present Mode

### **Table of Contents**

**Introduction**

**Examples**

**Reduction on GPU**

Sources

- [OpenCL.jl](#)
- [HandsOnOpenCL](#)
- [Optimizing Parallel Reduction in CUDA](#)
- [Parallel Computation Patterns \(Reduction\)](#)
- [Profiling, debugging and optimization](#)
- [How to use TAU for Performance Analysis](#)

# Introduction ⇄

## Context ⇄

### Critical feedback from ArgTools module

RequestError: HTTP/2 429 while requesting  
[https://upload.wikimedia.org/wikipedia/commons/2/25/WebGL\\_Logo.svg](https://upload.wikimedia.org/wikipedia/commons/2/25/WebGL_Logo.svg)

Drill down into execution genealogy...

```
1 hbox([
2     md""
3     * Most *dedicated* GPUs produced by
4     $(img("https://upload.wikimedia.org/wikipedia/commons/a/a4/NVIDIA_logo.svg",
5         :height => "15pt")) and
6     $(img("https://upload.wikimedia.org/wikipedia/commons/7/7c/AMD_Logo.svg", :height
7         => "15pt"))
8     * *Integrated* GPUs by
9     $(img("https://upload.wikimedia.org/wikipedia/commons/6/6a/Intel_logo_%282020%2C_dar
10         k_blue%29.svg", :height => "15pt")) used in laptops to reduce power consumption
11     * Designed for 3D rendering through ones of the APIs :
12     $(img("https://upload.wikimedia.org/wikipedia/commons/7/7f/Microsoft-DirectX-Logo-
13         wordmark.svg", :height => "20pt")),
14     $(img("https://upload.wikimedia.org/wikipedia/commons/2/21/OpenGL_logo.svg",
15         :height => "20pt")),
16     $(img("https://upload.wikimedia.org/wikipedia/commons/2/25/WebGL_Logo.svg", :height
17         => "20pt")),
18     $(img("https://upload.wikimedia.org/wikipedia/commons/2/2f/WebGPU_logo.svg",
19         :height => "25pt")),
20     $(img("https://upload.wikimedia.org/wikipedia/commons/f/fe/Vulkan_logo.svg",
21         :height => "20pt")) or Apple's Metal
22     $(img("https://upload.wikimedia.org/wikipedia/commons/8/8d/Metal_3_Logo.png",
23         :height => "20pt"))
24     * Illustration on the right is from [Charge's film]
25     (https://studio.blender.org/blog/charge-poster/?utm_medium=homepage), it shows how
```

```

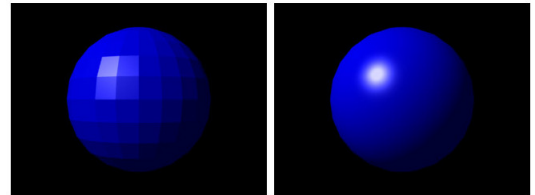
3D modeling works.
7  """
8    img("https://upload.wikimedia.org/wikipedia/commons/f/fd/Charge-
      movie_poster.jpg", :width => "120"),
9  ])

```

👁️ Analyzing stealth implementations

## General-Purpose computing on GPU (GPGPU) ⇄

Also known as *compute shader* as they abuse the programmable shading of GPUs by treating the data as texture maps.



## Critical feedback from ArgTools module

```

RequestError: HTTP/2 429 while requesting
https://upload.wikimedia.org/wikipedia/commons/b/b9/Nvidia_CUDA_Logo.jpg

```

Drill down into execution genealogy...

```

1  grid([
2    md"Hardware-specific"
    img("https://upload.wikimedia.org/wikipedia/commons/b/b9/Nvidia_CUDA_Logo.jpg",
      :height => "70pt")
    img("https://upload.wikimedia.org/wikipedia/commons/7/7b/ROCM_logo.png",
      :height => "60pt") md"```\`
    $(img("https://upload.wikimedia.org/wikipedia/commons/6/6a/Intel_logo_%282020%2C
      _dark_blue%29.svg", :height => "30pt"))
    $(img("https://upload.wikimedia.org/wikipedia/en/f/fa/OneAPI-rgb-3000.png",
      :height => "60pt"))""
3    md"Common interface"
    img("https://upload.wikimedia.org/wikipedia/commons/4/4d/OpenCL_logo.svg")
    img("https://upload.wikimedia.org/wikipedia/commons/1/12/SYCL_logo.svg")
    img("https://d29g4g2dyqv443.cloudfront.net/sites/default/files/akamai/designwork
      s/blog1/OpenACC-logo.png", :height => "50pt")
4  ])

```

👁️ Analyzing stealth implementations



```

6     - Number of processing elements equal to SIMD width
7 * GPUs:
8     - One device per GPU
9     "",
10    img("https://upload.wikimedia.org/wikipedia/de/9/96/Platform_architecture_2009-
11-08.svg", :width => "400pt"),
11 ])

```

👁️ Analyzing stealth implementations

compute device	compute unit	processing element
get_global_id	get_group_id	get_local_id
get_global_size	get_num_groups	get_local_size

## Memory ⇄

### Critical feedback from ArgTools module

RequestError: HTTP/2 429 while requesting  
[https://upload.wikimedia.org/wikipedia/de/d/d1/OpenCL\\_Memory\\_model.svg](https://upload.wikimedia.org/wikipedia/de/d/d1/OpenCL_Memory_model.svg)

Drill down into execution genealogy...

```

1 img("https://upload.wikimedia.org/wikipedia/de/d/d1/OpenCL_Memory_model.svg")

```

👁️ Analyzing stealth implementations

## OpenCL Platforms and Devices ⇄

- Platforms are OpenGL implementations, listed in `/etc/OpenCL/vendors`
- Devices are actual CPUs/GPUs
- ICD allows to change platform at runtime

```
1 OpenCL.versioninfo()
```

```
OpenCL.jl version 0.10.9

Toolchain:
- Julia v1.12.5
- OpenCL: 2024.10.24+1
- SPIRV_LLVM_Backend: 20.1.5+3

Julia packages:
- GPUArrays: 11.4.1
- GPUCompiler: 1.8.2
- KernelAbstractions: 0.9.40
- LLVM: 9.4.6
- SPIRVIntrinsics: 0.5.7

Available platforms: 1
- Portable Computing Language
  OpenCL 3.0, PoCL 7.1 Linux, Release, RELOC, SPIR-V, LLVM 20.1.2jl, SLEEF,
  DISTRO, POCL_DEBUG
  · cpu-haswell-AMD EPYC 7763 64-Core Processor (svm:c+f, usm:h+d, bda, fp64,
  il)
```

See also `clinfo` command line tool and `examples/OpenCL/common/device_info.c`.

## Tip

- ▶ **tl;dr** To refresh the list of platforms, you need to quit Julia and open a new session

## Important stats [↔](#)

- Platform
  - name: Portable Computing Language
  - profile: FULL\_PROFILE
  - vendor: The pocl project
  - version: PoCL 7.1 Linux, Release, RELOC, SPIR-V, LLVM 20.1.2jl, SLEEF, DISTRO, POCL\_DEBUG
- Device
  - name: cpu-haswell-AMD EPYC 7763 64-Core Processor
  - type: cpu

<u>clGetDeviceInfo</u>	<b>Value</b>
CL_DEVICE_GLOBAL_MEM_SIZE	11.62 GiB
CL_DEVICE_MAX_COMPUTE_UNITS	4
CL_DEVICE_LOCAL_MEM_SIZE	512.00 KiB
CL_DEVICE_MAX_WORK_GROUP_SIZE	4096
CL_DEVICE_NATIVE_VECTOR_WIDTH_HALF	0
CL_DEVICE_NATIVE_VECTOR_WIDTH_FLOAT	8
CL_DEVICE_NATIVE_VECTOR_WIDTH_DOUBLE	4
CL_DEVICE_MAX_CLOCK_FREQUENCY	3227 MHz
CL_DEVICE_PROFILING_TIMER_RESOLUTION	1.000 ns

Portable Computing Language ▼

cpu-haswell-AMD EPYC 7763 64-Core Processor ▼

# Examples

## Vectorized sum

```
__kernel void vadd(  
    __global const float *a,  
    __global const float *b,  
    __global float *c,  
    int verbose) {  
    int i = get_global_id(0);  
    c[i] = a[i] + b[i];  
}
```



vadd\_size =  512

vadd\_verbose =  0



```
5 warnings generated.  
CL_KERNEL_WORK_GROUP_SIZE          | 4096  
CL_KERNEL_COMPILE_WORK_GROUP_SIZE  | (0, 0, 0)  
CL_KERNEL_LOCAL_MEM_SIZE           | 0 bytes  
CL_KERNEL_PRIVATE_MEM_SIZE          | 0 bytes  
CL_KERNEL_PREFERRED_WORK_GROUP_SIZE_MULTIPLE | 8  
Send command from host to device    | 1.312 µs  
Including data transfer               | 40.322 ms  
Execution of kernel                  | 58.739 µs
```



Portable Computing Language ▼

cpu-haswell-AMD EPYC 7763 64-Core Processor ▼

vadd (generic function with 1 method)



# Mandelbrot ↻

```
__kernel void mandelbrot(__global float2 *q,
    __global ushort *output, ushort const maxit) {

    int gid = get_global_id(0), it;
    if (gid == 0)
        printf("%d\n", get_num_groups(0));
    float tmp, real = 0, imag = 0;
    output[gid] = 0;
    for(it = 0; it < maxit; it++) {
        tmp = real * real - imag * imag + q[gid].x;
        imag = 2 * real * imag + q[gid].y;
        real = tmp;
        if (real * real + imag * imag > 4.0f)
            output[gid] = it;
    }
}
```



mandel\_size =  512

maxiter =  100

```
1 q = [ComplexF32(r,i) for i=1:-((2.0/mandel_size):-1, r=-1.5:(3.0/mandel_size):0.5];
```

# Critical alert

ArgumentError: Illegal conversion of a OpenCL.cl.UnifiedDeviceMemory to a Ptr{ComplexF32}

## Execution genealogy

Here's the operational timeline, prioritized by recency:

1. `convert(T::Type{...}, mem::OpenCL.cl.UnifiedDeviceMemory) ...show types...`  
sourced from `memory.jl:8`
2. `unsafe_convert(P::Type{...}, mem::OpenCL.cl.UnifiedDeviceMemory) ...show types...`  
sourced from `memory.jl:15`
3. `unsafe_clconvert(typ::Type{...}, mem::OpenCL.cl.UnifiedDeviceMemory) ...show types...`  
sourced from `kernel.jl:440`
4. macro expansion  
sourced from `kernel.jl:352`
5. macro expansion  
sourced from
6. `convert_arguments(::OpenCL.cl.var"#276#277"{...}, ::Type{...}, ::OpenCL.CLArray{...}, ::OpenCL.CLArray{...}, ::Int64) ...show types...`  
sourced from
7. `#clcall#274(::OpenCL.cl.Kernel, ::Type{...}, ::OpenCL.CLArray{...}, ::OpenCL.CLArray{...}, ::Int64; kwargs::@Kwargs{...}) ...show types...`  
sourced from `kernel.jl:339`
8. `anonymous_function() ...show types...`  
sourced from `Adjacent statement: directive 16`

```
14 cl.queue!(::protile) do
15     for _ in 1:num_runs
16         evt = clcall(kernel, args...; kws...)
17         wait(evt)
18     end
end
```

cell preview
9. `Expand visibility...`

```
1 mandel_image = mandel(q, maxiter, mandel_device; global_size=length(q));
```



```
1 warning generated.  
CL_KERNEL_WORK_GROUP_SIZE          4096  
CL_KERNEL_COMPILE_WORK_GROUP_SIZE  (0, 0, 0)  
CL_KERNEL_LOCAL_MEM_SIZE           0 bytes  
CL_KERNEL_PRIVATE_MEM_SIZE         0 bytes  
CL_KERNEL_PREFERRED_WORK_GROUP_SIZE_MULTIPLE 8
```



Portable Computing Language ▼

cpu-haswell-AMD EPYC 7763 64-Core Processor ▼

## Critical alert

*Another statement defining `mandel_image` exhibits suboptimal performance indicators.*

```
1 aside(CairoMakie.image(CairoMakie.rotr90(mandel_image)), v_offset = -400)
```

Analyzing stealth implementations

mandel (generic function with 1 method)

```
1 function mandel(q::Array{ComplexF32}, maxiter::Int64, device; kws...)
2     cl.device!(device)
3     q = CLArray{q}
4     o = CLArray{Cushort}(undef, size(q))
5
6     prg = cl.Program(; source = mandel_source.code) |> cl.build!
7     k = cl.Kernel(prg, "mandelbrot")
8
9     timed_clcall(k, Tuple{Ptr{ComplexF32}, Ptr{Cushort}, Cushort},
10                q, o, maxiter; kws...)
11
12     return Array(o)
13 end
```

```
1 mandel_source = code(Example("OpenCL/mandelbrot/mandel.cl"));
```

# Compute $\pi$ ⇄

## Critical alert

ArgumentError: Illegal conversion of a OpenCL.cl.UnifiedDeviceMemory to a Ptr{Float32}

## Execution genealogy

Here's the operational timeline, prioritized by recency:

1. `convert(T::Type{...}, mem::OpenCL.cl.UnifiedDeviceMemory) ...show types...`  
sourced from `memory.jl:8`
2. `unsafe_convert(P::Type{...}, mem::OpenCL.cl.UnifiedDeviceMemory) ...show types...`  
sourced from `memory.jl:15`
3. `unsafe_clconvert(typ::Type{...}, mem::OpenCL.cl.UnifiedDeviceMemory) ...show types...`  
sourced from `kernel.jl:440`
4. macro expansion  
sourced from `kernel.jl:352`
5. macro expansion  
sourced from
6. `convert_arguments(::OpenCL.cl.var"#276#277"{...}, ::Type{...}, ::Int64, ::Float64, ::OpenCL.cl.LocalMem{...}, ::OpenCL.CLArray{...}) ...show types...`  
sourced from
7. `#clcall#274(::OpenCL.cl.Kernel, ::Type{...}, ::Int64, ::Float64, ::OpenCL.cl.LocalMem{...}, ::OpenCL.CLArray{...}; kwargs::@Kwargs{...}) ...show types...`  
sourced from `kernel.jl:339`
8. `anonymous function() ...show types...`  
sourced from `Adjacent statement: directive 16`

```
14 cl.queue!(::profile) do
15     for _ in 1:num_runs
16         evt = clcall(kernel, args...; kws...)
17         wait(evt)
18     end
end
```

cell preview

## 9. Expand visibility...

1 `mypi()`

```
CL_KERNEL_WORK_GROUP_SIZE      4096
CL_KERNEL_COMPILE_WORK_GROUP_SIZE (0, 0, 0)
CL_KERNEL_LOCAL_MEM_SIZE      0 bytes
CL_KERNEL_PRIVATE_MEM_SIZE     0 bytes
CL_KERNEL_PREFERRED_WORK_GROUP_SIZE_MULTIPLE 8
```

### ► How to compute $\pi$ with a kernel ?

Portable Computing Language ▼

cpu-haswell-AMD EPYC 7763 64-Core Processor ▼

`mypi` (generic function with 1 method)

## First element ⇄

Let's write a simple kernel that returns the first element of a vector in global memory.

```
--kernel void first_el(__global float* glob, __global float* result) {
    int item = get_local_id(0);
    if (item == 0)
        *result = glob[item];
}
```

0.9104708f0


1 `first_el(rand(Float32, first_el_len))`

```
CL_KERNEL_WORK_GROUP_SIZE      4096
CL_KERNEL_COMPILE_WORK_GROUP_SIZE (0, 0, 0)
CL_KERNEL_LOCAL_MEM_SIZE      0 bytes
CL_KERNEL_PRIVATE_MEM_SIZE     0 bytes
CL_KERNEL_PREFERRED_WORK_GROUP_SIZE_MULTIPLE 8
Send command from host to device 2.675 µs
Including data transfer          23.116 ms
Execution of kernel              19.005 µs
```

first\_el (generic function with 1 method)

Portable Computing Language ▼

cpu-haswell-AMD EPYC 7763 64-Core Processor ▼

first\_el\_len =  16

## Copy to local memory ⇄

```
__kernel void copy_to_local(__global float* glob, __local float* shared) {  
    int global_size = get_global_size(0);  
    int local_size = get_local_size(0);  
    int item = get_local_id(0);  
    shared[item] = 0;  
    for (int i = 0; i < global_size; i += local_size) {  
        shared[item] += glob[i + item];  
    }  
}
```


1 [copy\\_to\\_local](#)([copy\\_global\\_len](#), [copy\\_local\\_len](#))

```
CL_KERNEL_WORK_GROUP_SIZE          4096  
CL_KERNEL_COMPILE_WORK_GROUP_SIZE (0, 0, 0)  
CL_KERNEL_LOCAL_MEM_SIZE           0 bytes  
CL_KERNEL_PRIVATE_MEM_SIZE         0 bytes  
CL_KERNEL_PREFERRED_WORK_GROUP_SIZE_MULTIPLE 8  
Send command from host to device    1.863 µs  
Including data transfer              26.215 ms  
Execution of kernel                 25.829 µs
```

copy\_to\_local (generic function with 1 method)

Portable Computing Language ▼

cpu-haswell-AMD EPYC 7763 64-Core Processor ▼

copy\_global\_len =  16

copy\_local\_len =  16



# Reduction on GPU ↺

Many operations can be framed in terms of a MapReduce operation.

- Given a vector of data
- It first map each elements through a given function
- It then reduces the results into a single element

The mapping part is easily embarassingly parallel but the reduction is harder to parallelize. Let's see how this reduction step can be achieved using arguably the simplest example of mapreduce , the sum (corresponding to an identity map and a reduction with +).

## Sum ↺

► (OpenCL = 8.81811, Classical = 8.81811)

1 `local_sum(global_len, local_len, local_code, local_device)`



```
CL_KERNEL_WORK_GROUP_SIZE          4096
CL_KERNEL_COMPILE_WORK_GROUP_SIZE  (0, 0, 0)
CL_KERNEL_LOCAL_MEM_SIZE           0 bytes
CL_KERNEL_PRIVATE_MEM_SIZE         0 bytes
CL_KERNEL_PREFERRED_WORK_GROUP_SIZE_MULTIPLE 8
Send command from host to device    1.202 µs
Including data transfer              48.031 ms
Execution of kernel                 15.729 µs
```




► **How to compute the sum an array in local memory with a kernel ?**

local\_sum (generic function with 1 method)

Portable Computing Language ▼

cpu-haswell-AMD EPYC 7763 64-Core Processor ▼

global\_len =  16

local\_len =  16



# Blocked sum ↔

► (OpenCL = 8.81811, Classical = 8.81811)

1 `block_local_sum(block_global_len, block_local_len, factor)`

```
CL_KERNEL_WORK_GROUP_SIZE      4096
CL_KERNEL_COMPILE_WORK_GROUP_SIZE (0, 0, 0)
CL_KERNEL_LOCAL_MEM_SIZE      0 bytes
CL_KERNEL_PRIVATE_MEM_SIZE     0 bytes
CL_KERNEL_PREFERRED_WORK_GROUP_SIZE_MULTIPLE 8
Send command from host to device 1.523 µs
Including data transfer          35.418 ms
Execution of kernel             25.679 µs
```


► **How to reduce the amount of barrier synchronizations ?**

► **Was it beneficial in terms of performance for GPUs like in the case of OpenMP ?**

`block_local_sum` (generic function with 1 method)

Portable Computing Language ▼

cpu-haswell-AMD EPYC 7763 64-Core Processor ▼

`block_global_len` =  16

`block_local_len` =  16

`factor` =  16

## Back to SIMD ⇄

---

- Also called Single Instruction Multiple Threads (SIMT)
- CUDA Warp : width of 32 threads
- AMD wavefront : width of 64 threads
- In general : CL\_KERNEL\_PREFERRED\_WORK\_GROUP\_SIZE\_MULTIPLE
- Consecutive `get_local_id()` starting from 0
  - So the thread of local id from 0 to 31 are in the same CUDA warp.
- Threads execute the **same instruction** at the same time so no need for barrier .

## Warp divergence ⇄

---

Suppose a kernel is executed on a nvidia GPU with `global_size` threads. How much time will it take to execute it ?

```
__kernel void diverge(n)
{
    int item = get_local_id(0);
    if (item < n) {
        do_task_A(); // 'a' ns
    } else {
        do_task_B(); // 'b' ns
    }
}
```



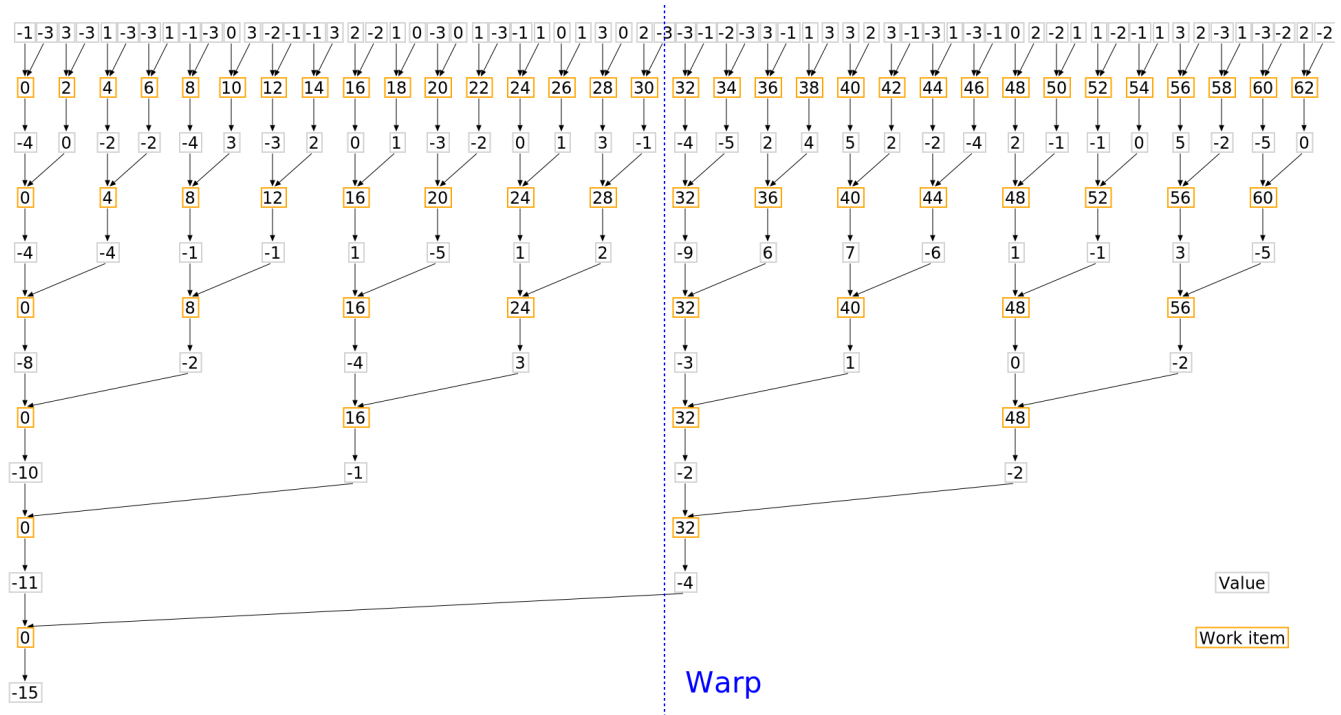
► **How much time will it take to execute it if `global_size` is 32 and `n` is 16 ?**

► **How much time will it take to execute it if `global_size` is 64 and `n` is 32 ?**

Are the threads that are still active in the same warp for you sum example ?

## Warp diversion for our sum ↻

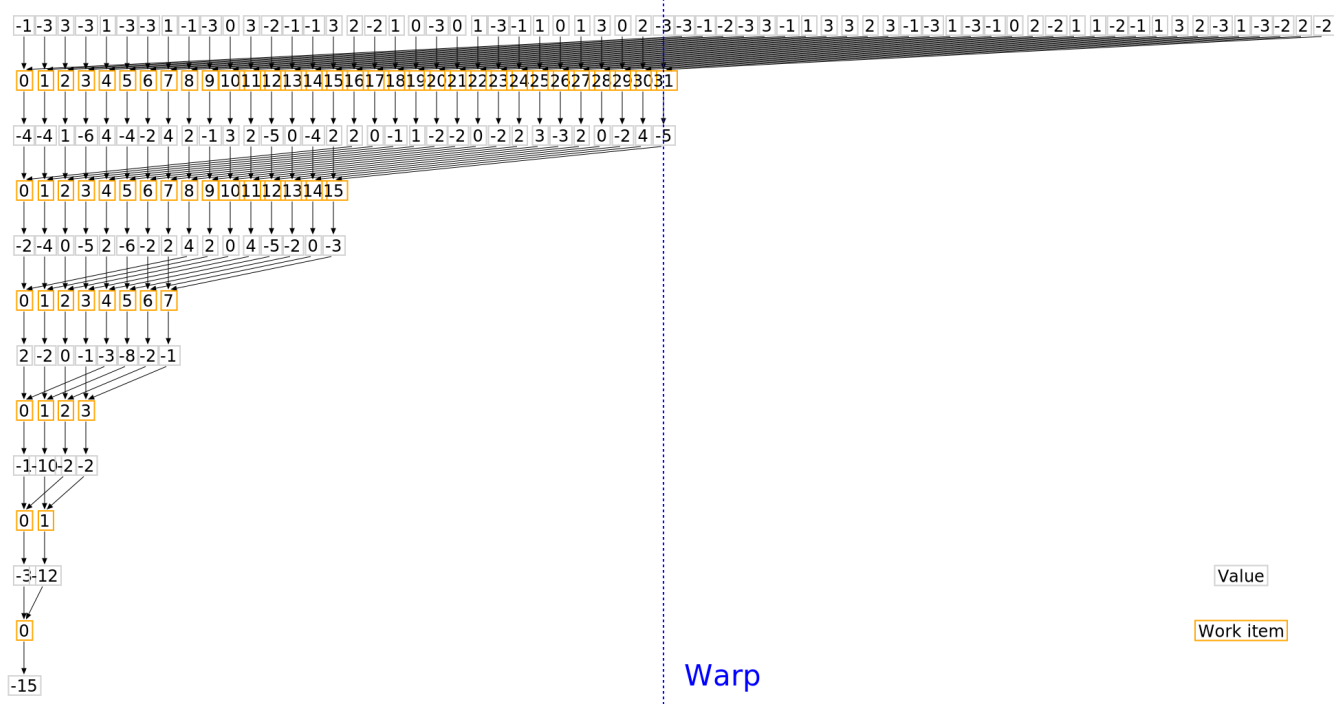
► We are still using different warps until the end. Is that a good thing ?



How should we change the sum to keep the working threads on the same warp ?

## No warp divergence ↻

Now the same warp is used for all threads so we don't need `barrier` and it frees other warps to stay idle (reducing power consumption) or do other tasks.



## Reordered local sum ⇄


```
__kernel void local_sum(__local float* shared)
{
    int items = get_local_size(0);
    int item = get_local_id(0);
    int stride = items / 2;
    float other_val = 0;
    while (stride > 0) {
        barrier(CLK_LOCAL_MEM_FENCE);
        if (item < stride) {
            other_val = 0;
            if (item + stride < items)
                other_val = shared[item+stride];
            shared[item] += other_val;
        }
        stride /= 2;
    }
}
```


► (OpenCL = 8.81811, Classical = 8.81811)

```
CL_KERNEL_WORK_GROUP_SIZE      4096
CL_KERNEL_COMPILE_WORK_GROUP_SIZE (0, 0, 0)
CL_KERNEL_LOCAL_MEM_SIZE      0 bytes
CL_KERNEL_PRIVATE_MEM_SIZE    0 bytes
CL_KERNEL_PREFERRED_WORK_GROUP_SIZE_MULTIPLE 8
Send command from host to device 531.000 ns
Including data transfer         42.802 ms
Execution of kernel            24.536 µs
```

Portable Computing Language ▼

cpu-haswell-AMD EPYC 7763 64-Core Processor ▼

reordered\_global\_size =  16

reordered\_local\_size =  16

## SIMT sum ⇄

```
__kernel void simt_sum(volatile __local float* shared)
{
    int items = get_local_size(0);
    int item = get_local_id(0);
    barrier(CLK_LOCAL_MEM_FENCE);
    while (items > 1) {
        items /= 2;
        shared[item] += shared[item + items];
    }
}
```

► (OpenCL = 2.52133, Classical = 8.81811)

```
1 local_sum(simt_global_size, simt_local_size, simt_code, simt_device)
```



```
CL_KERNEL_WORK_GROUP_SIZE      4096
CL_KERNEL_COMPILE_WORK_GROUP_SIZE (0, 0, 0)
CL_KERNEL_LOCAL_MEM_SIZE      0 bytes
CL_KERNEL_PRIVATE_MEM_SIZE    0 bytes
CL_KERNEL_PREFERRED_WORK_GROUP_SIZE_MULTIPLE 8
Send command from host to device 992.000 ns
Including data transfer        30.810 ms
Execution of kernel           23.353 µs
```



► **Why don't we check any condition on `item`, aren't some thread computing data that won't be used ?**

Portable Computing Language ▼

cpu-haswell-AMD EPYC 7763 64-Core Processor ▼

simt\_global\_size =  16

simt\_local\_size =  16

## Beware!

POCL does not synchronize, even for `simt_len <= 8`

► **Why do we need `volatile` ?**

## Unrolled sum ⇄

► **How to get even faster performance by assuming that `items` is a power of 2 smaller than 512 and that the SIMT width is 32 ?**

► (OpenCL = 2.52133, Classical = 8.81811)

```
1 local_sum(unrolled_global_size, unrolled_local_size, unrolled_code, unrolled_device)
```




```
CL_KERNEL_WORK_GROUP_SIZE          | 4096
CL_KERNEL_COMPILE_WORK_GROUP_SIZE  | (0, 0, 0)
CL_KERNEL_LOCAL_MEM_SIZE           | 0 bytes
CL_KERNEL_PRIVATE_MEM_SIZE         | 0 bytes
CL_KERNEL_PREFERRED_WORK_GROUP_SIZE_MULTIPLE | 8
Send command from host to device    | 862.000 ns
Including data transfer              | 34.929 ms
Execution of kernel                  | 23.775 µs
```




### ► How to have portable code using unrolling ?

Portable Computing Language ▼

cpu-haswell-AMD EPYC 7763 64-Core Processor ▼

unrolled\_global\_size =  16

unrolled\_local\_size =  16

## Utils

\_pretty\_time (generic function with 1 method)

```
1 _pretty_time(x) = BenchmarkTools.prettytime(minimum(x))
```

timed\_clcall (generic function with 1 method)

```
1 function timed_clcall(kernel, args...; kws...)
2   info = cl.work_group_info(kernel, cl.device())
3   # See
4   https://registry.khronos.org/OpenCL/sdk/3.0/docs/man/html/clGetKernelWorkGroupInfo.html
5   println("CL_KERNEL_WORK_GROUP_SIZE | ", info.size)
6   println("CL_KERNEL_COMPILE_WORK_GROUP_SIZE | ", info.compile_size)
7   println("CL_KERNEL_LOCAL_MEM_SIZE | ",
8   BenchmarkTools.prettymemory(info.local_mem_size))
9   println("CL_KERNEL_PRIVATE_MEM_SIZE | ",
10  BenchmarkTools.prettymemory(info.private_mem_size))
11  println("CL_KERNEL_PREFERRED_WORK_GROUP_SIZE_MULTIPLE | ",
12  info.preferred_size_multiple)
13
14  # ':profile' sets 'CL_QUEUE_PROFILING_ENABLE' to the command queue
15  queued_submit = Float64[]
16  submit_start = Float64[]
17  start_end = Float64[]
18  cl.queue!(:profile) do
19    for _ in 1:num_runs
20      evt = clcall(kernel, args...; kws...)
21      wait(evt)
22
23      # See
24      https://registry.khronos.org/OpenCL/sdk/3.0/docs/man/html/clGetEventProfilingInfo.html
25      push!(queued_submit, evt.profile_submit - evt.profile_queued)
26      push!(submit_start, evt.profile_start - evt.profile_submit)
27      push!(start_end, evt.profile_end - evt.profile_start)
28    end
29  end
30
31  println("Send command from host to device | $(_pretty_time(queued_submit))")
32  println("Including data transfer | $(_pretty_time(submit_start))")
33  println("Execution of kernel | $(_pretty_time(start_end))")
34 end
```

num\_runs =  1



Activating project at '~/work/LINMA2710/LINMA2710/Lectures'





```
1 using OpenCL, pocl_jll # 'pocl_jll' provides the POCL OpenCL platform for CPU devices
```