# SYCL memory models & SparseCCL

Master 2 internship -> PhD student

IJCLab - CNRS

David Chamont

Hadrien Grasland

SAMOVAR - Télécom SudParis

Elisabeth Brunet

Gaël Thomas

Sylvain Joube – 2021-09-17

ACTS Parallelization Meetings

# Outline

Why SYCL & SparseCCL

Investigation with a micro-benchmark
- Explicit memory copy
- In-place filling

Investigation with SparseCCL
- Graph pointer
- Flat arrays

Conclusion & future work

# Why SYCL ?

Similar to C++ standard

Compatible with many hardware chips

Based on an international consortium

Independent from any single manufacturer

*Multiple memory models*

# SYCL USM memory models

Unified Shared Memory [USM, SYCL 2020]

- "USM device" : located on GPU, explicit transfer
- "USM host" : host-pinned, device-accessible [main memory]
- "USM shared" : implementation decides data location

*Buffers and accessors [not tested here, SYCL 1.2]*

- *Dependency graph between kernels*

# Why SparseCCL?

Easy to understand and implement

First step in the track reconstruction

Only relies on csv input data

*First results hard to understand…*

# Why a microbenchmark ?

Test SYCL memory on a simple memory-bound program
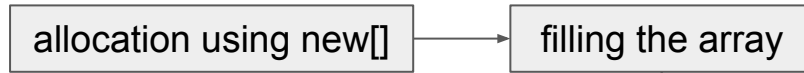
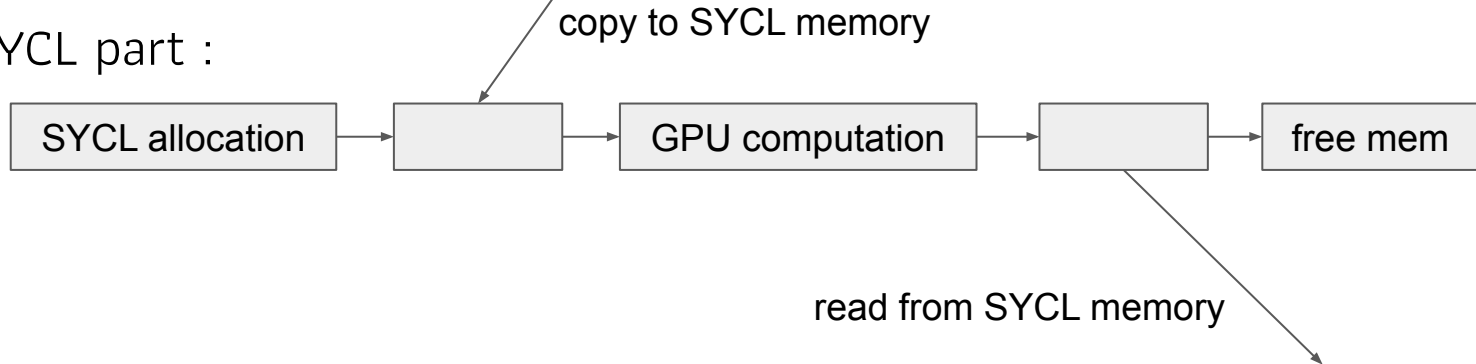Understand what to expect from SparseCCL on SYCL

*Data read only once*

# Microbenchmark : a simple reduction on GPU

Two main steps :

- Generic C++ code on CPU :

| allocation using new[] | → | filling the array |

copy to SYCL memory

- SYCL part :

| SYCL allocation | → | | → | GPU computation | → | | → | free mem |

read from SYCL memory

# What will be measured here ?

SYCL memory behaviour

- With all USM models [device, shared, host]
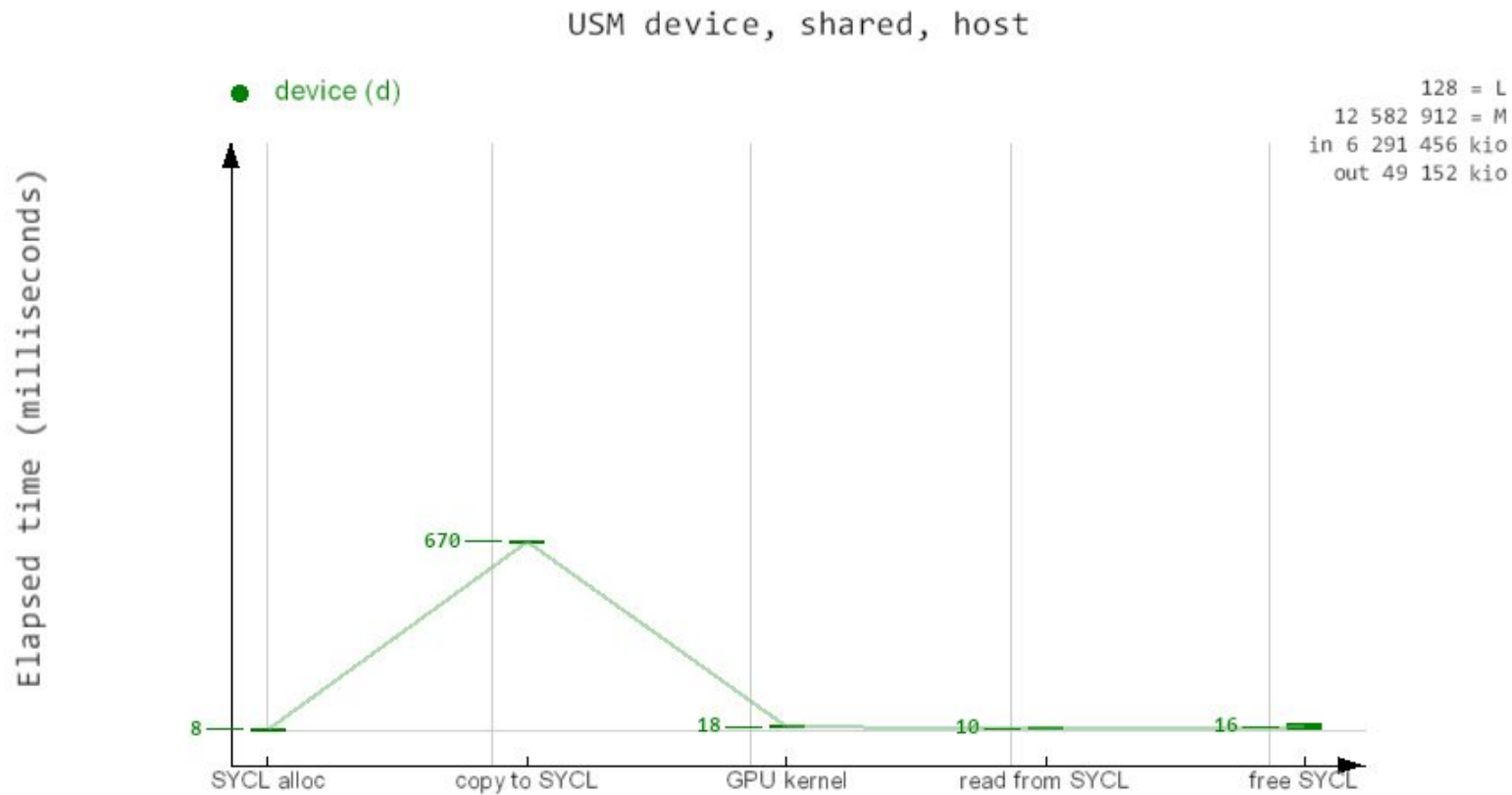- On a single server[1]
- Microbenchmark + SparseCCL

Dedicated server :

- 512 GB DDR4, 2400 MHz
- NVidia Quadro RTX 5000 : 16 GB GDDR6, 448 GB/s
- Debian 11.0 stable

[1]*singe server : but consistent quick cross-checked elsewhere with other hardware*

# USM models comparison



USM device, shared, host

device (d)

128 = L
12 582 912 = M
in 6 291 456 kio
out 49 152 kio

Elapsed time (milliseconds)

670

8          18          10          16

SYCL alloc    copy to SYCL    GPU kernel    read from SYCL    free SYCL

# USM models comparison



USM device, shared, host

● device (d)    ● shared (s)

128 = L
12 582 912 = M
in 6 291 456 kio
out 49 152 kio

Elapsed time (milliseconds)

s 960
d 670

d 8
s 1

s 24
d 18

s 23
d 10

s 166
d 16

SYCL alloc    copy to SYCL    GPU kernel    read from SYCL    free SYCL

# USM models comparison
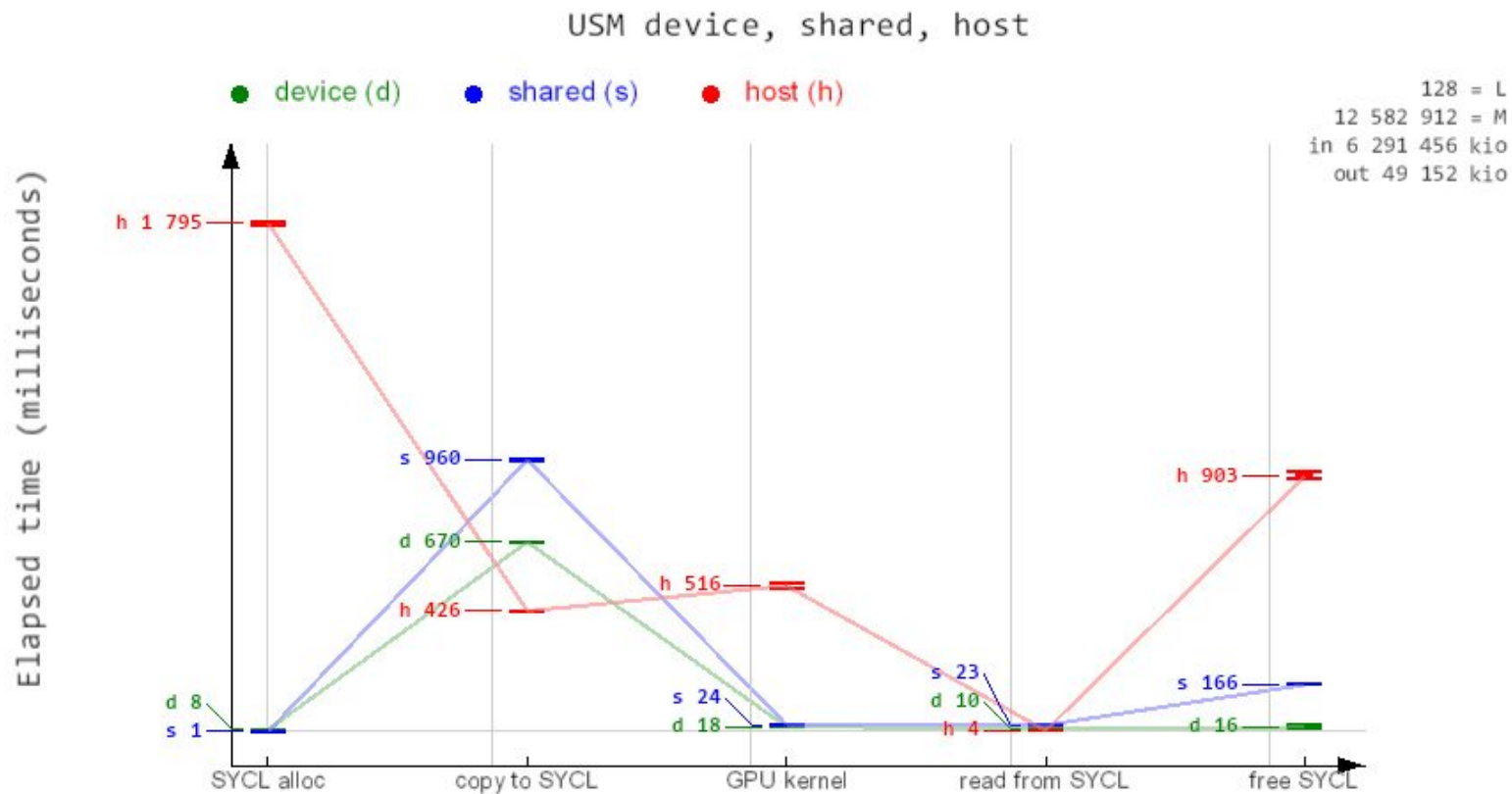


USM device, shared, host

# First findings

USM host :

- for unique memory accesses
- should reuse allocated memory

USM host and shared : also for data bigger than GPU memory

# In-place SYCL memory filling ?



USM device, shared, host

device (d)    shared (s)    host (h)

128 = L
12 582 912 = M
in 6 291 456 kio
out 49 152 kio

Elapsed time (milliseconds)

h 1 795

s 960

d 670

h 903

h 426

h 516

s 24
d 18

s 23
d 10

h 4

s 166

d 8
s 1

d 16

SYCL alloc    copy to SYCL    GPU kernel    read from SYCL    free SYCL

# In-place SYCL memory filling ?



USM device, shared, host + direct access to SYCL mem

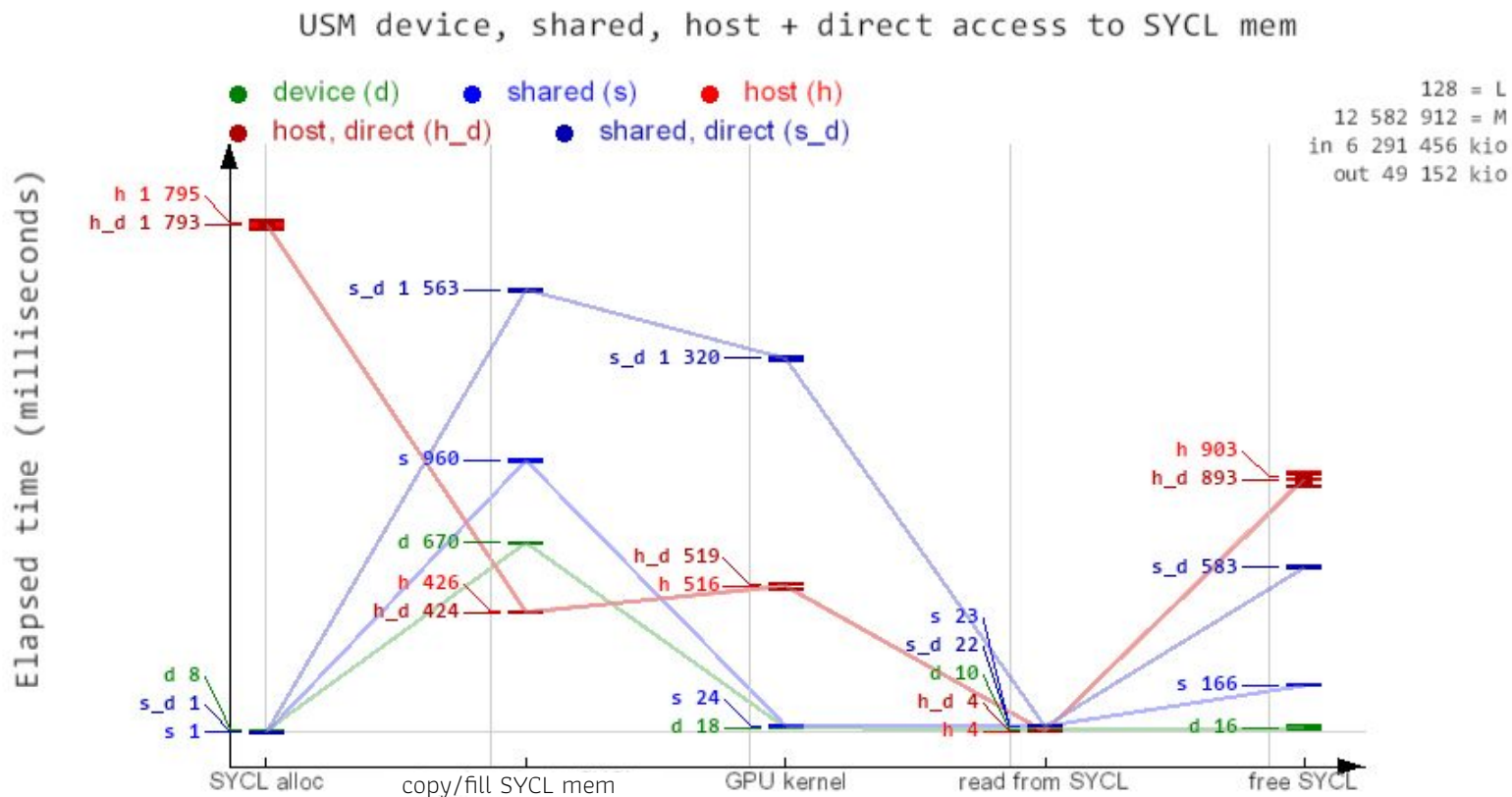# In-place SYCL memory filling



USM device, shared, host + direct access to SYCL mem

# More findings

USM shared :

- very expensive in-place filling

USM host :

- can be used as a local memory on host, with no extra cost

# SparseCCL

Questions :

- Previous results still applies ?
- How does the USMs compare ?
- How hard is adapting existing code to SYCL ?

Data structures options :

1. Pointer graph : usual structures
2. Flat arrays : made for GPUs

# SparseCCL, CPU only

Grouping timing of allocation & fill : to avoid lazy allocation bias



SparseCCL - CPU memory

CPU, pointer graph

in 165 540 kio
out 85 800 kio

Elapsed time μs

209 856

79 717

17 425

21 238

CPU alloc & fill        CPU compute        read from CPU mem        free CPU mem

# SparseCCL, CPU only : pointers vs flat

Similar between arrays and pointer graph. [10% on compute]



SparseCCL - CPU memory

● CPU, pointer graph       ● CPU, flat (f)

in 165 540 kio
out 85 800 kio

Elapsed time μs

209 856
f 182 303

f 83 531
79 717

17 425
f 5 017

21 238
f 814

CPU alloc & fill        CPU compute        read from CPU mem        free CPU mem

# SparseCCL on GPU - USM host



SparseCCL - host memory

● host, pointer graph

in 165 540 kio
out 85 800 kio

Elapsed time μs

9 659 879

5 640 651

2 034 793

23 543

SYCL alloc & fill        GPU kernel        read from SYCL mem        free SYCL mem

# SparseCCL on GPU - USM host

Allocation x80 ! Free x60 !

Global time x4



SparseCCL - host memory

● host, pointer graph     ● host, flat (f)

in 165 540 kio
out 85 800 kio

Elapsed time μs

9 659 879

5 640 651
f 4 243 436

2 034 793

23 543
f 5 189

f 126 027

f 33 126

SYCL alloc & fill     GPU kernel     read from SYCL mem     free SYCL mem

# SparseCCL on GPU - USM shared

[same scale than before]



SparseCCL - shared memory

● shared, pointer graph

in 165 540 kio
out 85 800 kio

Elapsed time μs

9 472 840

3 024 504

95 884          105 693

SYCL alloc & fill          GPU kernel          read from SYCL mem          free SYCL mem

# SparseCCL on GPU - USM shared

Faster kernel => much increased cost

Global time x50 !



SparseCCL - shared memory

shared, pointer graph    shared, flat (f)

in 165 540 kio
out 85 800 kio

Elapsed time μs

9 472 840

3 024 504

f 92 283     95 884     105 693     3 024 504
             f 62 855    f 39 942    f 44 293

SYCL alloc & fill     GPU kernel     read from SYCL mem     free SYCL mem

# Findings : pointer graph

Easy to replace malloc with malloc_[host/shared] but expensive
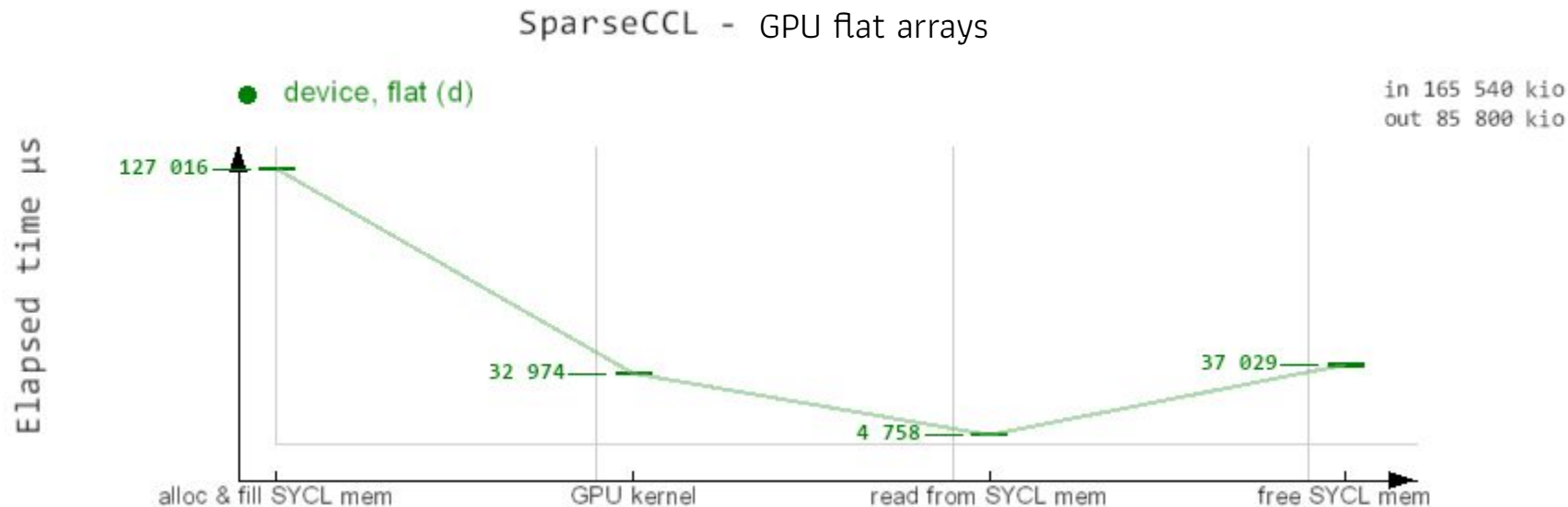
Needs structure adjustments to be efficient on GPU

USM device needs flat arrays


*Now, what about flat arrays ?*

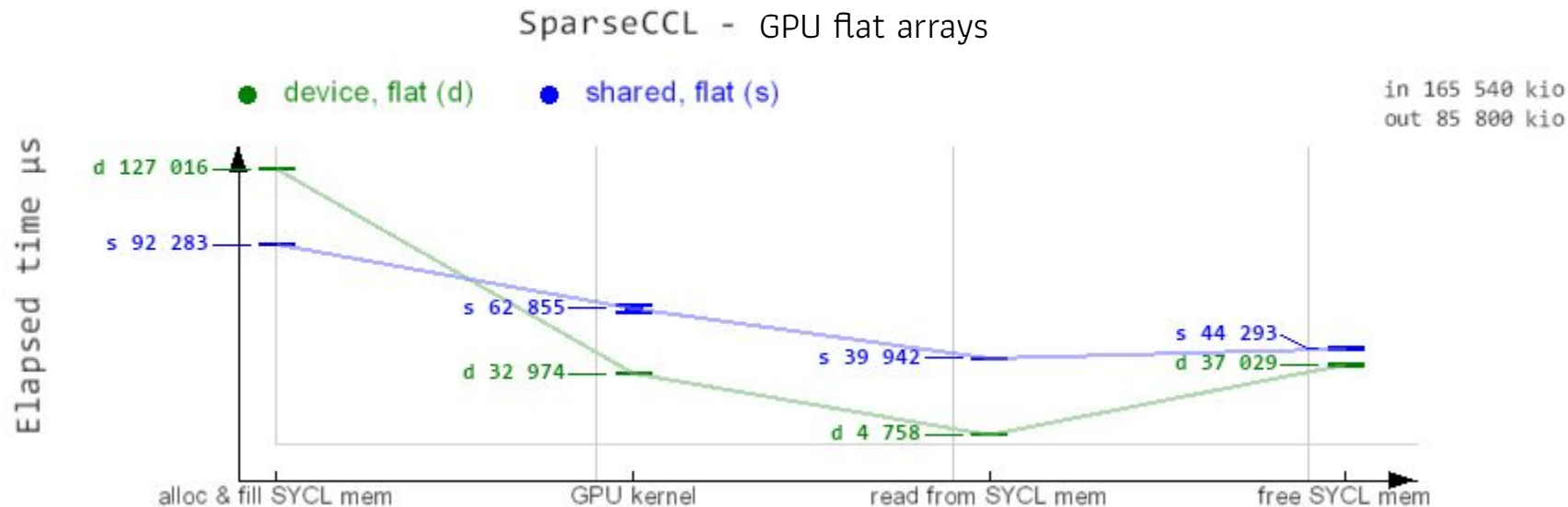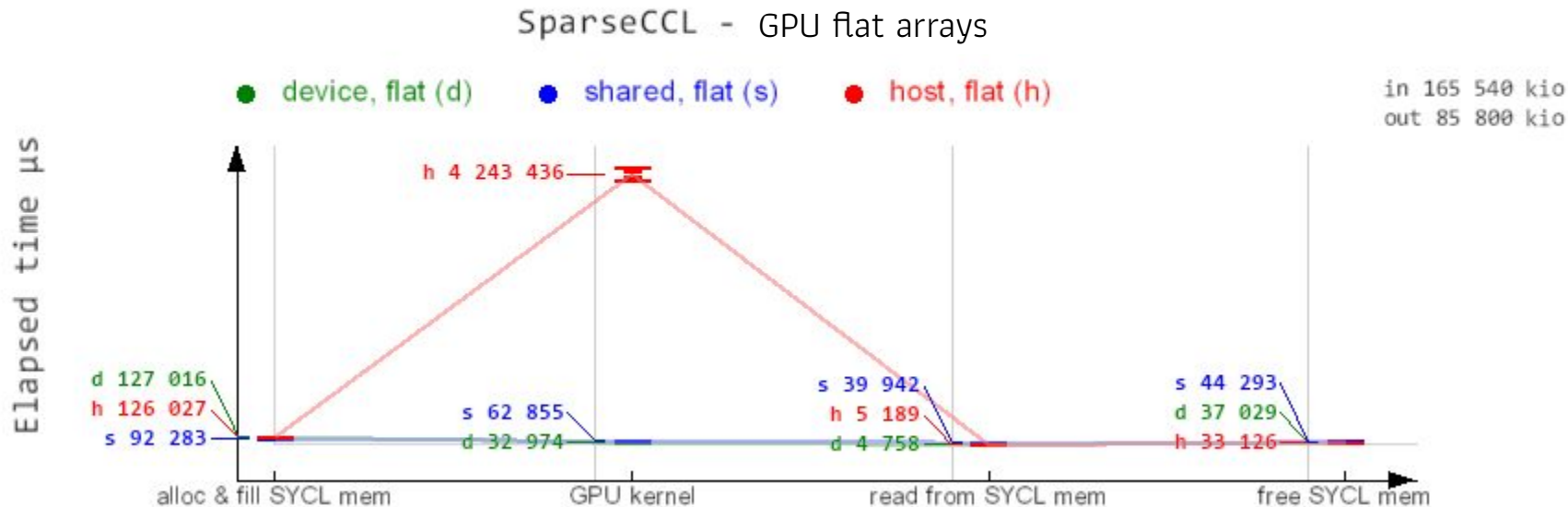# SparseCCL on GPU - flat arrays

USM device



SparseCCL - GPU flat arrays

device, flat (d)

in 165 540 kio
out 85 800 kio

Elapsed time µs

127 016

32 974

4 758

37 029

alloc & fill SYCL mem    GPU kernel    read from SYCL mem    free SYCL mem

# SparseCCL on GPU - flat arrays

USM device vs shared



SparseCCL - GPU flat arrays

device, flat (d)    shared, flat (s)

in 165 540 kio
out 85 800 kio

Elapsed time μs

d 127 016
s 92 283
s 62 855
s 44 293
s 39 942
d 37 029
d 32 974
d 4 758

alloc & fill SYCL mem    GPU kernel    read from SYCL mem    free SYCL mem

# SparseCCL on GPU - flat arrays

Multiple access to the same data => USM host kernel **very** expensive



SparseCCL - GPU flat arrays

- device, flat (d)     - shared, flat (s)     - host, flat (h)

in 165 540 kio
out 85 800 kio

Elapsed time µs

h 4 243 436

d 127 016
h 126 027
s 92 283

s 62 855
d 32 974

s 39 942
h 5 189
d 4 758

s 44 293
d 37 029
h 33 126

alloc & fill SYCL mem        GPU kernel        read from SYCL mem        free SYCL mem

# Findings : flat arrays

USM host only appropriate for single data access

USM shared is an option

USM device still faster

# General conclusions

Host & shared : Easy to reuse existing data structures but

*most likely* expensive.

USM host can *most likely* be used as a usual CPU-only memory.

USM device most efficient but needs flat data.

*Limitations :*

- *Only tried on SparseCCL and a microbenchmark.*
- *Only ran on a server and other smaller NVidia devices.*

# Future work

Prefetch USM shared

Buffers and accessors

Evaluate USMs models for seeding

More devices : NVidia, AMD, Intel…

A deeper understanding of SYCL

*Submitted abstract to ACAT'21[1]*

[1]ACAT'21 : https://indico.cern.ch/event/855454/

# Thank you ! Questions ?

## Links

My [WIP] github : https://github.com/SylvainJoube/SYCL_tests

hipSyCL : https://github.com/illuhad/hipSYCL

SYCL : https://www.khronos.org/sycl

Data Parallel C++ : https://software.intel.com/content/www/us/en/
                          develop/tools/oneapi/components/dpc-compiler.html