Migrating CUDA to oneAPI: A Smith-Waterman Case Study

Manuel Costanzo Enzo Rucci Carlos García-Sánchez Marcelo Naiouf Manuel Prieto-Matías









Agenda

- Motivation and Goal
- oneAPI
- Smith-Waterman
- Experimental Work
- Experimental Results
- Conclusions

Agenda

- Motivation and Goal
- oneAPI
- Smith-Waterman
- Experimental Work
- Experimental Results
- Conclusions

Heterogeneous computing and massively parallel architectures have proven to be an effective strategy for maximizing the performance and energy efficiency of computing system.



Heterogeneous computing and massively parallel architectures have proven to be an effective strategy for maximizing the performance and energy efficiency of computing system.



Because of that, programmers typically rely on a variety of hardware, such as CPUs, GPUs, FPGAs, and other kinds of accelerator

Heterogeneous computing and massively parallel architectures have proven to be an effective strategy for maximizing the performance and energy efficiency of computing system.



Because of that, programmers typically rely on a variety of hardware, such as CPUs, GPUs, FPGAs, and other kinds of accelerator



Heterogeneous computing and massively parallel architectures have proven to be an effective strategy for maximizing the performance and energy efficiency of computing system.







Heterogeneous computing and massively parallel architectures have proven to be an effective strategy for maximizing the performance and energy efficiency of computing system.

Because of that, programmers typically rely on a variety of hardware, such as CPUs, GPUs, FPGAs, and other kinds of accelerator



Heterogeneous computing and massively parallel architectures have proven to be an effective strategy for maximizing the performance and energy efficiency of computing system.

Because of that, programmers typically rely on a variety of hardware, such as CPUs, GPUs, FPGAs, and other kinds of accelerator



Motivation - oneAPI

Intel recently introduced the oneAPI programming ecosystem, which provides a unified programming model for a wide range of hardware architectures.



Motivation - oneAPI

Intel recently introduced the oneAPI programming ecosystem, which provides a unified programming model for a wide range of hardware architectures.

Data Parallel C++ (DPC++)

- C++
- SYCL



Motivation - oneAPI

Intel recently introduced the oneAPI programming ecosystem, which provides a unified programming model for a wide re

Same source code on different

- Data Parchitectures
 - SYCL



GPUs can be considered the dominant accelerator, while CUDA is the most popular programming language for them nowadays

GPUs can be considered the dominant accelerator, while CUDA is the most popular programming language for them nowadays



GPUs can be considered the dominant accelerator, while CUDA is the most popular programming language for them nowadays

GPUs can be considered the dominant accelerator, while CUDA is the most popular programming language for them nowadays



GPUs can be considered the dominant accelerator, while CUDA is the most popular programming language for them nowadays

GPUs can be considered the dominant accelerator, while CUDA is the most popular programming language for them nowadays

Bioinformatics and Computational Biology have been exploiting GPUs for more than two decades.

Molecular docking



GPUs can be considered the dominant accelerator, while CUDA is the most popular programming language for them nowadays

GPUs can be considered the dominant accelerator, while CUDA is the most popular programming language for them nowadays



GPUs can be considered the dominant accelerator, while CUDA is the most popular programming language for them nowadays

GPUs can be considered the dominant accelerator, while CUDA is the most popular programming language for them nowadays

Bioinformatics and Computational Biology have been exploiting GPUs for more than two decades.

Sequence alignment



GPUs can be considered the dominant accelerator, while CUDA is the most popular programming language for them nowadays

Bioinformatics and Computational Biology have been exploiting GPUs for more than two One of the most widely used implementations

Sequence alignment



GPUs can be considered the dominant accelerator, while CUDA is the most popular programming language for them nowadays

GPUs can be considered the dominant accelerator. while CUDA is nowaday Difficulty in porting the code to other architectures. Bioinform exploiting GPUs for more man two decades.

oneAPI facilitates the migration to the SYCL-based DPC++ programming language through the DPCT tool.

Intel® DPC++ Compatibility Tool Usage Flow



oneAPI facilitates the migration to the SYCL-based DPC++ programming language through the DPCT tool.

oneAPI facilitates the migration to the SYCL-based DPC++ programming language through the DPCT tool.

A few preliminary studies assessing the usefulness of dpct can be found

oneAPI facilitates the migration to the SYCL-based DPC++ programming language through the DPCT tool.

A few preliminary studies assessing the usefulness of dpct can be found

Simulation



oneAPI facilitates the migration to the SYCL-based DPC++ programming language through the DPCT tool.

A few preliminary studies assessing the usefulness of dpct can be found

oneAPI facilitates the migration to the SYCL-based DPC++ programming language through the DPCT tool.

A few preliminary studies assessing the usefulness of dpct can be found

Math

$$\frac{\partial}{\partial a} \ln f_{a,\sigma^2}(\xi_1) = \frac{(\xi_1 - a)}{\sigma^2} f_{a,\sigma^2}(\xi_1) = \frac{1}{\sqrt{2\pi\sigma}} \int_{\sigma^2} f_{a,\sigma^2}(\xi_1) = \frac{1}$$

oneAPI facilitates the migration to the SYCL-based DPC++ programming language through the DPCT tool.

A few preliminary studies assessing the usefulness of dpct can be found

oneAPI facilitates the migration to the SYCL-based DPC++ programming language through the DPCT tool.

A few preliminary studies assessing the usefulness of dpct can be found

Cryptography



oneAPI facilitates the migration to the SYCL-based DPC++ programming language through the DPCT tool.

A few preliminary studies assessing the usefulness of dpct can be found

oneAPI facilitates the migration to the SYCL-based DPC++ programming language through the DPCT tool.

A few preliminary studies assessing the usefulness of dpct can be found



Goal

To present our experiences porting a biological software tool to DPC++ using dpct

SW#db: a CUDA-based, memory efficient implementation of the Smith-Waterman algorithm
Goal

To present our experiences porting a biological software tool to DPC++ using dpct

SW#db: a CUDA-based, memory efficient implementation of the Smith-Waterman algorithm

- → An analysis of the dpct efectiveness for CUDA code migration
- → An analysis of the DPC++ code's portability, considering different target platforms and vendors (CPU and GPUs).
- → A comparison of the performance on different hardware architectures (CPU and GPUs).

Agenda

Motivation and Goal

oneAPI

- Smith-Waterman
- Experimental Work
- Experimental Results
- Conclusions

Technologies - oneAPI

oneAPI is a programming ecosystem, which provides a unified programming model for a wide range of hardware architectures



SYCL allows the programmer to write host code in C++ and make it compatible to run on different architectures.

SYCL allows the programmer to write host code in C++ and make it compatible to run on different architectures.



SYCL allows the programmer to write host code in C++ and make it compatible to run on different architectures.

DPC++ combines the SYCL language with modified C++.



SYCL allows the programmer to write host code in C++ and make it compatible to run on different architectures.

DPC++ combines the SYCL language with modified C++.





Agenda

- Motivation and Goal
- oneAPI
- Smith-Waterman
- Experimental Work
- Experimental Results
- Conclusions

Obtains the optimal local alignment between two biological sequences.



Obtains the optimal local alignment between two biological sequences.

It has been used as the basis for many subsequent algorithms



Obtains the optimal local alignment between two biological sequences.

It has been used as the basis for many subsequent algorithms

Is often employed as a benchmark when comparing different alignment techniques



Obtai

sequ

It has

jical

SW# is a CUDA tool for computing biological sequence alignments that works with both protein and DNA sequences.

orithms

Is often employed as a benchmark when comparing different alignment techniques



Agenda

- Motivation and Goal
- oneAPI
- Smith-Waterman
- Experimental Work
- Experimental Results
- Conclusions



The migration process consists of 4 stages:

1. Use the DPCT tool for generating a first version (Stage 1).



- 1. Use the DPCT tool for generating a first version (Stage 1).
- 2. Based on the DPCT alerts, modify the resulting code (Stage 2).



- 1. Use the DPCT tool for generating a first version (Stage 1).
- 2. Based on the DPCT alerts, modify the resulting code (Stage 2).
- **3**. Based on runtime errors, modify the code (Stage 3).



- 1. Use the DPCT tool for generating a first version (Stage 1).
- 2. Based on the DPCT alerts, modify the resulting code (Stage 2).
- **3**. Based on runtime errors, modify the code (Stage 3).
- 4. Verify that the results are correct (Stage 4).



Generate first version of DPC++ code using DPCT tool

Generate first version of DPC++ code using DPCT tool



Generate first version of DPC++ code using DPCT tool

1. intercept-build make



Generate first version of DPC++ code using DPCT tool

1. intercept-build make

2. dpct -p compile_commands.json -in-root=\$PROJ_DIR



DPCT1003: Migrated API does not return error code. (*, 0) is inserted. You may need to rewrite this code

DPCT1003: Migrated API does not return error code. (*, 0) is inserted. You may need to rewrite this code

```
size_t valuesSize =
1
        databaseLen * sizeof(double);
2
    double* valuesGpu;
3
4
    CUDA_SAFE_CALL(
5
        cudaMalloc(
6
            &valuesGpu, valuesSize
7
8
    ):
9
```

```
size_t valuesSize =
databaseLen * sizeof(double);
double* valuesGpu;
CUDA_SAFE_CALL((
 valuesGpu = (double *)sycl::malloc_device(
 valuesSize, dpct::get_default_queue()),
 0));
```

```
(b) oneAPI
```

```
(a) CUDA
```

(a) CUDA

DPCT1003: Migrated API does not return error code. (*, 0) is inserted. You may need to rewrite this code

DPCT1009: SYCL uses exceptions to report errors and does not use the error codes. The original code was commented out and a warning string was inserted. You need to rewrite this code.



DPCT1003: Migrated API does not return error code. (*, 0) is inserted. You may need to rewrite this code

```
size_t valuesSize =
1
        databaseLen * sizeof(double);
2
    double* valuesGpu;
3
4
    CUDA_SAFE_CALL(
5
        cudaMalloc(
6
            &valuesGpu, valuesSize
7
8
    ):
9
```

```
size_t valuesSize =
databaseLen * sizeof(double);
double* valuesGpu;
CUDA_SAFE_CALL((
 valuesGpu = (double *)sycl::malloc_device(
 valuesSize, dpct::get_default_queue()),
 0));
```

```
(b) oneAPI
```

```
(a) CUDA
```

DPCT1003: Migrated API does not return error code. (*, 0) is inserted. You may need to rewrite this code

DPCT1009: SYCL uses exceptions to report errors and does not use the error codes. The original code was commented out and a warning string was inserted. You need to rewrite this code.

```
size_t valuesSize =
1
                                                 1
        databaseLen * sizeof(double);
2
                                                 2
    double* valuesGpu;
3
                                                 3
4
                                                 4
    CUDA_SAFE_CALL(
5
                                                 5
        cudaMalloc(
6
                                                 6
             &valuesGpu, valuesSize
7
                                                 7
8
                                                      ());
                                                 8
    );
9
```

(a) CUDA

```
size_t valuesSize =
    databaseLen * sizeof(double);
double* valuesGpu;
CUDA_SAFE_CALL((
    valuesGpu = (double * sycl::malloc_device(
        valuesSize, dpct::get_default_queue()),
0));
```

```
(b) oneAPI
```

DPCT1003: Migrated API does not return error code. (*, 0) is inserted. You may need to rewrite this code

```
size_t valuesSize =
1
        databaseLen * sizeof(double);
2
    double* valuesGpu;
3
4
    CUDA_SAFE_CALL(
5
        cudaMalloc(
6
            &valuesGpu, valuesSize
7
8
    ):
9
```

```
size_t valuesSize =
databaseLen * sizeof(double);
double* valuesGpu;
CUDA_SAFE_CALL((
 valuesGpu = (double *)sycl::malloc_device(
 valuesSize, dpct::get_default_queue()),
 0));
```

```
(b) oneAPI
```

```
(a) CUDA
```

DPCT1003: Migrated API does not return error code. (*, 0) is inserted. You may need to rewrite this code



DPCT1003: Migrated API does not return error code. (*, 0) is inserted. You may need to rewrite this code

```
size_t valuesSize =
1
        databaseLen * sizeof(double);
2
    double* valuesGpu;
3
4
    CUDA_SAFE_CALL(
5
        cudaMalloc(
6
            &valuesGpu, valuesSize
7
8
    ):
9
```

```
size_t valuesSize =
databaseLen * sizeof(double);
double* valuesGpu;
CUDA_SAFE_CALL((
 valuesGpu = (double *)sycl::malloc_device(
 valuesSize, dpct::get_default_queue()),
 0));
```

```
(b) oneAPI
```

```
(a) CUDA
```

DPCT1005: The SYCL device version is different from CUDA Compute Compatibility. You may need to rewrite this code.

DPCT1005: The SYCL device version is different from CUDA Compute Compatibility. You may need to rewrite this code.

```
1 cudaDeviceProp properties;
2 cudaGetDeviceProperties(
3    &properties, card
4 );
5 
6 bool major = properties.major < 2;
7 int threads = major ? 64 : 128;
8 int blocks = major ? 360 : 480;
```

```
dpct::device_info properties;
1
2
    dpct::dev_mgr::instance()
3
        .get_device(card)
4
        .get_device_info(properties);
5
6
    bool major = false;
7
    int threads = major ? 64 : 128;
8
    int blocks = major ? 360 : 480;
9
```

```
(a) CUDA
```

DPCT1005: The SYCL device version is different from CUDA Compute Compatibility. You may need to rewrite this code.

```
dpct::device_info properties;
1
2
    dpct::dev_mgr::instance()
3
        .get_device(card)
4
        .get_device_info(properties);
5
6
    bool major = false;
7
    int threads = major ? 64 : 128;
8
    int blocks = major ? 360 : 480;
9
```

(a) CUDA

(b) oneAPI

DPCT1005: The SYCL device version is different from CUDA Compute Compatibility. You may need to rewrite this code.

```
1 cudaDeviceProp properties;
2 cudaGetDeviceProperties(
3    &properties, card
4 );
5 
6 bool major = properties.major < 2;
7 int threads = major ? 64 : 128;
8 int blocks = major ? 360 : 480;
```

```
dpct::device_info properties;
1
2
    dpct::dev_mgr::instance()
3
        .get_device(card)
4
        .get_device_info(properties);
5
6
    bool major = false;
7
    int threads = major ? 64 : 128;
8
    int blocks = major ? 360 : 480;
9
```

```
(a) CUDA
```

DPCT1005: The SYCL device version is different from CUDA Compute Compatibility. You may need to rewrite this code.



(a) CUDA

1	<pre>dpct::device_info properties;</pre>
2	
3	dpct::dev_mgr::instance()
4	.get_device(card)
5	<pre>.get_device_info(properties);</pre>
6	
7	<pre>bool major = false;</pre>
8	<pre>int threads = major ? 64 : 128;</pre>
9	<pre>int blocks = major ? 360 : 480;</pre>

DPCT1005: The SYCL device version is different from CUDA Compute Compatibility. You may need to rewrite this code.

```
1 cudaDeviceProp properties;
2 cudaGetDeviceProperties(
3    &properties, card
4 );
5 
6 bool major = properties.major < 2;
7 int threads = major ? 64 : 128;
8 int blocks = major ? 360 : 480;
```

```
dpct::device_info properties;
1
2
    dpct::dev_mgr::instance()
3
        .get_device(card)
4
        .get_device_info(properties);
5
6
    bool major = false;
7
    int threads = major ? 64 : 128;
8
    int blocks = major ? 360 : 480;
9
```

```
(a) CUDA
```
DPCT1049: The workgroup size passed to the SYCL kernel may exceed the limit. To get the device limit, query info::device::max_work_group_size. Adjust the workgroup size if needed.

DPCT1049: The workgroup size passed to the SYCL kernel may exceed the limit. To get the device limit, query info::device::max_work_group_size. Adjust the workgroup size if needed.

solveShort<<<blocks, threads>>>(...);

(a) CUDA

dpct::get_default_queue() 1 .submit([&](sycl::handler &cgh) { 2 3 . . . cgh.parallel_for(4 sycl::nd_range<3> 5 (sycl::range<3>(1, 1, blocks) * 6 sycl::range<3>(1, 1, threads), 7 sycl::range<3>(1, 1, threads)), 8 [=](sycl::nd_item<3> item_ct1) { 9 solveShort(...); 10 }); 11 }); 12

DPCT1049: The workgroup size passed to the SYCL kernel may exceed the limit. To get the device limit, query info::device::max_work_group_size. Adjust the workgroup size if needed.

1 solveShort<<<blocks, threads>>>(...);

(a) CUDA

dpct::get_default_queue() 1 .submit([&](sycl::handler &cgh) { 2 3 cgh.parallel_for(4 sycl::nd_range<3> 5 (sycl::range<3>(1, 1, blocks) * 6 sycl::range<3>(1, 1, threads), 7 sycl::range<3>(1, 1, threads)), 8 [=](sycl::nd_item<3> item_ct1) { 9 solveShort(...); 10 }); 11 12 });

DPCT1049: The workgroup size passed to the SYCL kernel may exceed the limit. To get the device limit, query info::device::max_work_group_size. Adjust the workgroup size if needed.

solveShort<<<blocks, threads>>>(...);

(a) CUDA

dpct::get_default_queue() 1 .submit([&](sycl::handler &cgh) { 2 3 . . . cgh.parallel_for(4 sycl::nd_range<3> 5 (sycl::range<3>(1, 1, blocks) * 6 sycl::range<3>(1, 1, threads), 7 sycl::range<3>(1, 1, threads)), 8 [=](sycl::nd_item<3> item_ct1) { 9 solveShort(...); 10 }); 11 }); 12

1		1	
2	<pre>syncthreads();</pre>	2	item_ct1.barrier();
3		3	
	(a) CUDA		(b) oneAPI



1		1	
2	<pre>syncthreads();</pre>	2	item_ct1.barrier();
3		3	
	(a) CUDA		(b) oneAPI

1

2

3

. . .

. . .



DPCT1059: SYCL only supports 4-channel image format. Adjust the code.

1

2

8

9

10

11

12

13

14

15

DPCT1059: SYCL only supports 4-channel image format. Adjust the code.

```
texture<char> colTexture;
1
2
     int colSize = colsGpu * sizeof(char);
3
     char *colGpu;
4
     cudaMalloc(&colGpu, colSize);
5
     cudaMemcpy(colGpu, colCpu,
6
         colSize, TO_GPU);
7
     cudaBindTexture(NULL, colTexture,
8
         colGpu, colSize);
9
10
     char v = tex1Dfetch(colTexture, 10);
11
                 (a) CUDA
```

dpct::image_wrapper<sycl::char4, 1> colTexture;

//dpct::image_wrapper<char, 1> colTexture;

```
dpct::get_default_queue()
   .memcpy(colGpu, colCpu, colSize).wait();
```

```
colTexture.attach(colGpu, colSize);
```

```
// DIV 4 y MOD 3
char v = colTexture.read(10 >> 2)[10 & 3];
```

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

DPCT1059: SYCL only supports 4-channel image format. Adjust the code.

```
texture<char> colTexture;
1
2
     int colSize = colsGpu * sizeof(char);
3
     char *colGpu;
4
     cudaMalloc(&colGpu, colSize);
5
     cudaMemcpy(colGpu, colCpu,
6
         colSize, TO_GPU);
7
     cudaBindTexture(NULL, colTexture,
8
         colGpu, colSize);
9
10
     char v = tex1Dfetch(colTexture, 10);
11
                 (a) CUDA
```

```
//dpct::image_wrapper<char, 1> colTexture;
dpct::image_wrapper<sycl::char4, 1> colTexture;
int colSize = colsGpu * sizeof(char);
char* colGpu;
colGpu = (char *)sycl::malloc_device(
    colSize, dpct::get_default_queue());
dpct::get_default_queue()
    .memcpy(colGpu, colCpu, colSize).wait();
colTexture.attach(colGpu, colSize);
// DIV 4 y MOD 3
char v = colTexture.read(10 >> 2)[10 & 3];
```

1

2

8

9

10

11

12

13

14

15

DPCT1059: SYCL only supports 4-channel image format. Adjust the code.

```
texture<char> colTexture;
1
2
     int colSize = colsGpu * sizeof(char);
3
     char *colGpu;
4
     cudaMalloc(&colGpu, colSize);
5
     cudaMemcpy(colGpu, colCpu,
6
         colSize, TO_GPU);
7
     cudaBindTexture(NULL, colTexture,
8
         colGpu, colSize);
9
10
     char v = tex1Dfetch(colTexture, 10);
11
                 (a) CUDA
```

dpct::image_wrapper<sycl::char4, 1> colTexture;

//dpct::image_wrapper<char, 1> colTexture;

```
dpct::get_default_queue()
   .memcpy(colGpu, colCpu, colSize).wait();
```

```
colTexture.attach(colGpu, colSize);
```

```
// DIV 4 y MOD 3
char v = colTexture.read(10 >> 2)[10 & 3];
```

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

DPCT1059: SYCL only supports 4-channel image format. Adjust the code.

```
texture<char> colTexture;
1
2
     int colSize = colsGpu * sizeof(char);
3
     char *colGpu;
4
     cudaMalloc(&colGpu, colSize);
5
     cudaMemcpy(colGpu, colCpu,
6
         colSize, TO_GPU);
7
     cudaBindTexture(NULL, colTexture,
8
         colGpu, colSize);
9
10
     char v = tex1Dfetch(colTexture, 10);
11
                 (a) CUDA
```

```
Int colsize = colsepu * sizeor(cnar);
char* colGpu;
colGpu = (char *)sycl::malloc_device(
    colSize, dpct::get_default_queue());
dpct::get_default_queue()
    .memcpy(colGpu, colCpu, colSize).wait();
colTexture.attach(colGpu, colSize);
// DIV 4 y MOD 3
```

dpct::image_wrapper<sycl::char4, 1> colTexture;

```
char v = colTexture.read(10 >> 2)[10 & 3];
```

1

2

8

9

10

11

12

13

14

15

DPCT1059: SYCL only supports 4-channel image format. Adjust the code.

```
texture<char> colTexture;
1
2
     int colSize = colsGpu * sizeof(char);
3
     char *colGpu;
4
     cudaMalloc(&colGpu, colSize);
5
     cudaMemcpy(colGpu, colCpu,
6
         colSize, TO_GPU);
7
     cudaBindTexture(NULL, colTexture,
8
         colGpu, colSize);
9
10
     char v = tex1Dfetch(colTexture, 10);
11
                 (a) CUDA
```

dpct::image_wrapper<sycl::char4, 1> colTexture;

//dpct::image_wrapper<char, 1> colTexture;

```
dpct::get_default_queue()
   .memcpy(colGpu, colCpu, colSize).wait();
```

```
colTexture.attach(colGpu, colSize);
```

```
// DIV 4 y MOD 3
char v = colTexture.read(10 >> 2)[10 & 3];
```

6

7

8

9

10

11

12

13

14

15

DPCT1059: SYCL only supports 4-channel image format. Adjust the code.

```
texture<char> colTexture;
1
2
    int colSize = colsGpu * sizeof(char);
3
    char *colGpu;
4
    cudaMalloc(&colGpu, colSize);
5
    cudaMemcpy(colGpu, colCpu,
6
         colSize, TO_GPU);
7
    cudaBindTexture(NULL, colTexture,
8
         colGpu, colSize);
10
    char v = tex1Dfetch(colTexture, 10);
11
                 (a) CUDA
```

```
1 //dpct::image_wrapper<char, 1> colTexture;
2 dpct::image_wrapper<sycl::char4, 1> colTexture;
3 int colSize = colsGpu * sizeof(char);
4 char* colGpu;
5
```

```
colGpu = (char *)sycl::malloc_device(
    colSize, dpct::get_default_queue());
```

```
dpct::get_default_queue()
    .memcpy(colGpu, colCpu, colSize).wait();
```

```
colTexture.attach(colGpu, colSize);
```

```
// DIV 4 y MOD 3
char v = colTexture.read(10 >> 2)[10 & 3];
```

1

2

8

9

10

11

12

13

14

15

DPCT1059: SYCL only supports 4-channel image format. Adjust the code.

```
texture<char> colTexture;
1
2
     int colSize = colsGpu * sizeof(char);
3
     char *colGpu;
4
     cudaMalloc(&colGpu, colSize);
5
     cudaMemcpy(colGpu, colCpu,
6
         colSize, TO_GPU);
7
     cudaBindTexture(NULL, colTexture,
8
         colGpu, colSize);
9
10
     char v = tex1Dfetch(colTexture, 10);
11
                 (a) CUDA
```

dpct::image_wrapper<sycl::char4, 1> colTexture;

//dpct::image_wrapper<char, 1> colTexture;

```
dpct::get_default_queue()
   .memcpy(colGpu, colCpu, colSize).wait();
```

```
colTexture.attach(colGpu, colSize);
```

```
// DIV 4 y MOD 3
char v = colTexture.read(10 >> 2)[10 & 3];
```

5

6

7

8

9

10

11

12

13

14

15

DPCT1059: SYCL only supports 4-channel image format. Adjust the code.

```
texture<char> colTexture;
1
2
    int colSize = colsGpu * sizeof(char);
3
    char *colGpu;
4
    cudaMalloc(&colGpu, colSize);
5
    cudaMemcpy(colGpu, colCpu,
6
         colSize, TO_GPU);
7
    cudaBindTexture(NULL, colTexture,
8
         colGpu, colSize);
9
10
    char v = tex1Dfetch(colTexture, 10);
11
                 (a) CUDA
```

```
//dpct::image_wrapper<char, 1> colTexture;
1
    dpct::image_wrapper<sycl::char4, 1> colTexture;
2
    int colSize = colsGpu * sizeof(char);
3
    char* colGpu;
4
```

```
colGpu = (char *)sycl::malloc_device(
    colSize, dpct::get_default_queue());
```

```
dpct::get_default_queue()
    .memcpy(colGpu, colCpu, colSize).wait();
```

```
colTexture.attach(colGpu, colSize);
```

```
// DIV 4 y MOD 3
char v = colTexture.read(10 >> 2)[10 \& 3];
```

1

2

8

9

10

11

12

13

14

15

DPCT1059: SYCL only supports 4-channel image format. Adjust the code.

```
texture<char> colTexture;
1
2
     int colSize = colsGpu * sizeof(char);
3
     char *colGpu;
4
     cudaMalloc(&colGpu, colSize);
5
     cudaMemcpy(colGpu, colCpu,
6
         colSize, TO_GPU);
7
     cudaBindTexture(NULL, colTexture,
8
         colGpu, colSize);
9
10
     char v = tex1Dfetch(colTexture, 10);
11
                 (a) CUDA
```

dpct::image_wrapper<sycl::char4, 1> colTexture;

//dpct::image_wrapper<char, 1> colTexture;

```
dpct::get_default_queue()
   .memcpy(colGpu, colCpu, colSize).wait();
```

```
colTexture.attach(colGpu, colSize);
```

```
// DIV 4 y MOD 3
char v = colTexture.read(10 >> 2)[10 & 3];
```

For a 1D/2D image/image array, the width must be a Value >= 1 and <= CL_DEVICE_IMAGE2D_MAX_WIDTH.

For a 1D/2D image/image array, the width must be a Value >= 1 and <= CL_DEVICE_IMAGE2D_MAX_WIDTH.

```
texture<int, 2,
1
       cudaReadModeElementType> seqsTexture;
                                                    static int *seqsGpu;
 2
                                              1
 3
                                                2
     cudaArray *sequencesGpu;
                                                    seqsGpu = (int *)sycl::malloc_device(
 4
                                                3
     cudaChannelFormatDesc channel =
                                                     sequencesCols * sequencesRows * sizeof(int),
5
                                                4
       seqsTexture.channelDesc;
                                                     dpct::get_default_queue());
 6
                                                5
     cudaMallocArray(&sequencesGpu,
 7
                                                6
         &channel, sequencesCols, sequencesRows;
8
                                                    dpct::get_default_queue()
     cudaMemcpyToArray(sequencesGpu, 0, 0,
                                                     .memcpy(seqsGpu, sequences, sequencesCols *
9
                                                8
         sequences, sequencesSize, TO_GPU);
                                                            sequencesRows * sizeof(int))
10
                                                9
     cudaBindTextureToArray(
                                                      .wait():
11
                                               10
         seqsTexture, sequencesGpu);
12
```

(a) CUDA

1int columnCodes = tex2D(1int columnCodes =2seqsTexture, colOff, j + rowOff);2seqsGpu[(j + rowOff) * sequencesCols + colOff];

```
(a) CUDA
```

For a 1D/2D image/image array, the width must be a Value >= 1 and <= CL_DEVICE_IMAGE2D_MAX_WIDTH.

1 2 3	<pre>texture<int, 2,<br="">cudaReadModeElementType> seqsTexture;</int,></pre>	1	<pre>static int *seqsGpu;</pre>
4	cudaArray *sequencesGpu;	3	<pre>seqsGpu = (int *)sycl::malloc_device(</pre>
5	cudaChannelFormatDesc channel =	4	<pre>sequencesCols * sequencesRows * sizeof(int),</pre>
6	<pre>seqsTexture.channelDesc;</pre>	5	<pre>dpct::get_default_queue());</pre>
7	cudaMallocArray(&sequencesGpu,	6	
8	&channel, sequencesCols, sequencesR	ows);	dpct::get_default_queue()
9	<pre>cudaMemcpyToArray(sequencesGpu, 0, 0,</pre>	8	<pre>.memcpy(seqsGpu, sequences, sequencesCols *</pre>
10	<pre>sequences, sequencesSize, TO_GPU);</pre>	9	<pre>sequencesRows * sizeof(int))</pre>
11	cudaBindTextureToArray(10	.wait();
12	<pre>seqsTexture, sequencesGpu);</pre>		

(a) CUDA

1int columnCodes = tex2D(1int columnCodes =2seqsTexture, colOff, j + rowOff);2seqsGpu[(j + rowOff) * sequencesCols + colOff];

(a) CUDA

For a 1D/2D image/image array, the width must be a Value >= 1 and <= CL_DEVICE_IMAGE2D_MAX_WIDTH.

```
texture<int, 2,
1
       cudaReadModeElementType> seqsTexture;
                                                    static int *seqsGpu;
 2
                                              1
 3
                                                2
     cudaArray *sequencesGpu;
                                                    seqsGpu = (int *)sycl::malloc_device(
 4
                                                3
     cudaChannelFormatDesc channel =
                                                     sequencesCols * sequencesRows * sizeof(int),
5
                                                4
       seqsTexture.channelDesc;
                                                     dpct::get_default_queue());
 6
                                                5
     cudaMallocArray(&sequencesGpu,
 7
                                                6
         &channel, sequencesCols, sequencesRows;
8
                                                    dpct::get_default_queue()
     cudaMemcpyToArray(sequencesGpu, 0, 0,
                                                     .memcpy(seqsGpu, sequences, sequencesCols *
9
                                                8
         sequences, sequencesSize, TO_GPU);
                                                            sequencesRows * sizeof(int))
10
                                                9
     cudaBindTextureToArray(
                                                      .wait():
11
                                               10
         seqsTexture, sequencesGpu);
12
```

(a) CUDA

1int columnCodes = tex2D(1int columnCodes =2seqsTexture, colOff, j + rowOff);2seqsGpu[(j + rowOff) * sequencesCols + colOff];

```
(a) CUDA
```

For a 1D/2D image/image array, the width must be a Value >= 1 and <= CL_DEVICE_IMAGE2D_MAX_WIDTH.

<pre>texture<int, 2,<br="">cudaReadModeElementType> seqsTexture;</int,></pre>	<pre>1 static int *seqsGpu;</pre>
<pre>cudaArray *sequencesGpu; cudaChannelFormatDesc channel = seqsTexture.channelDesc; cudaMallocArray(&sequencesGpu,</pre>	<pre>2 3 seqsGpu = (int *)sycl::malloc_device(4 sequencesCols * sequencesRows * sizeof(int), 5 dpct::get_default_queue()); 6</pre>
<pre>&channel, sequencesCols, sequencesRow cudaMemcpyToArray(sequencesGpu, 0, 0, sequences, sequencesSize, TO_GPU); cudaBindTextureToArray(seqsTexture, sequencesGpu);</pre>	<pre>sp); dpct::get_default_queue() 8 .memcpy(seqsGpu, sequences, sequencesCols * 9</pre>
	(b) oneAPI

(a) CUDA

```
1int columnCodes = tex2D(1int columnCodes =2seqsTexture, colOff, j + rowOff);2seqsGpu[(j + rowOff) * sequencesCols + colOff];
```

For a 1D/2D image/image array, the width must be a Value >= 1 and <= CL_DEVICE_IMAGE2D_MAX_WIDTH.

```
texture<int, 2,
1
       cudaReadModeElementType> seqsTexture;
                                                    static int *seqsGpu;
 2
                                              1
 3
                                                2
     cudaArray *sequencesGpu;
                                                    seqsGpu = (int *)sycl::malloc_device(
 4
                                                3
     cudaChannelFormatDesc channel =
                                                     sequencesCols * sequencesRows * sizeof(int),
5
                                                4
       seqsTexture.channelDesc;
                                                     dpct::get_default_queue());
 6
                                                5
     cudaMallocArray(&sequencesGpu,
 7
                                                6
         &channel, sequencesCols, sequencesRows;
8
                                                    dpct::get_default_queue()
     cudaMemcpyToArray(sequencesGpu, 0, 0,
                                                     .memcpy(seqsGpu, sequences, sequencesCols *
9
                                                8
         sequences, sequencesSize, TO_GPU);
                                                            sequencesRows * sizeof(int))
10
                                                9
     cudaBindTextureToArray(
                                                      .wait():
11
                                               10
         seqsTexture, sequencesGpu);
12
```

(a) CUDA

1int columnCodes = tex2D(1int columnCodes =2seqsTexture, colOff, j + rowOff);2seqsGpu[(j + rowOff) * sequencesCols + colOff];

```
(a) CUDA
```

For a 1D/2D image/image array, the width must be a Value >= 1 and <= CL_DEVICE_IMAGE2D_MAX_WIDTH.

1	texture <int, 2,<="" th=""><th></th><th></th></int,>		
2	<pre>cudaReadModeElementType> seqsTexture;</pre>	1	<pre>static int *seqsGpu;</pre>
3		2	
4	cudaArray *sequencesGpu;	3	<pre>seqsGpu = (int *)sycl::malloc_device(</pre>
5	cudaChannelFormatDesc channel =	4	<pre>sequencesCols * sequencesRows * sizeof(int),</pre>
6	<pre>seqsTexture.channelDesc;</pre>	5	<pre>dpct::get_default_queue());</pre>
7	cudaMallocArray(&sequencesGpu,	6	
8	&channel, sequencesCols, sequencesRe	; (rawc	dpct::get_default_queue()
9	<pre>cudaMemcpyToArray(sequencesGpu, 0, 0,</pre>	8	<pre>.memcpy(seqsGpu, sequences, sequencesCols *</pre>
10	<pre>sequences, sequencesSize, TO_GPU);</pre>	9	<pre>sequencesRows * sizeof(int))</pre>
11	cudaBindTextureToArray(10	.wait();
12	<pre>seqsTexture, sequencesGpu);</pre>		

(b) oneAPI

(a) CUDA



For a 1D/2D image/image array, the width must be a Value >= 1 and <= CL_DEVICE_IMAGE2D_MAX_WIDTH.

```
texture<int, 2,
1
       cudaReadModeElementType> seqsTexture;
                                                    static int *seqsGpu;
 2
                                              1
 3
                                                2
     cudaArray *sequencesGpu;
                                                    seqsGpu = (int *)sycl::malloc_device(
 4
                                                3
     cudaChannelFormatDesc channel =
                                                     sequencesCols * sequencesRows * sizeof(int),
5
                                                4
       seqsTexture.channelDesc;
                                                     dpct::get_default_queue());
 6
                                                5
     cudaMallocArray(&sequencesGpu,
 7
                                                6
         &channel, sequencesCols, sequencesRows;
8
                                                    dpct::get_default_queue()
     cudaMemcpyToArray(sequencesGpu, 0, 0,
                                                     .memcpy(seqsGpu, sequences, sequencesCols *
9
                                                8
         sequences, sequencesSize, TO_GPU);
                                                            sequencesRows * sizeof(int))
10
                                                9
     cudaBindTextureToArray(
                                                      .wait():
11
                                               10
         seqsTexture, sequencesGpu);
12
```

(a) CUDA

1int columnCodes = tex2D(1int columnCodes =2seqsTexture, colOff, j + rowOff);2seqsGpu[(j + rowOff) * sequencesCols + colOff];

```
(a) CUDA
```



Different tests were run to prove that the DPC++ code generates correct results

Agenda

- Motivation and Goal
- oneAPI
- Smith-Waterman
- Experimental Work
- Experimental Results
- Conclusions

Experimental Results - Hardware Platforms

CPU			GPU			
ID	Processor	RAM (Memory)	ID	Vedor (Type)	Model (Architecture)	GFLOPS peak (SP)
Core-i5	Intel Core i5- 7400	16 GB	Titan	NVIDIA (Discrete)	Titan X (Pascal)	10970
Core-i3	Intel Core i3- 4160	8 GB	RTX	NVIDIA (Discrete)	RTX 2070 (Turing)	7465
Core-i9	Intel Core i9- 10920X	32 GB	Iris XE	Intel (Discrete)	Iris Xe MAX Graphics (Gen 12.1)	2534
Xeon	Intel Xeon E- 2176G	65 GB	P630	NVIDIA (Integrated)	UHD Graphics P630 (Gen 9.5)	441.6

Experimental Results - Design

 UniProtKB/Swiss-Prot database (release 2021_04) contains 204173280 amino acid residues in 565928 sequences with a maximum length of 35213

Experimental Results - Design

- UniProtKB/Swiss-Prot database (release 2021_04) contains 204173280 amino acid residues in 565928 sequences with a maximum length of 35213
- 20 queries with length from 144 to 5478 (accession numbers: P02232, P05013, P14942, P07327, P01008, P03435, P42357, P21177, Q38941, P27895, P07756, P04775, P19096, P28167, P0C6B8, P20930, P08519, Q7TMA5, P33450, and Q9UKN1)

Experimental Results - Design

- UniProtKB/Swiss-Prot database (release 2021_04) contains 204173280 amino acid residues in 565928 sequences with a maximum length of 35213
- 20 queries with length from 144 to 5478 (accession numbers: P02232, P05013, P14942, P07327, P01008, P03435, P42357, P21177, Q38941, P27895, P07756, P04775, P19096, P28167, P0C6B8, P20930, P08519, Q7TMA5, P33450, and Q9UKN1)
- Scoring matrix: BLOSUM62 Gap insertion: 10 Gap extension: 2

Experimental Results - Configuration & tests

1. Single thread at the CPU level (using flag T=1).

Experimental Results - Configuration & tests

- 1. Single thread at the CPU level (using flag T=1).
- 2. Different work group sizes were configured for kernel execution.

Experimental Results - Configuration & tests

- 1. Single thread at the CPU level (using flag T=1).
- 2. Different work group sizes were configured for kernel execution.
- 3. Each test was run twenty times and the performance was calculated as the average.
Experimental Results - Configuration & tests

- 1. Single thread at the CPU level (using flag T=1).
- 2. Different work group sizes were configured for kernel execution.
- 3. Each test was run twenty times and the performance was calculated as the average.
- 4. GCUPS as the performance metric.



Q → Database length
D → Sequence length
t → Execution time in seconds































Experimental Results - Cross-GPU vendor portability



Experimental Results - Cross-GPU vendor portability



Experimental Results - Cross-GPU vendor portability













Agenda

- Motivation and Goal
- oneAPI
- Smith-Waterman
- Experimental Work
- Experimental Results
- Conclusions

This paper presents our experiences migrating a CUDA-based, biological software tool to DPC++ using the oneAPI framework.

The main contributions are:

This paper presents our experiences migrating a CUDA-based, biological software tool to DPC++ using the oneAPI framework.

The main contributions are:

1. DPCT proved to be an effective tool for code migration to DPC++

This paper presents our experiences migrating a CUDA-based, biological software tool to DPC++ using the oneAPI framework.

The main contributions are:

- 1. DPCT proved to be an effective tool for code migration to DPC++
- 2. The migrated code could be successfully executed on CPUs and also GPUs from different vendors, demonstrating its cross-architecture, cross-vendor GPU portability

This paper presents our experiences migrating a CUDA-based, biological software tool to DPC++ using the oneAPI framework.

The main contributions are:

- 1. DPCT proved to be an effective tool for code migration to DPC++
- 2. The migrated code could be successfully executed on CPUs and also GPUs from different vendors, demonstrating its cross-architecture, cross-vendor GPU portability
- 3. The performance results showed that the migrated DPC++ code is comparable to the original CUDA one.

Conclusions - Future Work

 Understanding the gap in performance between DPC++ and CUDA code, and optimizing DPC++ code to reach its maximum performance.

Conclusions - Future Work

 Understanding the gap in performance between DPC++ and CUDA code, and optimizing DPC++ code to reach its maximum performance.

2. Carrying out more exhaustive experimental work.

Conclusions - Future Work

 Understanding the gap in performance between DPC++ and CUDA code, and optimizing DPC++ code to reach its maximum performance.

2. Carrying out more exhaustive experimental work.

3. Running the DPC++ code on other architectures such as FPGAs, to verify its cross-architecture portability.

¡Thank you very much for your time!







