

Embedding a DSL in SYCL for Productive and Performant Tensor

Computing on Heterogeneous Devices



Abenezer Wudenhe

University of California, Riverside

Parallel Computing Lab, Intel

2022



Overview of T2S

• A constructive programming approach for spatial architectures (e.g. FPGAs) and vector architectures (e.g. GPUs)



Integrating T2S with SYCL





.



ínte



<pre>#define P #define P_cii_minus_1 #define P_kx_minus_1 #define P_kx_minus_1 #define P_ci_minus_1 #define P_cooo_minus_ #define P_out // Linearized address #define total_iy #define total_ix #define total_oy #define total_ox #define total_co #define total_ci // Type of the data t #define CTYPE float</pre>	cii, cooo, yyy, xxx, coo, yy, xx, ky, kx, ci, y, x, co, n cii+CII-1, cooo, yyy, xxx, coo, yy, xx, ky, kx, ci, y, x, co, n cii+CII-1, cooo, yyy, xxx, coo, yy, xx, ky+KY-1, kx-1, ci, y, x, co, n cii+CII-1, cooo, yyy, xxx, coo, yy, xx, ky+KY-1, kx-1, ci-1, y, x, co, n cii+CII-1, cooo, yyy, xxx, coo, yy, xx, ky, kx, ci, y, x, co, n cii, cooo, yyy, xxx, coo, yy, xx, ky, kx, ci, y, x, co, n cii, cooo, yyy, xxx, coo, yy, xx, ky, kx, ci, y, x, co, n coo, yyy, xxx, coo, yy, xx, ky, kx, ci, y, x, co, n (yyy + YYY*yY + YYY*YY*y + ky) (xxx + XXX*xx + XXX*X*x + kx) (yyy + YYY*yY + YYY*YY*y) (xxx + XXX*xx + XXX*X*x) (cooo + COO0*cco + COO0*CO0*co) (cii + CII*ci) to process in C and T25	Standard SYCL code Normal DPC++ compile fl Fat binary
#define TTYPE Float(3	2)	
<pre>// Inputs #i ImageParam I("I", TTY #define P_1 total #define P_K total #define P_0 total</pre>	<pre>/PE, 2), K("K", TTYPE, 2); _c1 + (IUTAL_C1) * n, total_iy + (TOTAL_IY) * total_ix L_co + (TOTAL_C0) * kx, total_ci + (TOTAL_CI) * ky L_co + (TOTAL_C0) * n, total_oy + (TOTAL_OY) * total_ox</pre>	
#define UN (I.di #else ImageParam I("I", TTY	.m(0).extent() / TOTAL_CI) /PE, 4), K("K", TTYPE, 4); Liv. total ix. total ci. n	

(intel)

Preprocessor (transparent to user)

Preprocessor (transparent to user) T2S spec **Unified programming model: T2S** Other code T2S compiler standard SYCL code using queue submit // URES 65 Standard SYCL code 66 Var cii("cii"), ci("ci"), cooo("coo"), co("coo"), ky("ky"), kx("kx"), yyy("yyy"), xxx("xxx"), yy("yy"), xx("xx"), y("y"), x("x"), n("n"); URE A("A", TTYPE, {P}), B("B", TTYPE, {P}), C("C", TTYPE, {P}), Out("Out"); 67 Normal DPC++ compile flow A(P) = select(cooo == 0, I(P_I), A(P_cooo_minus_1)); 68 B(P) = select(yyy == 0, K(P_K), B(P_yyy_minus_1)); 69 Fat binary C(P) = select(cii == 0 && ky == 0 && kx == 0 && ci == 0, 0, 70 71 select(cii == 0, select(ky == 0, select(kx == 0, C(P_ci_minus_1), C(P_kx_minus_1)), C(P_ky_minus_1)), C(P_cii_minus_1))) 72 + A(P) * B(P); **Uniform recurrence equations** Out(P Out) = select(cii == CII-1 && ky == KY-1 && kx == KX-1 && ci == CI-1, C(P)); 73 74 75 // Put all the UREs inside the same loop nest of X. A.merge_ures(B, C, Out); 76 77 // Explicitly set the loop bounds 78 A.set_bounds(cooo, 0, COO0, coo, 0, COO, co, 0, CO) 79 80 .set bounds(ky, 0, KY, kx, 0, KX) .set_bounds(cii, 0, CII, ci, 0, CI) 81 .set_bounds(yyy, 0, YYY, xxx, 0, XXX) 82 83 .set_bounds(yy, 0, YY, xx, 0, XX) .set_bounds(y, 0, Y, x, 0, X) 84 .set_bounds(n, 85 0, UN); 86 // Create a systolic array 87 **Space-time transform** A.space_time_transform(cooo, yyy); 88 89 // GPU can have many threads running in parallel. 90

SIMT

- #ifdef GPU 91
- A.gpu_blocks(x, co, n).gpu_threads(yy, xx); 92
- #endif 93





′inte



- const std::vector<Argument> &: A set of inputs to the Func or Stensor pipeline stage
- const std::string & fn_name: Name of the compiled function and file prefix
- const Target & target: Set of Halide Targets

i.e Target::IntelFPGA or Target::IntelGPU



Integrating T2S with OneAPI and SYCL

•••

```
1 /* gemm.cpp */
 2
3 // The only header file needed for including T2S.
4 #include "HalideBuffer.h"
5
 6 int main(){
 7
 8
      // A T2S specification that would create a matrix multiply accelerator on an FPGA.
9
      #pragma t2s spec start
10
11
          ImageParam A("A", TTYPE, 2), B("B", TTYPE, 2);
12
13
          // Synthesise the spec into a function "gemm", which is OneAPI/SYCL device code, for FPGA
14
          C.compile_to_oneapi( { A, B }, "gemm", IntelFPGA);
15
      #pragma t2s spec end
16
17
      // Initalize data
18
      const int TOTAL_I = III * II * I;
19
    const int TOTAL J = JJJ * JJ * J:
20
      const int TOTAL K = KKK * KK * K:
21 float *a = new float[TOTAL_K * TOTAL_I];
22
      float *b = new float[TOTAL J * TOTAL K]:
23
      float *c = new float[JJJ * III * JJ * II * J * I];
24
      for(unsigned int i = 0; i < (TOTAL_K * TOTAL_I); i++){ a[i] = random(); }</pre>
25
      for(unsigned int i = 0; i < (TOTAL_J * TOTAL_K); i++){ b[i] = random(); }</pre>
26
      for(unsigned int i = 0; i < (JJJ * III * JJ * II * J * I); i++){ c[i] = 0.0f; }
27
28
      // Dimensions of the data
29
      int a_dim[2] = {TOTAL_K, TOTAL_I};
30
      int b dim[2] = {TOTAL J, TOTAL K}:
31
      int c_dim[6] = {JJJ, III, JJ, II, J, I};
32
33
      // Below we invoke the generated matrix multiple accelerator
34
35
      // Call a function "gemm: that is to be generated by the embedded T2S specification above (C2)
36
      // Pass in the data and dimensions. For example, "A(A, a, a dim)" means that "A" in the
37
      // spec above (C2) corresponds to array "a" with dimension "a_dim".
38
      #pragma t2s_submit gemm (A, a, a_dim) (C, c, c_dim) (B, b, b_dim)
39
40
      return 0;
41 }
```







Integrating T2S

- 1. Wrapper
- 2. Host/Device Kernels
- 3. Memory Management
- 4. Performance Metrics

Preprocessor (transparent to user) T2S spec Other Host/Device Kernels code T2S compiler Memory Management standard SYCL code Performance metrics

Wrapper

1.

2.

3.

4.





Integrating T2S





- 3. Memory Management
- 4. Performance metrics



- Generated GEMM code takes in 2 types of arguments
 - sycl::ext::intel::fpga_selector device_selector
 - Used to either use FPGA emulation or synthesize to real hardware
 - Halide::Runtime::Buffer
 - T2SP's Halide Runtime Buffers. Used to manage device and host memory copies and management





• Host Kernels

Integrating T2S

- Host kernels are executed as normal C++ code on host memory
- Device Kernels

 Device kernels are submitted into device queues and kernel events are recorded in a std::vector for performance analysis at the completion of the application

```
1 /* gemm.sycl.h */
2 // Device Kernel
3 oneapi_kernel_events.push_back( q_device.submit([&](sycl::handler &h){
4    h.single_task<class kernel_1>([=](){
5         // ...
6    }); // h.single_task kernel_1
7 }) ); // q_device.submit
```



Host/Device Kernels

Integrating T2S

- Memory management
 - Device and Host memory are allocated though sycl::malloc_device() and std::malloc() respectively
- Memory copies (memcpy) are explicit and executed depending on the memory operation of that kernel
 - If kernel requires load operation, memcpy host->device before kernel submission
 - If kernel requires store operation, memcpy device->host after kernel execution is complete
- SYCL Pipes are used to communicate between data between kernels on the device
- Currently there is no implement of deviceto-device memory copies



Integrating T2S

7

8 9

10

11 12 }

14 15 } 16 }

1 /* gemm.sycl.h */ 2 // Performance

double tmp_start =

double $tmp_end =$



```
18 double events_time = (k_latest_end_time - k_earliest_start_time);
```

```
17 // Get time in ns
19 return events time:
```



Preprocessor (transparent to user) Other T2S spec T2S compiler standard SYCL code Using queue submit Standard SYCL code

Normal DPC++ compile flow

Fat binary

Recorded demo: T2S Full Compiling Flow GEMM



(intel) 19

Experiments/Results

- Benchmark Applications
 - General Matrix Multiply (GEMM)
 - 2-D convolution (CONV)
 - Capsule convolution (CAPSULE)

Experiments/Results

• GEMM/CONV/CAPSULE FPGA Synthesis Results for OpenCL and Generated OneAPI

Benchmark	ALUT	Registers	DSP blocks	RAM Blocks	Clock Freq.
GEMM	105151	180.874	89 / 1.518 (6 %)	462 / 2.713 (17 %)	264
CONV	98648	168.673	64 / 1.518 (4 %)	499 / 2,713 (18 %)	276
CAPSULE	119038	198,909	112 / 1,518 (7 %)	607 / 2,713 (22 %)	275

Experiments/Results

GEMM/CONV/CAPSULE FPGA Performance results for OpenCL and Generated OneAPI

Benchmark	GFLOPS
GEMM	302.402
CONV	285.301
CAPSULE	231.877

Conclusion

- Expanding Intel's T2SP to leverage SYCL implementation to extend the T2SP domain specific language to one uniform compiling flow for agnostic hardware acceleration
 - Developed Code generator to produce SYCL Device & Host Code
 - Implemented Clang source-to-source code for file preprocessing
 - FPGA implementation is able to non-trivial performance

Future Work

- Investigation as to performance bottlenecks to achieve faster tensor computing
- Extend to GPU, CPU, and other tensor computing accelerators

