



Towards performance portability of AI models using SYCL-DNN

Muhammad Tanvir – Staff Software Engineer – ML libraries

IWOCL – 10-12 May 2022


Outline

- Deep Learning Challenges
- Problem Statement – performance portability
- How does SYCL help?
- Performance Evaluations
- Deep Learning on RISC-V
- Conclusions

Deep Learning Challenges

- **Diversity** of Technologies and Techniques
- **Migrating** between Deep Learning frameworks
- Multiple implementations of same **NN Optimization algorithms**
- **Maintainability** of various version of low-level backend/libraries integrated in existing high-level frameworks
- **Hardcoded** implementation of inference engines for a restricted set of hardware
- Multiple implementations of same **NN algorithms** across various target **HW**

Deep Learning Challenges

- ✓ • **Diversity** of Technologies and Techniques.
 - ✓ • **Migrating** between DL framework.
 - ✓ • Multiple implementations of Same **NN Optimization algorithms**.
 - ✓ • **Maintainability** of various version of low-level libraries integrated in existing high-level frameworks.
 - ✓ • **Hardcoded** implementation of Inference engines for a restricted set of hardware.
 - Multiple implementations of same **NN algorithms** across various target **HW**
- 
- The diagram consists of two green curly braces on the right side of the list. The top brace groups the first two items, 'Diversity of Technologies and Techniques' and 'Migrating between DL framework', and is labeled 'ONNX' in green text. The bottom brace groups the next three items: 'Multiple implementations of Same NN Optimization algorithms', 'Maintainability of various version of low-level libraries integrated in existing high-level frameworks', and 'Hardcoded implementation of Inference engines for a restricted set of hardware'. This bottom brace is labeled 'ONNX Runtime' in green text.

Problem Statement

- Multiple implementations of same **NN algorithms** across various target **HW** – is there a way around this?
- An ideal solution(s) should:
 - Be **portable** across different platforms
 - Require **little (to no) changes** to the actual kernel code
 - Yield **acceptable performance** as compared to vendor optimized code
 - Have **backward compatibility (maintainability)**

SYCL

- Supports cross-platform portability
 - Different implementations of SYCL compilers provide a variety of targets
- Is maintainable
 - Open-source implementations
- Kernel Code modifications?
- Performance Portability?

SYCL

- **Kernel Code Modification?**
- **Possible Solution:** Use C++ template meta programming to write highly parametrized kernels (for the most compute intensive operations)
- **Benefits:**
 - Reuse the same kernel code
 - Modify (tune) the template parameters to maximize performance on target hardware

SYCL

- Performance Portability?

- Possible Solution:

- Tune the template parameters of the kernel to best match the underlying hardware*
- Expose maximum performance out of the tuned kernel

- Acceptable Performance:

- Within 70%-80% range of vendor optimized code's performance

*e.g. https://github.com/codeplaysoftware/sycl-blas/blob/master/tools/auto_tuner/gen/intel_gpu.json

SYCL

SYCL-DNN

- Conv
- Batchnorm
- Pool
- Softmax etc

<https://github.com/codeplaysoftware/SYCL-DNN>

SYCL-BLAS

- GEMM
- GEMV
- Reduction etc

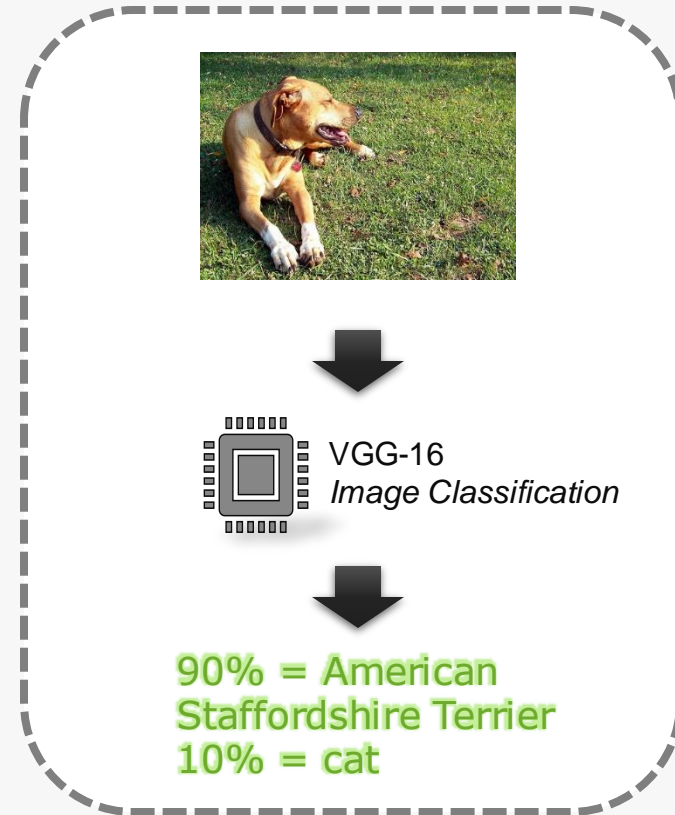
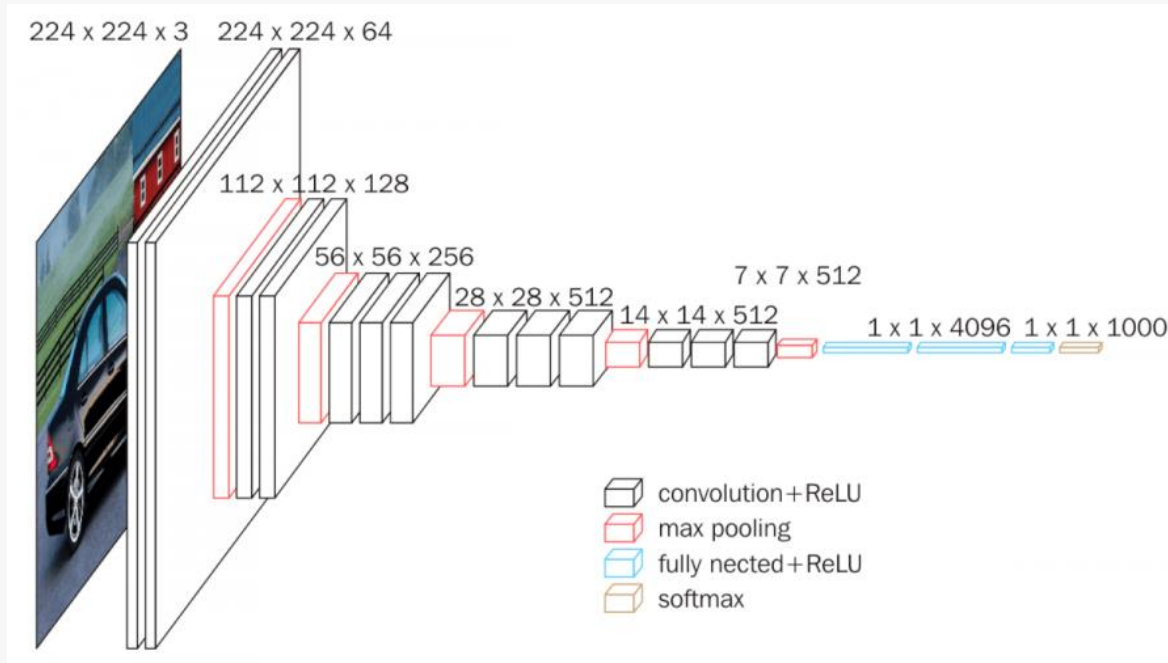
<https://github.com/codeplaysoftware/sycl-blas>

Performance Evaluations

1. Identify the most compute intensive operation(s)
2. Choose the target hardware
3. Tune SYCL code based on the target hardware
4. Evaluate performance
5. Repeat steps 2-4 for remaining target hardware

Performance Evaluations

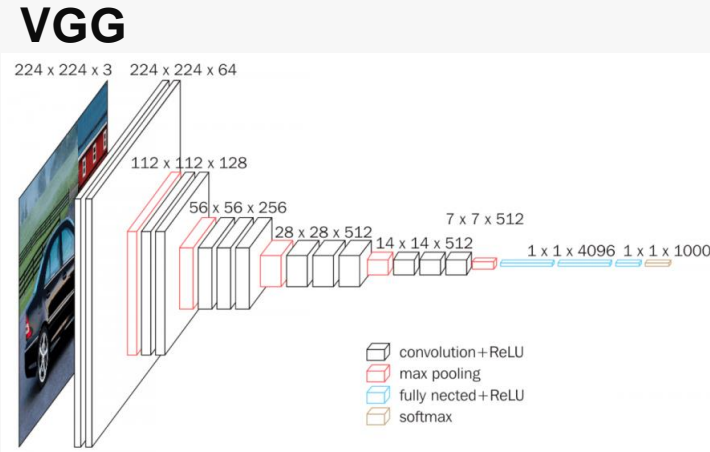
VGG



VGG Demo

Performance Evaluations

1. Identify the most compute intensive operation



Conv	->	sycl::conv2d
BiasAdd	->	sycl::biasAdd
Relu	->	sycl::relu
Pooling	->	sycl::pooling
Gemm	->	sycl::gemm
Softmax	->	sycl::softmax

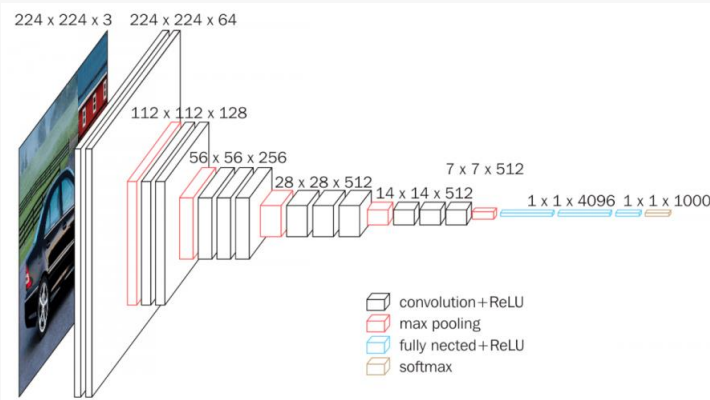
```
std::vector<float> output;
std::string data_dir(argv[1]);
auto input = read_image_data(argv[2], backend);
Network network(input, output, data_dir, backend, *selector);
network.add_layer<ConvolutionLayer, 3>({{3, 64, 224}});
network.add_layer<BiasAddLayer, 2>({64, 224});
network.add_layer<ReLULayer, 0>({});
network.add_layer<ConvolutionLayer, 3>({{64, 64, 224}});
network.add_layer<BiasAddLayer, 2>({64, 224});
network.add_layer<ReLULayer, 0>({});
network.add_layer<PoolingLayer, 2>({64, 224});
network.add_layer<ConvolutionLayer, 3>({{64, 128, 112}});
network.add_layer<BiasAddLayer, 2>({128, 112});
network.add_layer<ReLULayer, 0>({});
network.add_layer<ConvolutionLayer, 3>({{128, 128, 112}});
network.add_layer<BiasAddLayer, 2>({128, 112});
network.add_layer<ReLULayer, 0>({});
network.add_layer<PoolingLayer, 2>({128, 112});
network.add_layer<ConvolutionLayer, 3>({{128, 256, 56}});
network.add_layer<BiasAddLayer, 2>({256, 56});
network.add_layer<ReLULayer, 0>({});
network.add_layer<ConvolutionLayer, 3>({{256, 256, 56}});
network.add_layer<BiasAddLayer, 2>({256, 56});
network.add_layer<ReLULayer, 0>({});
network.add_layer<ConvolutionLayer, 3>({{256, 256, 56}});
network.add_layer<BiasAddLayer, 2>({256, 56});
network.add_layer<ReLULayer, 0>({});
network.add_layer<PoolingLayer, 2>({256, 56});
network.add_layer<ConvolutionLayer, 3>({{256, 512, 28}});
network.add_layer<BiasAddLayer, 2>({512, 28});
network.add_layer<ReLULayer, 0>({});
network.add_layer<ConvolutionLayer, 3>({{512, 512, 28}});
network.add_layer<BiasAddLayer, 2>({512, 28});
network.add_layer<ReLULayer, 0>({});
network.add_layer<PoolingLayer, 2>({512, 28});
network.add_layer<ConvolutionLayer, 3>({{512, 512, 14}});
network.add_layer<BiasAddLayer, 2>({512, 14});
network.add_layer<ReLULayer, 0>({});
network.add_layer<ConvolutionLayer, 3>({{512, 512, 14}});
network.add_layer<BiasAddLayer, 2>({512, 14});
network.add_layer<ReLULayer, 0>({});
network.add_layer<PoolingLayer, 2>({512, 14});
network.add_layer<FullyConnectedLayer, 1>({4096});
network.add_layer<BiasAddLayer, 2>({4096, 1});
network.add_layer<ReLULayer, 0>({});
network.add_layer<FullyConnectedLayer, 1>({4096});
network.add_layer<BiasAddLayer, 2>({4096, 1});
network.add_layer<ReLULayer, 0>({});
network.add_layer<FullyConnectedLayer, 1>({1000});
network.add_layer<BiasAddLayer, 2>({1000, 1});
network.add_layer<SoftmaxLayer, 0>({});
network.run();
```

Performance Evaluations

1. Identify the most compute intensive operation

GEMM

VGG



Conv	->	sycl::conv2d
BiasAdd	->	sycl::biasAdd
Relu	->	sycl::relu
Pooling	->	sycl::pooling
Gemm	->	sycl::gemm
Softmax	->	sycl::softmax

```
std::vector<float> output;
std::string data_dir(argv[1]);
auto input = read_image_data(argv[2], backend);
Network network(input, output, data_dir, backend, *selector);
network.add_layer<ConvolutionLayer, 3>({{3, 64, 224}});
network.add_layer<BiasAddLayer, 2>({{64, 224}});
network.add_layer<ReLULayer, 0>({});
network.add_layer<ConvolutionLayer, 3>({{64, 64, 224}});
network.add_layer<BiasAddLayer, 2>({{64, 224}});
network.add_layer<ReLULayer, 0>({});
network.add_layer<PoolingLayer, 2>({{64, 224}});
network.add_layer<ConvolutionLayer, 3>({{64, 128, 112}});
network.add_layer<BiasAddLayer, 2>({{128, 112}});
network.add_layer<ReLULayer, 0>({});
network.add_layer<ConvolutionLayer, 3>({{128, 128, 112}});
network.add_layer<BiasAddLayer, 2>({{128, 112}});
network.add_layer<ReLULayer, 0>({});
network.add_layer<PoolingLayer, 2>({{128, 112}});
network.add_layer<ConvolutionLayer, 3>({{128, 256, 56}});
network.add_layer<BiasAddLayer, 2>({{256, 56}});
network.add_layer<ReLULayer, 0>({});
network.add_layer<ConvolutionLayer, 3>({{256, 256, 56}});
network.add_layer<BiasAddLayer, 2>({{256, 56}});
network.add_layer<ReLULayer, 0>({});
network.add_layer<ConvolutionLayer, 3>({{256, 256, 56}});
network.add_layer<BiasAddLayer, 2>({{256, 56}});
network.add_layer<ReLULayer, 0>({});
network.add_layer<ConvolutionLayer, 3>({{256, 512, 28}});
network.add_layer<BiasAddLayer, 2>({{512, 28}});
network.add_layer<ReLULayer, 0>({});
network.add_layer<ConvolutionLayer, 3>({{512, 512, 28}});
network.add_layer<BiasAddLayer, 2>({{512, 28}});
network.add_layer<ReLULayer, 0>({});
network.add_layer<PoolingLayer, 2>({{512, 28}});
network.add_layer<ConvolutionLayer, 3>({{512, 512, 14}});
network.add_layer<BiasAddLayer, 2>({{512, 14}});
network.add_layer<ReLULayer, 0>({});
network.add_layer<ConvolutionLayer, 3>({{512, 512, 14}});
network.add_layer<BiasAddLayer, 2>({{512, 14}});
network.add_layer<ReLULayer, 0>({});
network.add_layer<PoolingLayer, 2>({{512, 14}});
network.add_layer<FullyConnectedLayer, 1>({{4096}});
network.add_layer<BiasAddLayer, 2>({{4096, 1}});
network.add_layer<ReLULayer, 0>({});
network.add_layer<FullyConnectedLayer, 1>({{4096}});
network.add_layer<BiasAddLayer, 2>({{4096, 1}});
network.add_layer<ReLULayer, 0>({});
network.add_layer<FullyConnectedLayer, 1>({{1000}});
network.add_layer<BiasAddLayer, 2>({{1000, 1}});
network.add_layer<SoftmaxLayer, 0>({});
network.run();
```

Performance Evaluations

2. Choose the target hardware

Intel CPU

Intel GPU

NVIDIA GPU

RISC-V

Performance Evaluations

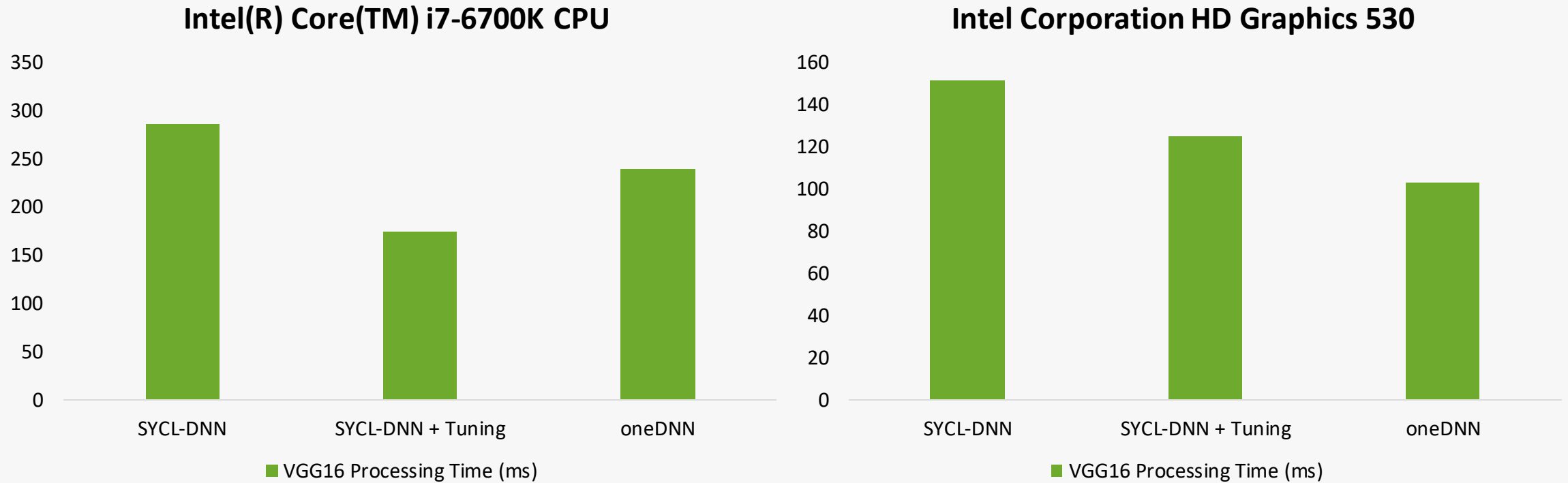
3. Tune SYCL code based on the target hardware

- Target Hardware
 - Intel(R) Core(TM) i7-6700K CPU @ 4 GHz
 - Intel Corporation HD Graphics 530
- Most Compute Intensive Operation
 - GEMM
- Tuning Methodology
 - Lawson, John, and Mehdi Goli. "Performance portability through machine learning guided kernel selection in SYCL libraries." *Parallel Computing* 107 (2021): 102813.

Tuning paradigm

- Every GEMM operation, in DNNs, has different compute intensity
- Most optimal solution
 - One tuned kernel per GEMM operation
 - Yields the best performance
 - Size of the library increases drastically
- Semi-optimal solution
 - One tuned kernel per DNN model
 - Yields semi-optimal performance
 - Limited no. of kernels in the library

Performance Evaluations



*With semi-optimal tuning regime

Performance Evaluations

2. Choose the target hardware

Intel CPU

Intel GPU

NVIDIA GPU

RISC-V

Performance Evaluations

3. Tune SYCL code based on the target hardware

- Target Hardware
 - NVIDIA Titan RTX GPU – 24 GB DDR6
- Most Compute Intensive Operation
 - GEMM
- Tuning tool
 - Modified version of SYCL-BLAS auto-tuner*

*https://github.com/codeplaysoftware/sycl-blas/tree/master/tools/auto_tuner

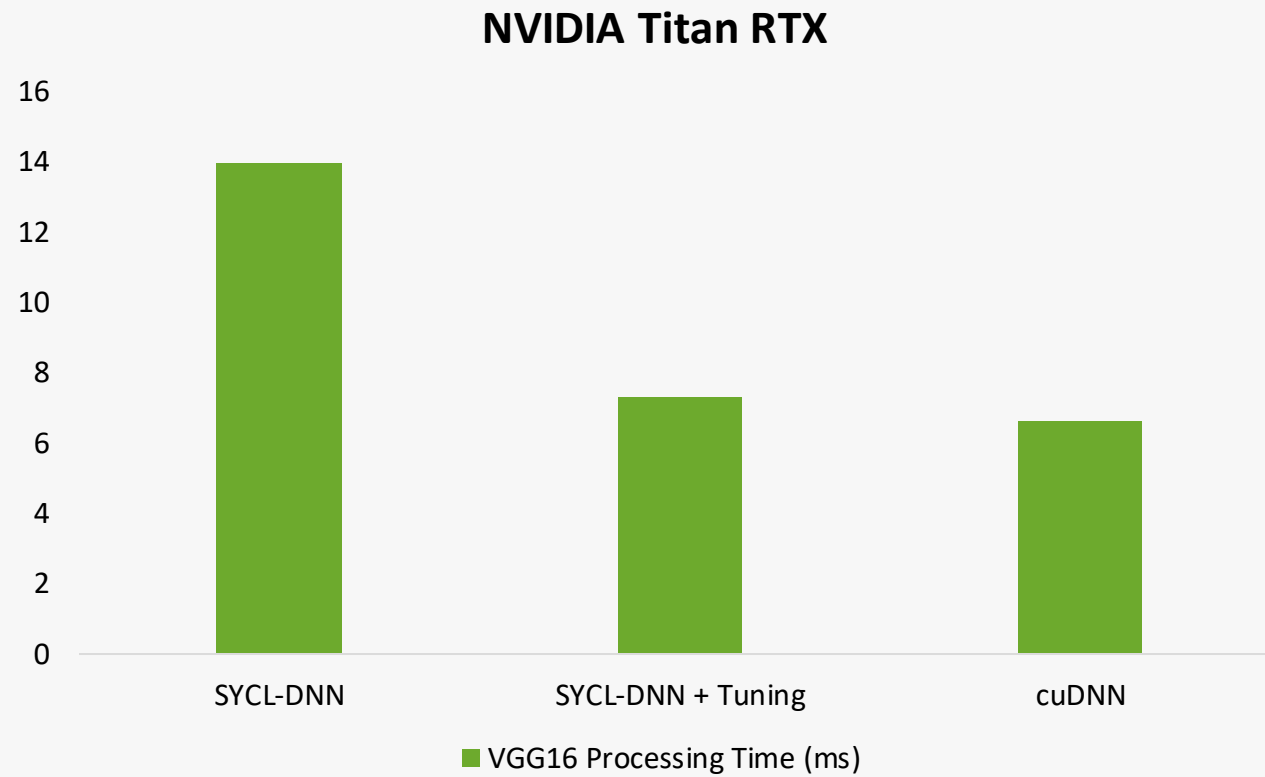
Performance Evaluations

3. Tune SYCL code based on the target hardware

- Search all possible GEMM configs* exhaustively
 - Choose a GEMM config
 - Run the entire VGG16 (DNN) model and measure performance
 - Record results for all possible GEMM configs
 - Choose the GEMM config which yields the best performance

*https://github.com/codeplaysoftware/sycl-blas/blob/master/tools/auto_tuner/gen/nvidia_gpu.json

Performance Evaluations



Performance Evaluations

2. Choose the target hardware

Intel CPU

Intel GPU

NVIDIA GPU

RISC-V

Deep Learning on RISC-V

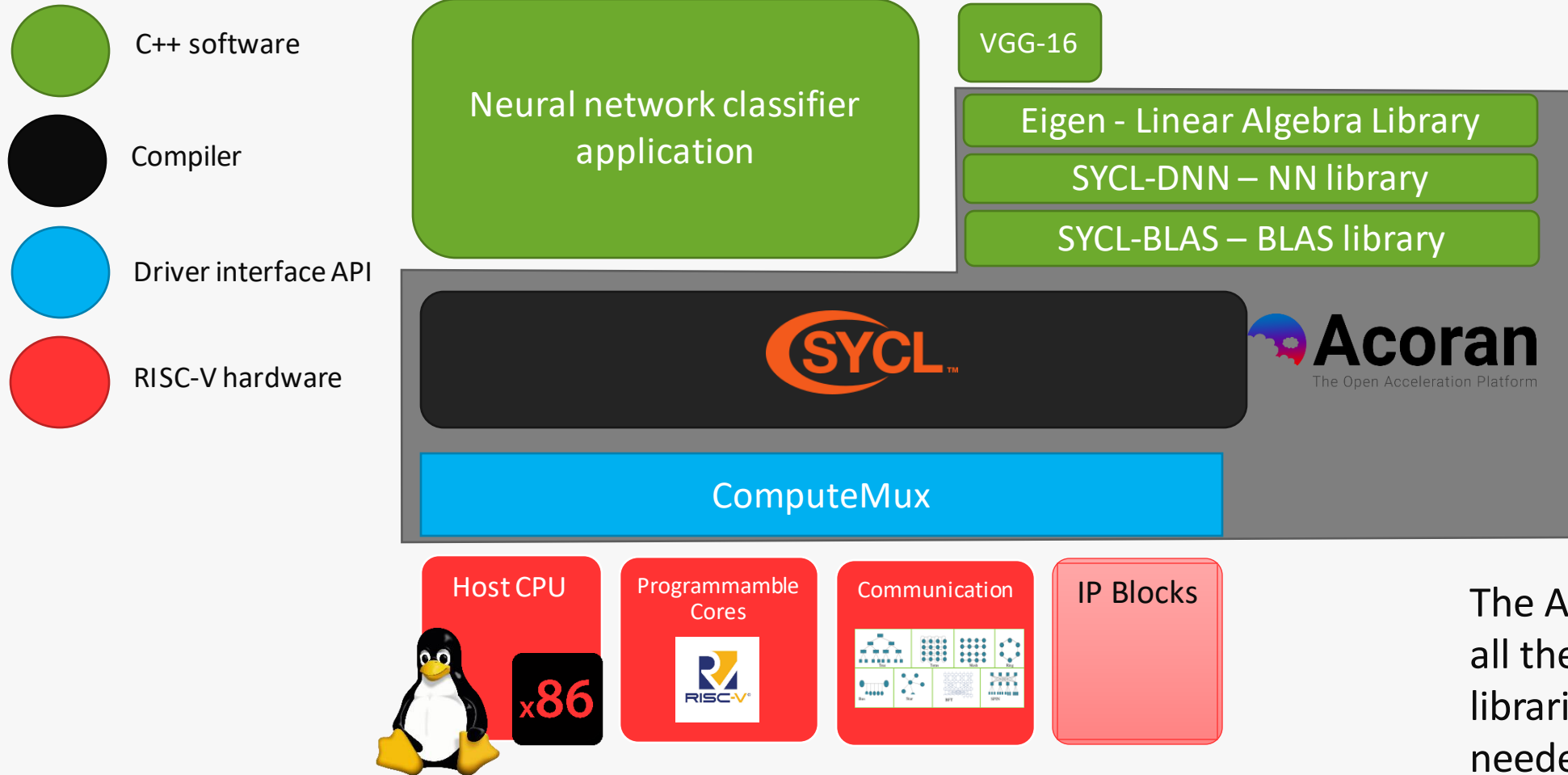
- Target Hardware
 - RISC-V spike simulator – single core
- Software Stack
 - Acoran – The Open Acceleration Platform*

*<https://www.codeplay.com/solutions/acoran/>

Why Neural networks on RISC-V?

- Domain specific accelerators are required to achieve **cost-effective performance** on-chip
- Cost effective performance requires **tuning the design** to the needs of the workload required
- RISC-V ISA has a minimalist base integer instruction set and provides custom extensions
 - An ideal starting point for creating special accelerators
- More companies are looking at RISC-V to enable AI software
- Designs can benefit from the RISC-V vector extension
 - Enables vectorization for various application
 - Helps achieve high compute density on chip

Acoran – The Open Acceleration Platform



The Acoran platform provides all the supporting open-source libraries and frameworks needed to build this neural network demonstration

GEMM Operator

Main Computation



```
...
for (int a = a_start, b = b_start; a <= a_end;
    a += blockSize, b += (blockSize * matSize)) {
    for (int k = 0; k < blockSize; k++) {
        tmp +=
            pBA[localY * blockSize + k] * pBB[localX * blockSize + k];
    }
    // The barrier ensures that all threads have written to local
    // memory before continuing
    it.barrier(access::fence_space::local_space);
}
...
pC[elemIndex] = tmp;
...
}
```

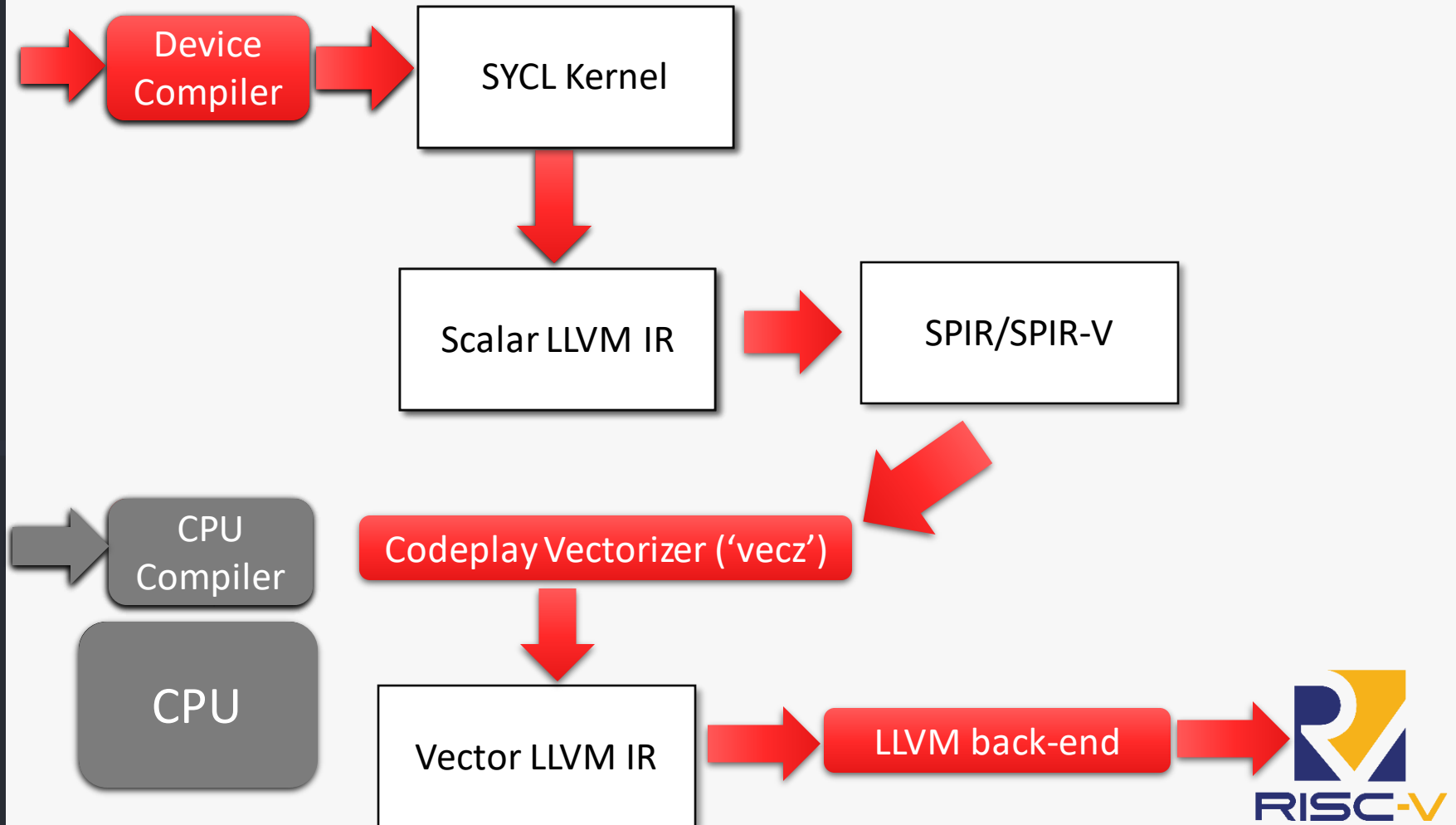
```

range<1> dimensions(matSize * matSize);
const property_list props = {property::buffer::use_host_ptr()};
4 buffer<T> bA(MA, dimensions, props);
5 buffer<T> bB(MB, dimensions, props);
6 buffer<T> bC(MC, dimensions, props);

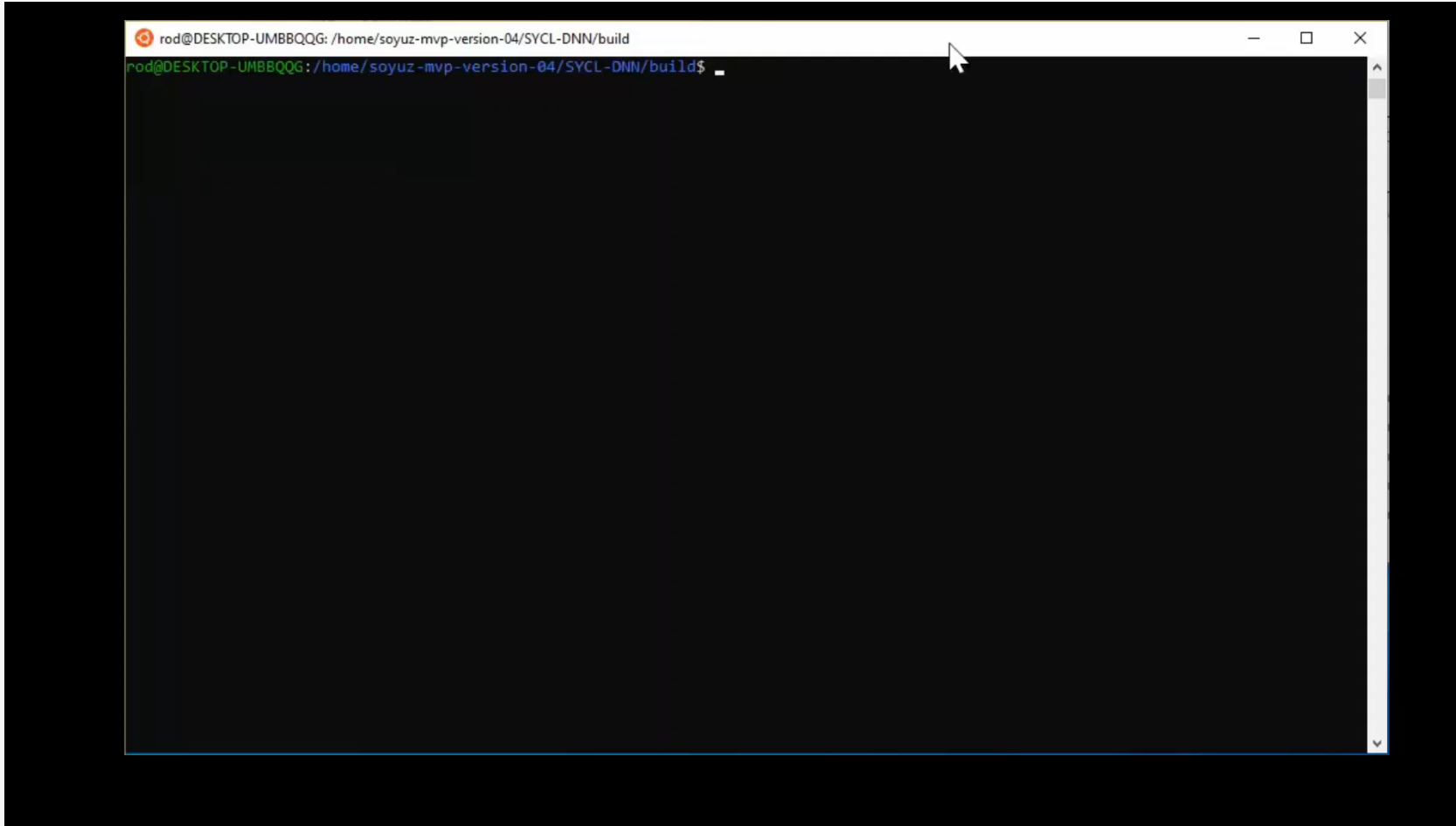
7
8 q.submit([&](handler& cgh) {
9     auto pA = bA.template get_access<access::mode::read>(cgh);
10    auto pB = bB.template get_access<access::mode::read>(cgh);
11    auto pC = bC.template get_access<access::mode::write>(cgh);
12    auto localRange = range<1>(blockSize * blockSize);
13
14    accessor<T, 1, access::mode::read_write, access::target::local> pBA(
15        localRange, cgh);
16    accessor<T, 1, access::mode::read_write, access::target::local> pBB(
17        localRange, cgh);
18
19    cgh.parallel_for<mxm_kernel>{
20        nd_range<2>{range<2>(matSize, matSize),
21                    range<2>(blockSize, blockSize)},
22        [=](nd_item<2> it) {
23            // Current block
24            int blockX = it.get_group(1);
25            int blockY = it.get_group(0);
26
27            // Current local item
28            int localX = it.get_local_id(1);
29            int localY = it.get_local_id(0);
30
31            // Start in the A matrix
32            int a_start = matSize * blockSize * blockY;
33            // End in the b matrix
34            int a_end = a_start + matSize - 1;
35            // Start in the b matrix
36            int b_start = blockSize * blockX;
37
38            // Result for the current C(i,j) element
39            T tmp = 0.0f;
40            // We go through all a, b blocks
41            for (int a = a_start, b = b_start; a <= a_end;
42                 a += blockSize, b += (blockSize * matSize)) {
43                // Copy the values in shared memory collectively
44                pBA[localY * blockSize + localX] =
45                    pA[a + matSize * localY + localX];
46                // Note the swap of X/Y to maintain contiguous access
47                pBB[localX * blockSize + localY] =
48                    pB[b + matSize * localY + localX];
49                it.barrier(access::fence_space::local_space);
50                // Now each thread adds the value of its sum
51                for (int k = 0; k < blockSize; k++) {
52                    tmp +=
53                        pBA[localY * blockSize + k] * pBB[localX * blockSize + k];
54                }
55                // The barrier ensures that all threads have written to local
56                // memory before continuing
57                it.barrier(access::fence_space::local_space);
58            }
59            auto elemIndex = it.get_global_id(0) * it.get_global_range()[1] +
60                it.get_global_id(1);
61            // Each thread updates its position
62            pC[elemIndex] = tmp;
63        });
64    });
65    return false;
66 }
67 }

```

RISC-V/RVV Kernel compilation flow FC



Deep Learning on RISC-V



VGG Demo

Conclusions

- SYCL-DNN / SYCL-BLAS have support for efficient acceleration of popular DNNs
- Acoran platform provides an end-to-end compute stack for accelerating DNNs on RISC-V
 - <https://developer.codeplay.com/products/acoran/pre-alpha>
- Recent Update: Adding SYCL to upstream ONNX Runtime
 - <https://github.com/codeplaysoftware/onnxruntime>

We're
Hiring!

codeplay.com/careers/



Enable AI & HPC to be Open, Safe and Accessible to All

Thank you

muhammad.tanvir@codeplay.com



[@codeplaysoft](https://twitter.com/codeplaysoft)



info@codeplay.com



codeplay.com