





THE UNIVERSITY *of* EDINBURGH **informatics**

- Largest Informatics Department in the UK:
 - > 500 academic and research staff
+ PhD students
- Overall 6 Research Institutes
 - 2 particular relevant for the topic of the talk:
- **ICSA** — Institute for Computing Systems Architecture
 - Compiler & Architecture
 - Parallel Computing
 - ...
- **LFCS** — Laboratory for Foundations of Computer Science
 - Programming Languages and Foundations
 - Software Engineering
 - ...



Structured Parallel Programming

From High-Level Functional Expressions
to High-Performance OpenCL Code

Michel Steuwer

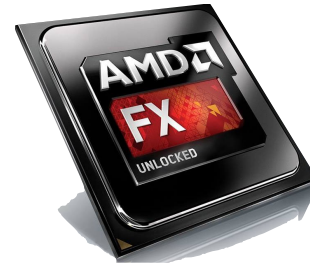
<http://homepages.inf.ed.ac.uk/msteuwer/>



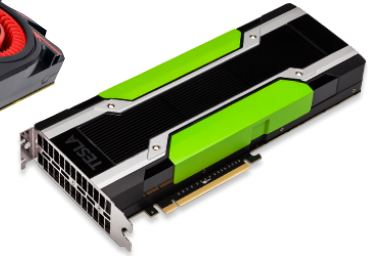
THE UNIVERSITY
of EDINBURGH

The Problem(s)

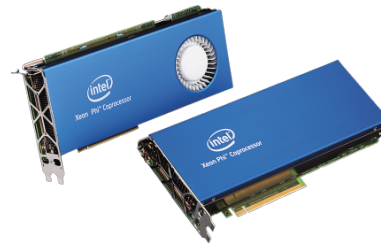
- Parallel processors everywhere
- Many different types: CPUs, GPUs, ...
- **First Major Challenge:**
Parallel programming is hard.
Optimising is even harder!
- **Second Major Challenges:**
No portability of performance!



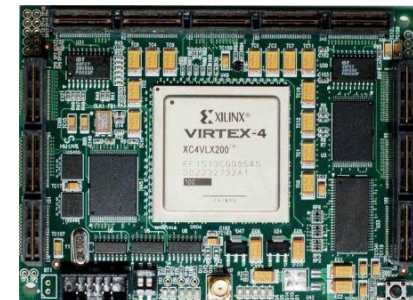
CPU



GPU



Accelerator



FPGA

Part I: Addressing the Programmability Challenge



Programming with OpenCL

- Case Study: Parallel reduction in OpenCL

```
kernel void reduce(global float* g_idata, global float* g_odata,
                  unsigned int n, local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i   = get_global_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    barrier(CLK_LOCAL_MEM_FENCE);
    // do reduction in local memory
    for (unsigned int s=1; s < get_local_size(0); s*= 2) {
        if ((tid % (2*s)) == 0) {
            l_data[tid] += l_data[tid + s];
            barrier(CLK_LOCAL_MEM_FENCE);
        }
    }
    // write result for this work-group to global memory
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0];
}
```

Programming with OpenCL

- Case Study: Parallel reduction in OpenCL

Kernel function executed in parallel by multiple **work-items**

```
kernel void reduce(global float* g_idata, global float* g_odata,
                  unsigned int n, local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i = get_global_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    barrier(CLK_LOCAL_MEM_FENCE);
    // do reduction in local memory
    for (unsigned int s=1; s < get_local_size(0); s*= 2) {
        if ((tid % (2*s)) == 0) {
            l_data[tid] += l_data[tid + s];
            barrier(CLK_LOCAL_MEM_FENCE);
        }
    }
    // write result for this work-group to global memory
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0];
}
```

Work-items are identified by a unique **global id**



Programming with OpenCL

- Case Study: Parallel reduction in OpenCL

Work-items are grouped into work-groups

Local id within work-group

```
kernel void reduce(global float* g_idata, global float* g_odata,
                  unsigned int n, local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i   = get_global_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    barrier(CLK_LOCAL_MEM_FENCE);
    // do reduction in local memory
    for (unsigned int s=1; s < get_local_size(0); s*= 2) {
        if ((tid % (2*s)) == 0) {
            l_data[tid] += l_data[tid + s];
            barrier(CLK_LOCAL_MEM_FENCE);
        }
    }
    // write result for this work-group to global memory
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0];
}
```


Programming with OpenCL

- Case Study: Parallel reduction in OpenCL

Big, but slow **global** memory

Small, but fast **local** memory

```
kernel void reduce(global float* g_idata, global float* g_odata,
                  unsigned int n, local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i   = get_global_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    barrier(CLK_LOCAL_MEM_FENCE);
    // do reduction in local memory
    for (unsigned int s=1; s < get_local_size(0); s*= 2) {
        if ((tid % (2*s)) == 0) {
            l_data[tid] += l_data[tid + s];
            barrier(CLK_LOCAL_MEM_FENCE);
        }
    }
    // write result for this work-group to global memory
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0];
}
```

Memory **barriers** for consistency



Programming with OpenCL

- Case Study: Parallel reduction in OpenCL

```
kernel void reduce(global float* g_idata, global float* g_odata,
                  unsigned int n, local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i   = get_global_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    barrier(CLK_LOCAL_MEM_FENCE);
    // do reduction in local memory
    for (unsigned int s=1; s < get_local_size(0); s*= 2) {
        if ((tid % (2*s)) == 0) {
            l_data[tid] += l_data[tid + s];
        }
        barrier(CLK_LOCAL_MEM_FENCE);
    }
    // write result for this work-group to global memory
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0];
}
```

Functionally correct implementations in OpenCL are hard!



The SkelCL Programming Model

Three high-level features added to OpenCL:

- **parallel container data types**

for unified memory management between CPU and (multiple) GPUs

- **implicit memory transfers between CPU and GPU**
- **lazy copying minimizes data transfers**

- **recurring patterns of parallelism**

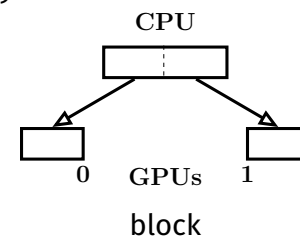
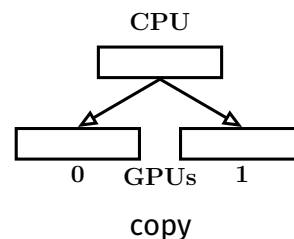
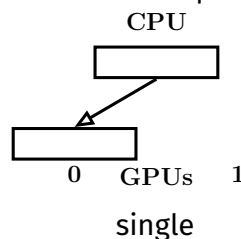
(a.k.a., algorithmic skeletons) for easily expressing parallel computation patterns;

$$\text{zip } (\oplus) [x_1, \dots, x_n] [y_1, \dots, y_n] = [x_1 \oplus y_1, \dots, x_n \oplus y_n]$$

$$\text{reduce } (\oplus) \oplus_{\text{id}} [x_1, \dots, x_n] = \oplus_{\text{id}} \oplus x_1 \oplus \dots \oplus x_n$$

- **data distribution and redistribution**

mechanisms for transparent data transfers in multi-GPU systems.



The SkelCL Library by Example

dotProduct A B = reduce (+) 0 ◦ zip (×) A B

```
#include <SkelCL/SkelCL.h>
#include <SkelCL/Zip.h>
#include <SkelCL/Reduce.h>
#include <SkelCL/Vector.h>

float dotProduct(const float* a, const float* b, int n) {
    using namespace skelcl;
    skelcl::init( 1_device.type(deviceType::ANY) );

    auto mult = zip([](float x, float y) { return x*y; });
    auto sum = reduce([](float x, float y) { return x+y; }, 0);

    Vector<float> A(a, a+n); Vector<float> B(b, b+n);

    Vector<float> C = sum( mult(A, B) );

    return C.front();
}
```

From SkelCL to OpenCL

```
1
#include <SkelCL/SkelCL.h>
#include <SkelCL/Zip.h>
#include <SkelCL/Reduce.h>
#include <SkelCL/Vector.h>

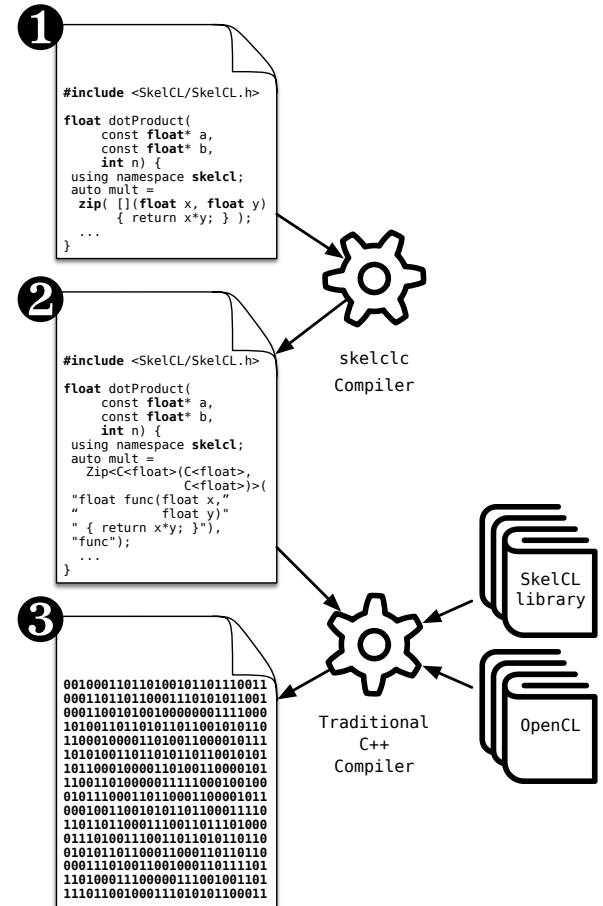
float dotProduct(const float* a, const float* b, int n) {
    using namespace skelcl;
    skelcl::init( 1_device.type(deviceType::ANY) );

    auto mult = zip([](float x, float y) { return x*y; });
    auto sum = reduce([](float x, float y) { return x+y; }, 0);

    Vector<float> A(a, a+n); Vector<float> B(b, b+n);

    Vector<float> C = sum( mult(A, B) );

    return C.front();
}
```



From SkelCL to OpenCL

2

```
#include <SkelCL/SkelCL.h>
#include <SkelCL/Zip.h>
#include <SkelCL/Reduce.h>
#include <SkelCL/Vector.h>

float dotProduct(const float* a, const float* b, int n) {
    using namespace skelcl;
    skelcl::init( 1_device.type(deviceType::ANY) );

    auto mult = Zip<Container<float>(Container<float>,
                                   Container<float>>(
        Source("float func(float x, float y) {return x*y;}"));
    auto sum = Reduce<Vector<float>(Vector<float>>(
        Source("float func(float x, float y) {return x+y;}"), "0"));

    Vector<float> A(a, a+n); Vector<float> B(b, b+n);

    Vector<float> C = sum( mult(A, B) );

    return C.front();
}
```

1

```
#include <SkelCL/SkelCL.h>

float dotProduct(
    const float* a,
    const float* b,
    int n) {
    using namespace skelcl;
    auto mult =
        zip( [](float x, float y)
            { return x*y; } );
    ...
}
```

2

```
#include <SkelCL/SkelCL.h>

float dotProduct(
    const float* a,
    const float* b,
    int n) {
    using namespace skelcl;
    auto mult =
        Zip<C<float>(C<float>>(
        "float func(float x,
                    float y)"
        " { return x*y; }"),
        "func");
    ...
}
```

3

```
00100011011010011101110011
0001101101100011101011001
00011001010010000000111000
101001101101011011001010110
11000100001101001100001011
1010100110110110110010101
1011000100001100010000101
11001101000001111000100100
01011000110110001100001011
00010011001010110110001110
11011011000110011011101000
0111010011100110110110110
01010110110001100011011010
00011101001100100011011101
110100011100000111001001101
1110110010001110101100011
```



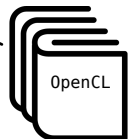
skelclc
Compiler



Traditional
C++
Compiler



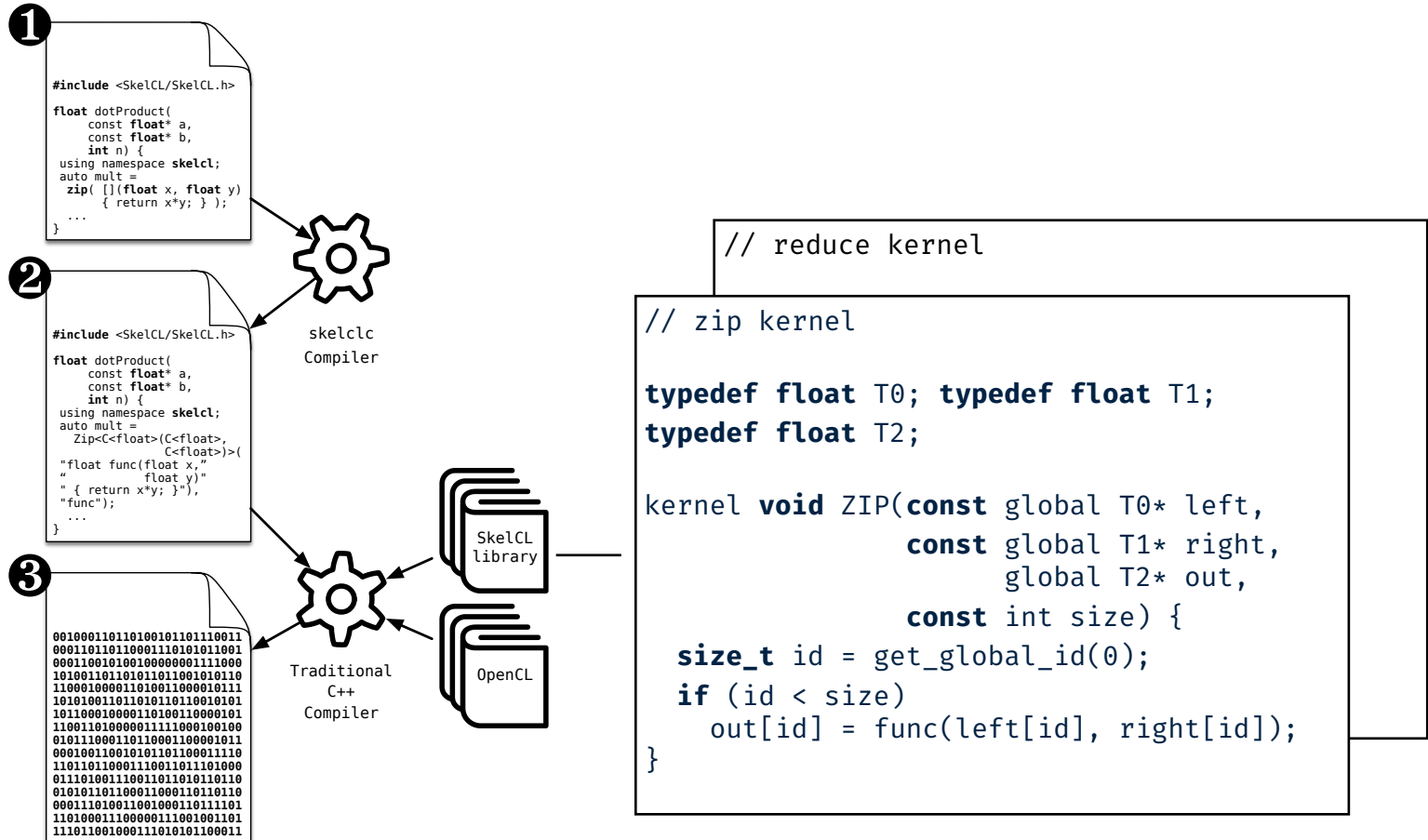
SkelCL
Library



OpenCL



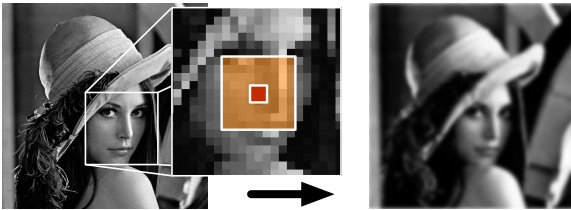
From SkelCL to OpenCL



Two Novel Algorithmic Skeletons

Stencil Computations

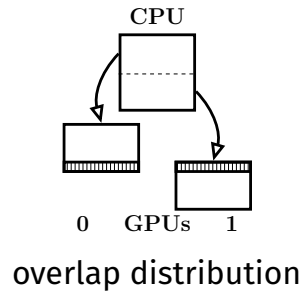
Example: Gaussian blur



$\text{gauss } M = \text{stencil } f \ 1 \ \bar{0} \ M$

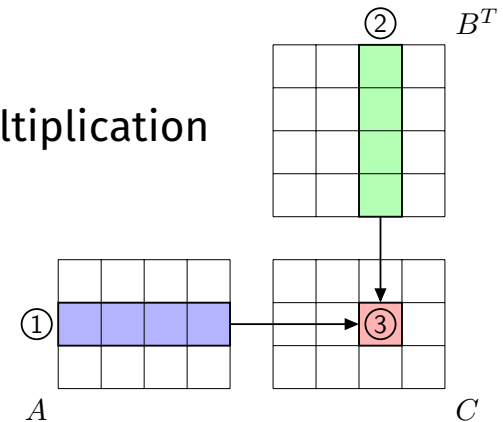
where f is the weighted gaussian kernel

Multi-GPU support:



Allpairs Computations

Example:
Matrix Multiplication



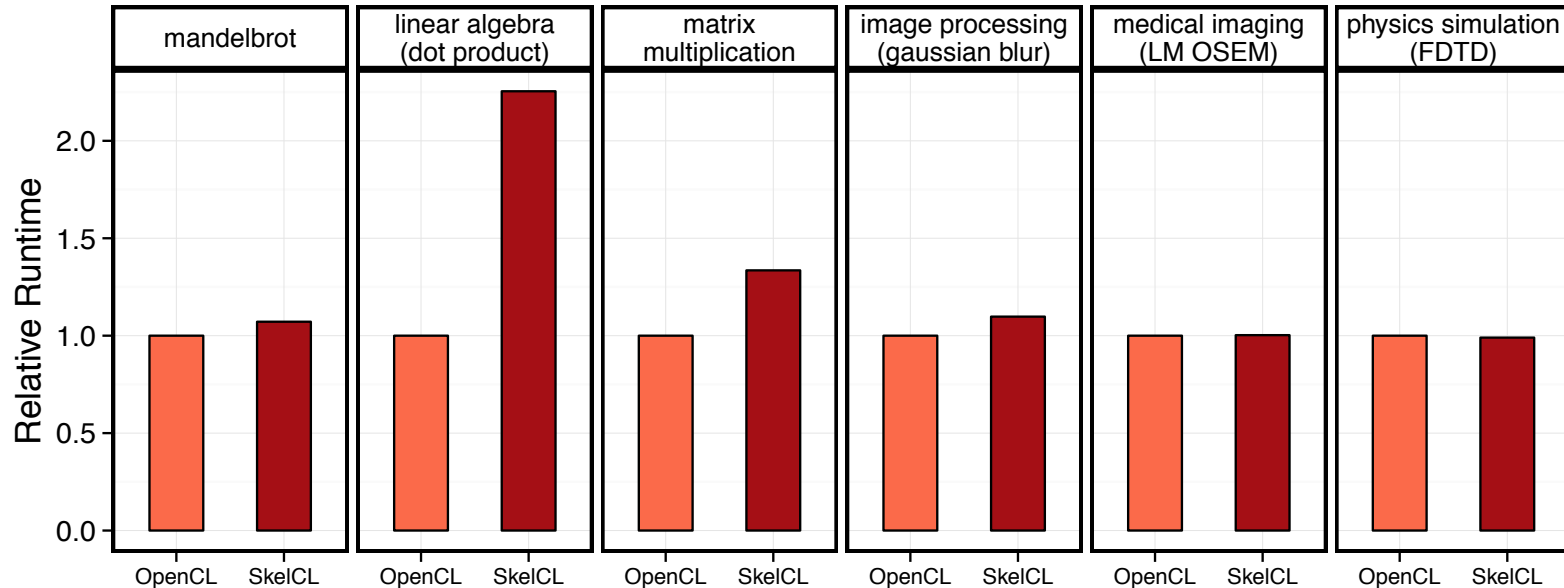
$A \times B = \text{allpairs dotProduct } A \ B^T$

Optimization for zipReduce patterns:

$\text{dotProduct } a \ b = \text{zipReduce } (+) \ 0 \ (\times) \ a \ b$

Multi-GPU support with
block and **copy** distribution

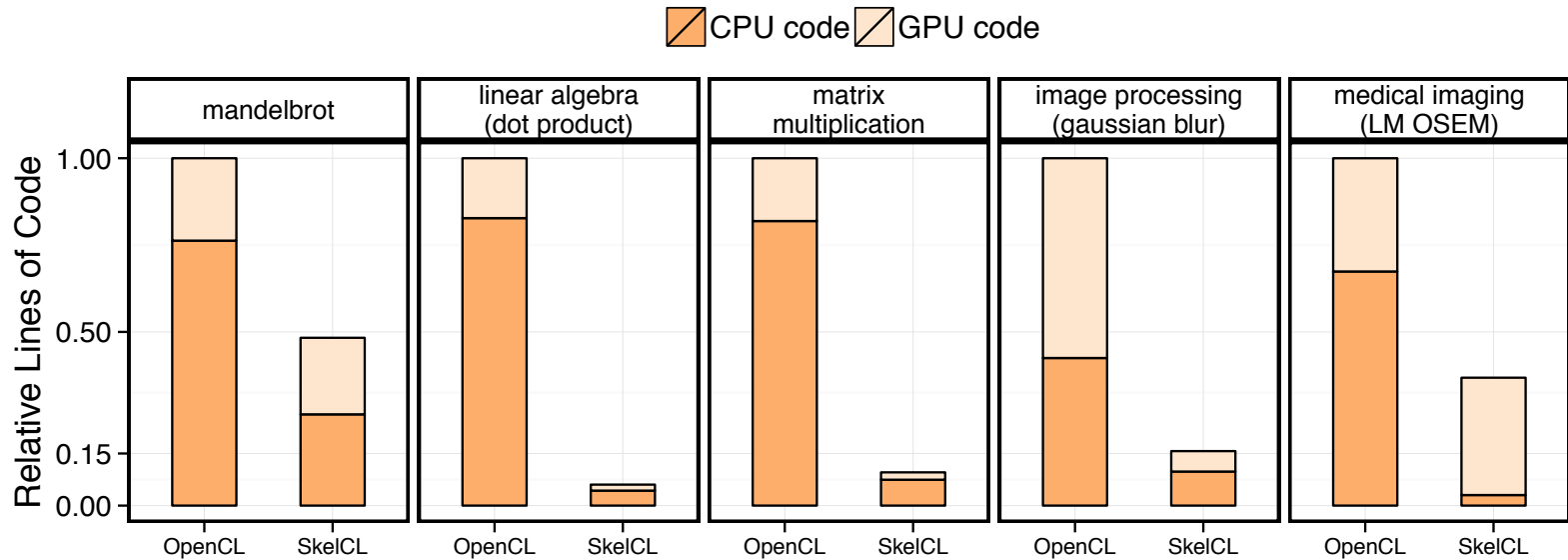
SkelCL Evaluation – Performance



SkelCL performance close to native OpenCL code!

(Exception: dot product ... wait for Part II)

SkelCL Evaluation — Productivity



SkelCL programs are significantly shorter!

SkelCL is open source software and available from <http://github.com/skelcl/skelcl>

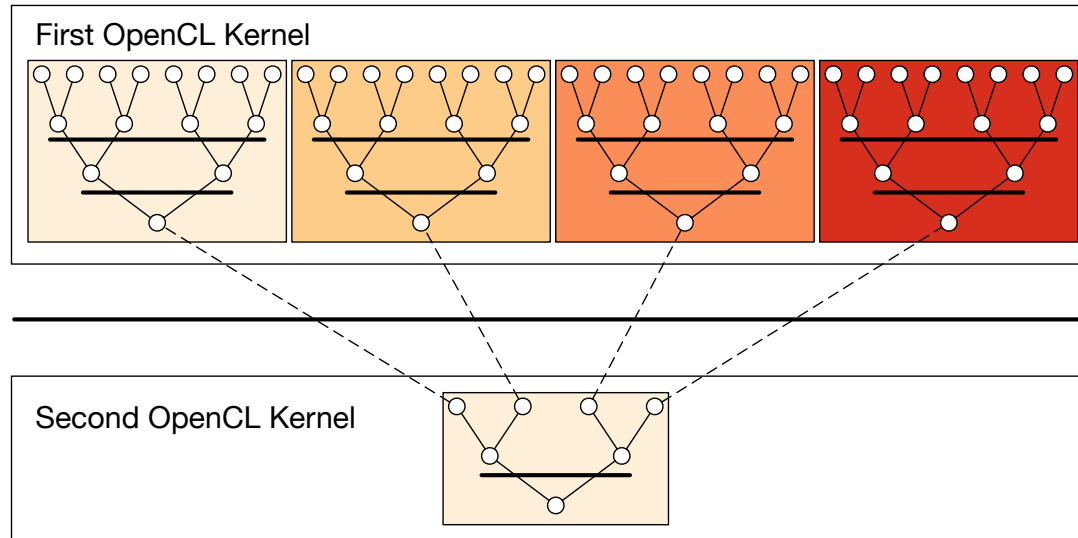


Part II: Addressing the Performance Portability Challenge



Case Study: Parallel Reduction in OpenCL

- Summing up all values of an array
- Comparison of 7 implementations by Nvidia
- Investigating complexity and efficiency of optimisations



Unoptimised Implementation Parallel Reduction

```
kernel void reduce0(global float* g_idata, global float* g_odata,
                   unsigned int n, local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i   = get_global_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    barrier(CLK_LOCAL_MEM_FENCE);
    // do reduction in local memory
    for (unsigned int s=1; s < get_local_size(0); s*= 2) {
        if ((tid % (2*s)) == 0) {
            l_data[tid] += l_data[tid + s];
            barrier(CLK_LOCAL_MEM_FENCE);
        }
    }
    // write result for this work-group to global memory
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0];
}
```

Avoid Divergent Branching

```
kernel void reduce1(global float* g_idata, global float* g_odata,
                    unsigned int n, local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i    = get_global_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    barrier(CLK_LOCAL_MEM_FENCE);

    for (unsigned int s=1; s < get_local_size(0); s*= 2) {
        // continuous work-items remain active
        int index = 2 * s * tid;
        if (index < get_local_size(0)) {
            l_data[index] += l_data[index + s];
        }
        barrier(CLK_LOCAL_MEM_FENCE);
    }
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0];
}
```

Avoid Interleaved Addressing

```
kernel void reduce2(global float* g_idata, global float* g_odata,
                   unsigned int n, local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i   = get_global_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    barrier(CLK_LOCAL_MEM_FENCE);

    // process elements in different order
    // requires commutativity
    for (unsigned int s=get_local_size(0)/2; s>0; s>>=1) {
        if (tid < s) {
            l_data[tid] += l_data[tid + s];
        }
        barrier(CLK_LOCAL_MEM_FENCE);
    }
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0];
}
```

Increase Computational Intensity per Work-Item

```
kernel void reduce3(global float* g_idata, global float* g_odata,
                   unsigned int n, local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i = get_group_id(0) * (get_local_size(0)*2)
                    + get_local_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    // performs first addition during loading
    if (i + get_local_size(0) < n)
        l_data[tid] += g_idata[i+get_local_size(0)];
    barrier(CLK_LOCAL_MEM_FENCE);

    for (unsigned int s=get_local_size(0)/2; s>0; s>>=1) {
        if (tid < s) {
            l_data[tid] += l_data[tid + s];
        }
        barrier(CLK_LOCAL_MEM_FENCE);
    }
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0];
}
```


Avoid Synchronisation inside a Warp

```
kernel void reduce4(global float* g_idata, global float* g_odata,
                   unsigned int n, local volatile float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i = get_group_id(0) * (get_local_size(0)*2)
                   + get_local_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    if (i + get_local_size(0) < n)
        l_data[tid] += g_idata[i+get_local_size(0)];
    barrier(CLK_LOCAL_MEM_FENCE);

    #pragma unroll 1
    for (unsigned int s=get_local_size(0)/2; s>32; s>>=1) {
        if (tid < s) { l_data[tid] += l_data[tid + s]; }
        barrier(CLK_LOCAL_MEM_FENCE); }

    // this is not portable OpenCL code!
    if (tid < 32) {
        if (WG_SIZE >= 64) { l_data[tid] += l_data[tid+32]; }
        if (WG_SIZE >= 32) { l_data[tid] += l_data[tid+16]; }
        if (WG_SIZE >= 16) { l_data[tid] += l_data[tid+ 8]; }
        if (WG_SIZE >=  8) { l_data[tid] += l_data[tid+ 4]; }
        if (WG_SIZE >=  4) { l_data[tid] += l_data[tid+ 2]; }
        if (WG_SIZE >=  2) { l_data[tid] += l_data[tid+ 1]; } }
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0]; }
```

Complete Loop Unrolling

```
kernel void reduce5(global float* g_idata, global float* g_odata,
                    unsigned int n, local volatile float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i = get_group_id(0) * (get_local_size(0)*2)
                    + get_local_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    if (i + get_local_size(0) < n)
        l_data[tid] += g_idata[i+get_local_size(0)];
    barrier(CLK_LOCAL_MEM_FENCE);

    if (WG_SIZE >= 256) {
        if (tid < 128) { l_data[tid] += l_data[tid+128]; }
        barrier(CLK_LOCAL_MEM_FENCE); }
    if (WG_SIZE >= 128) {
        if (tid < 64) { l_data[tid] += l_data[tid+ 64]; }
        barrier(CLK_LOCAL_MEM_FENCE); }
    if (tid < 32) {
        if (WG_SIZE >= 64) { l_data[tid] += l_data[tid+32]; }
        if (WG_SIZE >= 32) { l_data[tid] += l_data[tid+16]; }
        if (WG_SIZE >= 16) { l_data[tid] += l_data[tid+ 8]; }
        if (WG_SIZE >= 8)  { l_data[tid] += l_data[tid+ 4]; }
        if (WG_SIZE >= 4)  { l_data[tid] += l_data[tid+ 2]; }
        if (WG_SIZE >= 2)  { l_data[tid] += l_data[tid+ 1]; } }
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0]; }
```

Fully Optimised Implementation

```
kernel void reduce6(global float* g_idata, global float* g_odata,
                   unsigned int n, local volatile float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i = get_group_id(0) * (get_local_size(0)*2)
                   + get_local_id(0);
    unsigned int gridSize = WG_SIZE * get_num_groups(0);
    l_data[tid] = 0;
    while (i < n) { l_data[tid] += g_idata[i];
                   if (i + WG_SIZE < n)
                       l_data[tid] += g_idata[i+WG_SIZE];
                   i += gridSize; }
    barrier(CLK_LOCAL_MEM_FENCE);

    if (WG_SIZE >= 256) {
        if (tid < 128) { l_data[tid] += l_data[tid+128]; }
        barrier(CLK_LOCAL_MEM_FENCE); }
    if (WG_SIZE >= 128) {
        if (tid < 64) { l_data[tid] += l_data[tid+ 64]; }
        barrier(CLK_LOCAL_MEM_FENCE); }
    if (tid < 32) {
        if (WG_SIZE >= 64) { l_data[tid] += l_data[tid+32]; }
        if (WG_SIZE >= 32) { l_data[tid] += l_data[tid+16]; }
        if (WG_SIZE >= 16) { l_data[tid] += l_data[tid+ 8]; }
        if (WG_SIZE >=  8) { l_data[tid] += l_data[tid+ 4]; }
        if (WG_SIZE >=  4) { l_data[tid] += l_data[tid+ 2]; }
        if (WG_SIZE >=  2) { l_data[tid] += l_data[tid+ 1]; } }
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0]; }
```

Case Study Conclusions

- Optimising OpenCL is complex
 - Understanding of target hardware required
- Program changes not obvious
- Is it worth it? ...

```
kernel
void reduce0(global float* g_idata,
             global float* g_odata,
             unsigned int n,
             local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i = get_global_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    barrier(CLK_LOCAL_MEM_FENCE);

    for (unsigned int s=1;
         s < get_local_size(0); s*= 2) {
        if ((tid % (2*s)) == 0) {
            l_data[tid] += l_data[tid + s];
        }
        barrier(CLK_LOCAL_MEM_FENCE);
    }
    if (tid == 0)
        g_odata[get_group_id(0)] = l_data[0];
}
```

Unoptimized Implementation

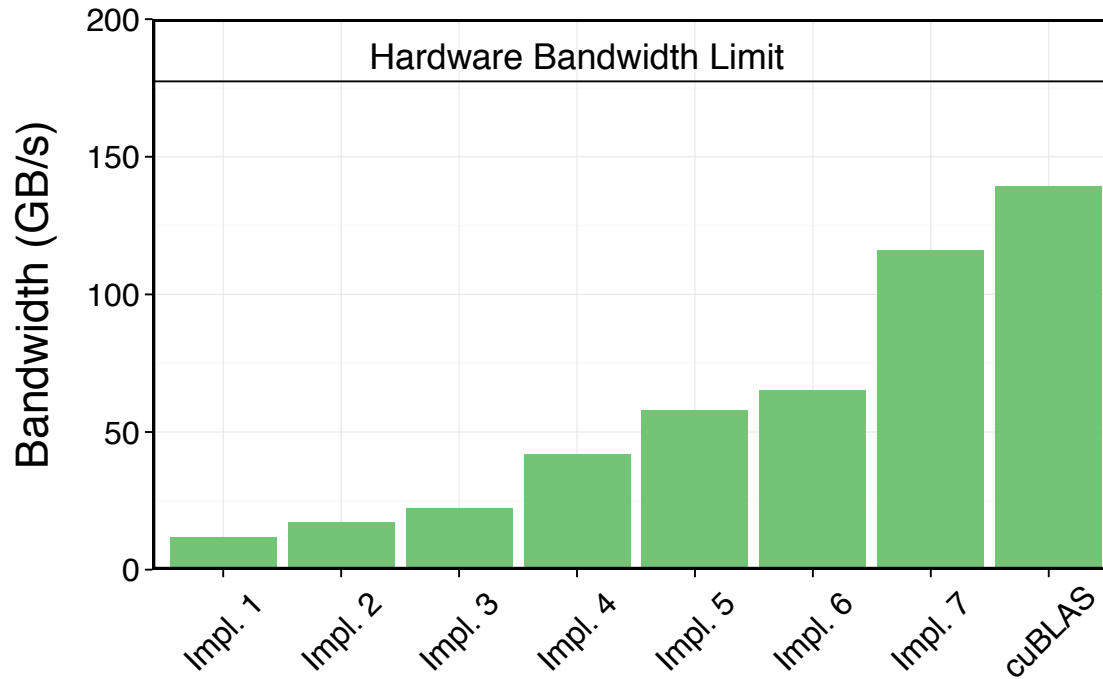
```
kernel
void reduce6(global float* g_idata,
             global float* g_odata,
             unsigned int n,
             local volatile float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i =
        get_group_id(0) * (get_local_size(0)*2)
        + get_local_id(0);

    unsigned int gridSize =
        WG_SIZE * get_num_groups(0);
    l_data[tid] = 0;
    while (i < n) {
        l_data[tid] += g_idata[i];
        if (i + WG_SIZE < n)
            l_data[tid] += g_idata[i+WG_SIZE];
        i += gridSize; }
    barrier(CLK_LOCAL_MEM_FENCE);

    if (WG_SIZE >= 256) {
        if (tid < 128) {
            l_data[tid] += l_data[tid+128]; }
        barrier(CLK_LOCAL_MEM_FENCE); }
    if (WG_SIZE >= 128) {
        if (tid < 64) {
            l_data[tid] += l_data[tid+ 64]; }
        barrier(CLK_LOCAL_MEM_FENCE); }
    if (tid < 32) {
        if (WG_SIZE >= 64) {
            l_data[tid] += l_data[tid+32]; }
        if (WG_SIZE >= 32) {
            l_data[tid] += l_data[tid+16]; }
        if (WG_SIZE >= 16) {
            l_data[tid] += l_data[tid+ 8]; }
        if (WG_SIZE >= 8) {
            l_data[tid] += l_data[tid+ 4]; }
        if (WG_SIZE >= 4) {
            l_data[tid] += l_data[tid+ 2]; }
        if (WG_SIZE >= 2) {
            l_data[tid] += l_data[tid+ 1]; } }
    if (tid == 0)
        g_odata[get_group_id(0)] = l_data[0];
}
```

Fully Optimized Implementation

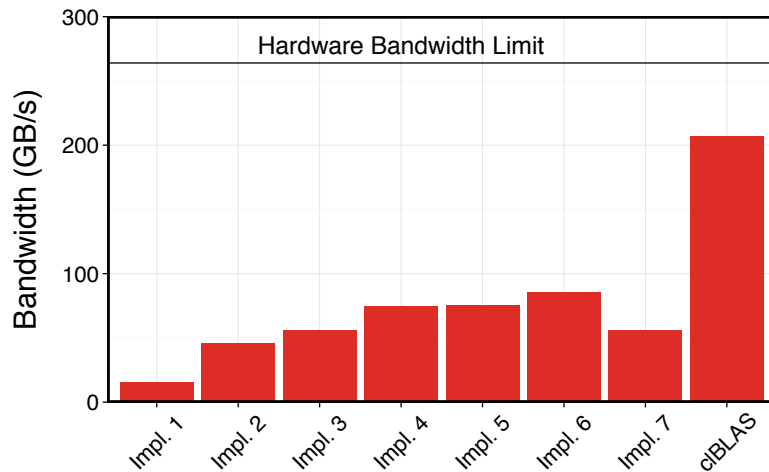
Performance Results Nvidia



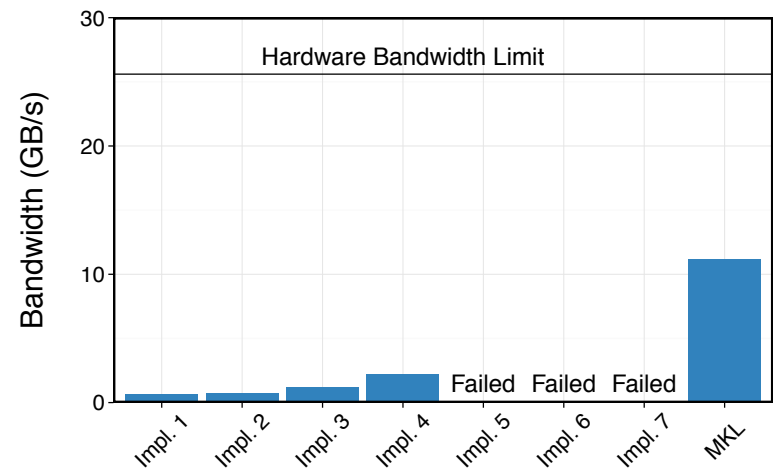
(a) Nvidia's GTX 480 GPU.

- ... Yes! Optimising improves performance by a factor of 10!
- Optimising is important, but ...

Performance Results AMD and Intel



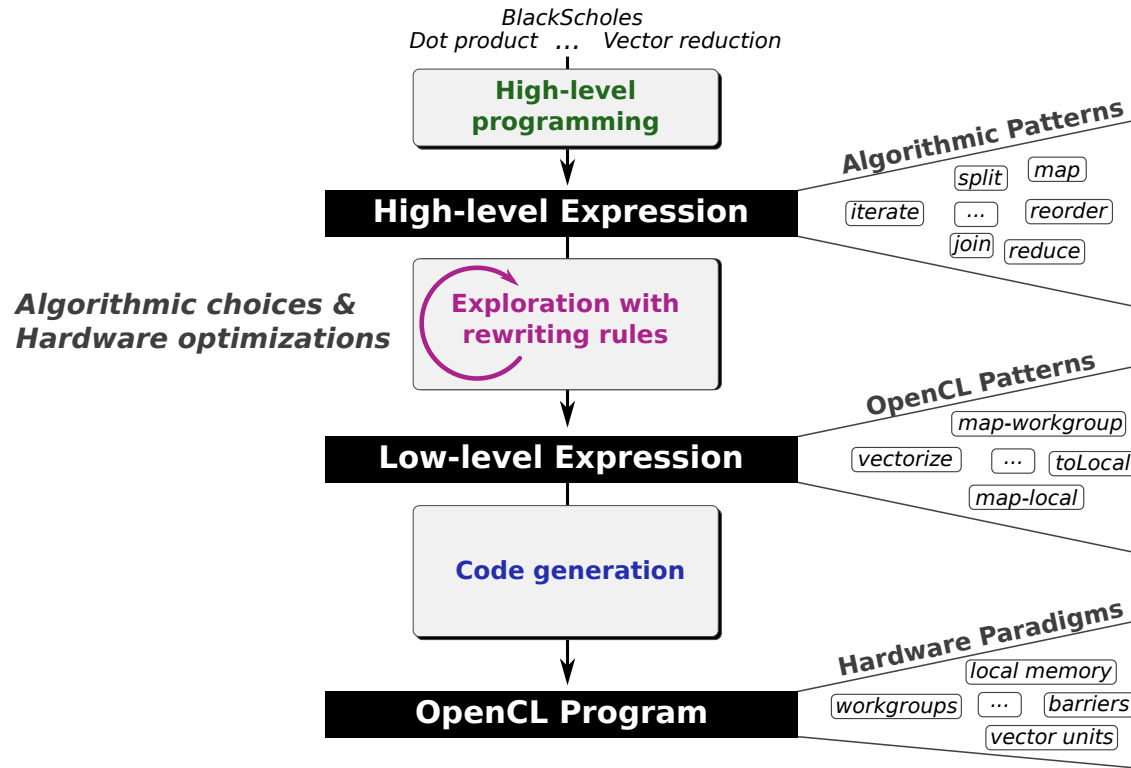
(b) AMD's HD 7970 GPU.



(c) Intel's E5530 dual-socket CPU.

- ... unfortunately, optimisations in OpenCL are not portable!
- **Challenge:** how to achieving portable performance?

Generating Performance Portable Code using Rewrite Rules



- **Ambition:** automatic generation of *Performance Portable* code

Walkthrough

① $\text{sum}(\text{vec}) = \text{reduce}(+, 0, \text{vec})$

↓
rewrite rules

↗
code generation

②

```
vecSum = reduce ◦ join ◦ map-workgroup (  
  join ◦ toGlobal (map-local (map-seq id)) ◦ split 1 ◦  
  join ◦ map-warp (  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 1 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 2 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 4 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 8 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 16 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 32  
  ) ◦ split 64 ◦  
  join ◦ map-local (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 64 ◦  
  join ◦ toLocal (map-local (reduce-seq (+) 0)) ◦  
  split (blockSize/128) ◦ reorder-stride 128  
) ◦ split blockSize
```

③

```
kernel  
void reduce6(global float* g_idata,  
             global float* g_odata,  
             unsigned int n,  
             local volatile float* l_data) {  
  unsigned int tid = get_local_id(0);  
  unsigned int i =  
    get_group_id(0) * (get_local_size(0)*2)  
    + get_local_id(0);  
  unsigned int gridSize =  
    WG_SIZE * get_num_groups(0);  
  l_data[tid] = 0;  
  while (i < n) {  
    l_data[tid] += g_idata[i];  
    if (i + WG_SIZE < n)  
      l_data[tid] += g_idata[i+WG_SIZE];  
    i += gridSize; }  
  barrier(CLK_LOCAL_MEM_FENCE);  
  
  if (WG_SIZE >= 256) {  
    if (tid < 128) {  
      l_data[tid] += l_data[tid+128]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (WG_SIZE >= 128) {  
    if (tid < 64) {  
      l_data[tid] += l_data[tid+ 64]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (tid < 32) {  
    if (WG_SIZE >= 64) {  
      l_data[tid] += l_data[tid+32]; }  
    if (WG_SIZE >= 32) {  
      l_data[tid] += l_data[tid+16]; }  
    if (WG_SIZE >= 16) {  
      l_data[tid] += l_data[tid+ 8]; }  
    if (WG_SIZE >= 8) {  
      l_data[tid] += l_data[tid+ 4]; }  
    if (WG_SIZE >= 4) {  
      l_data[tid] += l_data[tid+ 2]; }  
    if (WG_SIZE >= 2) {  
      l_data[tid] += l_data[tid+ 1]; } }  
  if (tid == 0)  
    g_odata[get_group_id(0)] = l_data[0];  
}
```



Walkthrough

① $\text{sum}(\text{vec}) = \text{reduce}(+, 0, \text{vec})$

rewrite rules

code generation

②

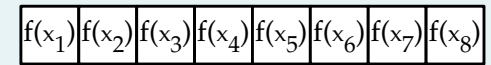
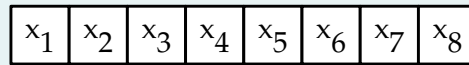
```
vecSum = reduce ◦ join ◦ map-workgroup (  
  join ◦ toGlobal (map-local (map-seq id)) ◦ split 1 ◦  
  join ◦ map-warp (  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 1 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 2 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 4 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 8 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 16 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 32  
  ) ◦ split 64 ◦  
  join ◦ map-local (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 64 ◦  
  join ◦ toLocal (map-local (reduce-seq (+) 0)) ◦  
  split (blockSize/128) ◦ reorder-stride 128  
) ◦ split blockSize
```

③

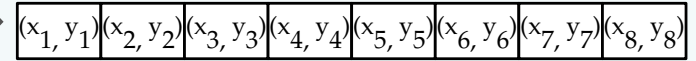
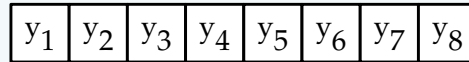
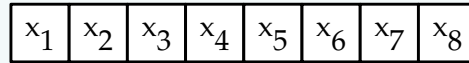
```
kernel  
void reduce6(global float* g_idata,  
             global float* g_odata,  
             unsigned int n,  
             local volatile float* l_data) {  
  unsigned int tid = get_local_id(0);  
  unsigned int i =  
    get_group_id(0) * (get_local_size(0)*2)  
    + get_local_id(0);  
  
  unsigned int gridSize =  
    WG_SIZE * get_num_groups(0);  
  l_data[tid] = 0;  
  while (i < n) {  
    l_data[tid] += g_idata[i];  
    if (i + WG_SIZE < n)  
      l_data[tid] += g_idata[i+WG_SIZE];  
    i += gridSize; }  
  barrier(CLK_LOCAL_MEM_FENCE);  
  
  if (WG_SIZE >= 256) {  
    if (tid < 128) {  
      l_data[tid] += l_data[tid+128]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (WG_SIZE >= 128) {  
    if (tid < 64) {  
      l_data[tid] += l_data[tid+ 64]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (tid < 32) {  
    if (WG_SIZE >= 64) {  
      l_data[tid] += l_data[tid+32]; }  
    if (WG_SIZE >= 32) {  
      l_data[tid] += l_data[tid+16]; }  
    if (WG_SIZE >= 16) {  
      l_data[tid] += l_data[tid+ 8]; }  
    if (WG_SIZE >= 8) {  
      l_data[tid] += l_data[tid+ 4]; }  
    if (WG_SIZE >= 4) {  
      l_data[tid] += l_data[tid+ 2]; }  
    if (WG_SIZE >= 2) {  
      l_data[tid] += l_data[tid+ 1]; } }  
  if (tid == 0)  
    g_odata[get_group_id(0)] = l_data[0];  
}
```

① Algorithmic Primitives (a.k.a. algorithmic skeletons)

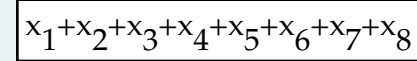
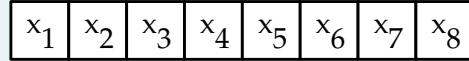
map(f, x):



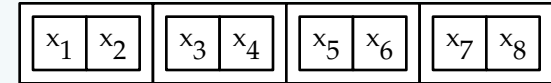
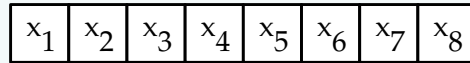
zip(x, y):



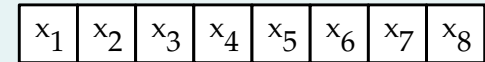
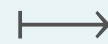
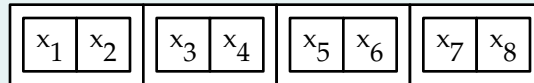
reduce($+, 0, x$):



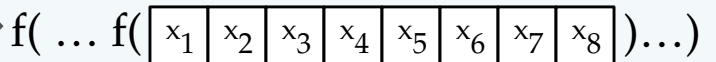
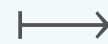
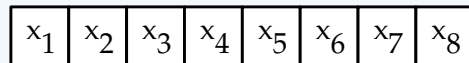
split(n, x):



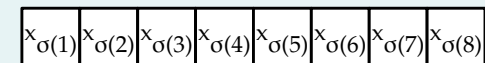
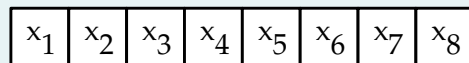
join(x):



iterate(f, n, x):



reorder(σ, x):



① High-Level Programs

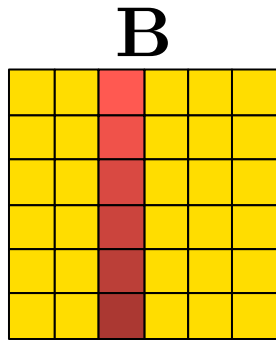
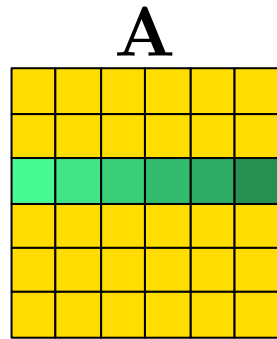
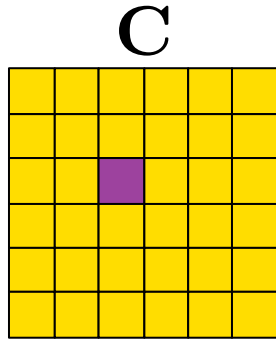
`scal(a, vec) = map($\lambda x \mapsto x*a$, vec)`

`asum(vec) = reduce(+, 0, map(abs, vec))`

`dotProduct(x, y) = reduce(+, 0, map(*, zip(x, y)))`

`gemv(mat, x, y, α , β) =
 map(+, zip(
 map(λ row \mapsto scal(α , dotProduct(row, x)), mat),
 scal(β , y)))`

① High-Level Programs



$A \times B =$
`map(λ rowA \mapsto
 map(λ colB \mapsto
 dotProduct(rowA, colB)
 , transpose(B))
 , A)`

Walkthrough

① $\text{sum}(\text{vec}) = \text{reduce}(+, 0, \text{vec})$

↓
rewrite rules

↗
code generation

②

```
vecSum = reduce ◦ join ◦ map-workgroup (  
  join ◦ toGlobal (map-local (map-seq id)) ◦ split 1 ◦  
  join ◦ map-warp (  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 1 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 2 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 4 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 8 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 16 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 32  
  ) ◦ split 64 ◦  
  join ◦ map-local (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 64 ◦  
  join ◦ toLocal (map-local (reduce-seq (+) 0)) ◦  
  split (blockSize/128) ◦ reorder-stride 128  
) ◦ split blockSize
```

③

```
kernel  
void reduce6(global float* g_idata,  
             global float* g_odata,  
             unsigned int n,  
             local volatile float* l_data) {  
  unsigned int tid = get_local_id(0);  
  unsigned int i =  
    get_group_id(0) * (get_local_size(0)*2)  
    + get_local_id(0);  
  unsigned int gridSize =  
    WG_SIZE * get_num_groups(0);  
  l_data[tid] = 0;  
  while (i < n) {  
    l_data[tid] += g_idata[i];  
    if (i + WG_SIZE < n)  
      l_data[tid] += g_idata[i+WG_SIZE];  
    i += gridSize; }  
  barrier(CLK_LOCAL_MEM_FENCE);  
  
  if (WG_SIZE >= 256) {  
    if (tid < 128) {  
      l_data[tid] += l_data[tid+128]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (WG_SIZE >= 128) {  
    if (tid < 64) {  
      l_data[tid] += l_data[tid+ 64]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (tid < 32) {  
    if (WG_SIZE >= 64) {  
      l_data[tid] += l_data[tid+32]; }  
    if (WG_SIZE >= 32) {  
      l_data[tid] += l_data[tid+16]; }  
    if (WG_SIZE >= 16) {  
      l_data[tid] += l_data[tid+ 8]; }  
    if (WG_SIZE >= 8) {  
      l_data[tid] += l_data[tid+ 4]; }  
    if (WG_SIZE >= 4) {  
      l_data[tid] += l_data[tid+ 2]; }  
    if (WG_SIZE >= 2) {  
      l_data[tid] += l_data[tid+ 1]; } }  
  if (tid == 0)  
    g_odata[get_group_id(0)] = l_data[0];  
}
```



Walkthrough

① $\text{sum}(\text{vec}) = \text{reduce}(+, 0, \text{vec})$

↓
rewrite rules

↗
code generation

②

```
vecSum = reduce ◦ join ◦ map-workgroup (  
  join ◦ toGlobal (map-local (map-seq id)) ◦ split 1 ◦  
  join ◦ map-warp (  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 1 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 2 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 4 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 8 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 16 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 32  
  ) ◦ split 64 ◦  
  join ◦ map-local (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 64 ◦  
  join ◦ toLocal (map-local (reduce-seq (+) 0)) ◦  
  split (blockSize/128) ◦ reorder-stride 128  
) ◦ split blockSize
```

③

```
kernel  
void reduce6(global float* g_idata,  
            global float* g_odata,  
            unsigned int n,  
            local volatile float* l_data) {  
  unsigned int tid = get_local_id(0);  
  unsigned int i =  
    get_group_id(0) * (get_local_size(0)*2)  
    + get_local_id(0);  
  unsigned int gridSize =  
    WG_SIZE * get_num_groups(0);  
  l_data[tid] = 0;  
  while (i < n) {  
    l_data[tid] += g_idata[i];  
    if (i + WG_SIZE < n)  
      l_data[tid] += g_idata[i+WG_SIZE];  
    i += gridSize; }  
  barrier(CLK_LOCAL_MEM_FENCE);  
  
  if (WG_SIZE >= 256) {  
    if (tid < 128) {  
      l_data[tid] += l_data[tid+128]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (WG_SIZE >= 128) {  
    if (tid < 64) {  
      l_data[tid] += l_data[tid+ 64]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (tid < 32) {  
    if (WG_SIZE >= 64) {  
      l_data[tid] += l_data[tid+32]; }  
    if (WG_SIZE >= 32) {  
      l_data[tid] += l_data[tid+16]; }  
    if (WG_SIZE >= 16) {  
      l_data[tid] += l_data[tid+ 8]; }  
    if (WG_SIZE >= 8) {  
      l_data[tid] += l_data[tid+ 4]; }  
    if (WG_SIZE >= 4) {  
      l_data[tid] += l_data[tid+ 2]; }  
    if (WG_SIZE >= 2) {  
      l_data[tid] += l_data[tid+ 1]; } }  
  if (tid == 0)  
    g_odata[get_group_id(0)] = l_data[0];  
}
```

② Algorithmic Rewrite Rules

- **Provably correct** rewrite rules
- Express algorithmic implementation choices

Split-join rule:

$$\text{map } f \rightarrow \text{join} \circ \text{map } (\text{map } f) \circ \text{split } n$$

Map fusion rule:

$$\text{map } f \circ \text{map } g \rightarrow \text{map } (f \circ g)$$

Reduce rules:

$$\text{reduce } f \ z \rightarrow \text{reduce } f \ z \circ \text{reducePart } f \ z$$

$$\text{reducePart } f \ z \rightarrow \text{reducePart } f \ z \circ \text{reorder}$$

$$\text{reducePart } f \ z \rightarrow \text{join} \circ \text{map } (\text{reducePart } f \ z) \circ \text{split } n$$

$$\text{reducePart } f \ z \rightarrow \text{iterate } n \ (\text{reducePart } f \ z)$$

② OpenCL Primitives

Primitive

mapGlobal

mapWorkgroup

mapLocal

mapSeq

reduceSeq

toLocal , *toGlobal*

mapVec ,

splitVec , *joinVec*

OpenCL concept

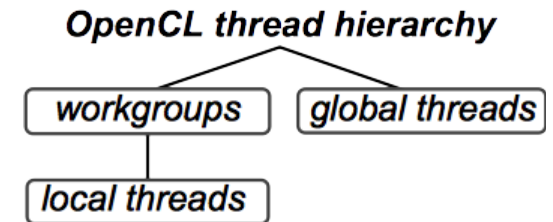
Work-items

Work-groups

Sequential implementations

Memory areas

Vectorisation



② OpenCL Rewrite Rules

- Express low-level implementation and optimisation choices

Map rules:

$$\text{map } f \rightarrow \text{mapWorkgroup } f \mid \text{mapLocal } f \mid \text{mapGlobal } f \mid \text{mapSeq } f$$

Local/ global memory rules:

$$\text{mapLocal } f \rightarrow \text{toLocal } (\text{mapLocal } f) \qquad \text{mapLocal } f \rightarrow \text{toGlobal } (\text{mapLocal } f)$$

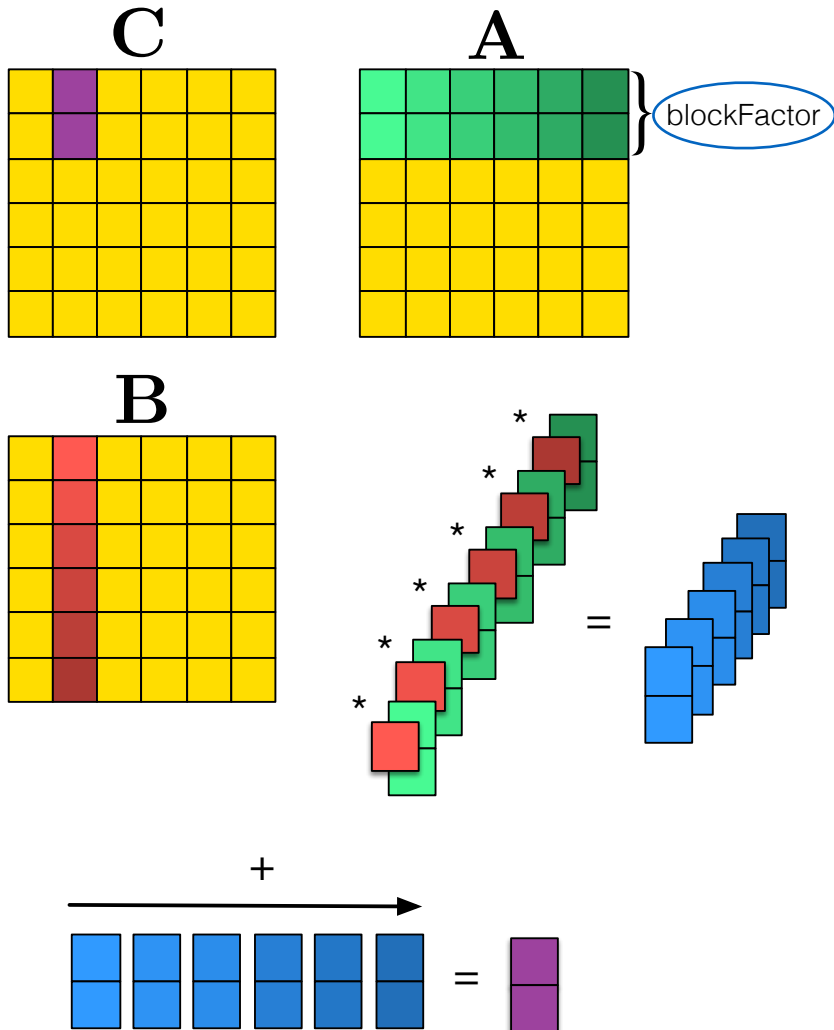
Vectorisation rule:

$$\text{map } f \rightarrow \text{joinVec} \circ \text{map } (\text{mapVec } f) \circ \text{splitVec } n$$

Fusion rule:

$$\text{reduceSeq } f \ z \circ \text{mapSeq } g \rightarrow \text{reduceSeq } (\lambda (acc, x). f (acc, g x)) \ z$$

② Optimisation Example: Register Blocking



```

1 kernel void KERNEL(
2   const global float* restrict A,
3   const global float* restrict B,
4   global float* C, int K, int M, int N)
5 {
6   float acc(blockFactor);
7
8   for (int glb_id_1 = get_global_id(1);
9        glb_id_1 < M / blockFactor;
10       glb_id_1 += get_global_size(1)) {
11     for (int glb_id_0 = get_global_id(0); glb_id_0 < N;
12         glb_id_0 += get_global_size(0)) {
13
14       for (int i = 0; i < K; i += 1)
15         float temp = B[i * N + glb_id_0];
16       for (int j = 0; j < blockFactor; j += 1)
17         acc[j] +=
18           A[blockFactor * glb_id_1 * K + j * K + i]
19           * temp;
20
21       for (int j = 0; j < blockFactor; j += 1)
22         C[blockFactor * glb_id_1 * N + j * N + glb_id_0]
23         = acc[j];
24     }
25 }
26 }

```

② Register Blocking as a Macro Rule

- Optimisations are expressed as *Macro Rules*:
 - Series of Rewrites applied to achieve an optimisation goal

registerBlocking =

$$\text{Map}(f) \Rightarrow \text{Join}() \circ \text{Map}(\text{Map}(f)) \circ \text{Split}(k)$$

$$\text{Map}(a \mapsto \text{Map}(b \mapsto f(a, b))) \Rightarrow \text{Transpose}() \circ \text{Map}(b \mapsto \text{Map}(a \mapsto f(a, b)))$$

$$\text{Map}(f \circ g) \Rightarrow \text{Map}(f) \circ \text{Map}(g)$$

$$\text{Map}(\text{Reduce}(f)) \Rightarrow \text{Transpose}() \circ \text{Reduce}((acc, x) \mapsto \text{Map}(f) \circ \text{Zip}(acc, x))$$

$$\text{Map}(\text{Map}(f)) \Rightarrow \text{Transpose}() \circ \text{Map}(\text{Map}(f)) \circ \text{Transpose}()$$

$$\text{Transpose}() \circ \text{Transpose}() \Rightarrow id$$

$$\text{Reduce}(f) \circ \text{Map}(g) \Rightarrow \text{Reduce}((acc, x) \mapsto f(acc, g(x)))$$

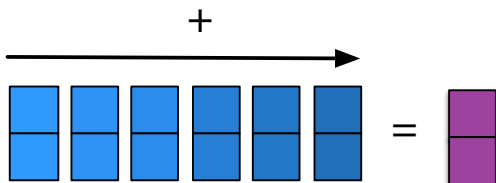
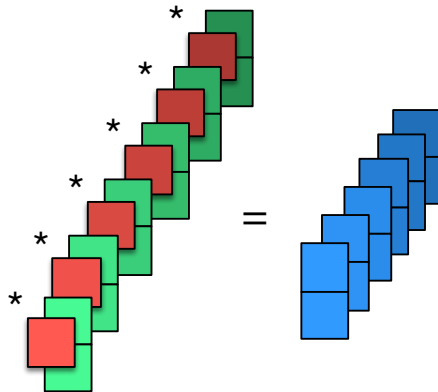
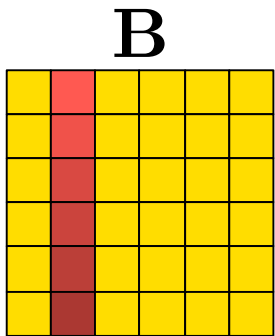
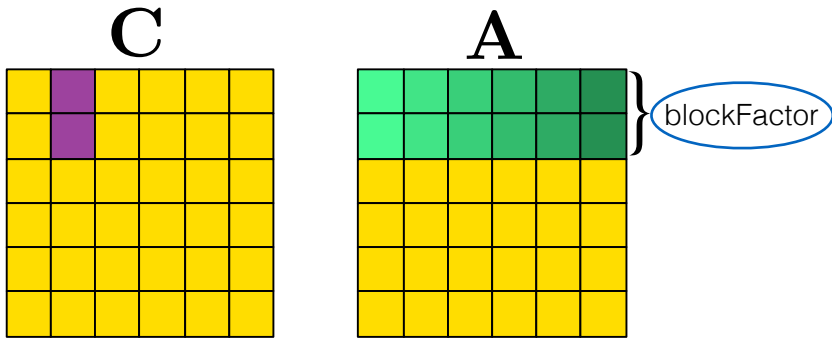
$$\text{Map}(f) \circ \text{Map}(g) \Rightarrow \text{Map}(f \circ g)$$

② Register Blocking as a Series of Rewrites

$$\begin{aligned}
 & \text{Map}(\overrightarrow{\text{rowA}} \mapsto \\
 & \quad \text{Map}(\overrightarrow{\text{colB}} \mapsto \\
 & \quad \quad \text{Reduce}(+) \circ \text{Map}(*)) \\
 & \quad \quad \$ \text{Zip}(\overrightarrow{\text{rowA}}, \overrightarrow{\text{colB}}) \\
 & \quad) \circ \text{Transpose}() \$ \mathbf{B} \\
 &) \$ \mathbf{A}
 \end{aligned}$$


$$\begin{aligned}
 & \text{Join}() \circ \text{Map}(\text{rowsA} \mapsto \\
 & \quad \text{Transpose}() \circ \text{Map}(\overrightarrow{\text{colB}} \mapsto \\
 & \quad \quad \text{Transpose}() \circ \text{Reduce}((\overrightarrow{\text{acc}}, \overrightarrow{\text{pair}}) \mapsto \\
 & \quad \quad \quad \text{Map}(x \mapsto x_0 + x_1 * \text{pair}._1) \\
 & \quad \quad \quad \quad \$ \text{Zip}(\overrightarrow{\text{acc}}, \text{pair}._0) \\
 & \quad \quad \quad) \$ \text{Zip}(\text{Transpose}() \$ \text{rowsA}, \overrightarrow{\text{colB}}) \\
 & \quad \quad) \circ \text{Transpose}() \$ \mathbf{B} \\
 &) \circ \text{Split}(\text{blockFactor}) \$ \mathbf{A}
 \end{aligned}$$

② Register Blocking Functionally Expressed



$$\begin{aligned}
 &Join() \circ Map(rowsA \mapsto \\
 &Transpose() \circ Map(\overrightarrow{colB} \mapsto \\
 &Transpose() \circ Reduce((\overrightarrow{acc}, \overrightarrow{pair}) \mapsto \\
 &Map(x \mapsto x._0 + x._1 * pair._1) \\
 &\$ Zip(\overrightarrow{acc}, pair._0) \\
 &)\$ Zip(Transpose() \$ rowsA, \overrightarrow{colB}) \\
 &)\circ Transpose() \$ B \\
 &)\circ Split(blockFactor) \$ A
 \end{aligned}$$

Walkthrough

① $vecSum = reduce (+) 0$

↓
rewrite rules

↗
code generation

②

```
vecSum = reduce ◦ join ◦ map-workgroup (  
  join ◦ toGlobal (map-local (map-seq id)) ◦ split 1 ◦  
  join ◦ map-warp (  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 1 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 2 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 4 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 8 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 16 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 32  
  ) ◦ split 64 ◦  
  join ◦ map-local (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 64 ◦  
  join ◦ toLocal (map-local (reduce-seq (+) 0)) ◦  
  split (blockSize/128) ◦ reorder-stride 128  
) ◦ split blockSize
```

③

```
kernel  
void reduce6(global float* g_idata,  
             global float* g_odata,  
             unsigned int n,  
             local volatile float* l_data) {  
  unsigned int tid = get_local_id(0);  
  unsigned int i =  
    get_group_id(0) * (get_local_size(0)*2)  
    + get_local_id(0);  
  unsigned int gridSize =  
    WG_SIZE * get_num_groups(0);  
  l_data[tid] = 0;  
  while (i < n) {  
    l_data[tid] += g_idata[i];  
    if (i + WG_SIZE < n)  
      l_data[tid] += g_idata[i+WG_SIZE];  
    i += gridSize; }  
  barrier(CLK_LOCAL_MEM_FENCE);  
  
  if (WG_SIZE >= 256) {  
    if (tid < 128) {  
      l_data[tid] += l_data[tid+128]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (WG_SIZE >= 128) {  
    if (tid < 64) {  
      l_data[tid] += l_data[tid+ 64]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (tid < 32) {  
    if (WG_SIZE >= 64) {  
      l_data[tid] += l_data[tid+32]; }  
    if (WG_SIZE >= 32) {  
      l_data[tid] += l_data[tid+16]; }  
    if (WG_SIZE >= 16) {  
      l_data[tid] += l_data[tid+ 8]; }  
    if (WG_SIZE >= 8) {  
      l_data[tid] += l_data[tid+ 4]; }  
    if (WG_SIZE >= 4) {  
      l_data[tid] += l_data[tid+ 2]; }  
    if (WG_SIZE >= 2) {  
      l_data[tid] += l_data[tid+ 1]; } }  
  if (tid == 0)  
    g_odata[get_group_id(0)] = l_data[0];  
}
```



Walkthrough

① $vecSum = reduce (+) 0$

rewrite rules

code generation

②

```
vecSum = reduce ◦ join ◦ map-workgroup (  
  join ◦ toGlobal (map-local (map-seq id)) ◦ split 1 ◦  
  join ◦ map-warp (  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 1 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 2 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 4 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 8 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 16 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 32  
  ) ◦ split 64 ◦  
  join ◦ map-local (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 64 ◦  
  join ◦ toLocal (map-local (reduce-seq (+) 0)) ◦  
  split (blockSize/128) ◦ reorder-stride 128  
) ◦ split blockSize
```

③

```
kernel  
void reduce6(global float* g_idata,  
             global float* g_odata,  
             unsigned int n,  
             local volatile float* l_data) {  
  unsigned int tid = get_local_id(0);  
  unsigned int i =  
    get_group_id(0) * (get_local_size(0)*2)  
    + get_local_id(0);  
  unsigned int gridSize =  
    WG_SIZE * get_num_groups(0);  
  l_data[tid] = 0;  
  while (i < n) {  
    l_data[tid] += g_idata[i];  
    if (i + WG_SIZE < n)  
      l_data[tid] += g_idata[i+WG_SIZE];  
    i += gridSize; }  
  barrier(CLK_LOCAL_MEM_FENCE);  
  
  if (WG_SIZE >= 256) {  
    if (tid < 128) {  
      l_data[tid] += l_data[tid+128]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (WG_SIZE >= 128) {  
    if (tid < 64) {  
      l_data[tid] += l_data[tid+ 64]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (tid < 32) {  
    if (WG_SIZE >= 64) {  
      l_data[tid] += l_data[tid+32]; }  
    if (WG_SIZE >= 32) {  
      l_data[tid] += l_data[tid+16]; }  
    if (WG_SIZE >= 16) {  
      l_data[tid] += l_data[tid+ 8]; }  
    if (WG_SIZE >= 8) {  
      l_data[tid] += l_data[tid+ 4]; }  
    if (WG_SIZE >= 4) {  
      l_data[tid] += l_data[tid+ 2]; }  
    if (WG_SIZE >= 2) {  
      l_data[tid] += l_data[tid+ 1]; } }  
  if (tid == 0)  
    g_odata[get_group_id(0)] = l_data[0];  
}
```



③ Pattern based OpenCL Code Generation

- Generate OpenCL code for each OpenCL primitive

mapGlobal f xs →

```
for (int g_id = get_global_id(0); g_id < n;  
     g_id += get_global_size(0)) {  
    output[g_id] = f(xs[g_id]);  
}
```

reduceSeq f z xs →

```
T acc = z;  
for (int i = 0; i < n; ++i) {  
    acc = f(acc, xs[i]);  
}
```

⋮

⋮

Rewrite rules define a space of possible implementations

$$\begin{array}{c} \textit{reduce (+) 0} \\ | \\ \textit{reduce (+) 0} \circ \textit{reducePart (+) 0} \end{array}$$



Rewrite rules define a space of possible implementations

reduce (+) 0

|

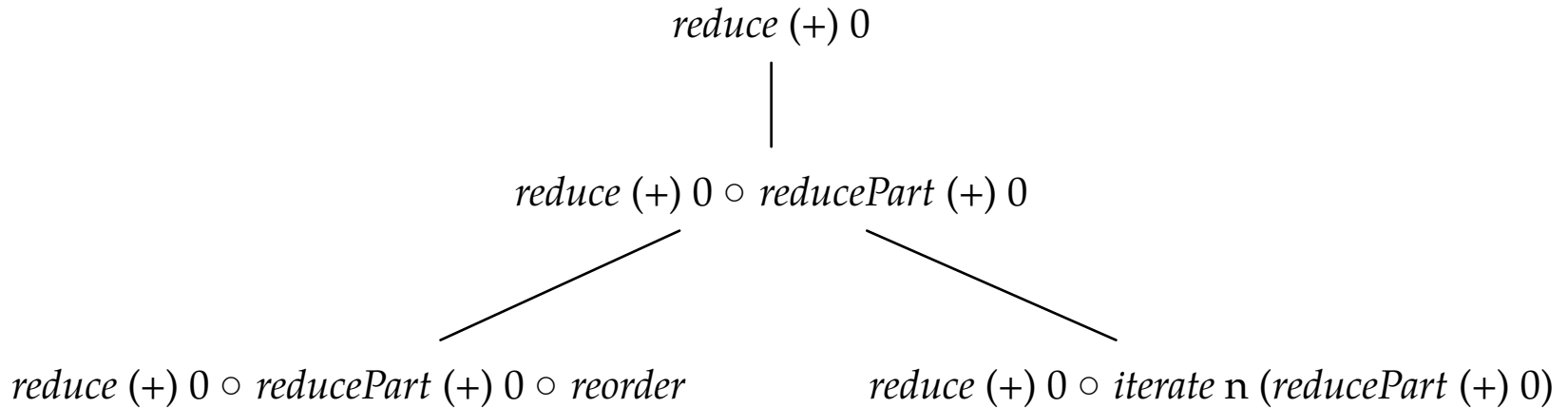
reduce (+) 0 ◦ *reducePart (+) 0*

↙

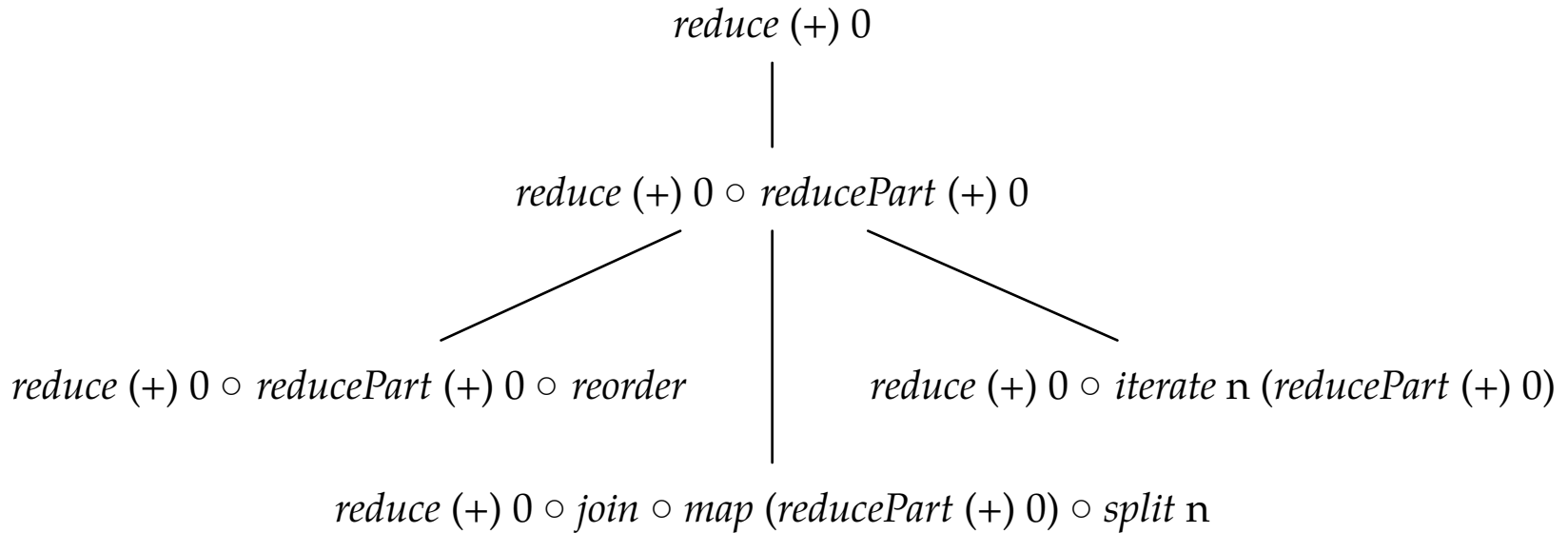
reduce (+) 0 ◦ *reducePart (+) 0* ◦ *reorder*



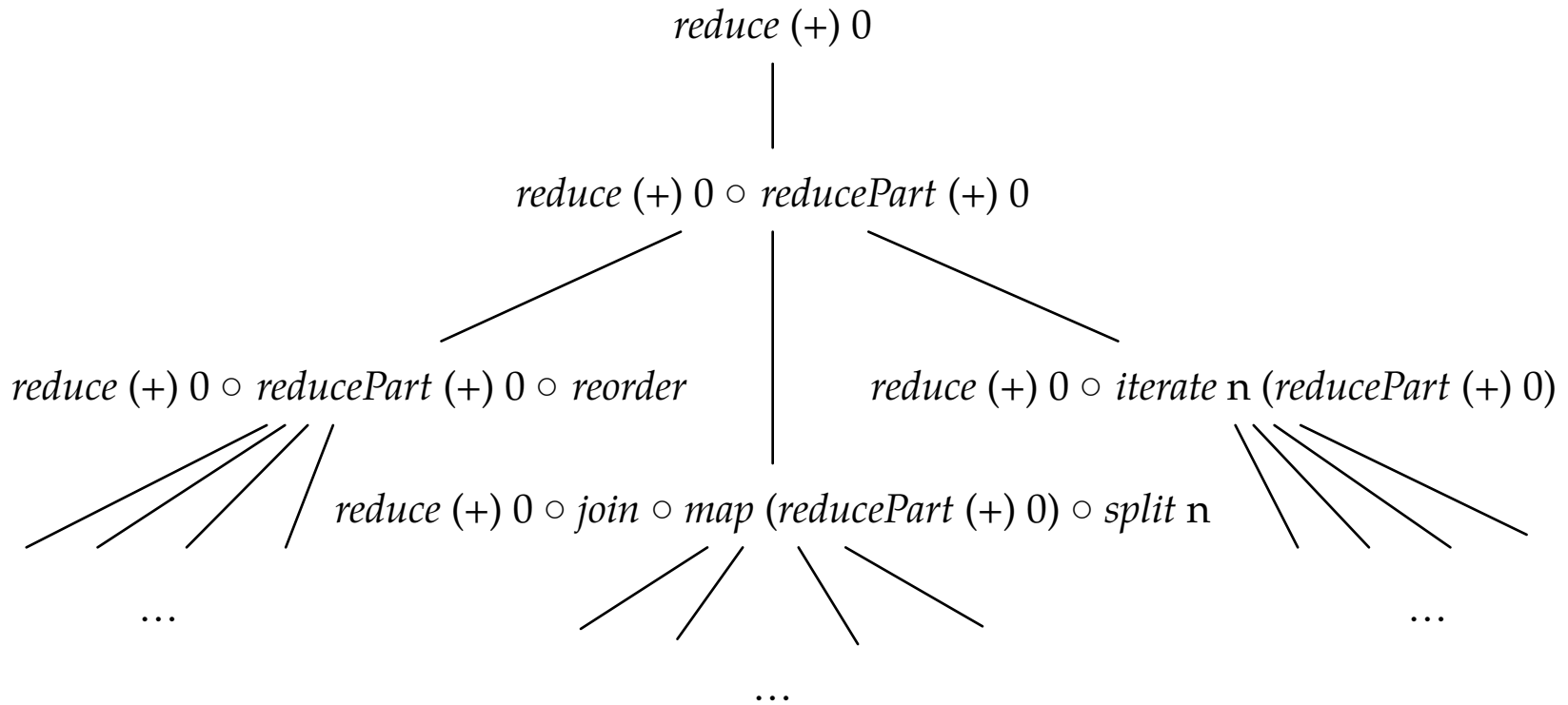
Rewrite rules define a space of possible implementations



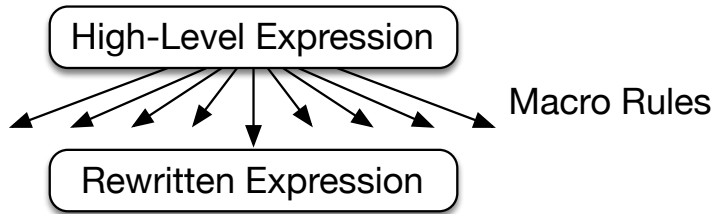
Rewrite rules define a space of possible implementations



Rewrite rules define a space of possible implementations



Exploration Strategy



1

$A * B =$

$Map(\overrightarrow{row A} \mapsto$

$Map(\overrightarrow{col B} \mapsto$

$DotProduct(\overrightarrow{row A}, \overrightarrow{col B}))$

$) \circ Transpose() \$ B$

$) \$ A$

1.1

```
TiledMultiply(A, B) =
  Untile() o
  Map(aRows ↦
    Map(bCols ↦
      Reduce((acc, pairOfTiles) ↦
        acc + pairOfTiles.0 * pairOfTiles.1
      ) $ Zip(aRows, bCols)
    ) o Transpose() o Tile(sizeN, sizeK) $ B
  ) o Tile(sizeM, sizeK) $ A
```

1.2

```
BlockedMultiply(A, B) =
  Join() o Map(Transpose()) o
  Map(aRows ↦
    Map(colB ↦
      Transpose() o
      Reduce((acc, rowElemPair) ↦
        acc + pairOfTiles.0 * pairOfTiles.1
      ) $ Zip(aRows, rowElemPair.1) $
      Zip(aRows, rowElemPair.0)
    ) o Transpose() $ B
  ) o Split(blockFactor) $ A
```

1.3

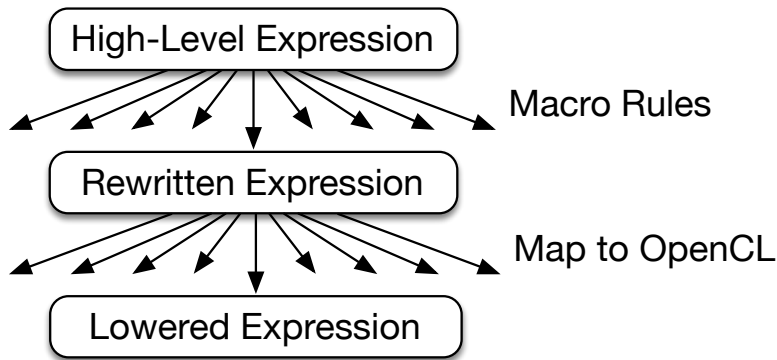
```
TiledMultiply(A, B) =
  Untile() o
  Map(aRows ↦
    Map(bCols ↦
      Reduce((acc, pairOfTiles) ↦
        acc + pairOfTiles.0 * pairOfTiles.1
      ) $ Zip(aRows, bCols)
    ) o Transpose() o Tile(sizeN, sizeK) $ B
  ) o Tile(sizeM, sizeK) $ A
```

1.4

```
BlockedMultiply(A, B) =
  Join() o Map(Transpose()) o
  Map(rowsA ↦
    Map(colB ↦
      Transpose() o
      Reduce((acc, rowElemPair) ↦
        acc + pairOfTiles.0 * pairOfTiles.1
      ) $ Zip(rowElemPair.0,
        rowElemPair.1) $
        Zip(rowsA, colB)
      ) o Transpose() $ B
    ) o Split(blockFactor) $ A
```



Exploration Strategy



1.3

TiledMultiply(A, B) =

Untile() ◦

1.3.1 *Map(aRows)* → 1.3.2 *Map(aRows)* → 1.3.3 *Map(aRows)*

1.3.1 *Map(bCols)* → 1.3.2 *Map(bCols)* → 1.3.3 *Map(bCols)*

1.3.1 *Reduce(acc, pairOfTiles)* → 1.3.2 *Reduce(acc, pairOfTiles)* → 1.3.3 *Reduce(acc, pairOfTiles)*

*acc + pairOfTiles..0 * pairOfTiles..1*

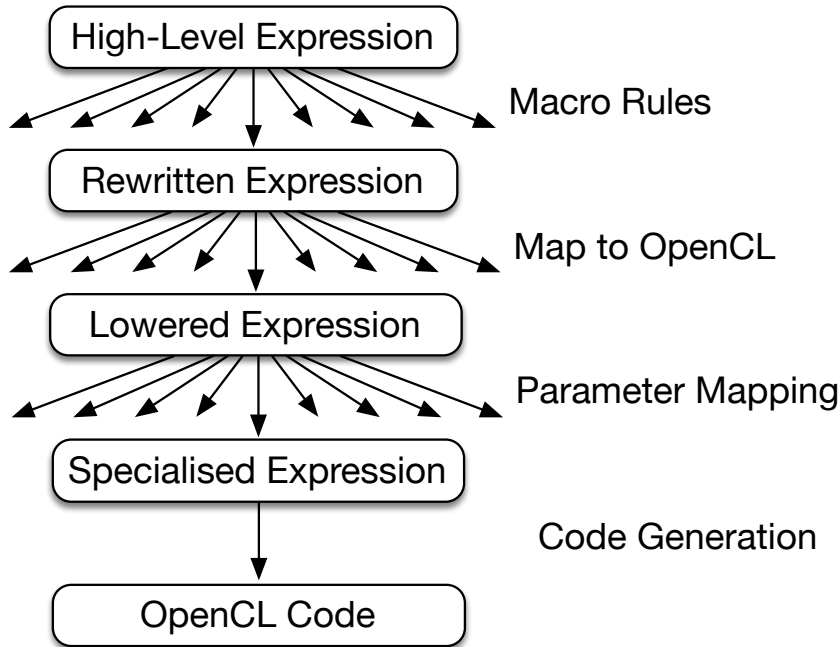
) \$ Zip(aRows, bCols)

) ◦ Transpose() ◦ Tile(sizeN, sizeK) \$ B

) ◦ Tile(sizeM, sizeK) \$ A



Exploration Strategy



1.3.2.5

```

1 kernel __min__(opt(global float *A, B, C,
2   TileMultiply(A, B) =
3   local float tileA[512]; tileB[512];
4
5   private float acc_0; ...; acc_31;
6   private float blockOfA_0; ...; blockOfA_7;
7   private float blockOfB_0; ...; blockOfB_3;
8
9   int lid0 = local_id(0); lid1 = local_id(1);
10  int w0 = group_id(0); w1 = group_id(1);
11  MapWrg(1)(aRows) ↦
12  for (int w1=wid1; w1<M/64; w1+=num_grps(1)) {
13    for (int w0=wid0; w0<N/64; w0+=num_grps(0)) {
14      MapWrg(0)(bCols) ↦
15      acc_0 = 0.0f; ...; acc_31 = 0.0f;
16      for (int i=0; i<K/8; i++) {
17        vstore4(vload4(lid1*M/4+2*i*M+16*w1+lid0,A), 16*lid1+lid0, tileA);
18        vstore4(vload4(lid1*N/4+2*i*N+16*w0+lid0,B), 16*lid1+lid0, tileB);
19        barrier (...);
20
21      }
22      ReduceSeq((acc, pairOfTiles) ↦
23      blockOfA_0 = tileA[0+64*i*8]; ...; blockOfA_7 = tileA[48+64*i*8];
24      blockOfB_0 = tileB[0+64*j+hd0]; ...; blockOfB_3 = tileB[48+64*j+hd0];
25
26      acc_0 += blockOfA_0 * blockOfB_0; ...; acc_31 += blockOfA_7 * blockOfB_3;
27      acc_1 += blockOfA_0 * blockOfB_1; ...; acc_29 += blockOfA_7 * blockOfB_3;
28      acc_2 += blockOfA_0 * blockOfB_2; ...; acc_30 += blockOfA_7 * blockOfB_2;
29      acc_3 += blockOfA_0 * blockOfB_3; ...; acc_31 += blockOfA_7 * blockOfB_3;
30    }
31    barrier (...);
32  }
33  C[0+8*lid1*N+64*w0+64*w1*N+0*N+lid0]=acc_0; ...; C[0+8*lid1*N+64*w0+64*w1*N+7*N+lid0]=acc_28;
34  C[16+8*lid1*N+64*w0+64*w1*N+0*N+lid0]=acc_1; ...; C[16+8*lid1*N+64*w0+64*w1*N+7*N+lid0]=acc_29;
35  C[32+8*lid1*N+64*w0+64*w1*N+0*N+lid0]=acc_2; ...; C[32+8*lid1*N+64*w0+64*w1*N+7*N+lid0]=acc_30;
36  C[48+8*lid1*N+64*w0+64*w1*N+0*N+lid0]=acc_3; ...; C[48+8*lid1*N+64*w0+64*w1*N+7*N+lid0]=acc_31;
37  } } }
  
```

Tile(128, 16) \$ B

Zip(aRows, bCols)

Tile(128, 16) \$ A

Heuristics for Matrix Multiplication

For Macro Rules:

- Nesting depth
- Distance of addition and multiplication
- Number of times rules are applied

For Map to OpenCL:

- Fixed parallelism mapping
- Limited choices for mapping to local and global memory
- Follows best practice

For Parameter Mapping:

- Amount of memory used
 - Global
 - Local
 - Registers
- Amount of parallelism
 - Work-items
 - Workgroup



Exploration in Numbers for Matrix Multiplication

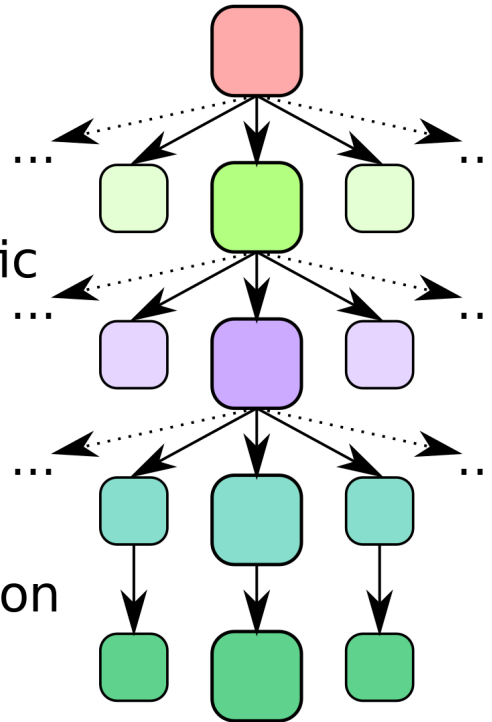
Phases:

Algorithmic
Exploration

OpenCL specific
Exploration

Parameter
Exploration

Code Generation



Program Variants:

High-Level Program 1

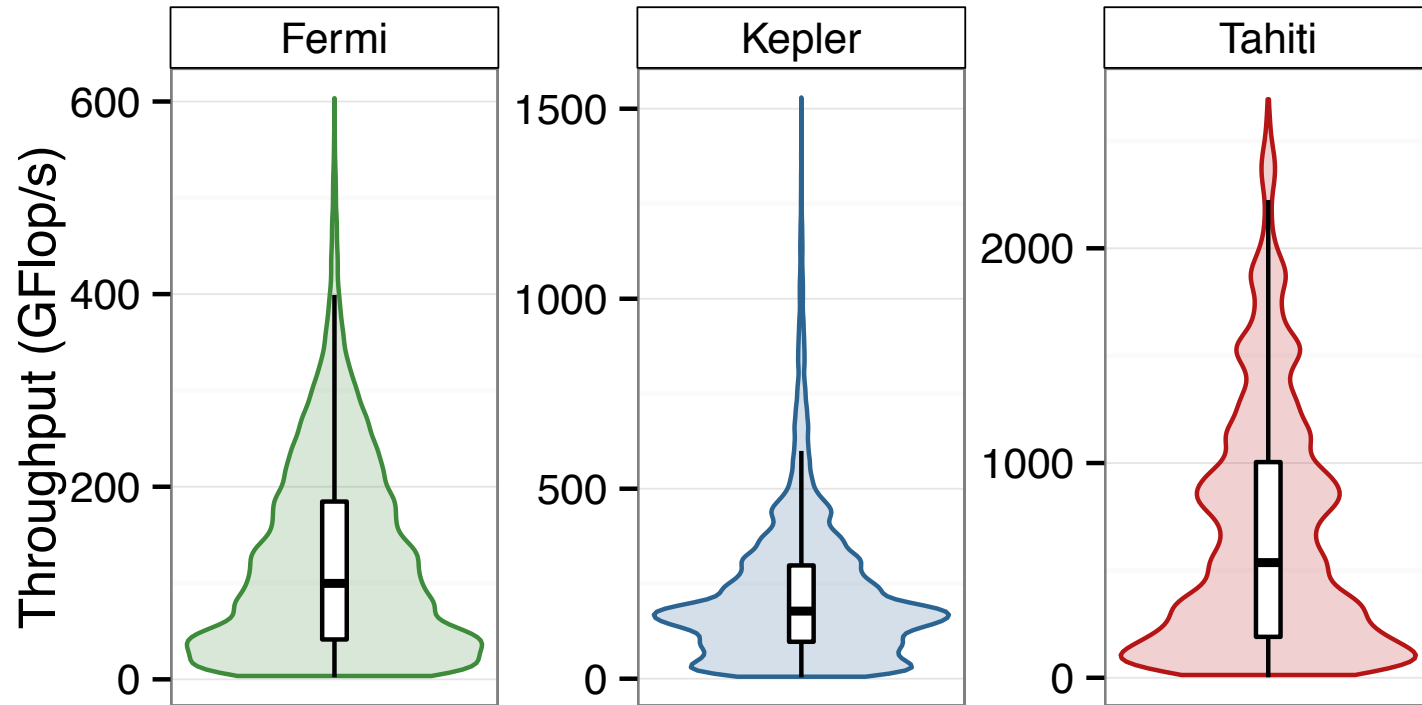
Algorithmic
Rewritten Program 8

OpenCL Specific
Program 760

Fully Specialized
Program 46,000

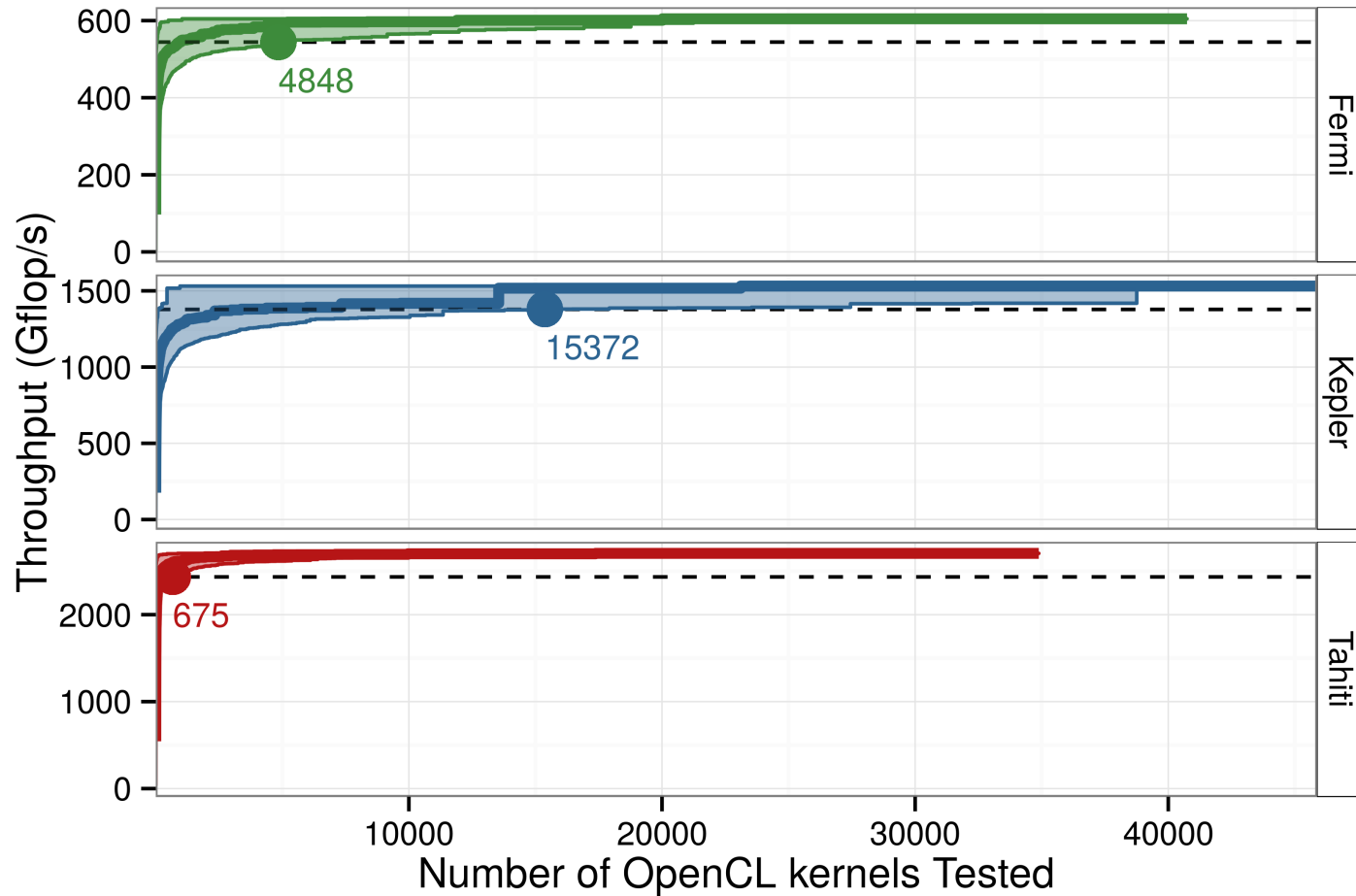
OpenCL Code 46,000

Exploration Space for Matrix Multiplication



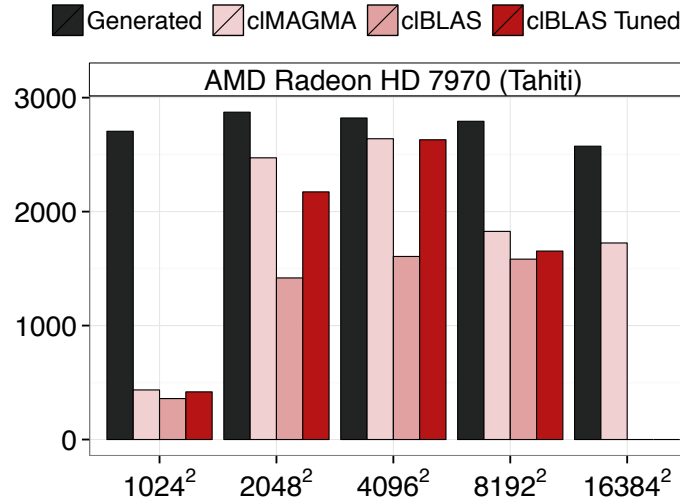
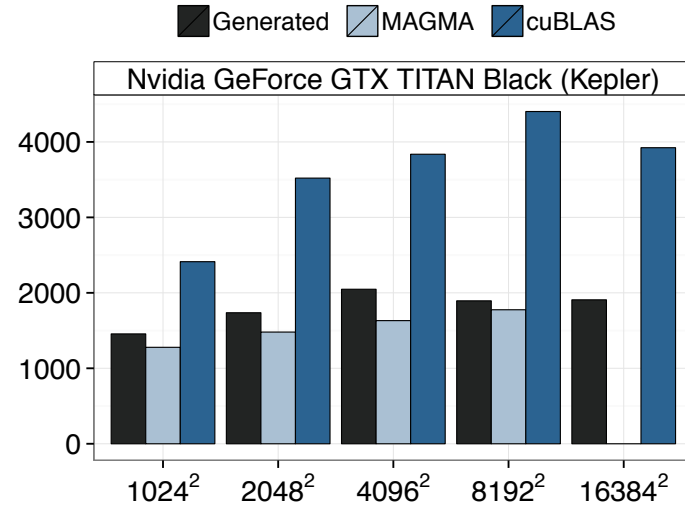
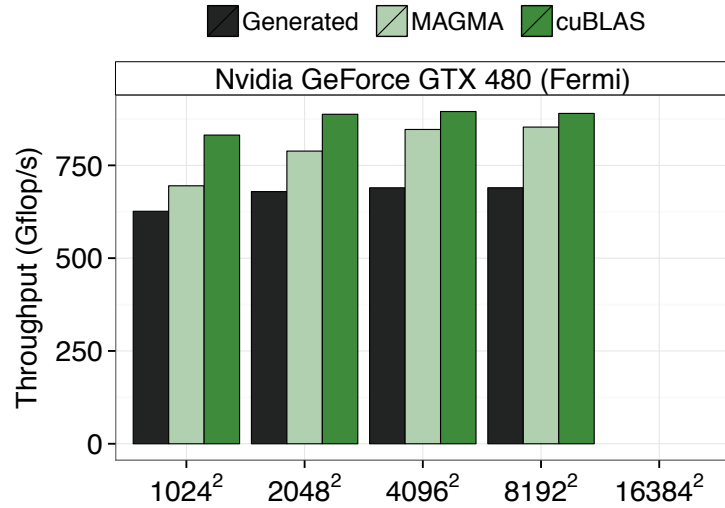
Only few OpenCL kernel with very good performance

Performance Evolution for Randomised Search



Even with a simple random search strategy one can expect to find a good performing kernel quickly

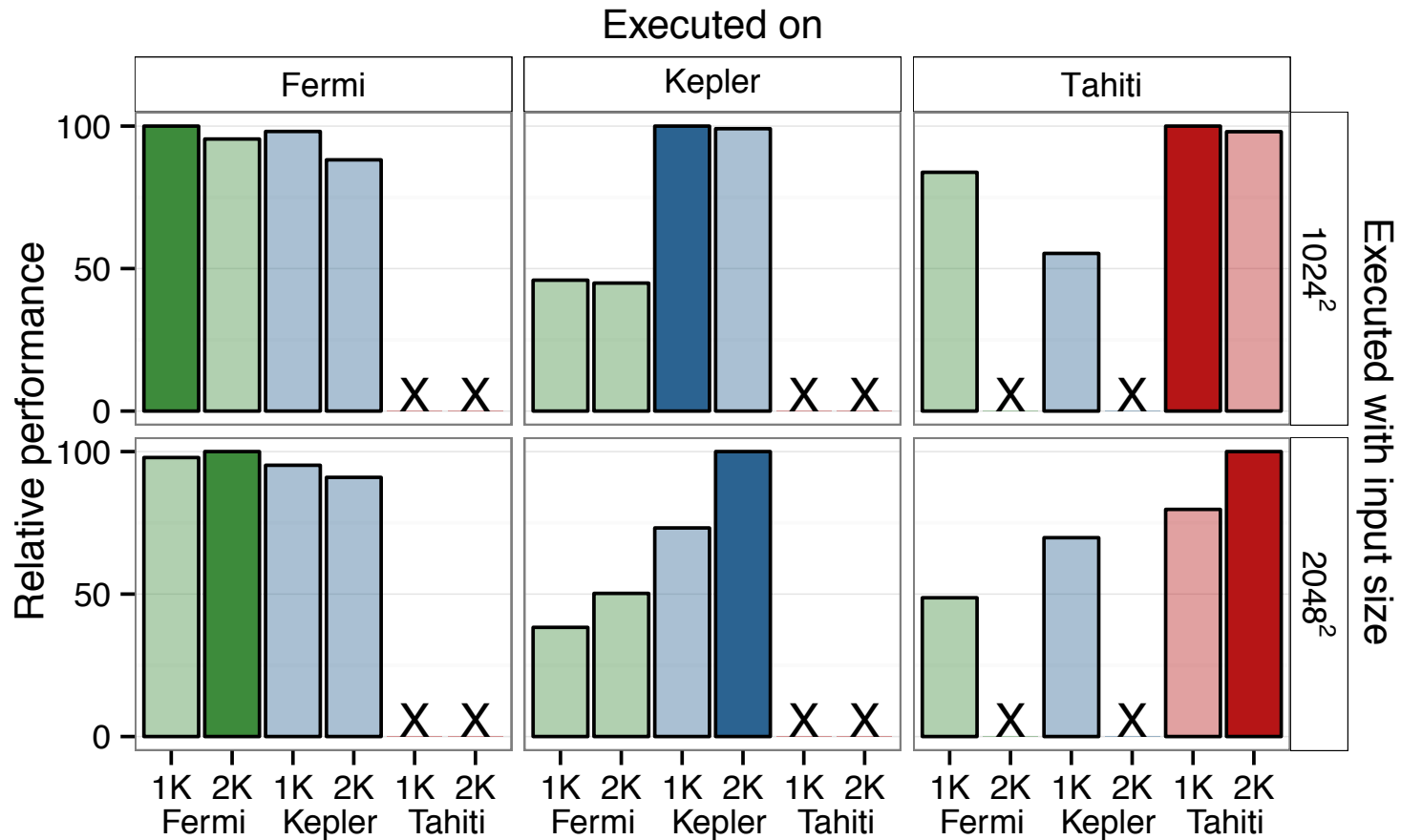
Performance Results Matrix Multiplication



Performance close or better than hand-tuned MAGMA library



Performance Portability Matrix Matrix Multiplication



The six specialized OpenCL kernels

Generated kernels are specialised for device and input size

Summary

- OpenCL code is *hard to write* and not *performance portable*
- Our approach uses
 - *portable* and functional **high-level primitives**,
 - **OpenCL-specific low-level primitives**, and
 - **rewrite-rules** to generate high *performance code*.
- Rewrite-rules define a space of possible implementations
- Performance on par with specialised, highly-tuned code



Christophe Dubach
christophe.dubach@ed.ac.uk



Michel Steuwer
michel.steuwer@ed.ac.uk



Thibaut Lutz
Now with Nvidia



Toomas Remmelg
toomas.remmelg@ed.ac.uk

More details in the **ICFP 2015**, **GPGPU 2016**, **CASES 2016** papers available at:
<http://www.lift-project.org>



THE UNIVERSITY of EDINBURGH
informatics

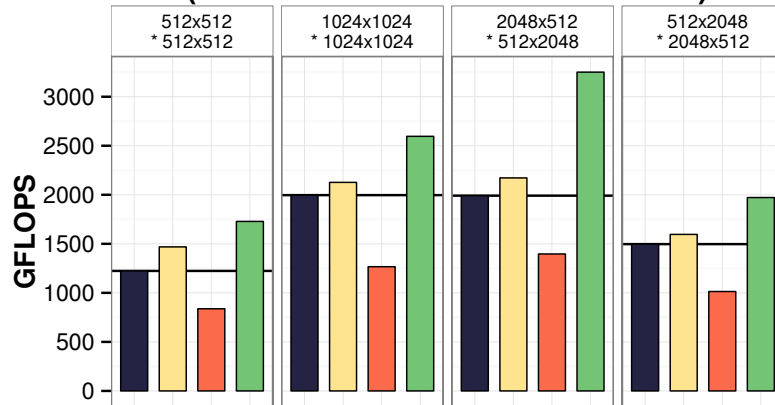
supported by:



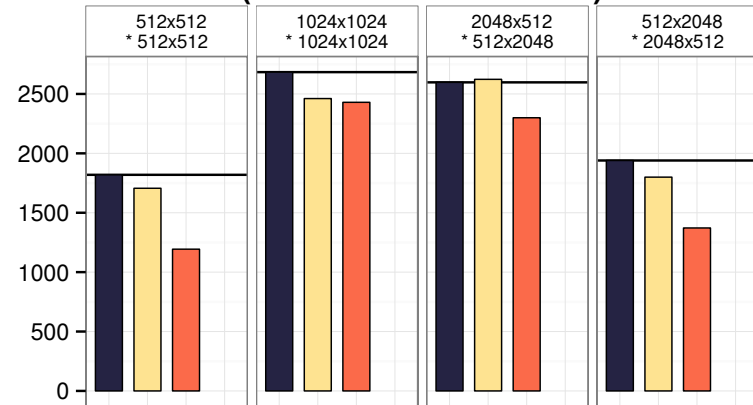
Oracle Labs

More Results Matrix Multiplication

**Desktop GPU
(Nvidia GeForce GTX Titan Black)**



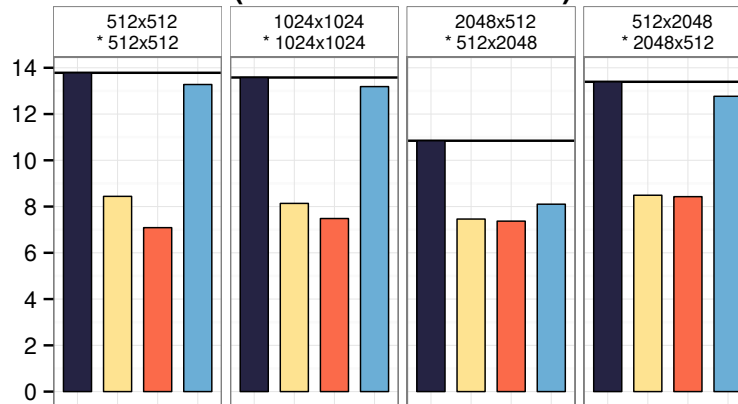
**Desktop GPU
(AMD Radeon HD 7970)**



Rewrite-based
 CLBlast + CLTune
 cBLAS
 cuBLAS

Rewrite-based
 CLBlast + CLTune
 cBLAS

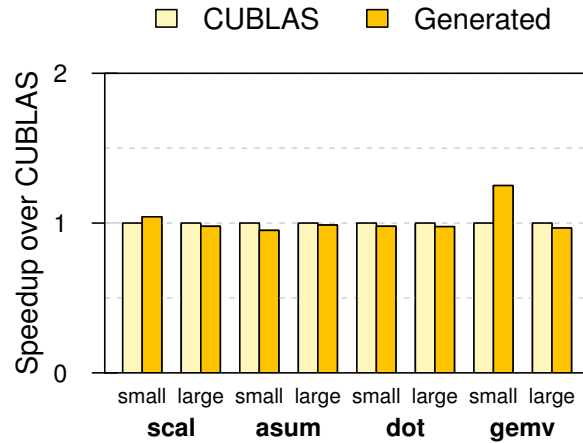
**Mobile GPU
(ARM Mali-T628 MP6)**



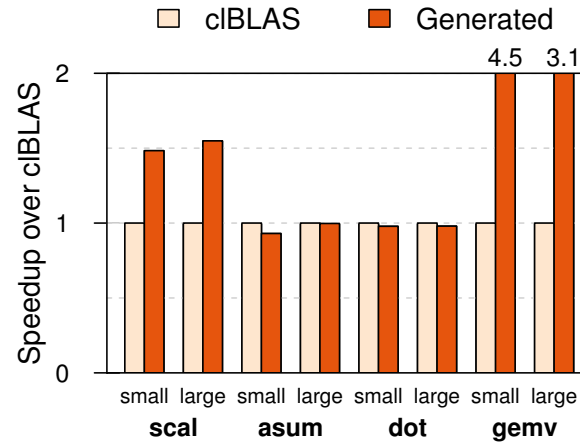
Rewrite-based
 CLBlast + CLTune
 cBLAS
 Hand optimized



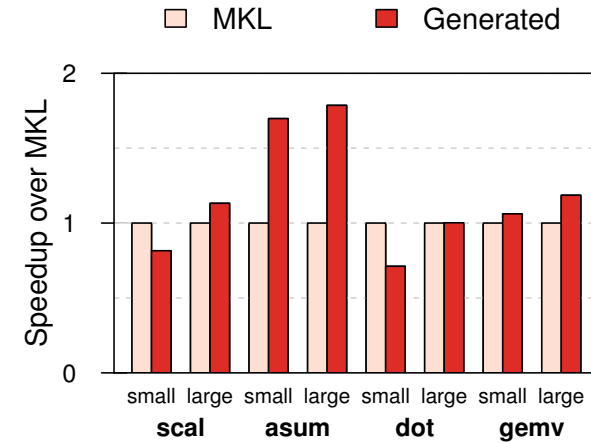
Performance Results more Benchmarks vs. Hardware-Specific Implementations



(a) Nvidia GPU



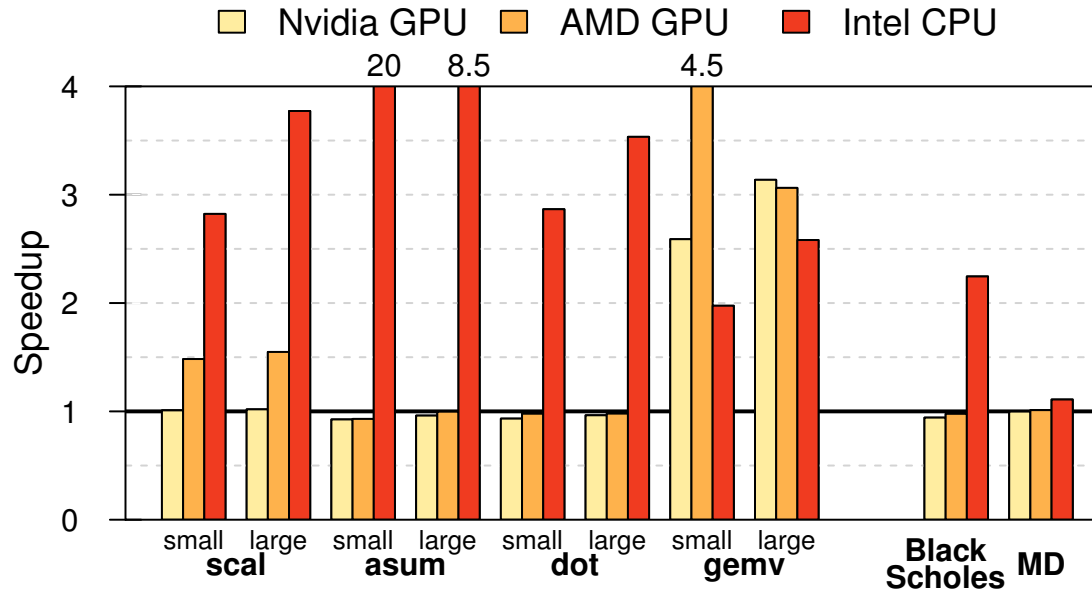
(b) AMD GPU



(c) Intel CPU

- Automatically generated code vs. expert written code
- Competitive performance vs. highly optimised implementations
- Up to **4.5x** speedup for *gemv* on AMD

Performance Results more Benchmarks vs. Portable Implementation



- Up to **20x** speedup on fairly simple benchmarks vs. portable cBLAS implementation