

Generating Performance Portable Code using Rewrite Rules

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

Michel Steuwer

Christian Fensch

Sam Lindley

Christophe Dubach

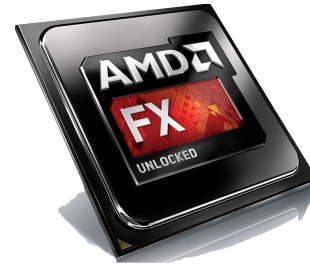


THE UNIVERSITY
of EDINBURGH



The Problem(s)

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



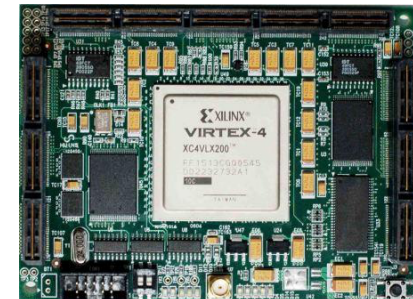
CPU



GPU



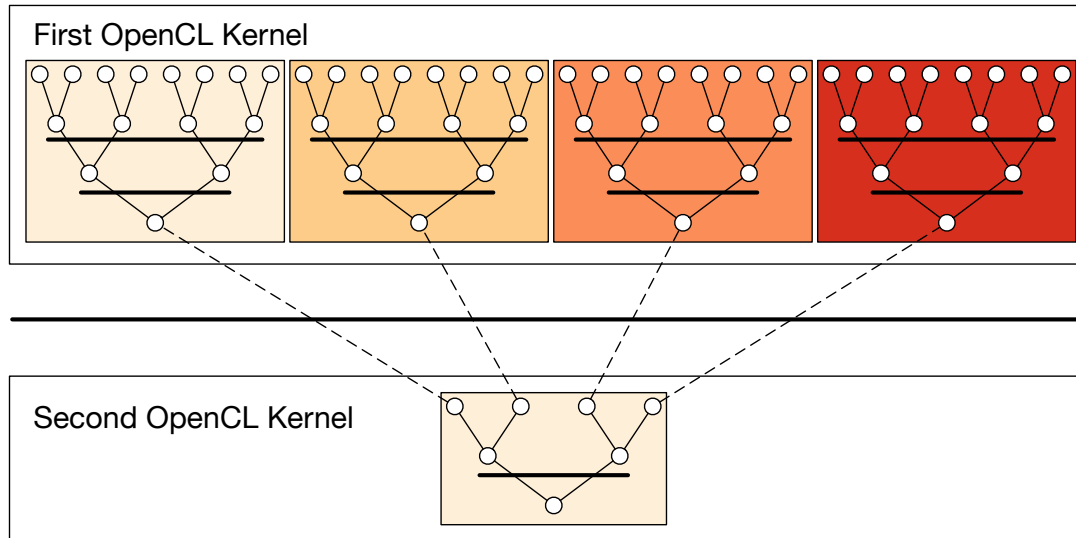
Accelerator



FPGA

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];
}
```



OpenCL Programming Model

```
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];
}
```

- Multiple *work-items* (threads) execute the same *kernel* function
- *Work-items* are organised for execution in *work-groups*



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];
}
```

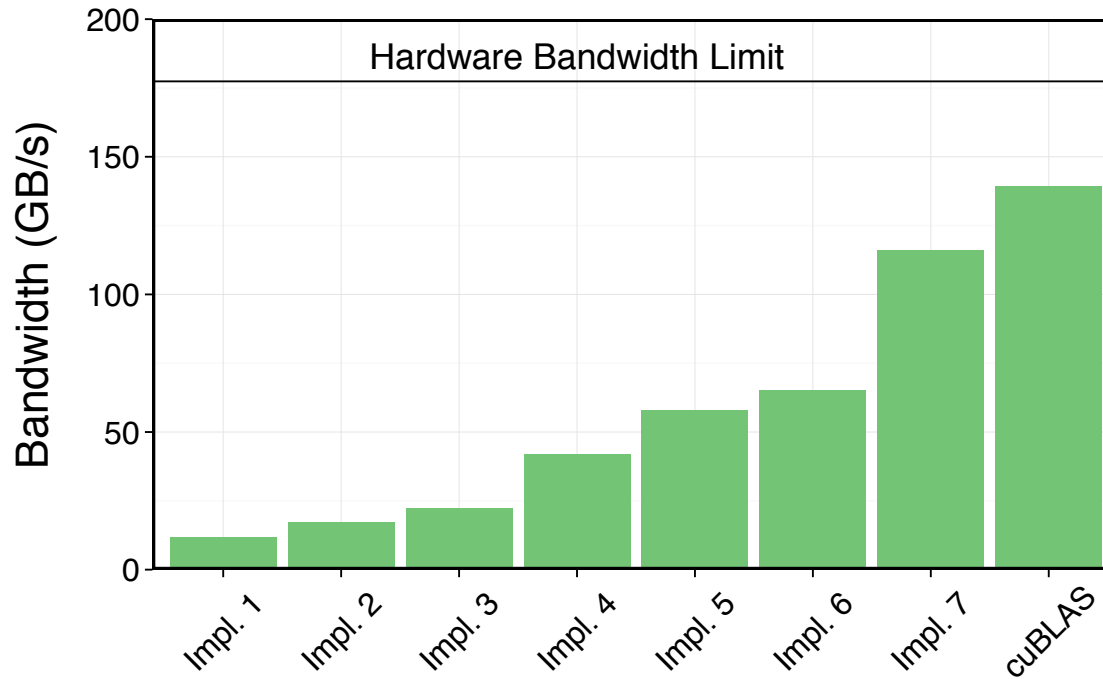
Unoptimized Implementation

```
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

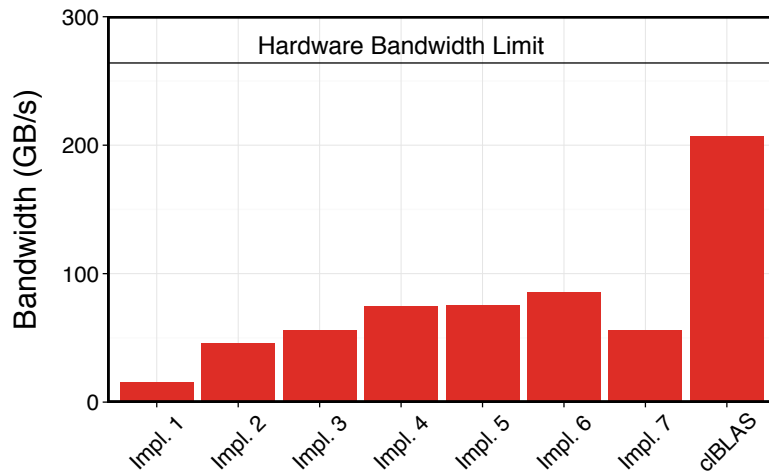
Performance Results Nvidia



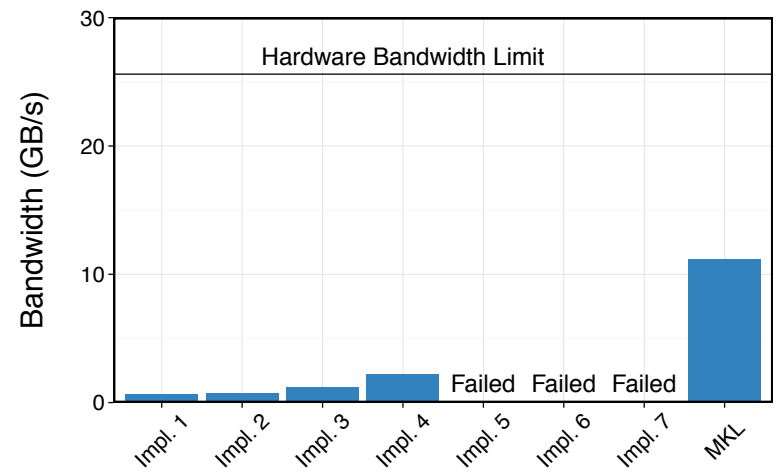
(a) Nvidia's GTX 480 GPU.

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

Performance Results AMD and Intel



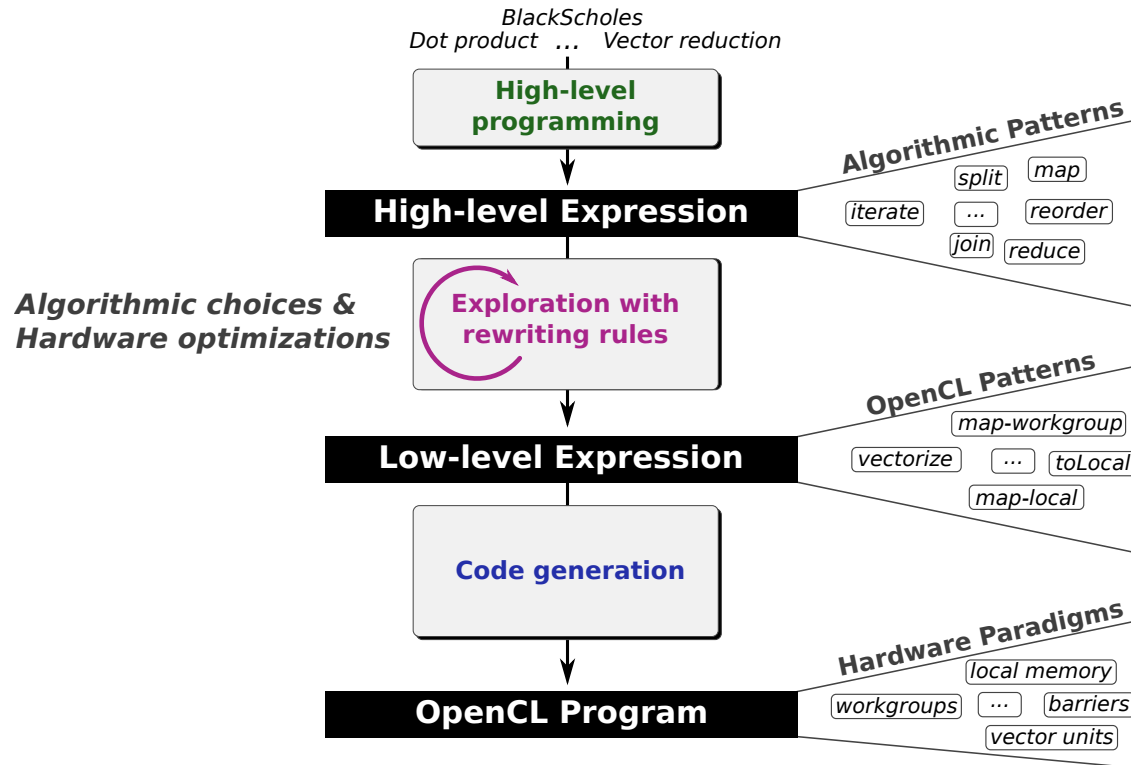
(b) AMD's HD 7970 GPU.



(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



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

Example Parallel Reduction ③

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
```

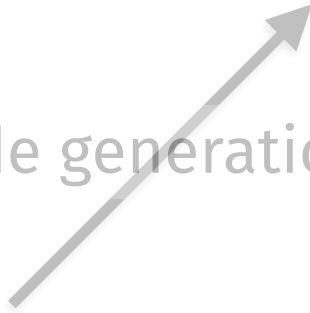
```
kernel  
void reduce6(global float* g_idata,  
             global float* g_odata,  
             unsigned int n,  
             local volatile float* l_data) {  
  unsigned int tid = get_local_id(0);  
  unsigned int i =  
    get_group_id(0) * (get_local_size(0)*2)  
    + get_local_id(0);  
 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


Example Parallel Reduction ③

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
```

```
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

$$\mathit{map}_{A,B,I} : (A \rightarrow B) \rightarrow [A]_I \rightarrow [B]_I$$

$$\mathit{zip}_{A,B,I} : [A]_I \rightarrow [B]_I \rightarrow [A \times B]_I$$

$$\mathit{reduce}_{A,I} : ((A \times A) \rightarrow A) \rightarrow A \rightarrow [A]_I \rightarrow [A]_1$$

$$\mathit{split}_{A,I} : (n : \mathit{size}) \rightarrow [A]_{n \times I} \rightarrow [[A]_n]_I$$

$$\mathit{join}_{A,I,J} : [[A]_I]_J \rightarrow [A]_{I \times J}$$

$$\mathit{iterate}_{A,I,J} : (n : \mathit{size}) \rightarrow ((m : \mathit{size}) \rightarrow [A]_{I \times m} \rightarrow [A]_m) \rightarrow [A]_{I^n \times J} \rightarrow [A]_J$$

$$\mathit{reorder}_{A,I} : [A]_I \rightarrow [A]_I$$

① High-Level Programs

$scal = \lambda a. map (*a)$

$asum = reduce (+) 0 \circ map abs$

$dot = \lambda xs\ ys. (reduce (+) 0 \circ map (*)) (zip\ xs\ ys)$

$gemv = \lambda mat\ xs\ ys\ \alpha\ \beta. map (+) (\$
 $zip (map (scal\ \alpha \circ dot\ xs) mat) (scal\ \beta\ ys))$

Example Parallel Reduction ③

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
```

```
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];  
}
```



Example Parallel Reduction ③

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
```

```
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

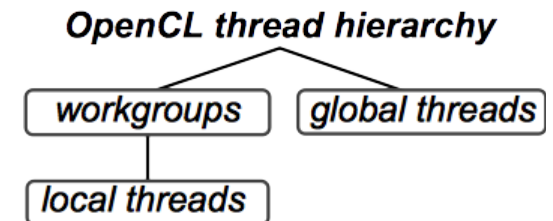
Work-items

Work-groups

Sequential implementations

Memory areas

Vectorization



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:

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

Example Parallel Reduction ③

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
```

```
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];  
}
```



Example Parallel Reduction ③

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
```

```
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

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

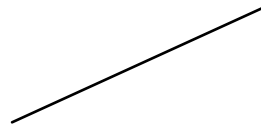


Rewrite rules define a space of possible implementations

reduce (+) 0



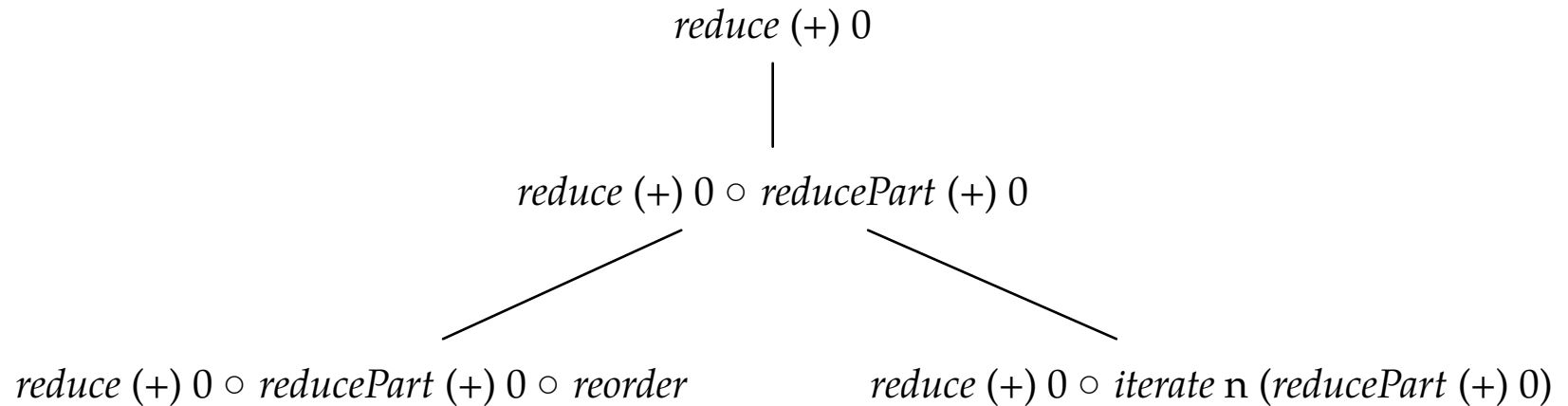
reduce (+) 0 \circ *reducePart (+) 0*



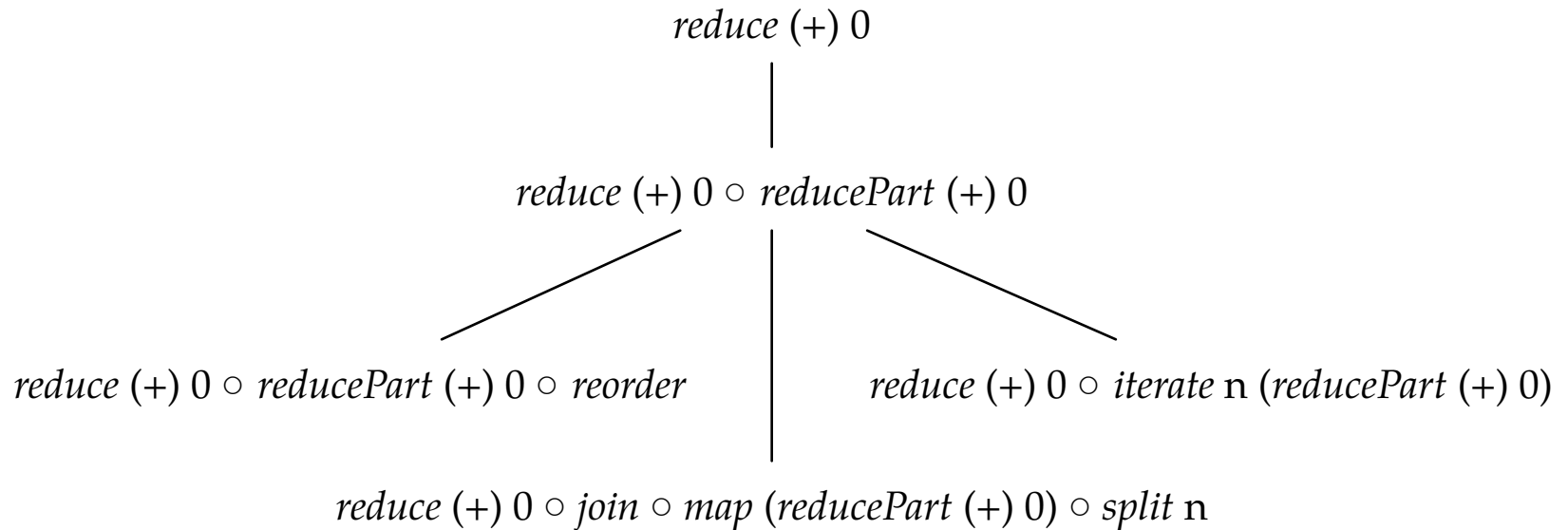
reduce (+) 0 \circ *reducePart (+) 0* \circ *reorder*



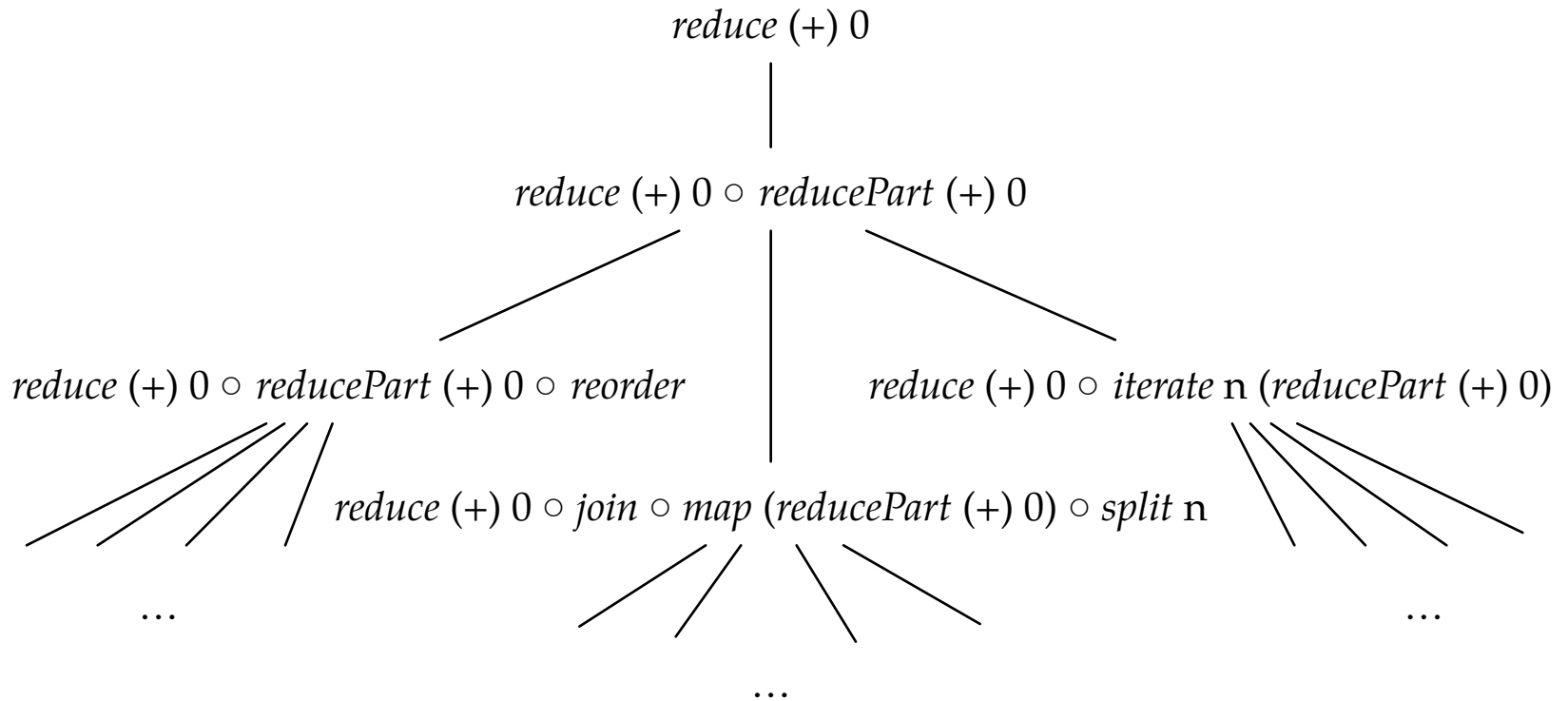
Rewrite rules define a space of possible implementations



Rewrite rules define a space of possible implementations



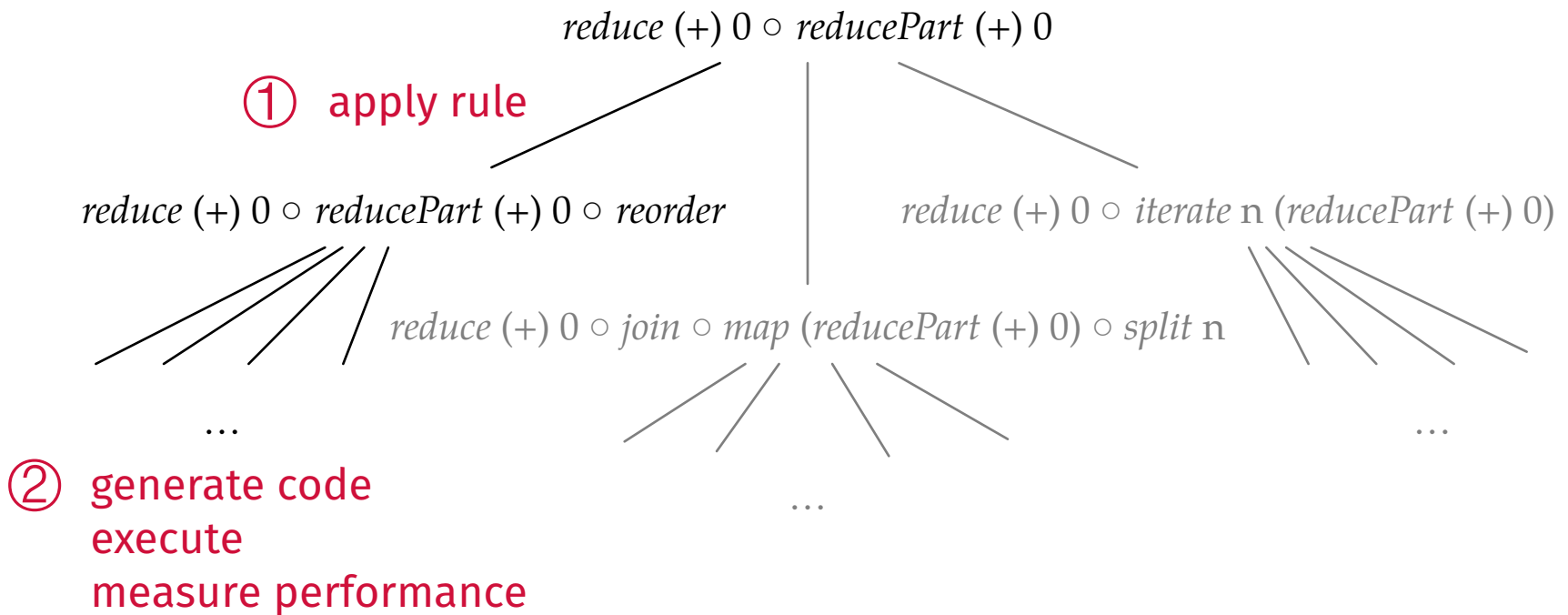
Rewrite rules define a space of possible implementations



- Fully automated search for good implementations possible

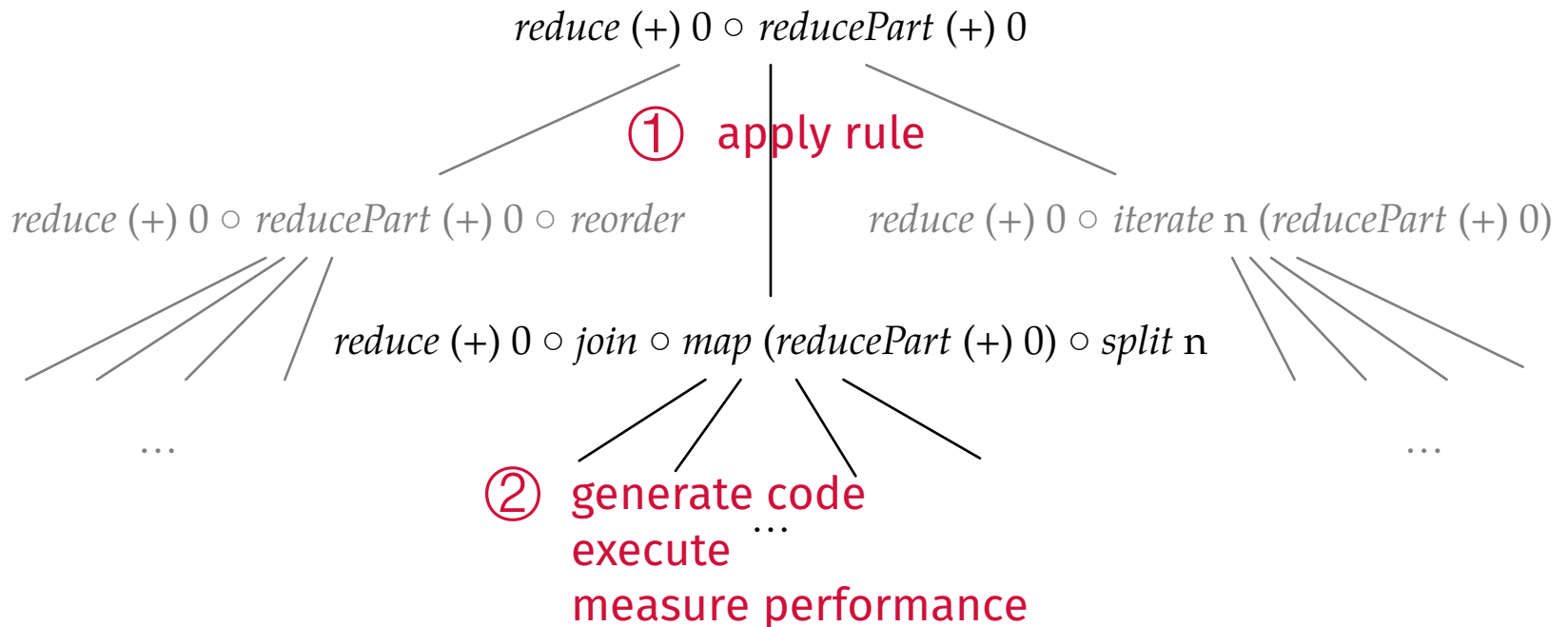
Search Strategy

- For each node in the tree:
 - Apply one rule and randomly sample subtree
 - Repeat for node with best performing subtree



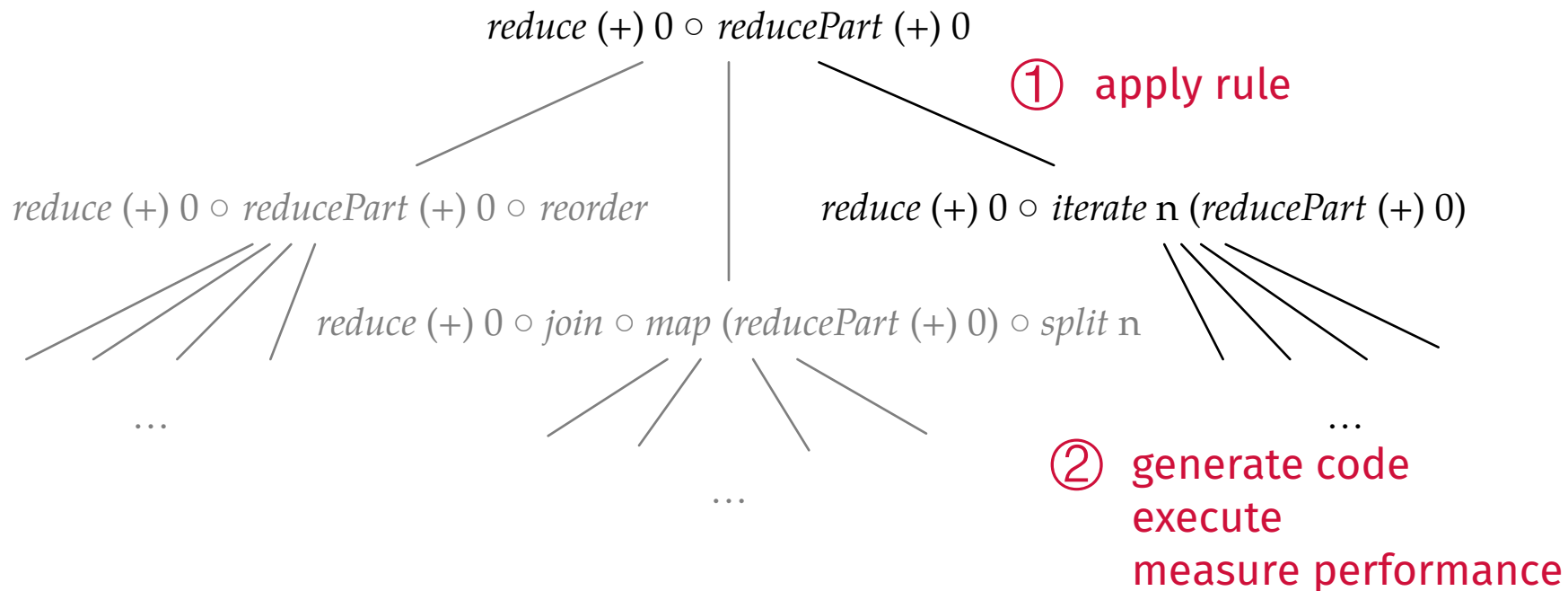
Search Strategy

- For each node in the tree:
 - Apply one rule and randomly sample subtree
 - Repeat for node with best performing subtree



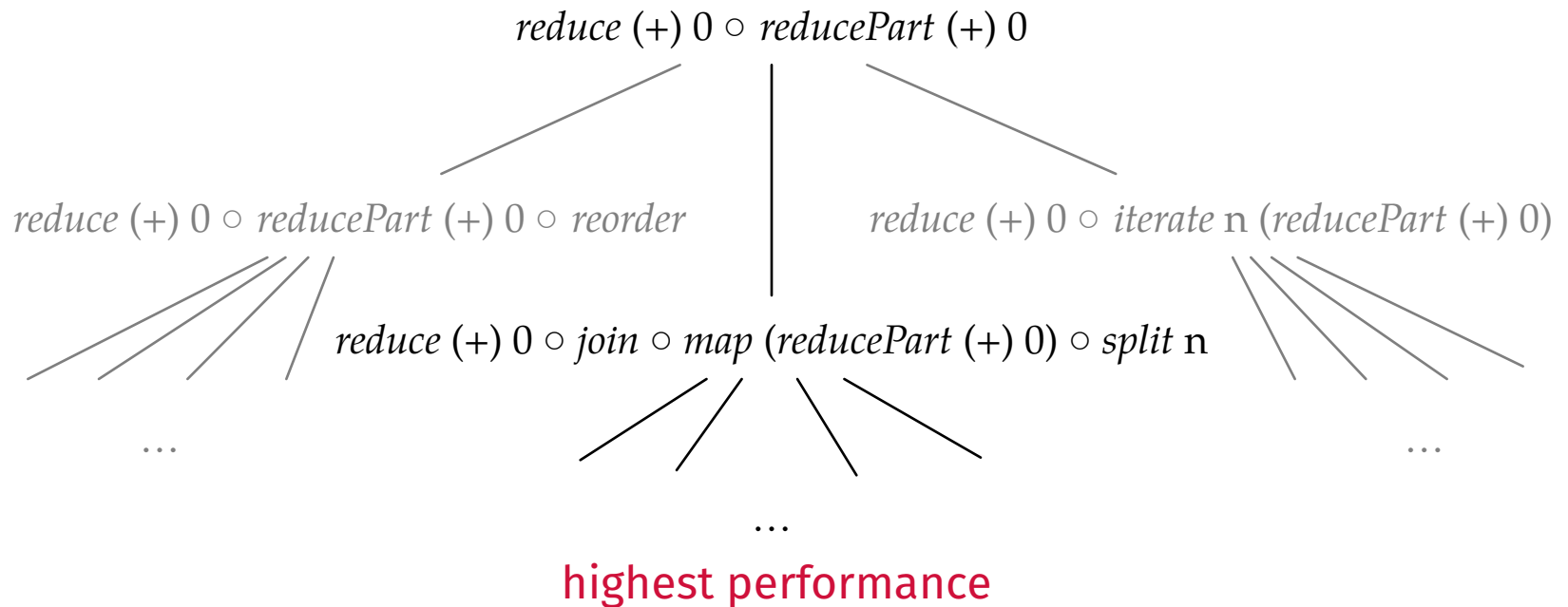
Search Strategy

- For each node in the tree:
 - Apply one rule and randomly sample subtree
 - Repeat for node with best performing subtree



Search Strategy

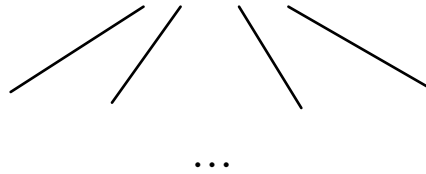
- For each node in the tree:
 - Apply one rule and randomly sample subtree
 - Repeat for node with best performing subtree



Search Strategy

- For each node in the tree:
 - Apply one rule and randomly sample subtree
 - Repeat for node with best performing subtree

reduce (+) 0 \circ *join* \circ *map* (*reducePart (+) 0*) \circ *split n*



③ repeat process

Search Results

Automatically Found Expressions

$asum = reduce (+) 0 \circ map\ abs$



Nvidia
GPU

```
 $\lambda x. (reduceSeq \circ join \circ join \circ mapWorkgroup ($   
     $toGlobal (mapLocal (reduceSeq (\lambda(a, b). a + (abs\ b)) 0)) \circ reorderStride\ 2048$   
     $) \circ split\ 128 \circ split\ 2048) x$ 
```

AMD
GPU

```
 $\lambda x. (reduceSeq \circ join \circ joinVec \circ join \circ mapWorkgroup ($   
     $mapLocal (reduceSeq (mapVec\ 2 (\lambda(a, b). a + (abs\ b))) 0 \circ reorderStride\ 2048$   
     $) \circ split\ 128 \circ splitVec\ 2 \circ split\ 4096) x$ 
```

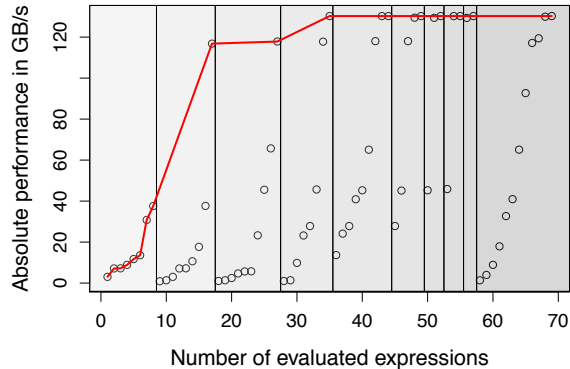
Intel
CPU

```
 $\lambda x. (reduceSeq \circ join \circ mapWorkgroup (join \circ joinVec \circ mapLocal ($   
     $reduceSeq (mapVec\ 4 (\lambda(a, b). a + (abs\ b))) 0$   
     $) \circ splitVec\ 4 \circ split\ 32768) \circ split\ 32768) x$ 
```

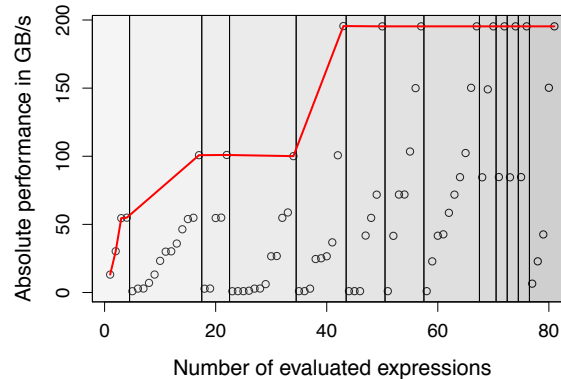
- Search on: **Nvidia** GTX 480 GPU, **AMD** Radeon HD 7970 GPU, **Intel** Xeon E5530 CPU

Search Results

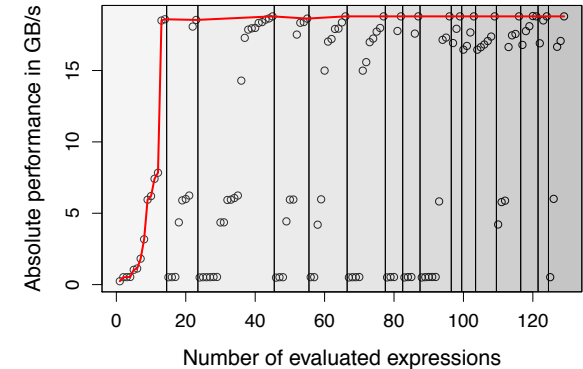
Search Efficiency



(a) Nvidia GPU



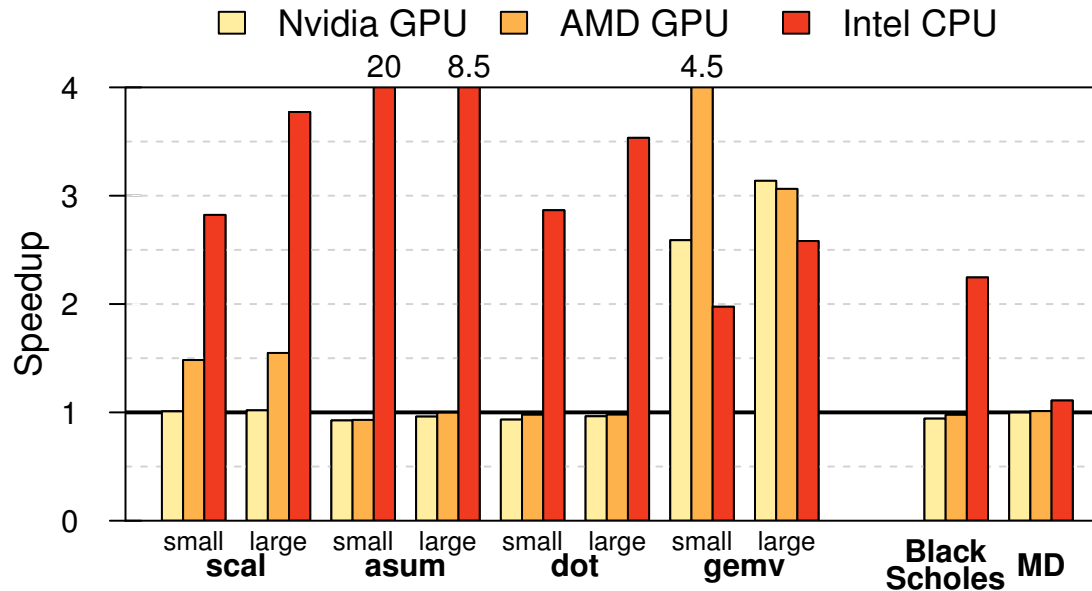
(b) AMD GPU



(c) Intel CPU

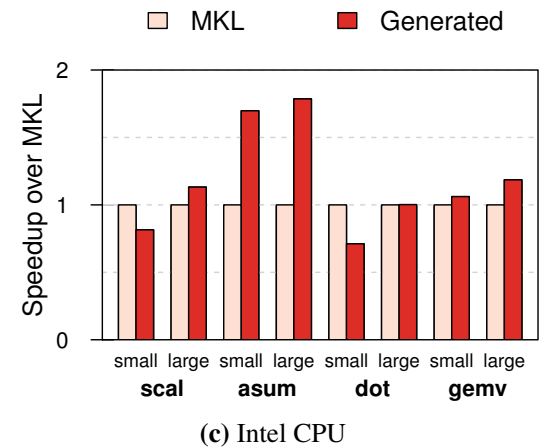
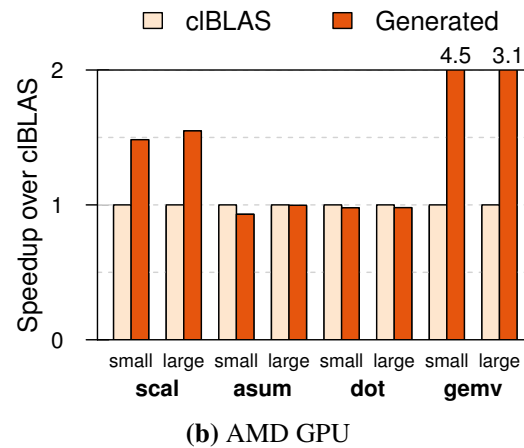
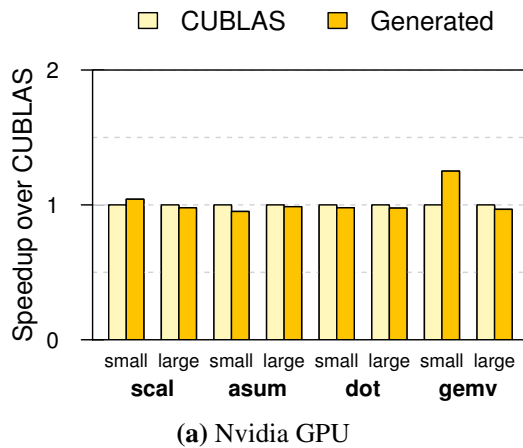
- Overall search on each platform took < 1 hour
- Average execution time per tested expression < 1/2 second

Performance Results vs. Portable Implementation



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

Performance Results vs. Hardware-Specific Implementations



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

Summary

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

Michel Steuwer — michel.steuwer@ed.ac.uk

