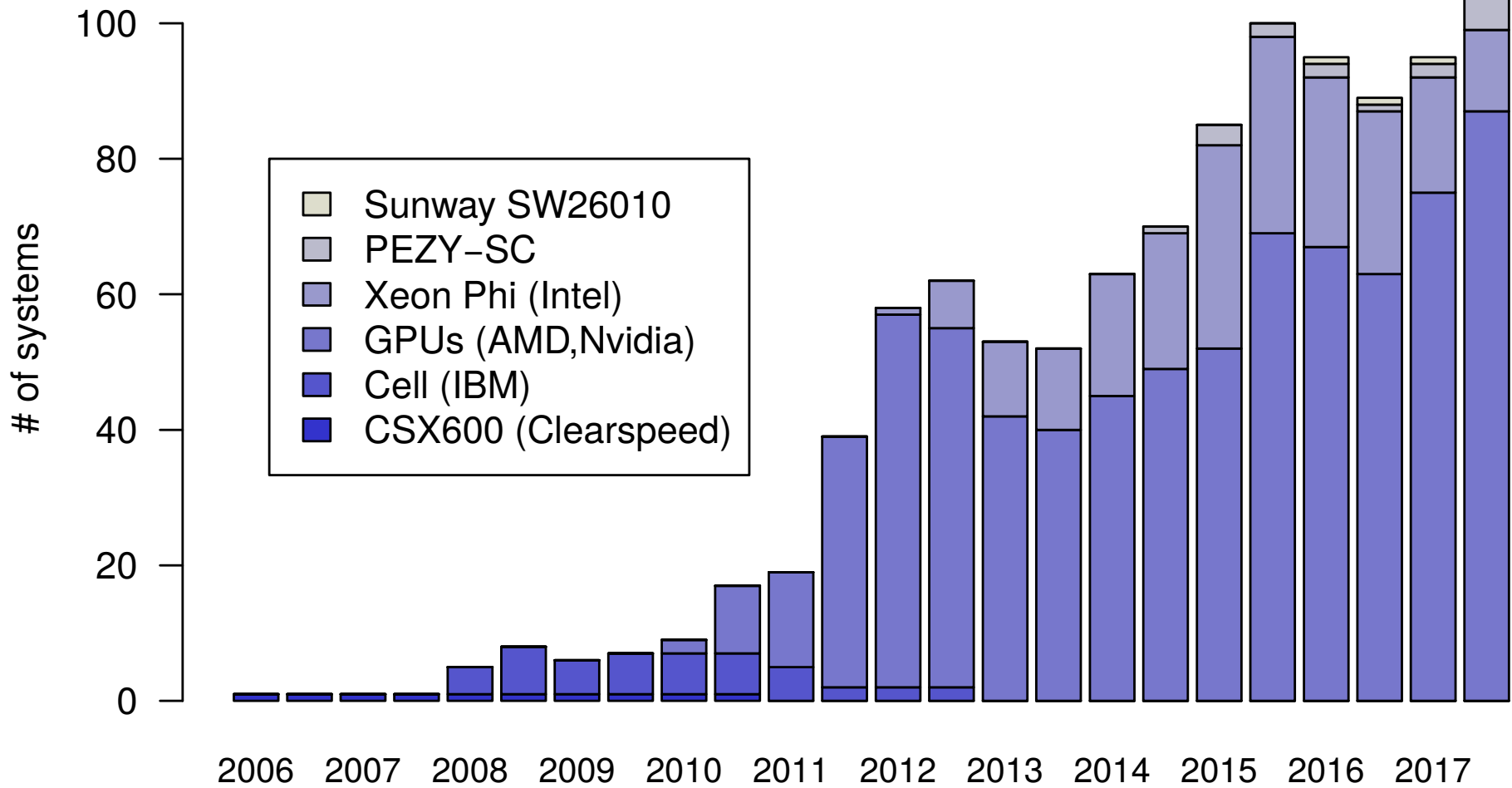# Lift
# Performance Portable Code Generation on Parallel Accelerators
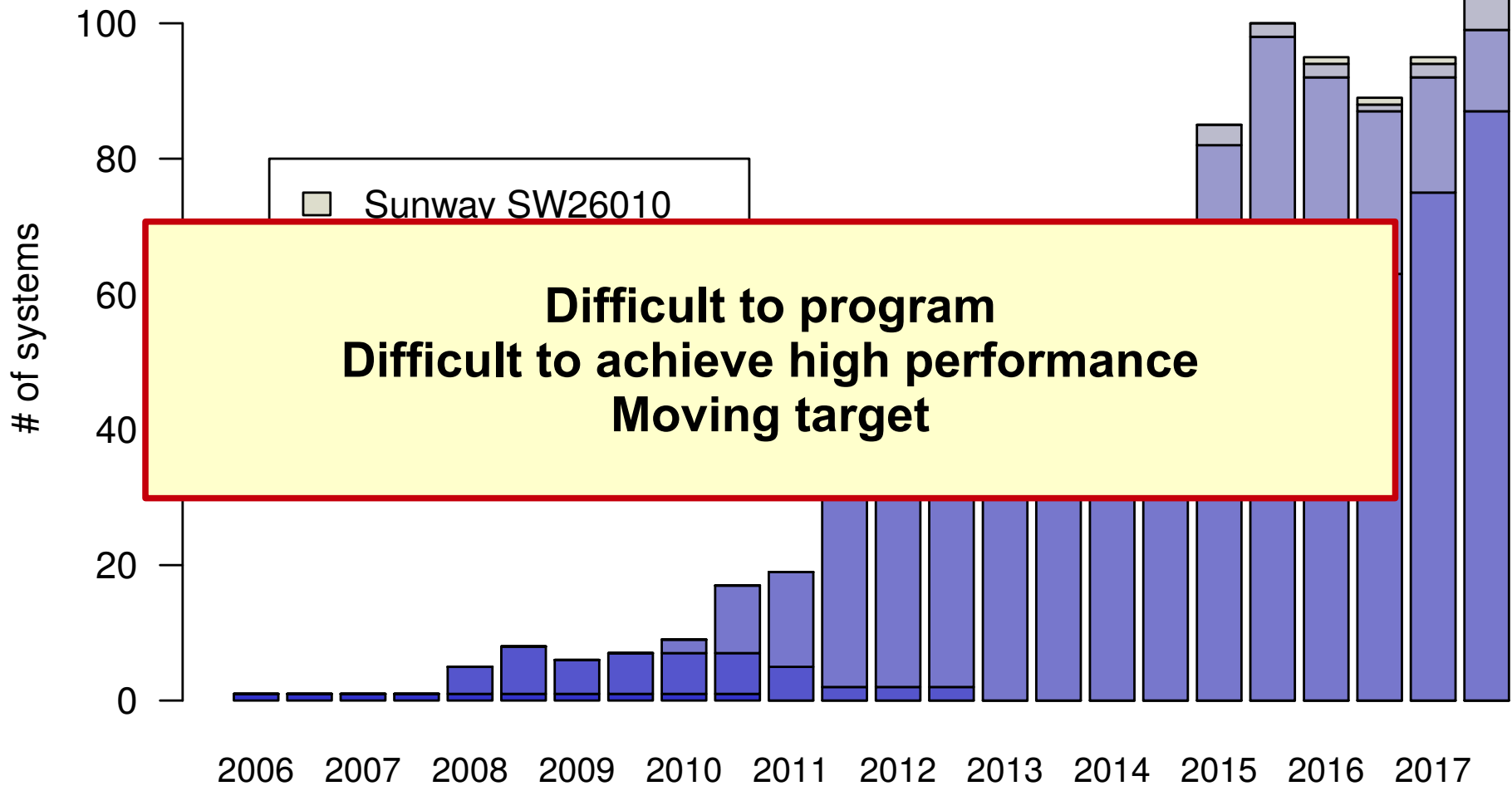
ISPASS tutorial

2 Apr. 2018

# Top 500 with parallel accelerators

# Top 500 with parallel accelerators



# of systems

Sunway SW26010

Difficult to program
Difficult to achieve high performance
Moving target

100

80

60

40

20

0

2006  2007  2008  2009  2010  2011  2012  2013  2014  2015  2016  2017

# How to sum an array?

# How to sum an array?

```
float acc = 0;
for (int i=0; i<N; i++)
  acc += input[i];
```

# How to really sum an array:

```
kernel void sum(global float* g_in, global float* g_out,
                unsigned int n, local volatile float* l_data) {

  unsigned int tid = get_local_id(0);
  unsigned int i   = get_group_id(0) * 256 + get_local_id(0);

  l_data[tid] = 0;
  while (i < n) {
    l_data[tid] += g_in[i];
    i += 256 * get_num_groups(0);
  }
  barrier(CLK_LOCAL_MEM_FENCE);

  if (tid < 128)
    l_data[tid] += l_data[tid+128];
  barrier(CLK_LOCAL_MEM_FENCE);
  if (tid <  64)
    l_data[tid] += l_data[tid+ 64];
  barrier(CLK_LOCAL_MEM_FENCE)

  if (tid < 32) {
    l_data[tid] += l_data[tid+32]; l_data[tid] += l_data[tid+16];
    l_data[tid] += l_data[tid+ 8]; l_data[tid] += l_data[tid+ 4];
    l_data[tid] += l_data[tid+ 2]; l_data[tid] += l_data[tid+ 1];
  }
  if (tid == 0)
    g_out[get_group_id(0)] = l_data[0];
}
```
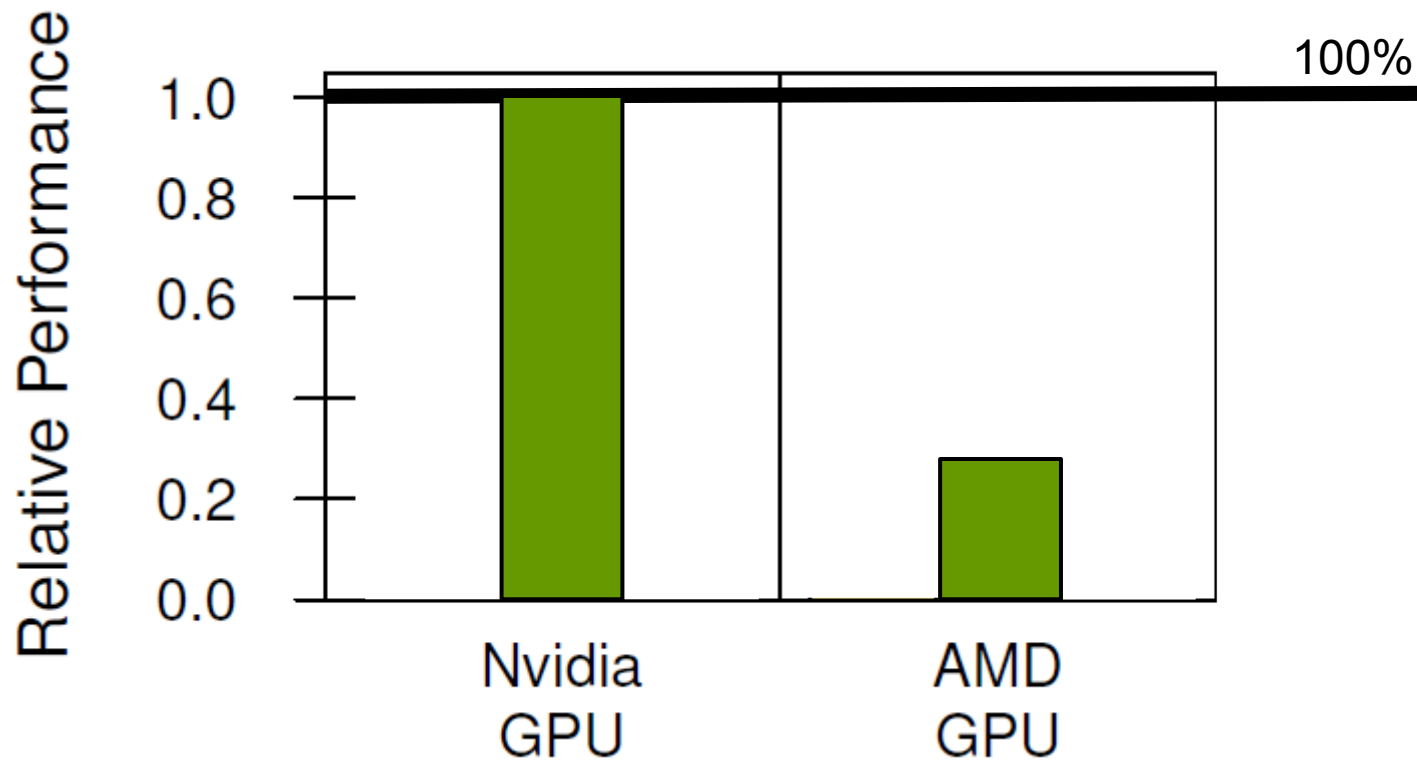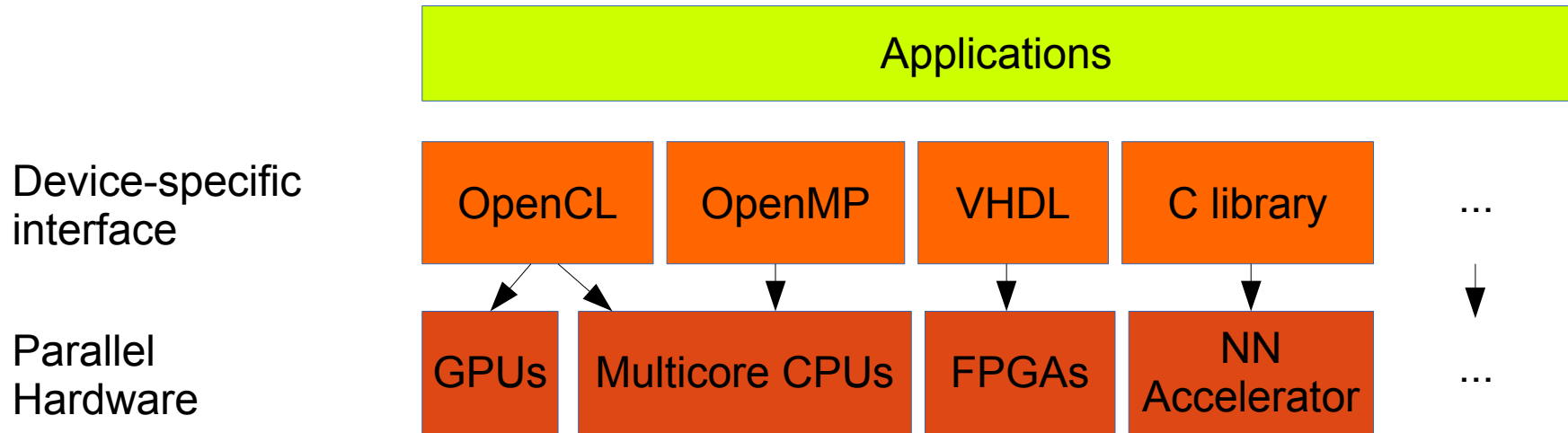
# Performance is not portable

# Current landscape
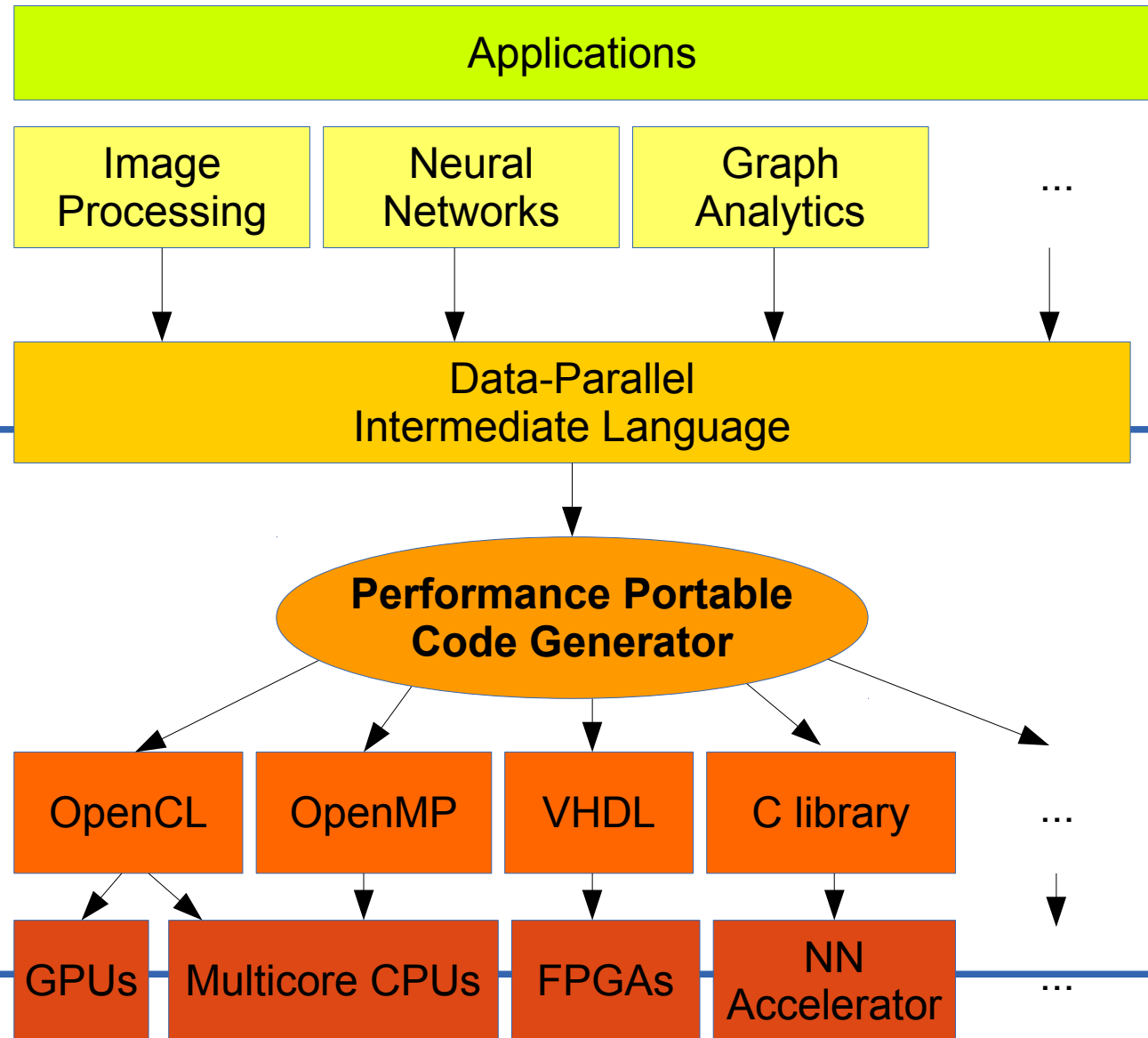
| Applications |
|:---:|

**Device-specific interface**

| OpenCL | OpenMP | VHDL | C library | ... |
|:---:|:---:|:---:|:---:|:---:|

**Parallel Hardware**

| GPUs | Multicore CPUs | FPGAs | NN Accelerator | ... |
|:---:|:---:|:---:|:---:|:---:|

# What we need

**Applications**

**Domain Specific Languages (DSLs)**

| Image Processing | Neural Networks | Graph Analytics | ... |

**Language for Parallelism**

Data-Parallel Intermediate Language

**Compiler Technology**

**Performance Portable Code Generator**

| OpenCL | OpenMP | VHDL | C library | ... |

| GPUs | Multicore CPUs | FPGAs | NN Accelerator | ... |

# Bottom up Approach

Applications

| Image Processing | Neural Networks | Graph Analytics | ... |

Data-Parallel Intermediate Language

**Performance Portable Code Generator**

| OpenCL | OpenMP | VHDL | C library | ... |

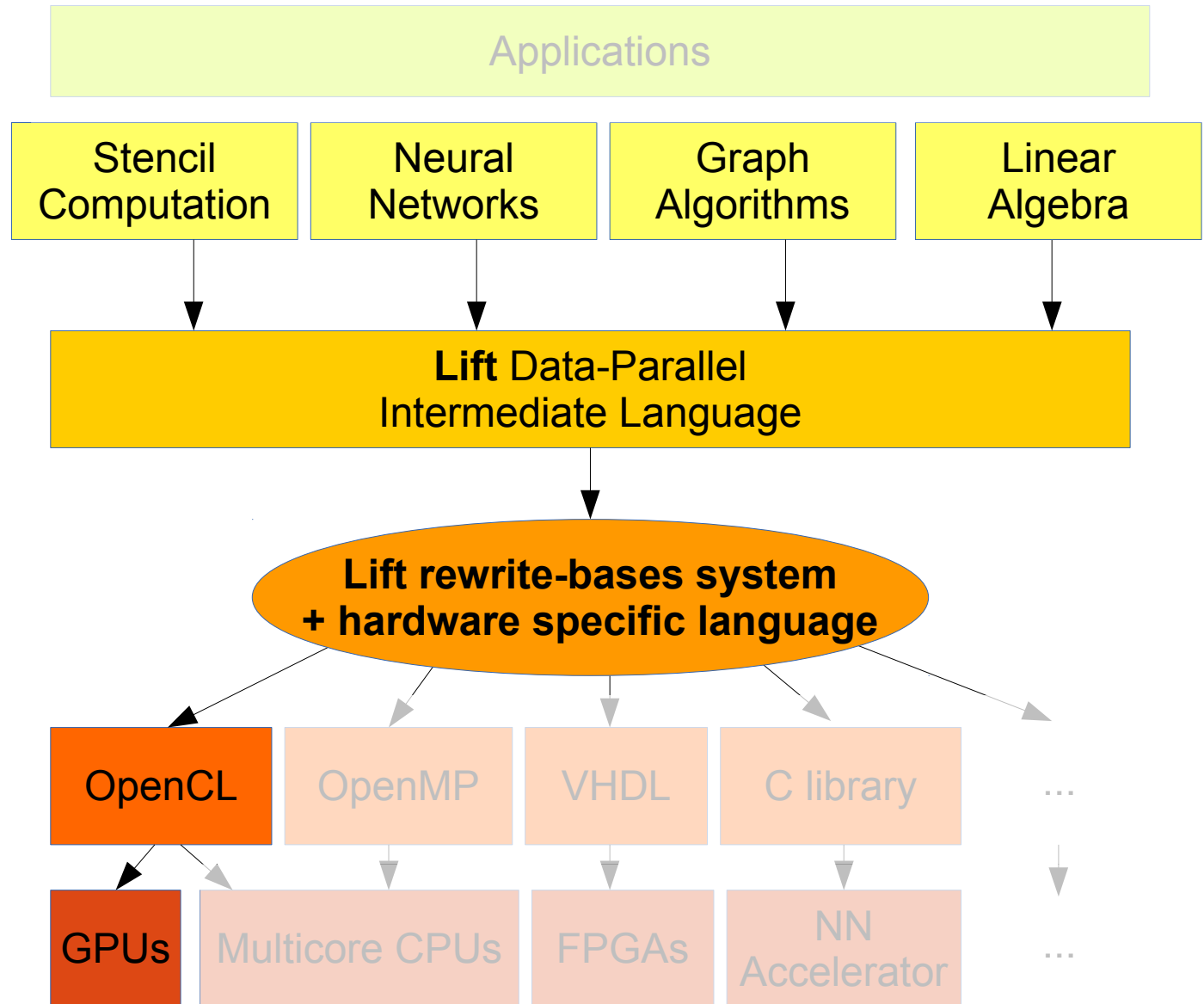| GPUs | Multicore CPUs | FPGAs | NN Accelerator | ... |

**Lift (bottom-up)**

# Current work

# Evolution of programming

# Evolution of programming

Assembly

```
            la    $s0, input;
            li    $t1, 0
    sum:    beq   $t1, $t0, endSun
            lw    $t7, 0(s0)
            addi  $t4, $t4, $t7
            addi  $s0, $s0, 4
            addi  $t3, $t3, 1
            jump  sum
```

# Evolution of programming

**Assembly**

```
        la    $s0, input;
        li    $t1, 0
sum:    beq   $t1, $t0, endSun
        lw    $t7, 0(s0)
        addi  $t4, $t4, $t7
        addi  $s0, $s0, 4
        addi  $t3, $t3, 1
        jump  sum
```

**Structured**

```
float acc = 0;
  for (int i=0; i<N; i++)
    acc += input[i];
```

# Evolution of programming

**Assembly**

```
        la    $s0, input;
        li    $t1, 0
sum:    beq   $t1, $t0, endSun
        lw    $t7, 0(s0)
        addi  $t4, $t4, $t7
        addi  $s0, $s0, 4
        addi  $t3, $t3, 1
        jump  sum
```
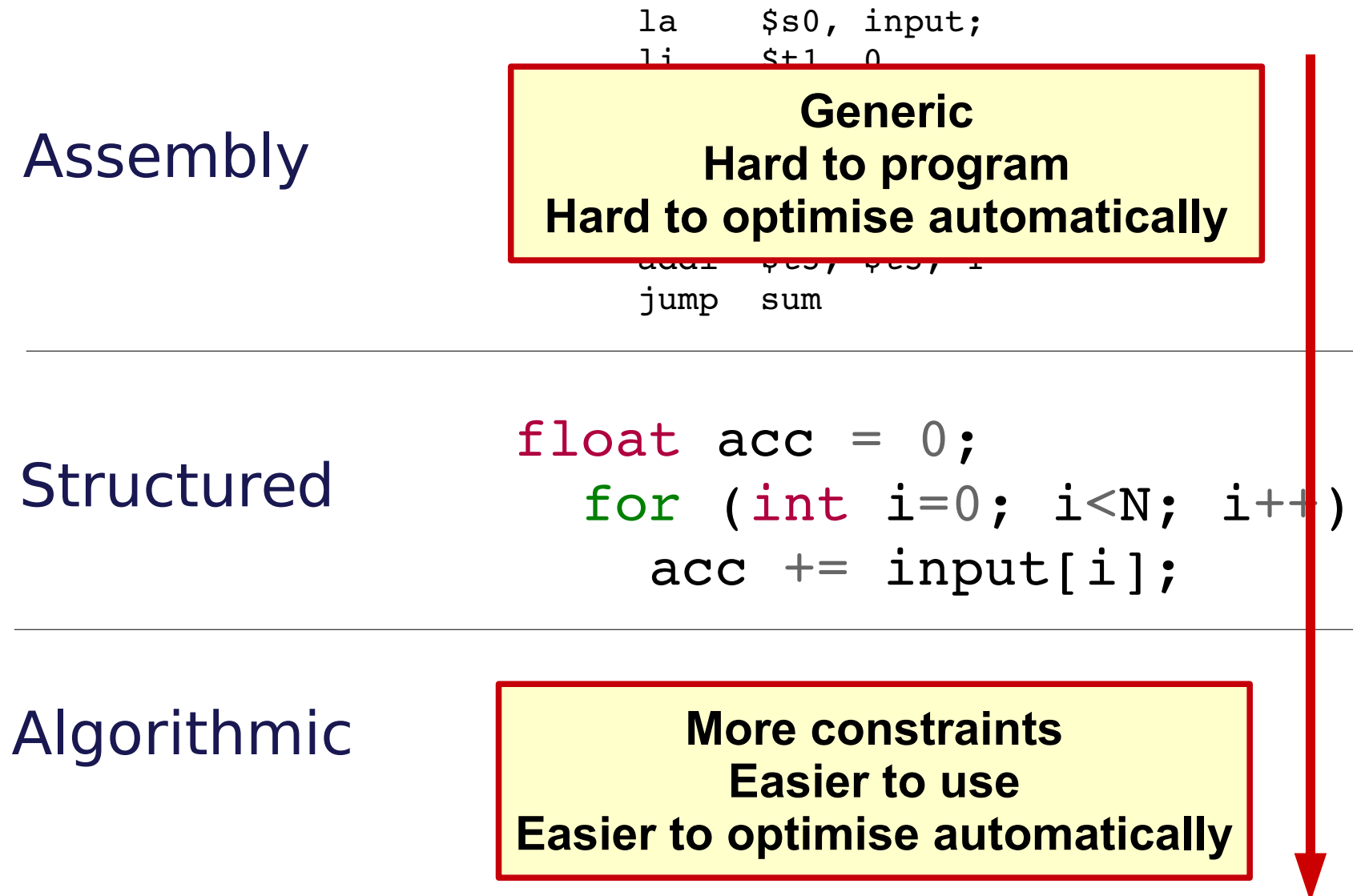
**Structured**

```
float acc = 0;
  for (int i=0; i<N; i++)
    acc += input[i];
```
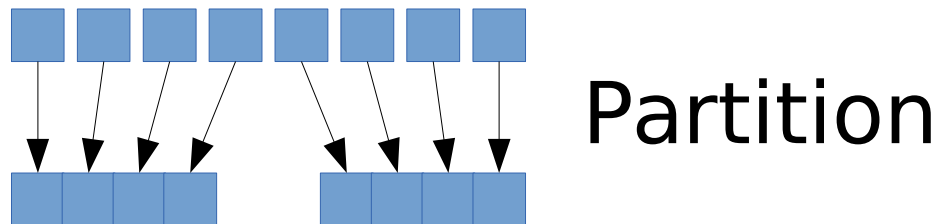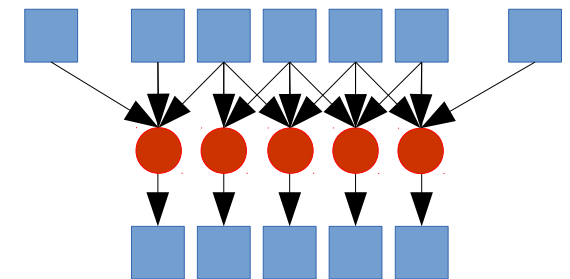
**Algorithmic**

```
reduce(0,+,input)
```

# Evolution of programming

## Assembly

```
la      $s0, input;
li      $t1, 0
```

**Generic**
**Hard to program**
**Hard to optimise automatically**
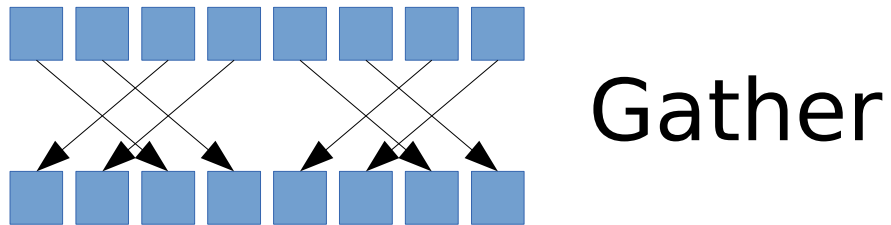
```
addi  $t3, $t3, 1
jump  sum
```
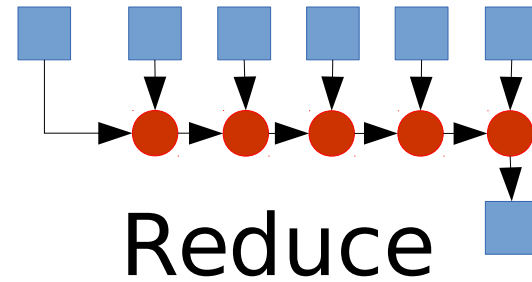
## Structured

```
float acc = 0;
  for (int i=0; i<N; i++)
    acc += input[i];
```

## Algorithmic

**More constraints**
**Easier to use**
**Easier to optimise automatically**

# Algorithmic Patterns



Map

Reduce

Gather

Stencil

Partition

...

# Separation of concern

**Programming abstraction**

Reduce

**Implementation**

Tree-based reduction

# Back to high performance code

```
kernel void sum(global float* g_in, global float* g_out,
                unsigned int n, local volatile float* l_data) {

  unsigned int tid = get_local_id(0);
  unsigned int i   = get_group_id(0) * 256 + get_local_id(0);

  l_data[tid] = 0;
  while (i < n) {
    l_data[tid] += g_in[i];
    i += 256 * get_num_groups(0);
  }
  barrier(CLK_LOCAL_MEM_FENCE);

  if (tid < 128)
    l_data[tid] += l_data[tid+128];
  barrier(CLK_LOCAL_MEM_FENCE);
  if (tid <  64)
    l_data[tid] += l_data[tid+ 64];
  barrier(CLK_LOCAL_MEM_FENCE)

  if (tid < 32) {
    l_data[tid] += l_data[tid+32]; l_data[tid] += l_data[tid+16];
    l_data[tid] += l_data[tid+ 8]; l_data[tid] += l_data[tid+ 4];
    l_data[tid] += l_data[tid+ 2]; l_data[tid] += l_data[tid+ 1];
  }
  if (tid == 0)
    g_out[get_group_id(0)] = l_data[0];
}
```

## How did we get there?

- ► Built step by step by hand
- ► Human "knows" good patterns of optimisations

# Back to high performance code

```
kernel void sum(global float* g_in, global float* g_out,
                unsigned int n, local volatile float* l_data) {

  unsigned int tid = get_local_id(0);
  unsigned int i   = get_group_id(0) * 256 + get_local_id(0);

  l_data[tid] = 0;
  while (i < n) {
    l_data[tid] += g_in[i];
    i += 256 * get_num_groups(0);
  }
  barrier(CLK_LOCAL_MEM_FENCE);

  if (tid < 128)
    l_data[tid] += l_data[tid+128];
  barrier(CLK_LOCAL_MEM_FENCE);
  if (tid <  64)
    l_data[tid] += l_data[tid+ 64];
  barrier(CLK_LOCAL_MEM_FENCE)

  if (tid < 32) {
    l_data[tid] += l_data[tid+32]; l_data[tid] += l_data[tid+16];
    l_data[tid] += l_data[tid+ 8]; l_data[tid] += l_data[tid+ 4];
    l_data[tid] += l_data[tid+ 2]; l_data[tid] += l_data[tid+ 1];
  }
  if (tid == 0)
    g_out[get_group_id(0)] = l_data[0];
}
```
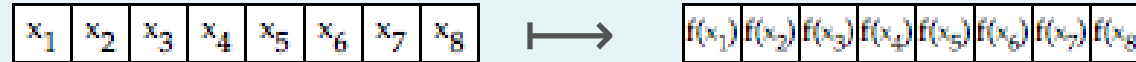
## How did we get there?
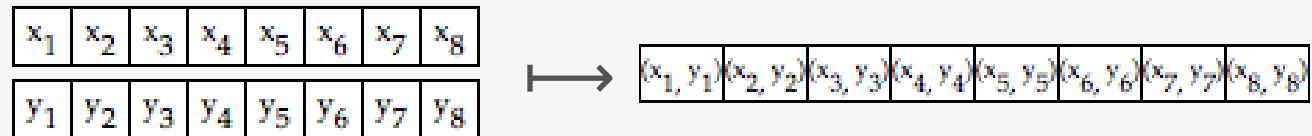
## How to get there automatically?
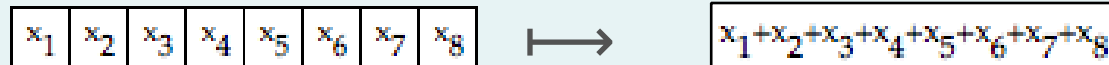
# The Lift Approach

# Lift Intermediate Language

**map**(f) : $\boxed{x_1\ x_2\ x_3\ x_4\ x_5\ x_6\ x_7\ x_8}$ $\longmapsto$ $\boxed{f(x_1)\ f(x_2)\ f(x_3)\ f(x_4)\ f(x_5)\ f(x_6)\ f(x_7)\ f(x_8)}$

**zip**: $\boxed{\begin{array}{c} x_1\ x_2\ x_3\ x_4\ x_5\ x_6\ x_7\ x_8 \\ y_1\ y_2\ y_3\ y_4\ y_5\ y_6\ y_7\ y_8 \end{array}}$ $\longmapsto$ $\boxed{(x_1, y_1)\ (x_2, y_2)\ (x_3, y_3)\ (x_4, y_4)\ (x_5, y_5)\ (x_6, y_6)\ (x_7, y_7)\ (x_8, y_8)}$

**reduce**(0,+): $\boxed{x_1\ x_2\ x_3\ x_4\ x_5\ x_6\ x_7\ x_8}$ $\longmapsto$ $\boxed{x_1+x_2+x_3+x_4+x_5+x_6+x_7+x_8}$

**split**(n): $\boxed{x_1\ x_2\ x_3\ x_4\ x_5\ x_6\ x_7\ x_8}$ $\longmapsto$ $\boxed{\boxed{x_1\ x_2}\ \boxed{x_3\ x_4}\ \boxed{x_5\ x_6}\ \boxed{x_7\ x_8}}$

**join**: $\boxed{\boxed{x_1\ x_2}\ \boxed{x_3\ x_4}\ \boxed{x_5\ x_6}\ \boxed{x_7\ x_8}}$ $\longmapsto$ $\boxed{x_1\ x_2\ x_3\ x_4\ x_5\ x_6\ x_7\ x_8}$

**iterate**(f, n): $\boxed{x_1\ x_2\ x_3\ x_4\ x_5\ x_6\ x_7\ x_8}$ $\longmapsto$ $f(\,\ldots\,f(\boxed{x_1\ x_2\ x_3\ x_4\ x_5\ x_6\ x_7\ x_8}\,)\ldots)$

**reorder**(σ): $\boxed{x_1\ x_2\ x_3\ x_4\ x_5\ x_6\ x_7\ x_8}$ $\longmapsto$ $\boxed{x_{\sigma(1)}\ x_{\sigma(2)}\ x_{\sigma(3)}\ x_{\sigma(4)}\ x_{\sigma(5)}\ x_{\sigma(6)}\ x_{\sigma(7)}\ x_{\sigma(8)}}$
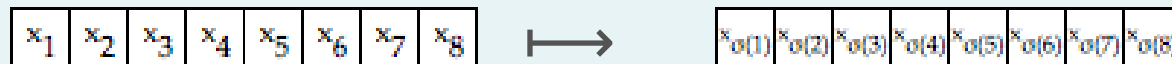
Focus on the **what** rather than **how**

# Abstracting the implementation

```c
kernel void sum(global float* g_in, global float* g_out,
                unsigned int n, local volatile float* l_data) {

  unsigned int tid = get_local_id(0);
  unsigned int i   = get_group_id(0) * 256 + get_local_id(0);

  l_data[tid] = 0;
  while (i < n) {
    l_data[tid] += g_in[i];
    i += 256 * get_num_groups(0);
  }
  barrier(CLK_LOCAL_MEM_FENCE);

  if (tid < 128)
    l_data[tid] += l_data[tid+128];
  barrier(CLK_LOCAL_MEM_FENCE);
  if (tid <  64)
    l_data[tid] += l_data[tid+ 64];
  barrier(CLK_LOCAL_MEM_FENCE)

  if (tid < 32) {
    l_data[tid] += l_data[tid+32]; l_data[tid] += l_data[tid+16];
    l_data[tid] += l_data[tid+ 8]; l_data[tid] += l_data[tid+ 4];
    l_data[tid] += l_data[tid+ 2]; l_data[tid] += l_data[tid+ 1];
  }
  if (tid == 0)
    g_out[get_group_id(0)] = l_data[0];
}
```

**OpenCL**

```
split(256) >>
reorder >>
mapWorkGroup(
    reduceSeq(
        toLocal(
            mapLocalThread(0)),
        zip(input) >>
    mapLocalThread(
        +)) >>
    iterate(
        reorder  >>
        split(2) >>
    mapLocalThread(
        reduceSeq(0,+))) >>
    iterate(
        reorder  >>
        split(2) >>
    mapWarp(
        reduceSeq(0,+))) >>
toGlobal(
    mapLocaThread(id)))
```

**Lift IR**

# Abstracting the implementation

```
kernel void sum(global float* g_in, global float* g_out,
                unsigned int n, local volatile float* l_data) {

  unsigned int tid = get_local_id(0);
  unsigned int i   = get_group_id(0) * 256 + get_local_id(0);

  l_data[tid] = 0;
  while (i < n) {
    l_data[tid] += g_in[i];
    i += 256 * get_num_groups(0);
  }
  barrier(CLK_LOCAL_MEM_FENCE);

  if (tid < 128)
    l_data[tid] += l_data[tid+128];
  barrier(CLK_LOCAL_MEM_FENCE);
  if (tid <  64)
    l_data[tid] += l_data[tid+ 64];
  barrier(CLK_LOCAL_MEM_FENCE)

  if (tid < 32) {
    l_data[tid] += l_data[tid+32]; l_data[tid] += l_data[tid+16];
    l_data[tid] += l_data[tid+ 8]; l_data[tid] += l_data[tid+ 4];
    l_data[tid] += l_data[tid+ 2]; l_data[tid] += l_data[tid+ 1];
  }
  if (tid == 0)
    g_out[get_group_id(0)] = l_data[0];
}
```

**OpenCL**

```
split(256) >>
reorder >>
mapWorkGroup(
    reduceSeq(
        toLocal(
            mapLocalThread(0)),
        zip(input) >>
        mapLocalThread(
            +)) >>
    iterate(
        reorder  >>
        split(2) >>
        mapLocalThread(
            reduceSeq(0,+))) >>
    iterate(
        reorder  >>
        split(2) >>
        mapWarp(
            reduceSeq(0,+))) >>
    toGlobal(
        mapLocaThread(id)))
```

**Lift IR**

# Abstracting the implementation

```c
kernel void sum(global float* g_in, global float* g_out,
                unsigned int n, local volatile float* l_data) {

  unsigned int tid = get_local_id(0);
  unsigned int i   = get_group_id(0) * 256 + get_local_id(0);

  l_data[tid] = 0;
  while (i < n) {
    l_data[tid] += g_in[i];
    i += 256 * get_num_groups(0);
  }
  barrier(CLK_LOCAL_MEM_FENCE);

  if (tid < 128)
    l_data[tid] += l_data[tid+128];
  barrier(CLK_LOCAL_MEM_FENCE);
  if (tid <  64)
    l_data[tid] += l_data[tid+ 64];
  barrier(CLK_LOCAL_MEM_FENCE)

  if (tid < 32) {
    l_data[tid] += l_data[tid+32]; l_data[tid] += l_data[tid+16];
    l_data[tid] += l_data[tid+ 8]; l_data[tid] += l_data[tid+ 4];
    l_data[tid] += l_data[tid+ 2]; l_data[tid] += l_data[tid+ 1];
  }
  if (tid == 0)
    g_out[get_group_id(0)] = l_data[0];
}
```

**OpenCL**

```
split(256) >>
reorder >>
mapWorkGroup(
    reduceSeq(
        toLocal(
            mapLocalThread(0)),
        zip(input) >>
        mapLocalThread(
            +)) >>
    iterate(
        reorder  >>
        split(2) >>
        mapLocalThread(
            reduceSeq(0,+))) >>
    iterate(
        reorder  >>
        split(2) >>
        mapWarp(
            reduceSeq(0,+))) >>
    toGlobal(
        mapLocaThread(id)))
```

**Lift IR**

# Abstracting the implementation

```
kernel void sum(global float* g_in, global float* g_out,
                unsigned int n, local volatile float* l_data) {

  unsigned int tid = get_local_id(0);
  unsigned int i   = get_group_id(0) * 256 + get_local_id(0);

  l_data[tid] = 0;
  while (i < n) {
    l_data[tid] += g_in[i];
    i += 256 * get_num_groups(0);
  }
  barrier(CLK_LOCAL_MEM_FENCE);

  if (tid < 128)
    l_data[tid] += l_data[tid+128];
  barrier(CLK_LOCAL_MEM_FENCE);
  if (tid <  64)
    l_data[tid] += l_data[tid+ 64];
  barrier(CLK_LOCAL_MEM_FENCE)

  if (tid < 32) {
    l_data[tid] += l_data[tid+32]; l_data[tid] += l_data[tid+16];
    l_data[tid] += l_data[tid+ 8]; l_data[tid] += l_data[tid+ 4];
    l_data[tid] += l_data[tid+ 2]; l_data[tid] += l_data[tid+ 1];
  }
  if (tid == 0)
    g_out[get_group_id(0)] = l_data[0];
}
```

**OpenCL**

```
split(256) >>
reorder >>
mapWorkGroup(
    reduceSeq(
        toLocal(
            mapLocalThread(0)),
        zip(input) >>
    mapLocalThread(
        +)) >>
iterate(
    reorder   >>
    split(2) >>
    mapLocalThread(
        reduceSeq(0,+))) >>
iterate(
    reorder   >>
    split(2) >>
    mapWarp(
        reduceSeq(0,+))) >>
toGlobal(
    mapLocaThread(id)))
```

**Lift IR**

# Abstracting the implementation

```
kernel void sum(global float* g_in, global float* g_out,
                unsigned int n, local volatile float* l_data) {

  unsigned int tid = get_local_id(0);
  unsigned int i   = get_group_id(0) * 256 + get_local_id(0);

  l_data[tid] = 0;
  while (i < n) {
    l_data[tid] += g_in[i];
    i += 256 * get_num_groups(0);
  }
  barrier(CLK_LOCAL_MEM_FENCE);

  if (tid < 128)
    l_data[tid] += l_data[tid+128];
  barrier(CLK_LOCAL_MEM_FENCE);
  if (tid <  64)
    l_data[tid] += l_data[tid+ 64];
  barrier(CLK_LOCAL_MEM_FENCE)

  if (tid < 32) {
    l_data[tid] += l_data[tid+32]; l_data[tid] += l_data[tid+16];
    l_data[tid] += l_data[tid+ 8]; l_data[tid] += l_data[tid+ 4];
    l_data[tid] += l_data[tid+ 2]; l_data[tid] += l_data[tid+ 1];
  }
  if (tid == 0)
    g_out[get_group_id(0)] = l_data[0];
}
```

**OpenCL**

```
split(256) >>
reorder >>
mapWorkGroup(
    reduceSeq(
        toLocal(
            mapLocalThread(0)),
        zip(input) >>
    mapLocalThread(
        +)) >>
    iterate(
        reorder  >>
        split(2) >>
    mapLocalThread(
        reduceSeq(0,+))) >>
    iterate(
        reorder  >>
        split(2) >>
    mapWarp(
        reduceSeq(0,+))) >>
    toGlobal(
        mapLocaThread(id)))
```

**Lift IR**

# Abstracting the implementation

```
kernel void sum(global float* g_in, global float* g_out,
                unsigned int n, local volatile float* l_data) {

  unsigned int tid = get_local_id(0);
  unsigned int i   = get_group_id(0) * 256 + get_local_id(0);

  l_data[tid] = 0;
  while (i < n) {
    l_data[tid] += g_in[i];
    i += 256 * get_num_groups(0);
  }
  barrier(CLK_LOCAL_MEM_FENCE);

  if (tid < 128)
    l_data[tid] += l_data[tid+128];
  barrier(CLK_LOCAL_MEM_FENCE);
  if (tid <  64)
    l_data[tid] += l_data[tid+ 64];
  barrier(CLK_LOCAL_MEM_FENCE)

  if (tid < 32) {
    l_data[tid] += l_data[tid+32]; l_data[tid] += l_data[tid+16];
    l_data[tid] += l_data[tid+ 8]; l_data[tid] += l_data[tid+ 4];
    l_data[tid] += l_data[tid+ 2]; l_data[tid] += l_data[tid+ 1];
  }
  if (tid == 0)
    g_out[get_group_id(0)] = l_data[0];
}
```

**OpenCL**

```
split(256) >>
reorder >>
mapWorkGroup(
    reduceSeq(
        toLocal(
            mapLocalThread(0)),
        zip(input) >>
    mapLocalThread(
            +)) >>
    iterate(
        reorder  >>
        split(2) >>
        mapLocalThread(
            reduceSeq(0,+))) >>
    iterate(
        reorder  >>
        split(2) >>
        mapWarp(
            reduceSeq(0,+))) >>
toGlobal(
    mapLocaThread(id)))
```

**Lift IR**

# Observations

- Can describe implementation as composition of primitives
  - map, reduce, zip, …
  - same ones we use at the high level!

- Can express complex optimisations
  - e.g. shared memory use

# How to generate optimised version?

## Programming abstraction

```
reduce(0,+,input)
```

→

## Implementation

```
split(256) >>
reorder >>
mapWorkGroup(
    reduceSeq(
        toLocal(
            mapLocalThread(0)),
        zip(input) >>
        mapLocalThread(
            +)) >>
    iterate(
        reorder  >>
        split(2) >>
        mapLocalThread(
            reduceSeq(0,+))) >>
    iterate(
        reorder  >>
        split(2) >>
        mapWarp(
            reduceSeq(0,+))) >>
    toGlobal(
        mapLocaThread(id)))
```

# Algorithmic Rewrite Rules

- Provably correct rewrite rules
- Express algorithmic implementation choices

**Split-join rule:**

**map**(f) → **split**(n) >> **map**(**map**(f)) >> **join**

**Map fusion rule:**

**map**(f) >> **map**(g) → **map**(f >> g)

**Reduce rules:**

**reduce**(z,f)      → reducePart(z,f) >> **reduce**(f)
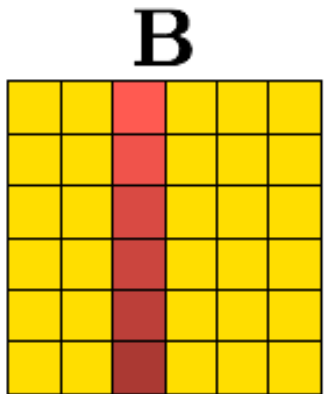
reducePart(z,f) → **reorder** >> reducePart(z,f)
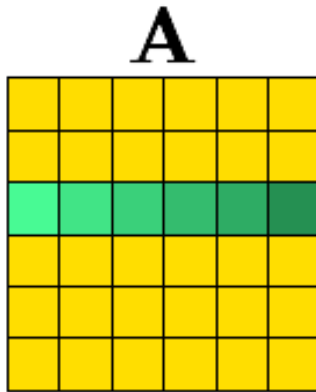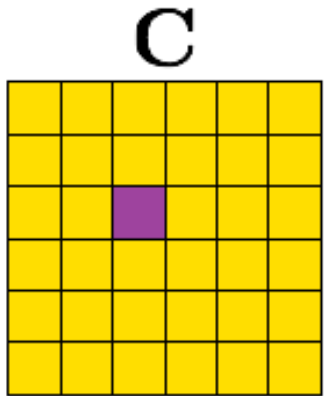reducePart(z,f) → **split**(n) >> **map**(reducePart(z,f)) >> **join**
reducePart(z,f) → **iterate**(reducePart(z,f))
reducePart(z,f) → **reduceSeq**(z,f)

# Matrix-multiplication example

# Matrix Multiplication in Lift



$$A \gg \mathrm{map}(\lambda \; \boxed{rowOfA} \mapsto$$
$$B \gg \mathrm{map}(\lambda \; \boxed{colOfB} \mapsto$$
$$\boxed{\mathrm{zip}\,(rowOfA,\; colOfB) \gg \\ \mathrm{map}(\mathrm{mult}) \gg \mathrm{reduce}(0.0\,\mathrm{f}, \mathrm{add})}$$
$$)$$
$$)$$

```
A >> map(λ rowOfA ↦
  B >> map(λ colOfB ↦
    zip(rowOfA, colOfB) >>
    map(mult) >> reduce(0.0f, add)
  )
)
```

$$A \gg \mathrm{map}(\lambda \ rowOfA \mapsto$$
$$\ B \gg \mathrm{map}(\lambda \ colOfB \mapsto$$
$$\ \ \mathrm{zip}(rowOfA, \ colOfB) \gg$$
$$\ \ \mathrm{map}(\mathrm{mult}) \gg \mathrm{reduce}(0.0\,\mathrm{f}, \mathrm{add})$$
$$\ )$$
$$)$$

```
1  for (int i = 0; i<M; i++) {
2    for (int j = 0; j<N; j++) {
3      for (int k = 0; k<K; k++) {
4        temp[k + K*N i + K*j] =
5          mult(A[k + K*i], B[k + K*j]);
6      }
7      for (int k = 0;k<K;k++) {
8        C[j + N*i] +=
9          temp[k + K*N*i + K*j];
10     }
11   }
12 }
```

$$A \gg \text{map} \ (\lambda \ rowOfA \ \mapsto$$
$$B \gg \text{map}(\lambda \ colOfB \ \mapsto$$
$$\text{zip} \ (rowOfA, \ colOfB) \gg$$
$$\text{map}(\text{mult}) \gg \text{reduce}(0.0\,f, \text{add})$$
$$)$$
$$)$$

```
1  for (int i = 0; i<M; i++) {
2    for (int j = 0; j<N; j++) {
3      for (int k = 0; k<K; k++) {
4        temp[k + K*N i + K*j] =
5          mult(A[k + K*i], B[k + K*j]);
6      }
7      for (int k = 0;k<K;k++) {
8        C[j + N*i] +=
9          temp[k + K*N*i + K*j];
10     }
11   }
12 }
```

$$Map(f) \Rightarrow Join() \circ Map(Map(f)) \circ Split(k)$$

```
A >> map (λ rowOfA ↦
 B >> map(λ colOfB ↦
   zip(rowOfA, colOfB) >>
   map(mult) >> reduce(0.0f, add)
 )
)
```

```
1  for (int i = 0; i<M; i++) {
2    for (int j = 0; j<N; j++) {
3      for (int k = 0; k<K; k++) {
4        temp[k + K*N i + K*j] =
5          mult(A[k + K*i], B[k + K*j]);
6      }
7      for (int k = 0;k<K;k++) {
8        C[j + N*i] +=
9          temp[k + K*N*i + K*j];
10     }
11   }
12 }
```

$$Map(f) \Rightarrow Join() \circ Map(Map(f)) \circ Split(k)$$

```
A >>  split(m)  >> map (λ rowsOfA ↦
  rowsOfA >> map (λ rowOfA ↦
   B >> map(λ colOfB ↦
     zip(rowOfA, colOfB) >>
     map(mult) >> reduce(0.0f, add)
   )
  )
 ) >> join
```

```
A >> map (λ rowOfA ↦
 B >> map(λ colOfB ↦
   zip (rowOfA , colOfB) >>
  map(mult) >> reduce (0.0f , add)
 )
)
```

```
1  for (int i = 0; i<M; i++) {
2    for (int j = 0; j<N; j++) {
3      for (int k = 0; k<K; k++) {
4        temp[k + K*N i + K*j] =
5          mult(A[k + K*i], B[k + K*j]);
6      }
7      for (int k = 0;k<K;k++) {
8        C[j + N*i] +=
9          temp[k + K*N*i + K*j];
10     }
11   }
12 }
```

$$Map(f) \Rightarrow Join() \circ Map(Map(f)) \circ Split(k)$$

```
A >>   split (m)  >> map (λ rowsOfA ↦
 rowsOfA >> map (λ rowOfA ↦
  B >> map(λ colOfB ↦
    zip (rowOfA , colOfB) >>
   map(mult) >> reduce (0.0f , add)
   )
  )
 ) >> join
```

```
1  for (int i = 0; i<M/m; i++) {
2    for (int l = 0 ; l < m ; l++ ) {
3      for (int j = 0; j<N; j++) {
4        for (int k = 0; k<K; k++) {
5          temp[k + 2*K*N*i + K*N*l + K*j] =
6            mult(A[k + K*l + 2*K*i], B[k + K*j]);
7        }
8        for (int k = 0;k<K;k++) {
9          C[j + N*l + 2*N*i] +=
10           temp[k + 2*K*N*i + K*N*l + K*j];
11       }
12     }
13   }
14 }
```

# Map interchanged

$$A \gg \text{split}(m) \gg \text{map}(\lambda \ rowsOfA \mapsto$$
$$B \gg \text{map}(\lambda \ colOfB \mapsto$$
$$rowsOfA \gg \text{map}(\lambda \ rowOfA \mapsto$$
$$\text{zip}(rowOfA, \ colOfB) \gg$$
$$\text{map}(\text{mult}) \gg \text{reduce}(0.0f, \text{add})$$
$$)$$
$$) \gg \text{transpose}$$
$$) \gg \text{join}$$

```
1  for (int i = 0; i<M/2; i++) {
2    for (int j = 0; j<N; j++) {
3      for (int l = 0; l<2; l++) {
4        for (int k = 0; k<K; k++) {
5          temp[k + 2*K*N*i + K*N*l + K*j] =
6            mult(A[k + K*l + 2*K*i], B[k + K*j]);
7        }
8        for (int k = 0;k<K;k++) {
9          C[j + N*l + 2*N*i] +=
10           temp[k + 2*K*N*i K*N*l + K*j];
11       }
12     }
13   }
14 }
```

# Split-join rule

$$A \gg split(m) \gg map(\lambda\ rowsOfA \mapsto$$
$$B \gg split(n) \gg map(\lambda\ colsOfB \mapsto$$
$$colsOfB \gg map(\lambda\ colOfB \mapsto$$
$$rowsOfA \gg map(\lambda\ rowOfA \mapsto$$
$$zip(rowOfA,\ colOfB) \gg$$
$$map(mult) \gg reduce(0.0f, add)$$
$$)$$
$$)$$
$$) \gg join \gg transpose$$
$$) \gg join$$

$\longrightarrow$

```
1  for (int i = 0; i<M/2; i++) {
2    for (int j = 0; j<N/2; j++) {
3      for (int m = 0; m<2; m++) {
4        for (int l = 0; l<2; l++) {
5          for (int k = 0; k<K; k++) {
6            temp[k + 4*K*N*i + 2*K*N*l + 2*K*j
                 + K*m] =
7            mult(A[k + K*l + 2*K*i], B[k + K*
                 m + 2*K*j]);
8          }
9          for (int k = 0;k<K;k++) {
10           C[m + 2*j + 2*N*l + 4*N*i] +=
11             temp[k + 4*K*N*i + 2*K*N*l + 2*
                  K*j + K*m];
12         }
13       }
14     }
15   }
16 }
```

# Map interchanged

$$A \gg \text{split}(m) \gg \text{map}(\lambda \; rowsOfA \mapsto$$
$$B \gg \text{split}(n) \gg \text{map}(\lambda \; colsOfB \mapsto$$
$$rowsOfA \gg \text{map}(\lambda \; rowOfA \mapsto$$
$$colsOfB \gg \text{map}(\lambda \; colOfB \mapsto$$
$$\text{zip}(rowOfA, \; colOfB) \gg$$
$$\text{map}(\text{mult}) \gg \text{reduce}(0.0\,\text{f}, \text{add})$$
$$)$$
$$) \gg \text{transpose}$$
$$) \gg \text{join} \gg \text{transpose}$$
$$) \gg \text{join}$$

$\longrightarrow$

```
1  for (int i = 0; i<M/2; i++) {
2    for (int j = 0; j<N/2; j++) {
3      for (int l = 0; l<2; l++) {
4        for (int m = 0; m<2; m++) {
5          for (int k = 0; k<K; k++) {
6            temp[k + 4*K*N*i + 2*K*N*l + 2*K*j
                 + K*m] =
7            mult(A[k + K*l + 2*K*i], B[k + K*
                 m + 2*K*j]);
8          }
9          for (int k = 0;k<K;k++) {
10           C[m + 2*j + 2*N*l + 4*N*i] +=
11           temp[k + 4*K*N*i + 2*K*N*l + 2*
                K*j + K*m];
12         }
13       }
14     }
15   }
16 }
```

# A few rewrite steps later...

# Tiled version

```
λ (A, B) ↦
  A >> split(m) >> map(λ nRowsOfA ↦
    B >> split(n) >> map(λ mColsOfB ↦
      zip( transpose(nRowsOfA) >> split(k),
           transpose(mColsOfB) >> split(k) ) >>
      reduceSeq(init = make2DArray(n,m, 0.0f),
       λ (accTile, (tileOfA, tileOfB)) ↦
        zip(accTile, transpose(tileOfA)) >>
        map(λ (accRow, rowOfTileOfA) ↦
         zip(accRow, transpose(tileOfB)) >>
         map(λ (acc, colOfTileOfB) ↦
          zip(rowOfTileOfA, colOfTileOfB) >>
          map(mult) >> reduce(acc, add)
         ) >> join
        )
      ) >> transpose() >>
      map(transpose) >> transpose
    ) >> join >> transpose
  ) >> join
```

```
 1  for (int i = 0;i<M/2; i++) {
 2    for (int j = 0;j<N/2; j++) {
 3      for (int k = 0;k<K/4; k++) {
 4        for (int l = 0;l<2; l++) {
 5          for (int m = 0;m<2; m++) {
 6            for (int n = 0;n<4; n++) {
 7              temp[n + 4*m + 8*N*i + 16*j + 8*l] =
 8                mult(
 9                  A[n + 2*K*i + 4*k + K*l],
10                  B[n + 2*K*j + 4*k + K*m]
11                );
12            }
13            for (int n = 0;n<4; n++) {
14              C[m + 2*N*i + 2*j + N*l] +=
15                temp[n + 4*m + 8*N*i + 16*j + 8*l];
16            }
17          }
18        }
19      }
20    }
21  }
```

# Vectorisation

$\lambda$ (A, B) $\mapsto$
  A >> split(m) >> map($\lambda$ $nRowsOfA$ $\mapsto$
    B >> split(n) >> map($\lambda$ $mColsOfB$ $\mapsto$
      zip( transpose($nRowsOfA$) >> split(k),
          transpose($mColsOfB$) >> split(k) ) >>
      reduceSeq(init = make2DArray(n,m, 0.0f),
       $\lambda$ ($accTile$, ($tileOfA$, $tileOfB$)) $\mapsto$
        zip($accTile$, transpose($tileOfA$)) >>
       map($\lambda$ ($accRow$, $rowOfTileOfA$) $\mapsto$
        zip($accRow$, transpose($tileOfB$)) >>
       map($\lambda$ ($acc$, $colOfTileOfB$) $\mapsto$
        zip($rowOfTileOfA$ >> asVector(k),
          $colOfTileOfB$ >> asVector(k)) >>
        map(mult4) >> asScalar >>
        reduce($acc$, add)
       ) >> join
      )
     ) >> transpose() >>
    map(transpose) >> transpose
   ) >> join >> transpose
 ) >> join

```
1  for (int i = 0;i<M/2; i++) {
2    for (int j = 0;j<N/2; j++) {
3      for (int k = 0;k<K/4; k++) {
4        for (int l = 0;l<2; l++) {
5          for (int m = 0;m<2; m++) {
6            float4 t = mult4(
7              vload4(A, K*i/2 + k + K*l/4),
8              vload4(B, K*j/2 + k + K*m/4)
9            );
10           vstore4(t, temp, m + 2*N*i + 4*j + 2*l);
11           for (int n = 0;n<4; n++) {
12             C[m + 2*N*i + 2*j + N*l] +=
13               temp[n + 4*m + 8*N*i + 16*j + 8*l];
14           }
15         }
16       }
17     }
18   }
19 }
```

# Mapping parallelism to global threads

$$
\begin{aligned}
&\lambda\ (A,\ B)\ \mapsto \\
&\ A \gg \mathtt{split}(m) \gg \mathrm{mapGlb}_0(\lambda\ nRowsOfA \mapsto \\
&\quad B \gg \mathtt{split}(n) \gg \mathrm{mapGlb}_1(\lambda\ mColsOfB \mapsto \\
&\quad\quad \mathtt{zip}(\ \mathtt{transpose}(nRowsOfA) \gg \mathtt{split}(k), \\
&\quad\quad\quad\quad \mathtt{transpose}(mColsOfB) \gg \mathtt{split}(k)\ ) \gg \\
&\quad\quad \mathtt{reduceSeq}(\mathtt{init} = \mathtt{make2DArray}(n,m,\ 0.0f), \\
&\quad\quad\ \lambda\ (accTile,\ (tileOfA,\ tileOfB)) \mapsto \\
&\quad\quad\ \ \mathtt{zip}(accTile,\ \mathtt{transpose}(tileOfA)) \gg \\
&\quad\quad\ \ \mathrm{mapSeq}(\lambda\ (accRow,\ rowOfTileOfA) \mapsto \\
&\quad\quad\ \ \ \mathtt{zip}(accRow,\ \mathtt{transpose}(tileOfB)) \gg \\
&\quad\quad\ \ \ \mathrm{mapSeq}(\lambda\ (acc,\ colOfTileOfB) \mapsto \\
&\quad\quad\ \ \ \mathtt{zip}(rowOfTileOfA \gg \mathtt{asVector}(k), \\
&\quad\quad\ \ \ \ \ colOfTileOfB \gg \mathtt{asVector}(k)) \gg \\
&\quad\quad\ \ \ \ \mathrm{mapSeq}(\mathrm{mult4}) \gg \mathtt{asScalar} \gg \\
&\quad\quad\ \ \ \ \mathtt{reduceSeq}(acc,\ \mathrm{add}) \\
&\quad\quad\ \ \ ) \gg \mathtt{join} \\
&\quad\quad\ \ ) \\
&\quad\quad\ ) \gg \mathtt{transpose}() \gg \\
&\quad\quad \mathrm{map}(\mathtt{transpose}) \gg \mathtt{transpose} \\
&\quad\ ) \gg \mathtt{join} \gg \mathtt{transpose} \\
&\ ) \gg \mathtt{join}
\end{aligned}
$$

$\longrightarrow$

```
1  int i = get_global_id(0);
2  int j = get_global_id(1);
3  for (int k = 0;k<K/4; k++) {
4    for (int l = 0;l<2; l++) {
5      for (int m = 0;m<2; m++) {
6        float4 t = mult4(
7          vload4(A, K*i/2 + k + K*l/4),
8          vload4(B, K*j/2 + k + K*m/4)
9        );
10       vstore4(t, temp, m + 2*N*i + 4*j + 2*l);
11       for (int n = 0;n<4; n++) {
12         C[m + 2*N*i + 2*j + N*l] +=
13           temp[n + 4*m + 8*N*i + 16*j + 8*l];
14       }
15     }
16   }
17 }
```

# Accumulating in private memory

$$\lambda\ (A,\ B) \mapsto$$
$$A \gg \text{split}(m) \gg \text{mapGlb}_0(\lambda\ nRowsOfA \mapsto$$
$$B \gg \text{split}(n) \gg \text{mapGlb}_1(\lambda\ mColsOfB \mapsto$$
$$\text{zip}(\ \text{transpose}(nRowsOfA) \gg \text{split}(k),$$
$$\text{transpose}(mColsOfB) \gg \text{split}(k)\ ) \gg$$
$$\text{reduceSeq}(\text{init} = \text{make2DArray}(n, m,\ 0.0f) \gg$$
$$\text{toPrivate}(\text{mapSeq}(\text{mapSeq}(id))),$$
$$\lambda\ (accTile,\ (tileOfA,\ tileOfB)) \mapsto$$
$$\text{zip}(accTile,\ \text{transpose}(tileOfA)) \gg$$
$$\text{mapSeq}(\lambda\ (accRow,\ rowOfTileOfA) \mapsto$$
$$\text{zip}(accRow,\ \text{transpose}(tileOfB)) \gg$$
$$\text{mapSeq}(\lambda\ (acc,\ colOfTileOfB) \mapsto$$
$$\text{zip}(rowOfTileOfA \gg \text{asVector}(k),$$
$$colOfTileOfB \gg \text{asVector}(k)) \gg$$
$$\text{mapSeq}(mult4) \gg \text{asScalar} \gg$$
$$\text{reduceSeq}(acc,\ add)$$
$$) \gg \text{join}$$
$$)$$
$$) \gg \text{toGlobal}(\text{mapSeq}(\text{mapSeq}(\text{mapSeq}(id)))$$
$$\gg \text{transpose}() \gg$$
$$\text{map}(\text{transpose}) \gg \text{transpose}$$
$$) \gg \text{join} \gg \text{transpose}$$
$$) \gg \text{join}$$

$\longrightarrow$

```
1  int i = get_global_id(0);
2  int j = get_global_id(1);
3
4  float4 temp_0; float4 temp_1;
5  float4 temp_2; float4 temp_3;
6  float acc_0; float acc_1;
7  float acc_2; float acc_3;
8
9  for (int k = 0;k<K/4; k++) {
10
11    temp_0 = mult4(vload4(k + K*i/2,A),
12       vload4(k + K*j/2,B));
13    acc_0 += temp_0.s0 + temp_0.s1 +
14       temp_0.s2 + temp_0.s3;
15
16    temp_1 = mult4(vload4(k + K*i/2,A,
17       vload4(k + K + 2*K*j/4,B));
18    acc_1 += temp_1.s0 + temp_1.s1 +
19       temp_1.s2 + temp_1.s3;
20
21    temp_2 = mult4(vload4(k + K + 2*K*i/4,A),
22       vload4(k + K*j/2,B));
23    acc_2 += temp_2.s0 + temp_2.s1 +
24       temp_2.s2 + temp_2.s3;
25
26    temp_3 = mult4(vload4(k + K + 2*K*i/4, A),
27       vload4(k + K + 2*K*j/4, B));
28    acc_3 += temp_3.s0 + temp_3.s1 +
29       temp_3.s2 + temp_3.s3;
30  }
31  C[2*N*i + 2*j] = id(acc_0);
32  C[1 + 2*N*i + 2*j] = id(acc_1);
33  C[N + 2*N*i + 2*j] = id(acc_2);
34  C[1 + N + 2*N*i + 2*j] = id(acc_3);
```
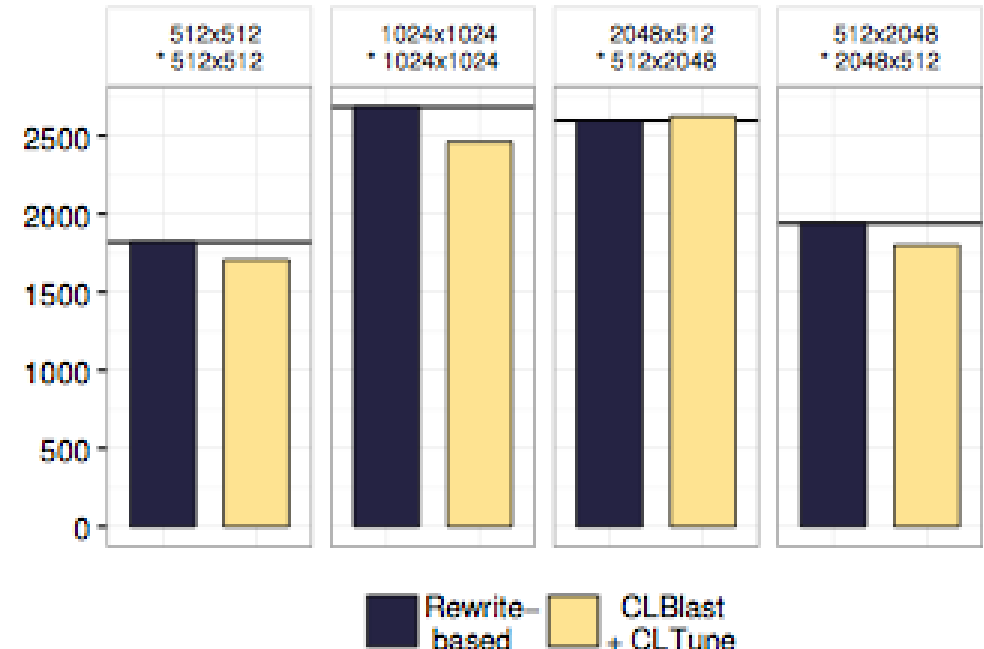
# Performance Portability Achieved

**Compiler input:**

$$
\begin{aligned}
A &>> \text{map}(\lambda \ rowOfA \mapsto \\
B &>> \text{map}(\lambda \ colOfB \mapsto \\
&\quad \text{zip}(rowOfA, \ colOfB) >> \\
&\quad \text{map}(\text{mult}) >> \text{reduce}(0.0\,\text{f}, \text{add}) \\
&) \\
)
\end{aligned}
$$



Desktop GPU (Nvidia GeForce GTX Titan Black) and Desktop GPU (AMD Radeon HD 7970) — GFLOPS comparison of Rewrite-based vs CLBlast + CLTune for matrix sizes 512x512 * 512x512, 1024x1024 * 1024x1024, 2048x512 * 512x2048, and 512x2048 * 2048x512.

# Also works on mobile GPUs

# Easily Extensible

- ► New rules can be added
- ► E.g. dot built-in

# Easily Extensible

- ‣ New rules can be added
- ‣ E.g. dot built-in

```
zip(x,y) >> mapSeq(mult4) >> reduceSeq(z, add4) → dot(x, y)
```

# Easily Extensible
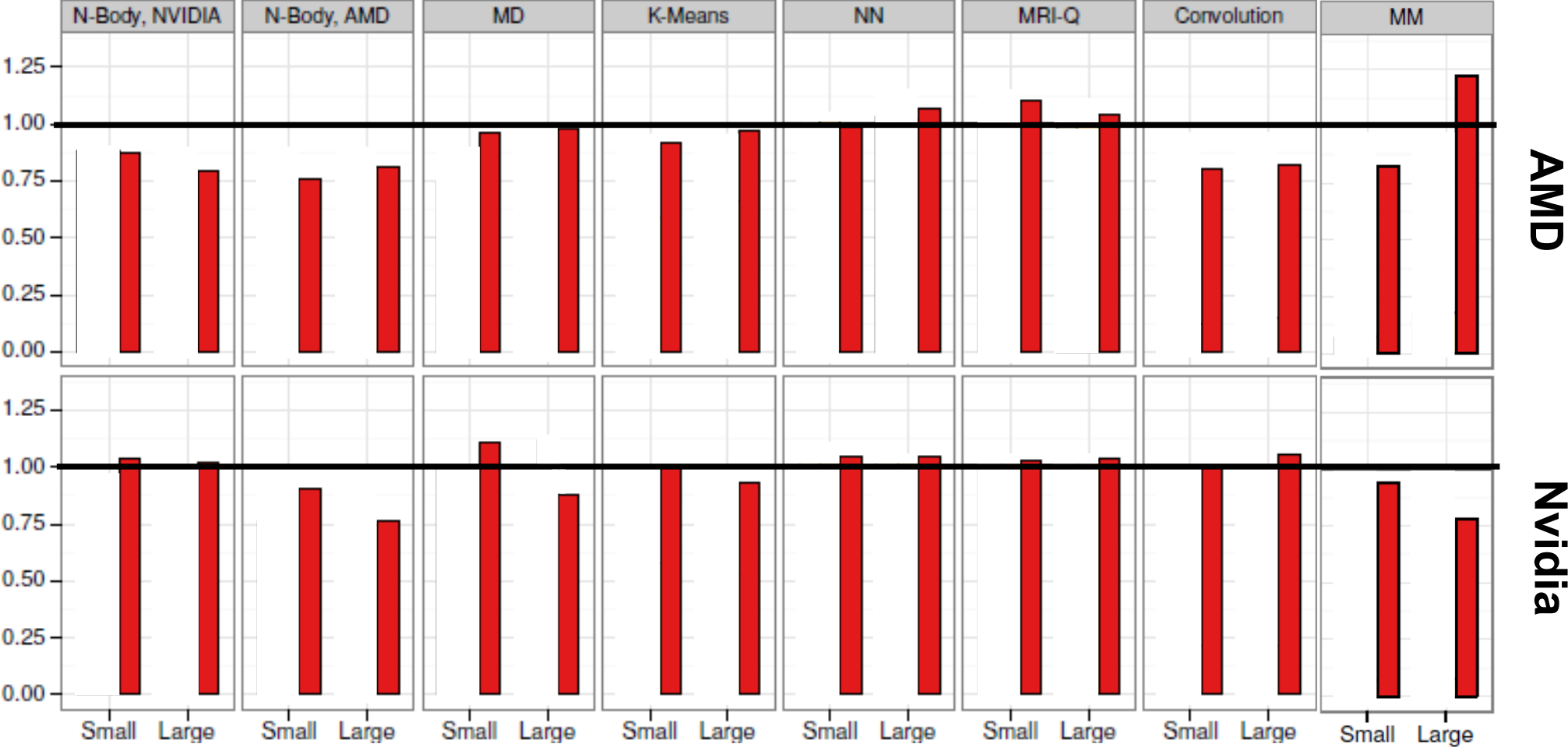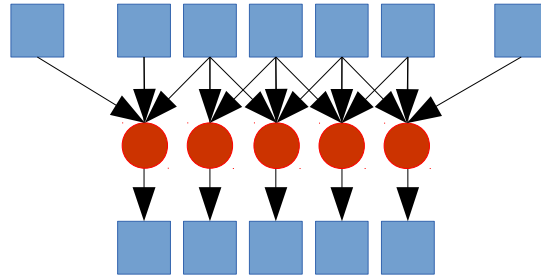
- ‣ New rules can be added
- ‣ E.g. dot built-in

```
zip(x,y) >> mapSeq(mult4) >> reduceSeq(z, add4) → dot(x, y)
```

**Mali GPU**

10.46    + dot    13.58

| 0 | 2 | 4 | 6 | 8 | 10 | 12 | 14 |

GFLOPS

**No need to modify original application!**

**=> Performance portable**

# Works for other programs too

# Stencil Computation



- ► important application
  - · image processing
  - · physic simulations
  - · neural networks
  - · pde solvers

- ► large opportunity for parallelism
  - · GPUs perfect fit
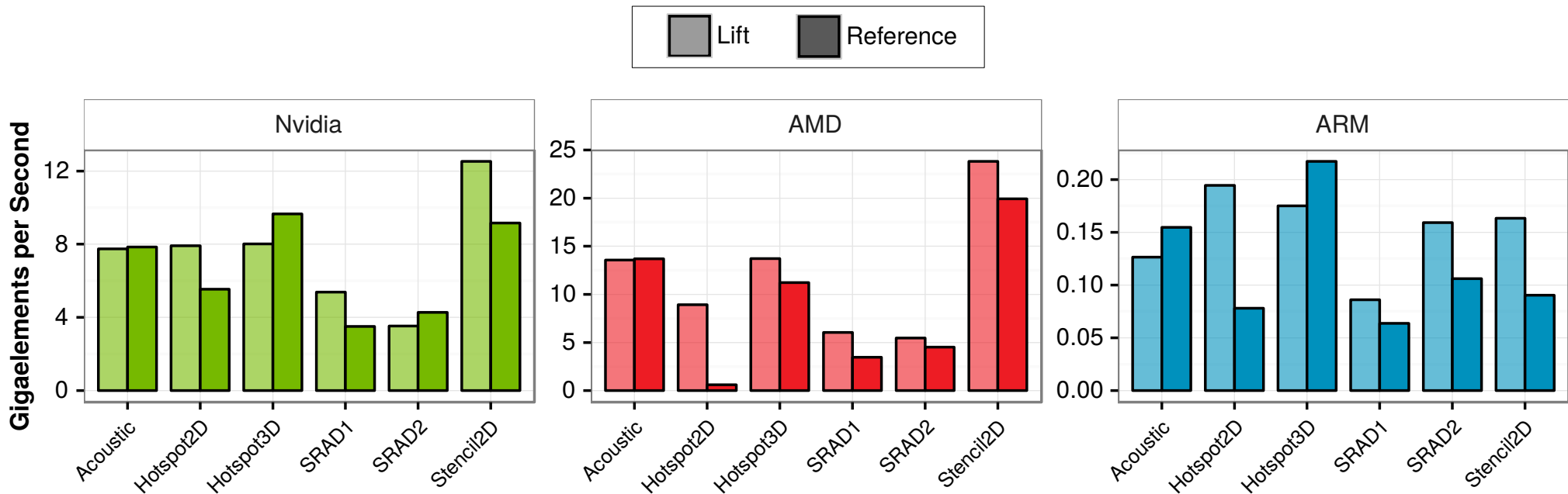
# Stencil Computation in Lift

► use a "slide" primitive

$$\textbf{stencil} = \textbf{slide}(\texttt{step,size}) >> \textbf{map}(\texttt{f})$$



► can work in 2D too

► leverage all existing rewrites

# Stencil Computation Results

- ► No conceptual changes in the compiler
  - just two new primitives (slide,pad)
  - one extra rewrite rules

# Summary

- Rewrite rules define a search space

  - formalisation of algorithmic and optimisation choices

- High performance

  - on par with highly-tuned code

# Future Directions

- How to search efficiently the space
  - machine-learning
- Look at other applications
  - e.g. neural network (convolution, recurrent)
- Support for different target programming models
  - MPI, FPGAs, …

# Purpose of the Tutorial

- ▶ Convince you to take up some of the ideas with you
  - · functional IR
  - · rewriting for optimisation

- ▶ Show you what Lift can do
  - · OpenCL code generation
  - · Rewriting exploration

- ▶ Find new potential collaboration
  - · FPGA backend
  - · MPI work
  - · new application domain
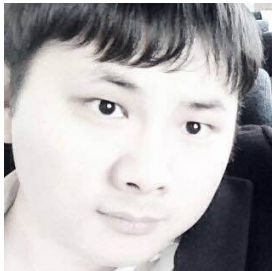
# www.lift-project.org

**Michel Steuwer**
Lectuer
Glasgow Uni.

**Toomas Remmelg**
PhD student
Edinburgh Uni.

**Adam Harris**
PhD student
Edinburgh Uni.

**Bastian Hagedorn**
PhD student
Muenster Uni

**Lu Li**
Postdoc
Edinburgh Uni.

**Federico Pizzuti**
PhD student
Edinburgh Uni.

**Larisa Stolztfus**
PhD student
Edinburgh Uni.

**Naums Mogers**
PhD student
Edinburgh Uni.