

Lift: a Functional Approach to Generating High Performance GPU Code using Rewrite Rules



Toomas
Remmelg



Michel
Steuwer



**Christophe
Dubach**



THE UNIVERSITY *of* EDINBURGH
informatics

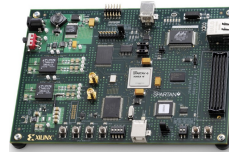
The 4th South of England Regional Programming Language Seminar

27th September 2016

Heterogeneity Everywhere



GPU



FPGA



CPU/GPU



Mobile devices



Game consoles

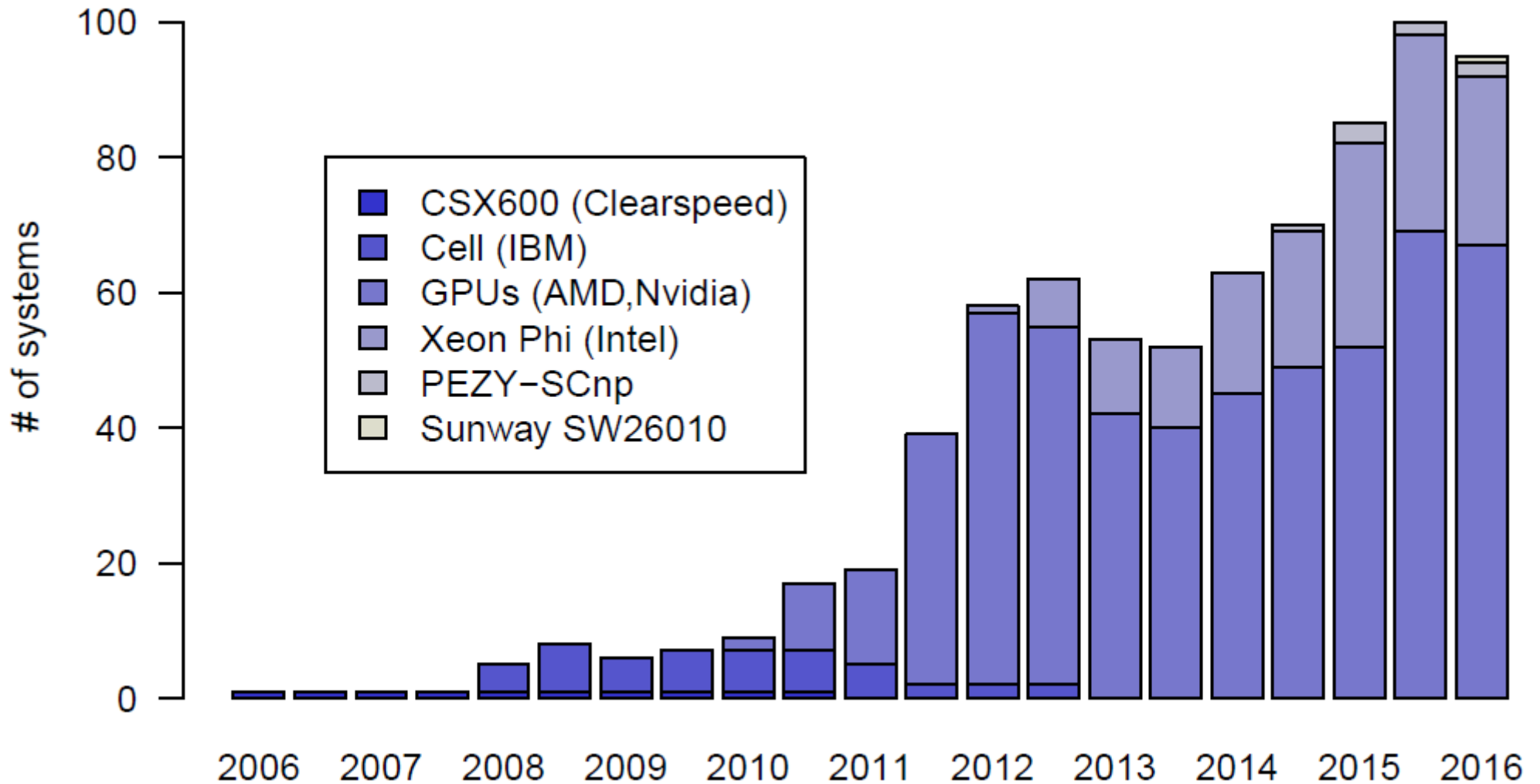


PCs

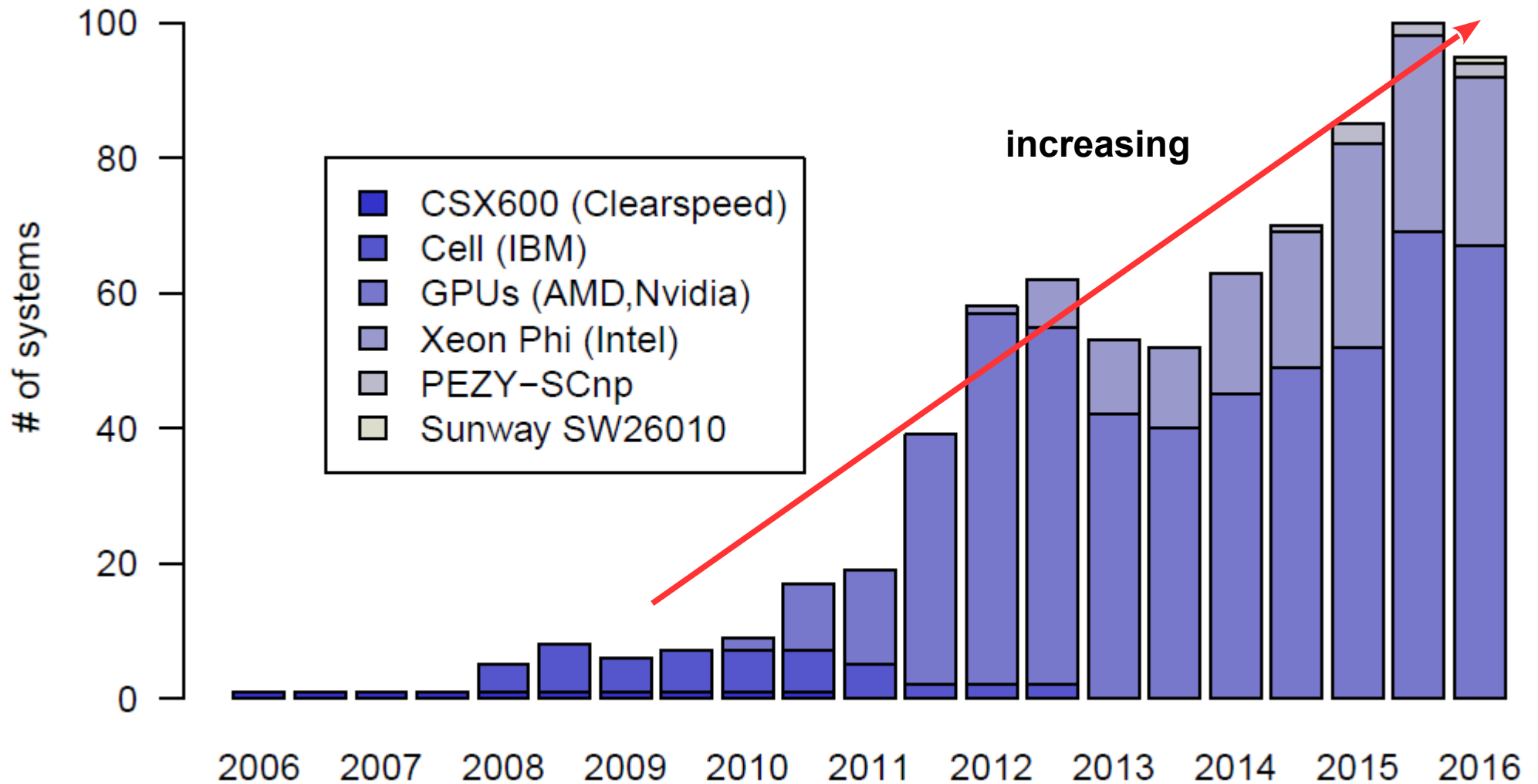


**Data centres
Supercomputers**

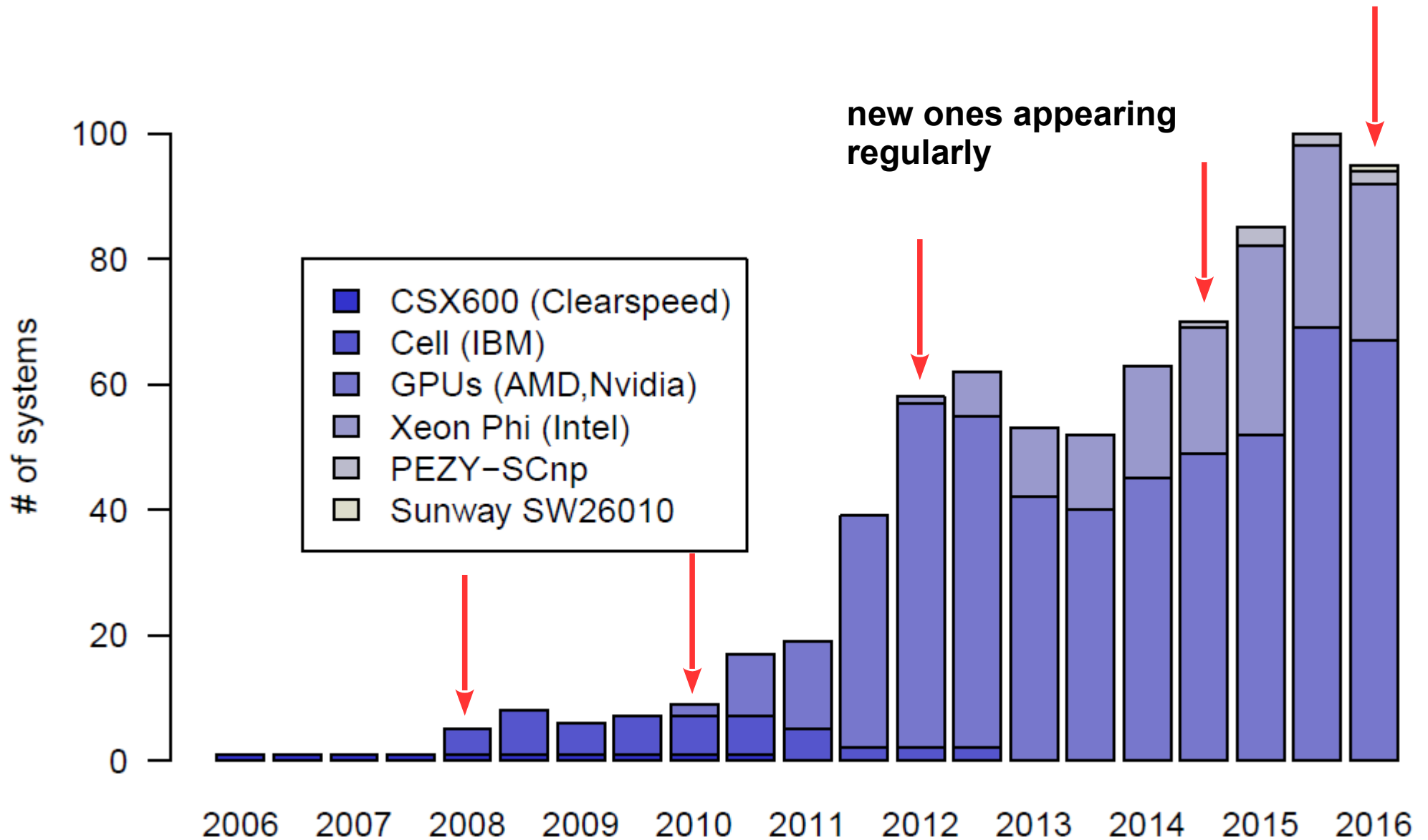
Top 500 with parallel accelerators



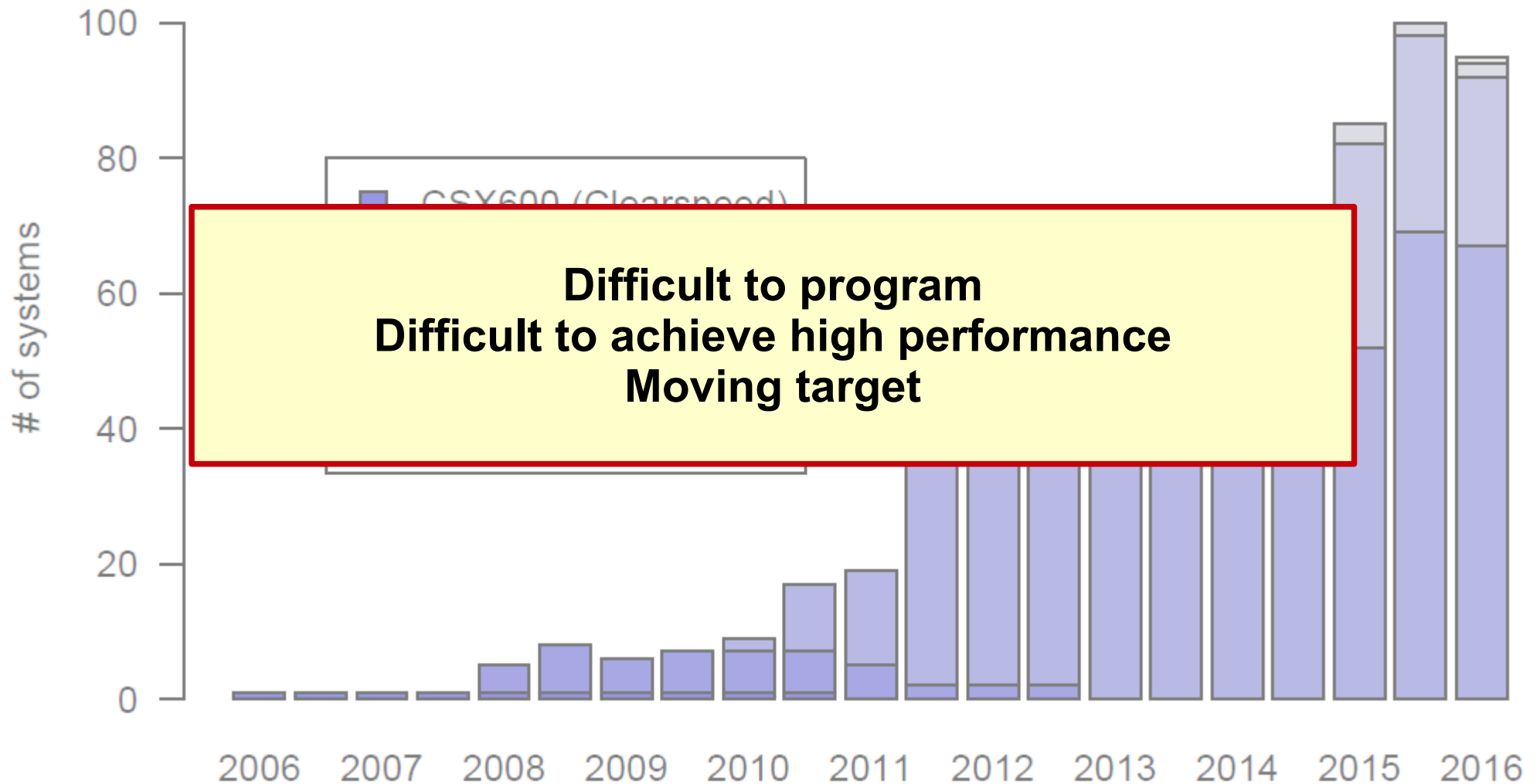
Top 500 with parallel accelerators



Top 500 with parallel accelerators



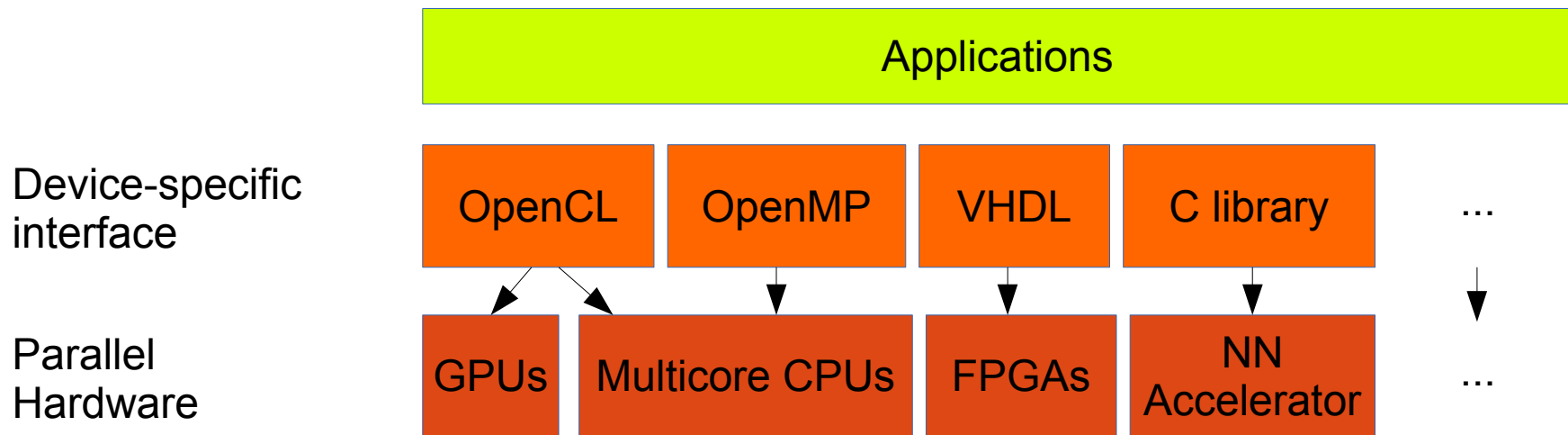
Top 500 with parallel accelerators



But also domain-specific accelerators

- ▶ E.g. Google neural network accelerators
 - TPU (Tensor Processor Units)
- ▶ E.g. Movidius vision accelerators
 - VPU (Vision Processing Units)

Current situation



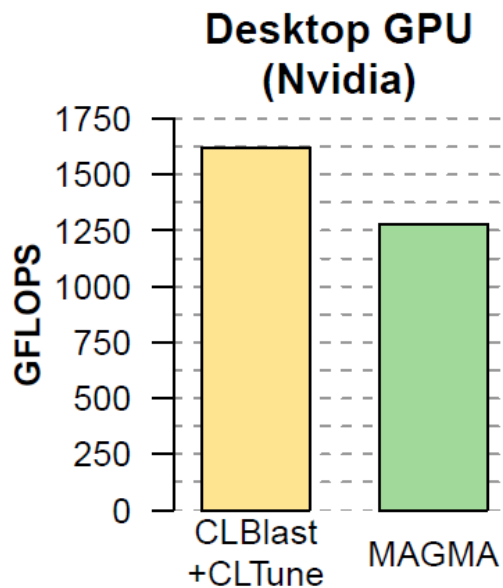
Programmers have to be expert in high performance computing!

+ Performance is not portable!

- ▶ hand-written implementations for each device
 - parametric auto-tuner

+ Performance is not portable even on same device class

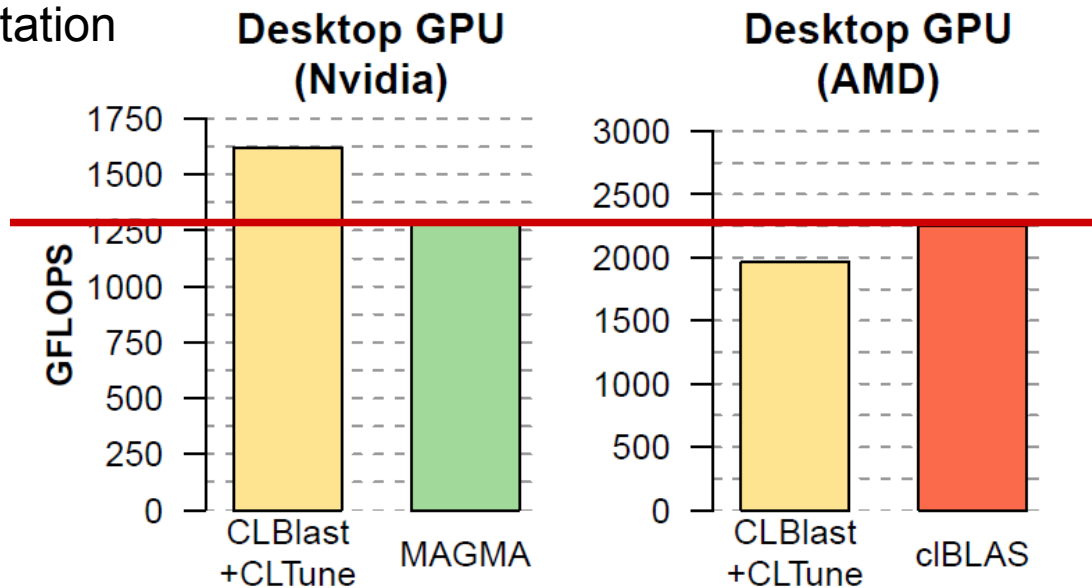
Matrix-matrix multiplication auto-tuner



+ Performance is not portable even on same device class

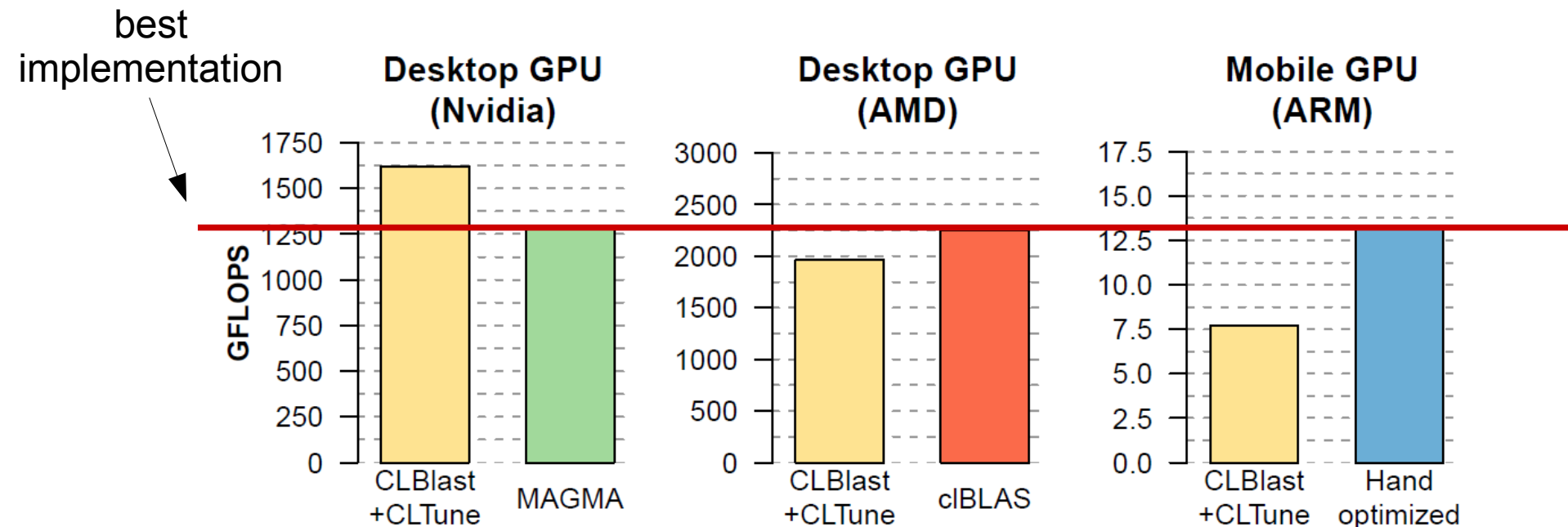
Matrix-matrix multiplication auto-tuner

best implementation



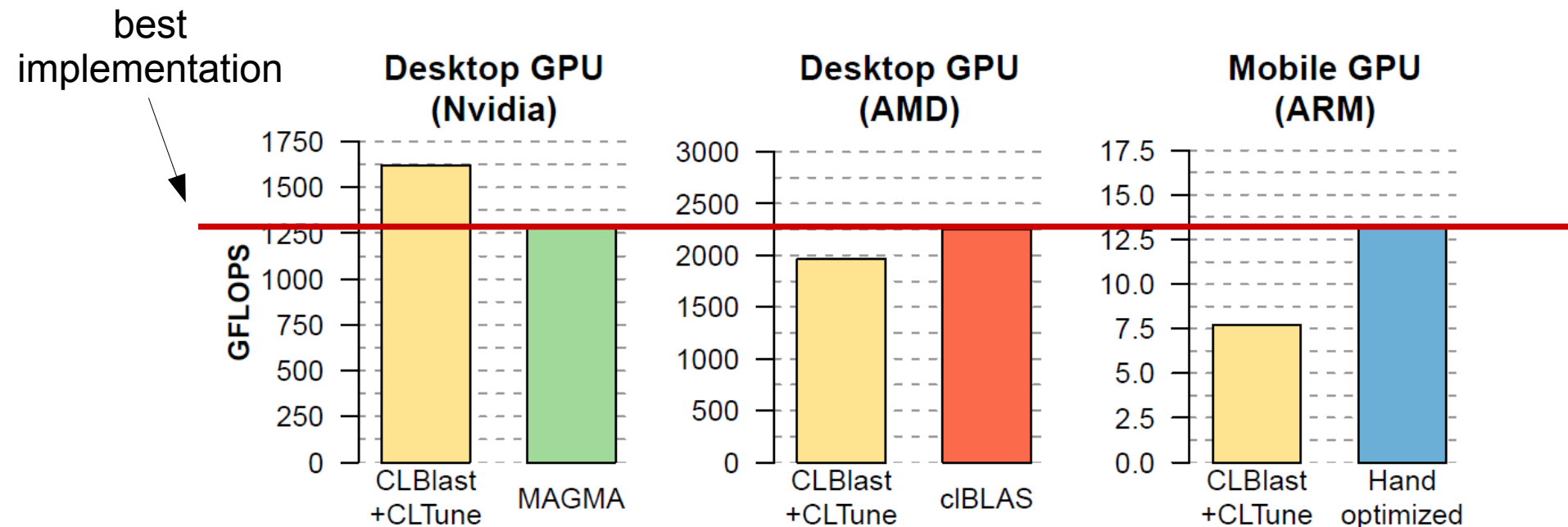
+ Performance is not portable even on same device class

Matrix-matrix multiplication auto-tuner



+ Performance is not portable even on same device class

Matrix-matrix multiplication auto-tuner



Auto-tuning alone fails to achieve portable performance

**Need for high-level programming
+
Code generation**

High-Level programming approaches for heterogeneous devices

- ▶ StreamIt (MIT)
 - Dataflow programming, architecture Independent
 - Multiple backend (tile architecture, C, FPGA, GPU)
- ▶ LiquidMetal (IBM)
 - Java + dataflow programming + (map / reduce)
 - C, OpenCL and FPGA backend
- ▶ Green-Marl (Stanford / Oracle)
 - DSL for graph analysis
 - operations described on graph's node or edges
- ▶ Halide (MIT)
 - DSL for image processing, functional in nature
 - GPU code generator
- ▶ TensorFlow (Google)
 - DSL library for AI applications
 - data flow model

Data flow & functional approaches:

- Hardware complexity is hidden away
 - → code portability
- Empowers the compiler
 - no need for complicated analysis
 - explicit data movement, no global state
 - can be mapped on various hardware
 - scheduling is easier

General purpose framework for data parallelism

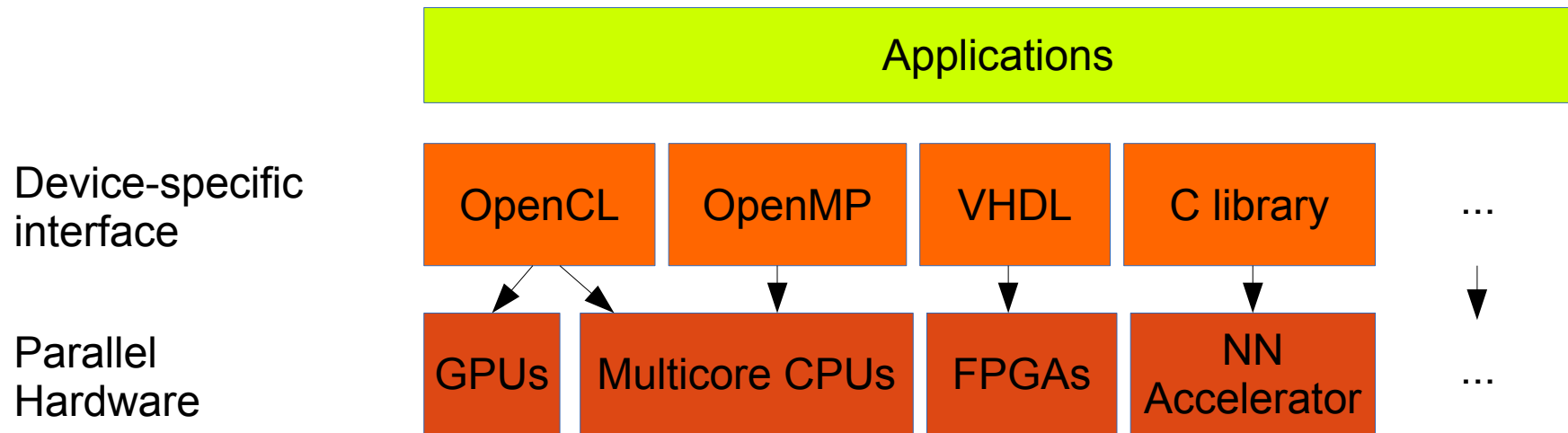
▸ Delite (Stanford + EPFL)

- framework for developing high-performance DSLs
- provide built-in parallel primitives
- Multicore + GPU code generation

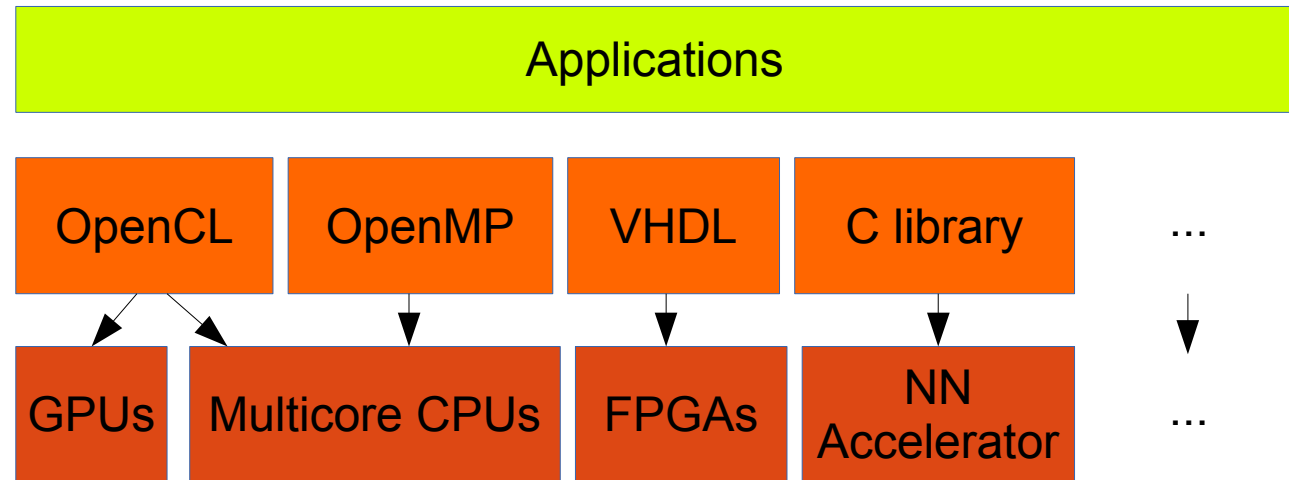
▸ Lift (Edinburgh)

- intermediate data-parallel language
- compiler optimisations expressed as rewrite rules
- Multicore + GPU code generation
- + more to come (MPI, OpenMP, VHDL)

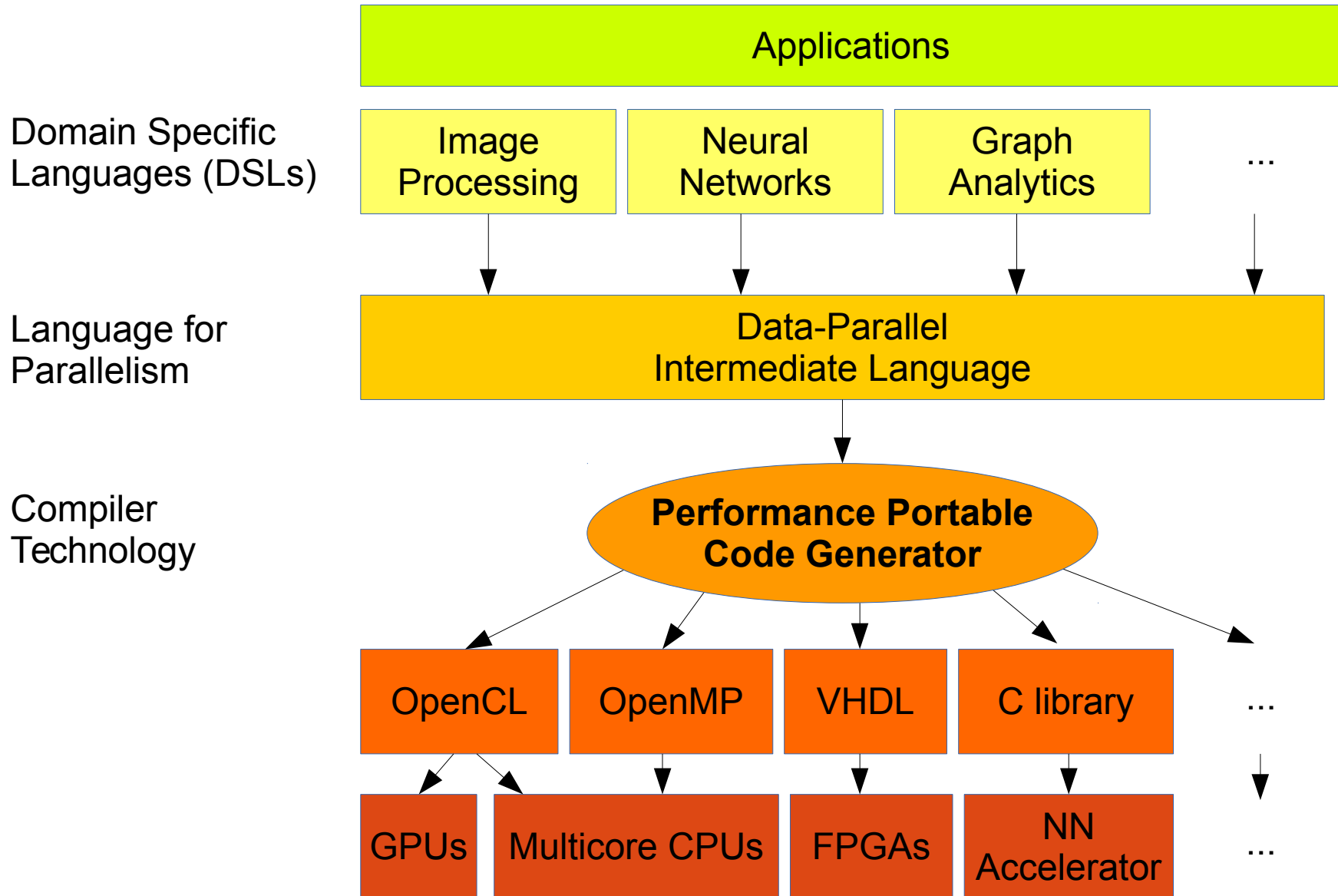
What we have



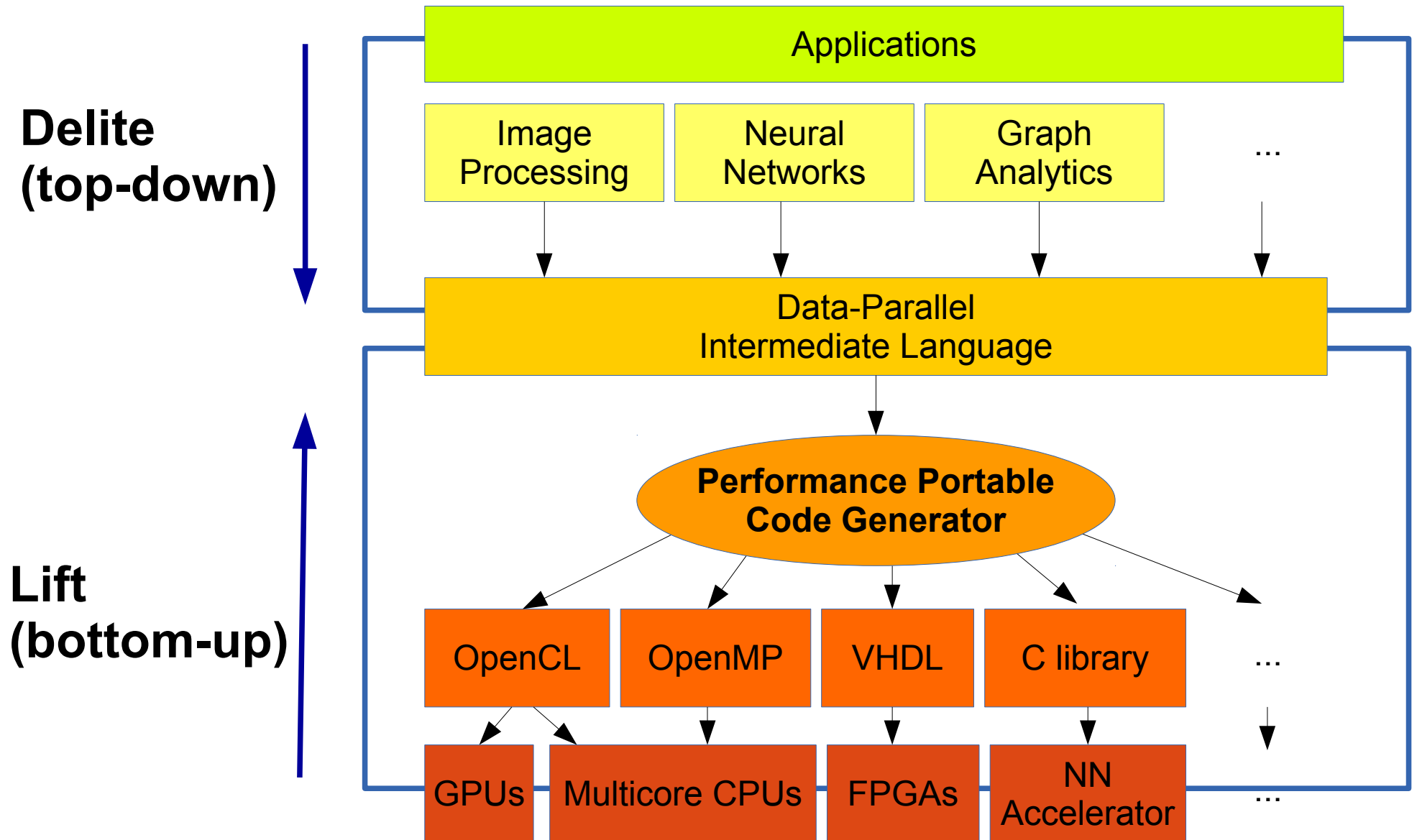
What we need



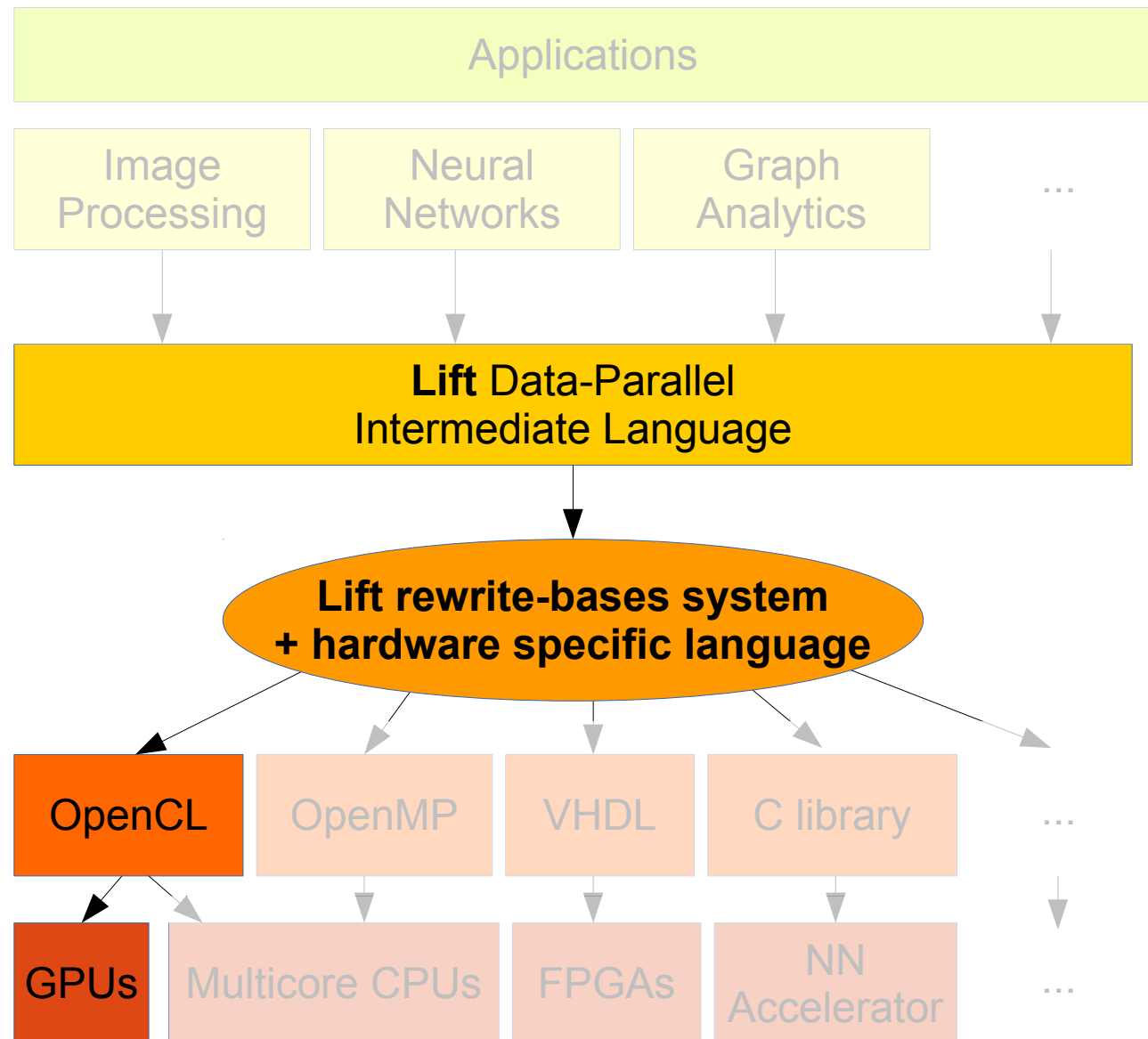
What we need



Different starting point



What this talk is about



Rest of the talk

- ▶ Part I: Lift data parallel language
- ▶ Part II: Optimisations as rewrite rules
- ▶ Part III: Code generation details

Part I

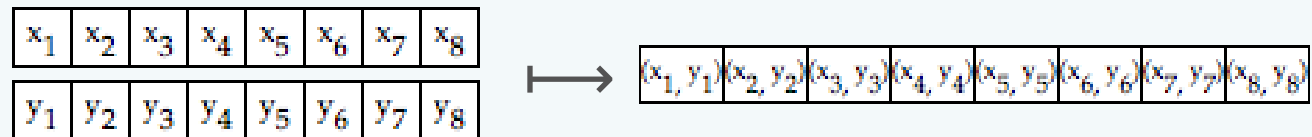
Lift: a Functional Data Parallel Language

Lift Intermediate Language

map(f) :



zip:



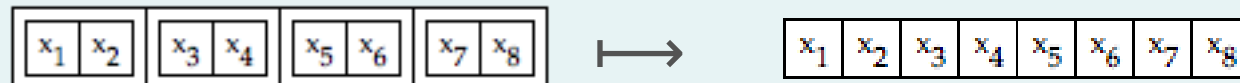
reduce(0,+):



split(n):



join:



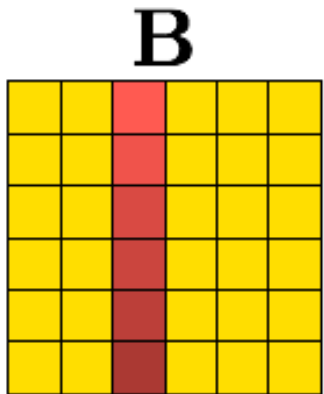
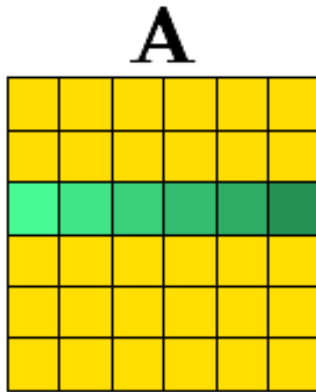
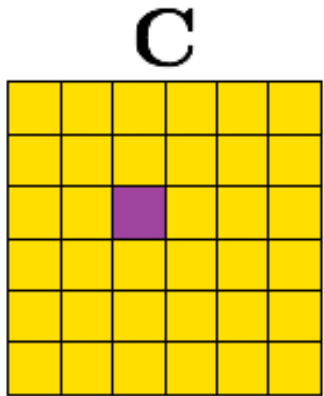
iterate(f, n):



reorder(σ):



Matrix Multiplication in Lift



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

Lift high-level matrix multiplication

```
A >> map( $\lambda$  rowOfA  $\mapsto$ 
  B >> map( $\lambda$  colOfB  $\mapsto$ 
    zip(rowOfA, colOfB) >>
    map(mult) >> reduce(0.0f, add)
  )
)
```



Naive OpenCL version

```
1 kernel void KERNEL(
2   const global float* restrict A,
3   const global float* restrict B
4   global float* C,
5   int M, int K, int N)
6 {
7   float acc = 0.0f;
8
9   for (int i = 0; i < K; i += 1)
10    acc = acc + A[id_A(glb_id_1, i)]
11    * B[id_B(i, glb_id_0)];
12
13   C[(id_C(glb_id_0, glb_id_1)] = acc;
14 }
```

Lift high-level matrix multiplication

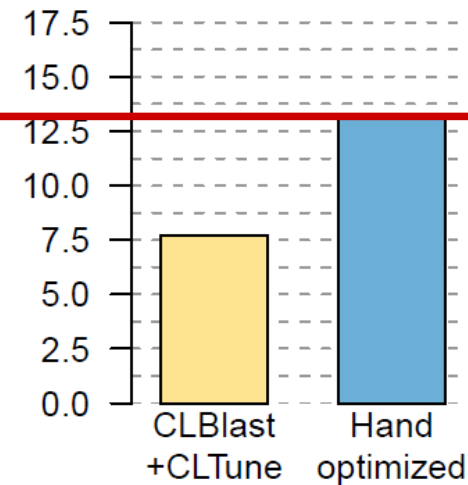
```
A >> map(λ rowOfA ↦  
  B >> map(λ colOfB ↦  
    zip(rowOfA, colOfB) >>  
    map(mult) >> reduce(0.0f, add)  
  )  
)
```

Naive OpenCL version

```
1 kernel void KERNEL(  
2   const global float* restrict A,  
3   const global float* restrict B  
4   global float* C,  
5   int M, int K, int N)  
6 {  
7   float acc = 0.0f;  
8  
9   for (int i = 0; i < K; i += 1)  
10    acc = acc + A[id_A(glb_id_1, i)]  
11    * B[id_B(i, glb_id_0)];  
12  
13   C[id_C(glb_id_0, glb_id_1)] = acc;  
14 }
```

?

Mobile GPU (ARM)



How to achieve high performance?

Lift high-level matrix multiplication

```
A >> map( $\lambda$  rowOfA  $\mapsto$ 
  B >> map( $\lambda$  colOfB  $\mapsto$ 
    zip(rowOfA, colOfB) >>
    map(mult) >> reduce(0.0f, add)
  )
)
```



Naive OpenCL version

```
1 kernel void KERNEL(
2   const global float* restrict A,
3   const global float* restrict B
4   global float* C,
5   int M, int K, int N)
6 {
7   float acc = 0.0f;
8
9   for (int i = 0; i < K; i += 1)
10    acc = acc + A[id_A(glb_id_1, i)]
11    * B[id_B(i, glb_id_0)];
12
13   C[id_C(glb_id_0, glb_id_1)] = acc;
14 }
```



ARM Mali optimised version

```
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);
```

How to achieve high performance?

Part II

Encoding optimisations choices as Rewrite Rules

Algorithmic Rewrite Rules

- Provably correct rewrite rules
- Express algorithmic and optimisations 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:

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

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

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

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

...

```
A >> map( $\lambda$  rowOfA  $\mapsto$ 
  B >> map( $\lambda$  colOfB  $\mapsto$ 
    zip(rowOfA, colOfB) >>
    map(mult) >> reduce(0.0f, add)
  )
)
```



```
A >> map( $\lambda$  rowOfA  $\mapsto$ 
  B >> map( $\lambda$  colOfB  $\mapsto$ 
    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 }
```

```

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) ⇒ Join() ∘ Map(Map(f)) ∘ 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) ⇒ Join() ∘ Map(Map(f)) ∘ 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) ⇒ Join() ∘ Map(Map(f)) ∘ 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 >> split(m) >> map( $\lambda$  rowsOfA  $\mapsto$ 
  B >> map( $\lambda$  colOfB  $\mapsto$ 
    rowsOfA >> map( $\lambda$  rowOfA  $\mapsto$ 
      zip(rowOfA, colOfB) >>
        map(mult) >> reduce(0.0f, add)
    )
  ) >> transpose
) >> 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 >> split(m) >> map(λ rowsOfA ↦  
  B >> split(n) >> map(λ colsOfB ↦  
    colsOfB >> map(λ colOfB ↦  
      rowsOfA >> map(λ rowOfA ↦  
        zip(rowOfA, colOfB) >>  
          map(mult) >> reduce(0.0 f, add)  
      )  
    )  
  ) >> join >> transpose  
) >> join
```



```
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  
7             + K*m] =  
8             mult(A[k + K*l + 2*K*i], B[k + K*  
9               m + 2*K*j]);  
10          }  
11          C[m + 2*j + 2*N*l + 4*N*i] +=  
12            temp[k + 4*K*N*i + 2*K*N*l + 2*  
13              K*j + K*m];  
14        }  
15      }  
16    }
```

Map interchanged

```
A >> split(m) >> map(λ rowsOfA ↦
  B >> split(n) >> map(λ colsOfB ↦
    rowsOfA >> map(λ rowOfA ↦
      colsOfB >> map(λ colOfB ↦
        zip(rowOfA, colOfB) >>
          map(mult) >> reduce(0.0f, add)
      )
    ) >> 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 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
7             + K*m] =
8             mult(A[k + K*l + 2*K*i], B[k + K*
9               m + 2*K*j]);
10          }
11          C[m + 2*j + 2*N*l + 4*N*i] +=
12            temp[k + 4*K*N*i + 2*K*N*l + 2*
13              K*j + K*m];
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, 0.0 f),
        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 }
```

Low-Level OpenCL-specific Extensions

Lift OpenCL-Specific Primitives

map-global(f)

map-workgroup(f)

map-local(f)

map-sequential(f)

reduce-sequential(z,f)

toLocal(f)

toGlobal(f)

vect-n(f)

asScalar

asVector-n

OpenCL thread hierarchy

Sequential implementations

Memory address spaces

Vectorisation

OpenCL-specific rewrites (examples)

OpenCL thread hierarchy:

$\text{map}(f) \implies \text{mapGlb}_{\{0,1,2\}}(f)$

$\text{map}(f) \implies \text{mapLcl}_{\{0,1,2\}}(f)$

OpenCL memory hierarchy:

$f \implies \text{toPrivate}(f)$

$f \implies \text{toLocal}(f)$

$f \implies \text{toGlobal}(f)$

OpenCL vector types and operations:

$\text{map}(f) \implies \text{asVector}(n, b)$
 $\gg \text{map}(\text{vectorize}(n, f)) \gg \text{asScalar}$

**Let's start applying
OpenCL-specific rules**

Vectorisation

```
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 >> 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

```
 $\lambda$  (A, B)  $\mapsto$   
A >> split(m) >> mapGlb0( $\lambda$  nRowsOfA  $\mapsto$   
B >> split(n) >> mapGlb1( $\lambda$  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)) >>  
mapSeq( $\lambda$  (accRow, rowOfTileOfA)  $\mapsto$   
zip(accRow, transpose(tileOfB)) >>  
mapSeq( $\lambda$  (acc, colOfTileOfB)  $\mapsto$   
zip(rowOfTileOfA >> asVector(k),  
    colOfTileOfB >> asVector(k)) >>  
mapSeq(mult4) >> asScalar >>  
reduceSeq(acc, add)  
    ) >> join  
  )  
    ) >> transpose() >>  
  map(transpose) >> transpose  
    ) >> join >> transpose  
  ) >> join
```



```
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

```
λ (A, B) ↦
  A >> split(m) >> mapGlb0(λ nRowsOfA ↦
    B >> split(n) >> mapGlb1(λ mColsOfB ↦
      zip( transpose(nRowsOfA) >> split(k),
          transpose(mColsOfB) >> split(k) ) >>
      reduceSeq( init = make2DArray(n,m, 0.0 f) >>
                  toPrivate( mapSeq( mapSeq( id ) ) ),
        λ (accTile, (tileOfA, tileOfB)) ↦
          zip(accTile, transpose(tileOfA)) >>
          mapSeq(λ (accRow, rowOfTileOfA) ↦
            zip(accRow, transpose(tileOfB)) >>
            mapSeq(λ (acc, colOfTileOfB) ↦
              zip(rowOfTileOfA >> asVector(k),
                  colOfTileOfB >> asVector(k)) >>
              mapSeq(mult4) >> asScalar >>
              reduceSeq(acc, add)
            ) >> join
          )
        ) >> toGlobal( mapSeq( mapSeq( mapSeq( id ) ) )
      >> transpose() >>
      map(transpose) >> transpose
    ) >> join >> transpose
  ) >> join
```



```
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);
```


Result: high performance

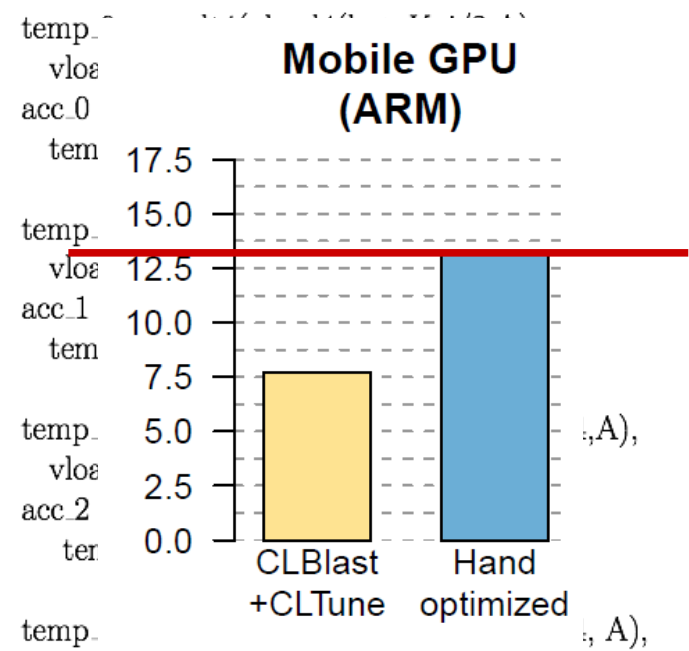
```

λ (A, B) ↦
  A >> split(m) >> mapGlb0(λ nRowsOfA ↦
    B >> split(n) >> mapGlb1(λ mColsOfB ↦
      zip( transpose(nRowsOfA) >> split(k),
          transpose(mColsOfB) >> split(k) ) >>
      reduceSeq( init = make2DArray(n,m, 0.0 f) >>
                  toPrivate( mapSeq( mapSeq( id ) ) ) ,
        λ (accTile, (tileOfA, tileOfB)) ↦
          zip( accTile, transpose(tileOfA) ) >>
          mapSeq( λ (accRow, rowOfTileOfA) ↦
            zip( accRow, transpose(tileOfB) ) >>
            mapSeq( λ (acc, colOfTileOfB) ↦
              zip( rowOfTileOfA >> asVector(k),
                  colOfTileOfB >> asVector(k) ) >>
              mapSeq( mult4 ) >> asScalar >>
              reduceSeq( acc, add )
            ) >> join
          )
        ) >> toGlobal( mapSeq( mapSeq( mapSeq( id ) ) )
      >> transpose() >>
      map( transpose ) >> transpose
    ) >> join >> transpose
  ) >> join
  
```



```

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 = ...;
12   vload ...
13   acc_0 = ...
14   temp_1 = ...
15   temp_2 = ...
16   vload ...
17   acc_1 = ...
18   temp_3 = ...
19   vload ...
20   temp_0 = ...
21   temp_1 = ...
22   vload ...
23   acc_2 = ...
24   temp_2 = ...
25
26   temp_3 = ...
27   vload ...
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);
  
```



Now just need to search the space

...

100,000 runs later

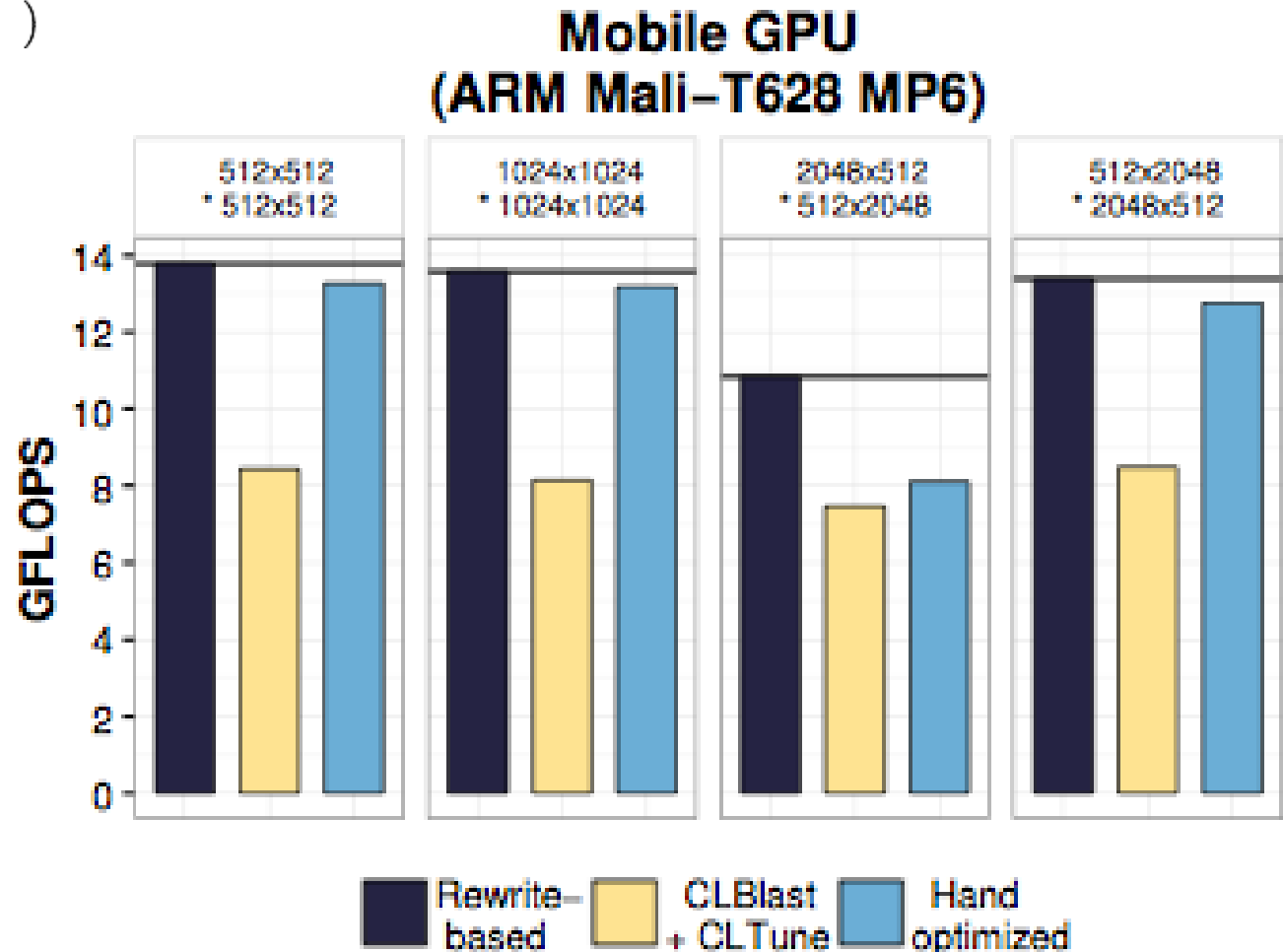
...

Performance Portability Achieved

Compiler input:

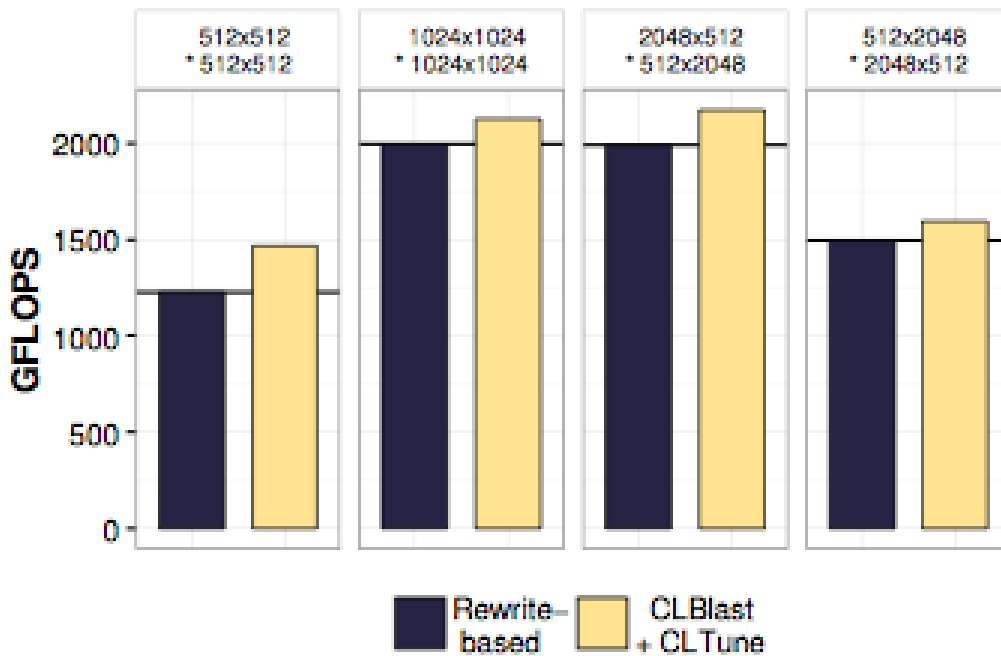
```
A >> map( $\lambda$  rowOfA  $\mapsto$   
B >> map( $\lambda$  colOfB  $\mapsto$   
zip(rowOfA, colOfB) >>  
map(mult) >> reduce(0.0f, add)  
)  
)
```

Search result:

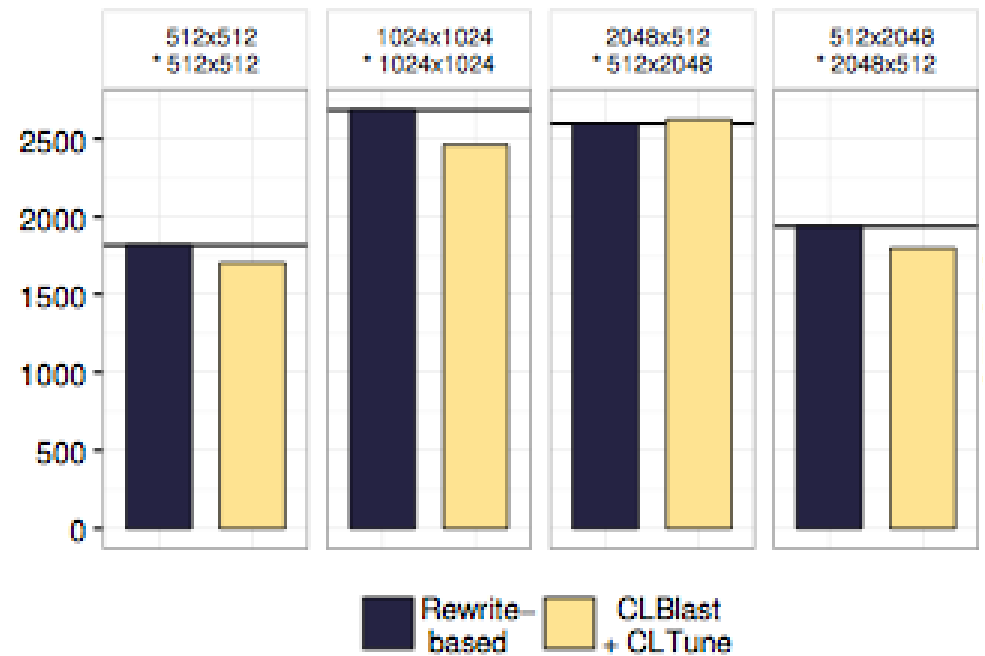


Desktop GPUs

Desktop GPU (Nvidia GeForce GTX Titan Black)



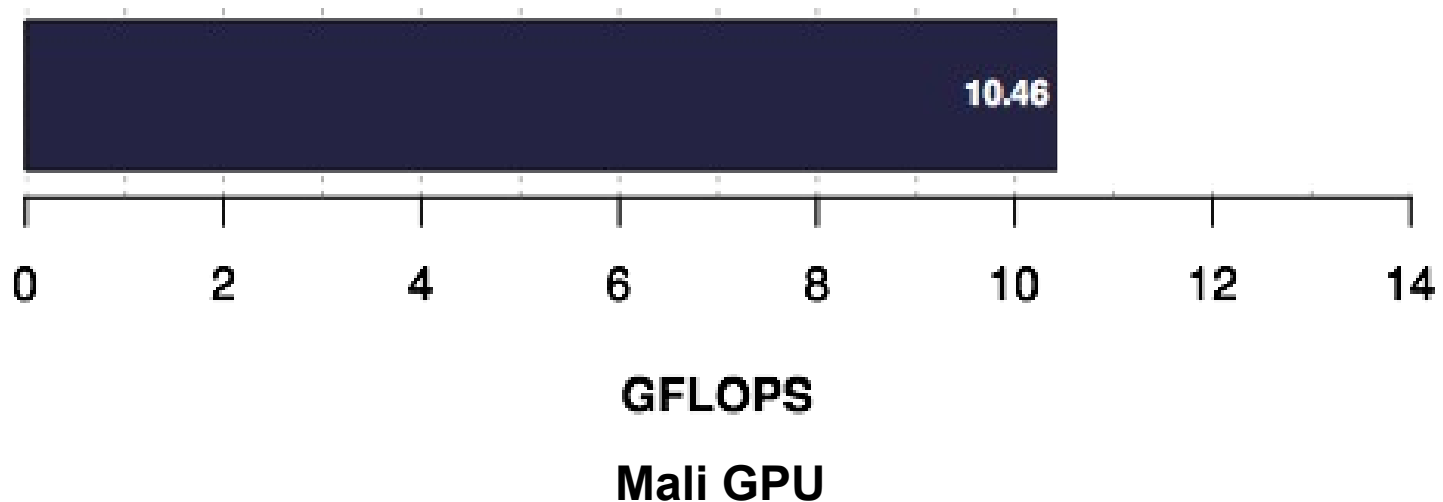
Desktop GPU (AMD Radeon HD 7970)



Easily Extensible

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

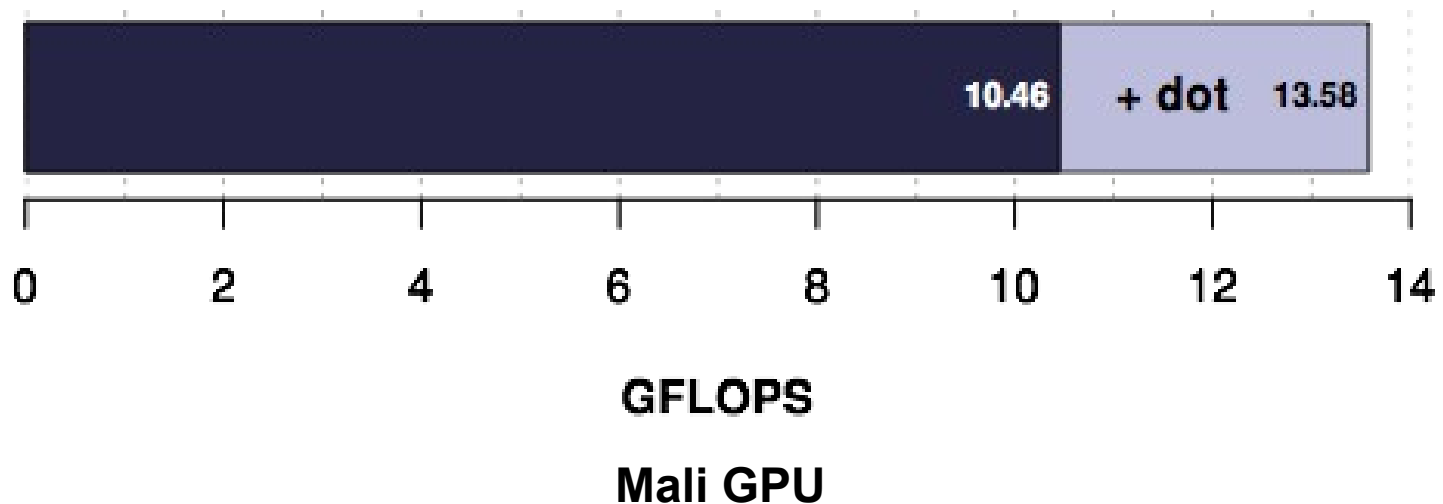
```
reduceSeq(z, add4) o mapSeq(mult4) o zip(x,y)  
=>  
dot(x, y)
```



Easily Extensible

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

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



Part III

Efficient Code Generation

Dot-product example

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


Dot-product example

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

after rewrite:

```
(join ◦ mapWrg0(
  join ◦ toGlobal(mapLcl0(mapSeq(id))) ◦ split1 ◦
  iterate6( join ◦
    mapLcl0( toLocal(mapSeq(id)) ◦
      reduceSeq(add, 0) ) ◦
    split2 ) ◦
  join ◦ mapLcl0( toLocal(mapSeq(id)) ◦
    reduceSeq(multAndSumUp, 0) ) ◦ split2
) ◦ split128)( zip(x, y) )
```

Dot-product example

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

after rewrite:

```
(join ◦ mapWrg0(  
  join ◦ toGlobal(mapLcl0(mapSeq(id))) ◦ split1 ◦  
  iterate6( join ◦  
    mapLcl0( toLocal(mapSeq(id)) ◦  
      reduceSeq(add, 0) ) ◦  
    split2 ) ◦
```

```
  join ◦ mapLcl0( toLocal(mapSeq(id)) ◦  
    reduceSeq(multAndSumUp, 0) ) ◦ split2
```

**copy to local
memory**

```
) ◦ split128)( zip(x, y) )
```

Dot-product example

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

after rewrite:

```
(join ◦ mapWrg0(  
  join ◦ toGlobal(mapLcl0(mapSeq(id))) ◦ split1 ◦  
  iterate6( join ◦  
    mapLcl0( toLocal(mapSeq(id)) ◦  
      reduceSeq(add, 0) ) ◦  
    split2 ) ◦  
  join ◦ mapLcl0( toLocal(mapSeq(id)) ◦  
    reduceSeq(multAndSumUp, 0) ) ◦ split2  
) ◦ split128)( zip(x, y) )
```

**iterative
reduction in
local memory**

Dot-product example

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

after rewrite:

```
(join ◦ mapWrg0(  
  join ◦ toGlobal(mapLcl0(mapSeq(id))) ◦ split1 ◦ write back to  
  iterate6( join ◦ global mem.  
    mapLcl0( toLocal(mapSeq(id)) ◦  
      reduceSeq(add, 0) ) ◦  
    split2 ) ◦  
  join ◦ mapLcl0( toLocal(mapSeq(id)) ◦  
    reduceSeq(multAndSumUp, 0) ) ◦ split2  
) ◦ split128)( zip(x, y) )
```

How to generate efficient OpenCL code?

```
(join ◦ mapWrg0(
  join ◦ toGlobal(mapLcl0(mapSeq(id))) ◦ split1 ◦
  iterate6( join ◦
    mapLcl0( toLocal(mapSeq(id)) ◦
      reduceSeq(add, 0) ) ◦
    split2 ) ◦
  join ◦ mapLcl0( toLocal(mapSeq(id)) ◦
    reduceSeq(multAndSumUp, 0) ) ◦ split2
) ◦ split128)( zip(x, y) )
```

→ Pattern based code generator

map-global (f,input)

```
for (uint i=get_global_id;  
     i<n;  
     i+= get_global_size) {  
    output[i] = f(input[i]);  
}
```

...

map-sequential (f,input)

```
for (uint i=0; i<n; i++) {  
    output[i] = f(input[i]);  
}
```

reduce-sequential (f,z,input)

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

What about split, zip, join, ... ?

```
(join ◦ mapWrg0 (
  join ◦ toGlobal (mapLcl0 (mapSeq(id))) ◦ split1 ◦
  iterate6 ( join ◦
    mapLcl0 ( toLocal (mapSeq(id)) ◦
      reduceSeq(add, 0) ) ◦
    split2 ) ◦
  join ◦ mapLcl0 ( toLocal (mapSeq(id)) ◦
    reduceSeq(multAndSumUp, 0) ) ◦ split2
) ◦ split128) ( zip(x, y) )
```

- ▶ Need to avoid temporary results
- ▶ split, zip, ... merely just change data layout
 - → “lazy evaluation”

dot-product (a tiny bit of it):

```
(join ◦ mapWrg0( ...  
  join ◦ mapLcl0( ...  
    reduceSeq(λ(a, xy) ↦ a + (xy0 * xy1), 0)) ◦ split2  
  ) ◦ split128)( zip(x, y) )
```

desired OpenCL output:

```
int wg_id = get_group_id(0) ;  
int l_id = get_local_id(0) ;  
...  
acc = 0.0f;  
for (int i = 0 ; i < 2 ; I += 1) {  
  acc = acc +  
    x [ 2 * l_id + 128 * wg_id + I ] *  
    y [ 2 * l_id + 128 * wg_id + I ] ;  
}
```


Construct a representation of the effects of data layout functions:

Views

Consumes the views to generate correct array indices

View Construction

```
(join ◦ mapWrg0( ...  
  join ◦ mapLcl0( ...  
    reduceSeq(λ(a, xy) ↦ a + (xy0 * xy1), 0)) ◦ split2  
  ) ◦ split128)( zip(x, y) )
```

View Construction

```
(join ◦ mapWrg0( ...  
  join ◦ mapLcl0( ...  
    reduceSeq(λ(a, xy) ↦ a + (xy0 * xy1), 0)) ◦ split2  
  ) ◦ split128)( zip(x, y) )
```

TupleAccessView(0)



View Construction

```
(join ◦ mapWrg0( ...  
  join ◦ mapLcl0( ...  
    reduceSeq(λ(a, xy) ↦ a + (xy0 * xy1), 0)) ◦ split2  
  ) ◦ split128)( zip(x, y) )
```

TupleAccessView(0)

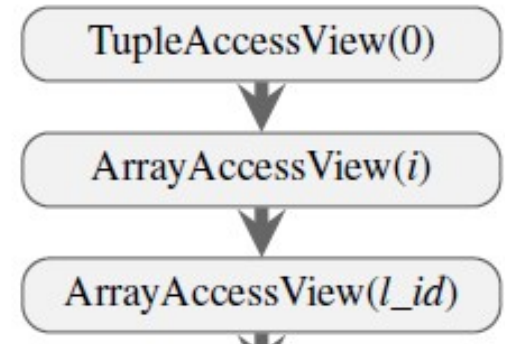


ArrayAccessView(i)



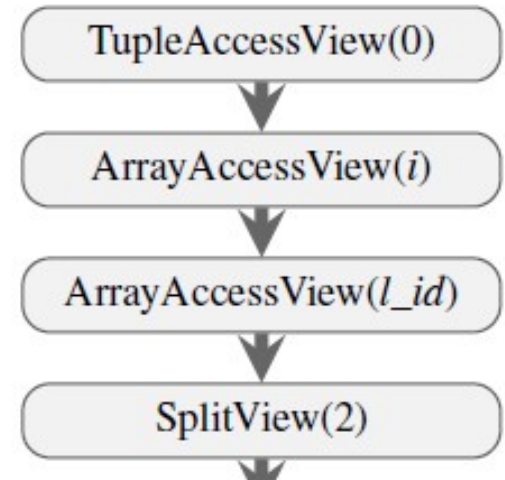
View Construction

```
(join ◦ mapWrg0( ...  
  join ◦ mapLcl0( ...  
    reduceSeq(λ(a, xy) ↦ a + (xy0 * xy1), 0)) ◦ split2  
  ) ◦ split128)( zip(x, y) )
```



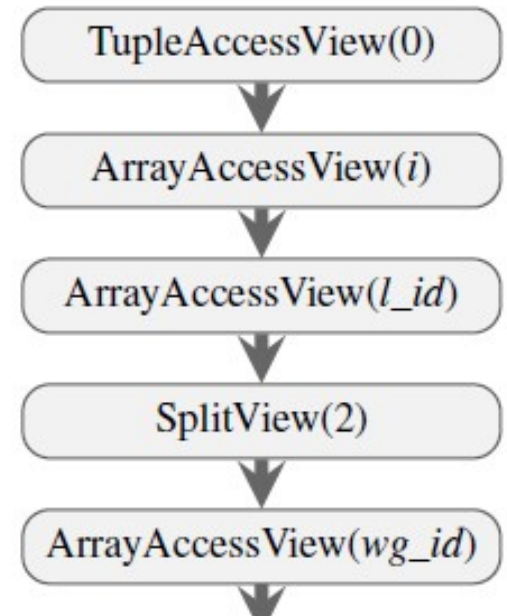
View Construction

```
(join ◦ mapWrg0( ...  
  join ◦ mapLcl0( ...  
    reduceSeq(λ(a, xy) ↦ a + (xy0 * xy1), 0)) ◦ split2  
  ) ◦ split128)( zip(x, y) )
```



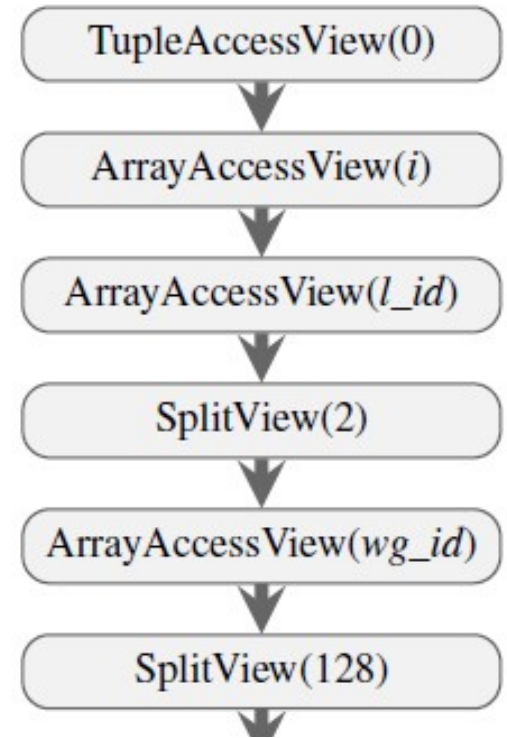
View Construction

```
(join ◦ mapWrg0( ...  
  join ◦ mapLcl0( ...  
    reduceSeq( $\lambda(a, xy) \mapsto a + (xy_0 * xy_1), \emptyset$ ) ◦ split2  
  ) ◦ split128)( zip(x, y) )
```



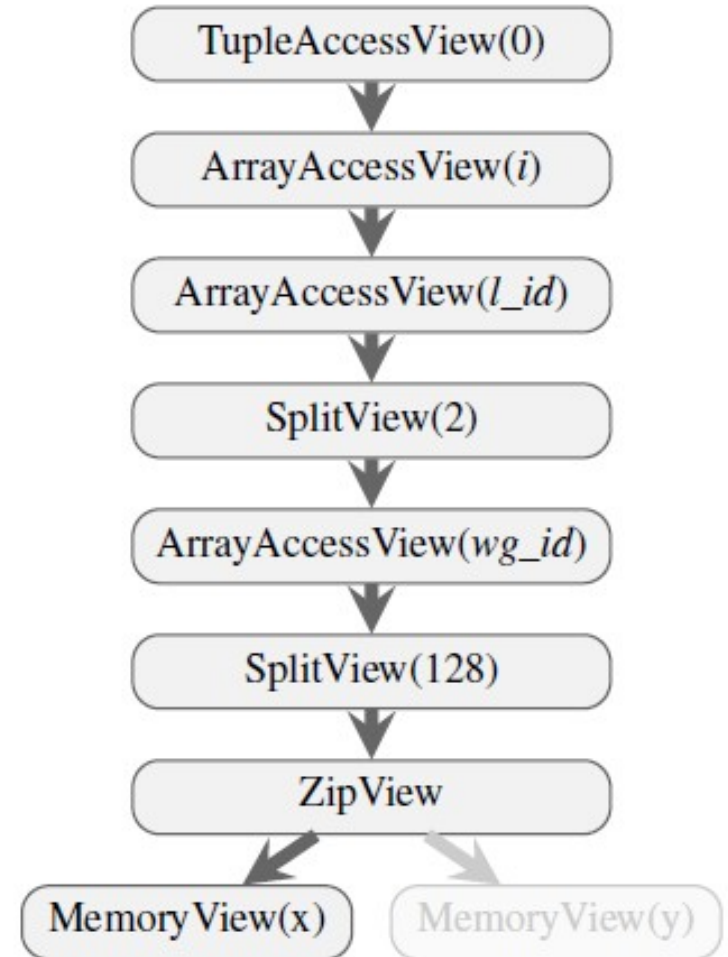
View Construction

```
(join ◦ mapWrg0( ...  
  join ◦ mapLcl0( ...  
    reduceSeq(λ(a, xy) ↦ a + (xy0 * xy1), 0)) ◦ split2  
  ) ◦ split128)( zip(x, y) )
```



View Construction

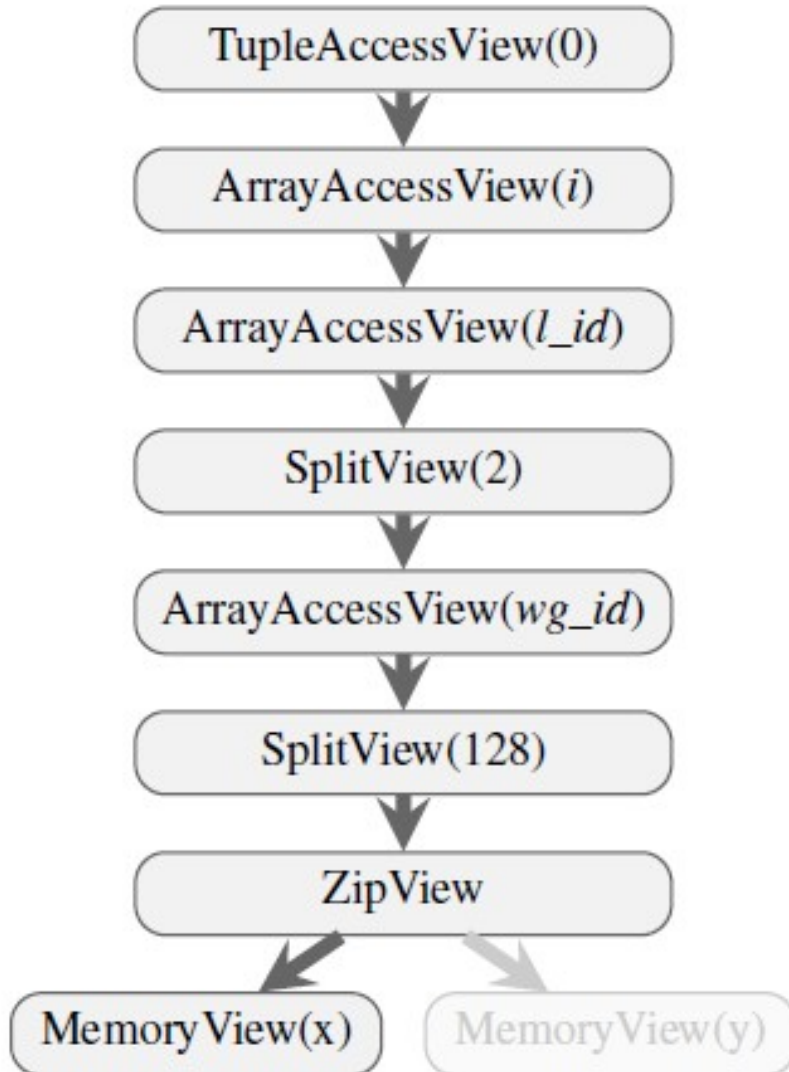
```
(join ◦ mapWrg0( ...  
  join ◦ mapLcl0( ...  
    reduceSeq(λ(a, xy) ↦ a + (xy0 * xy1), 0)) ◦ split2  
  ) ◦ split128)( zip(x, y) )
```



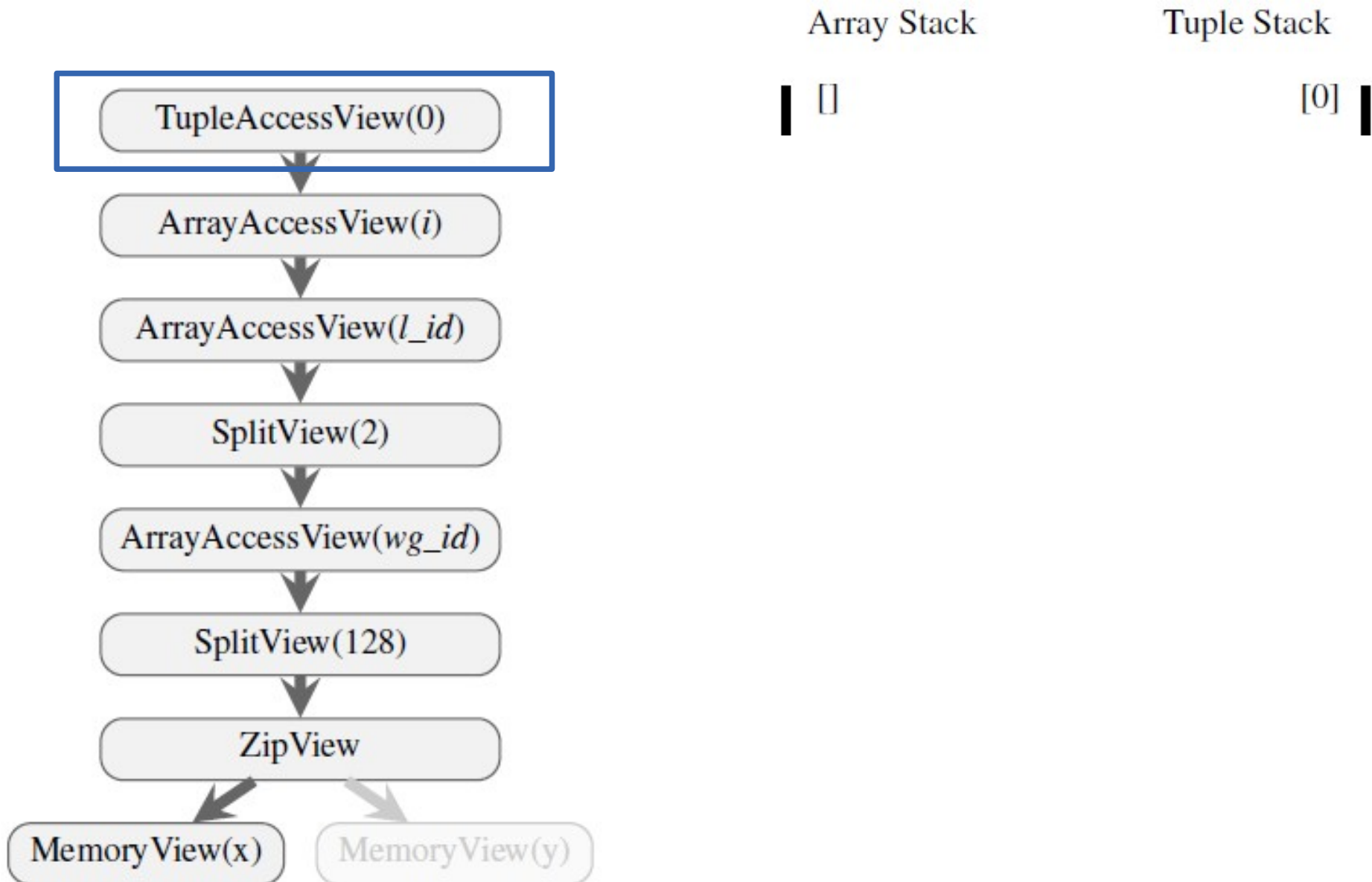
View Consumption

Array Stack

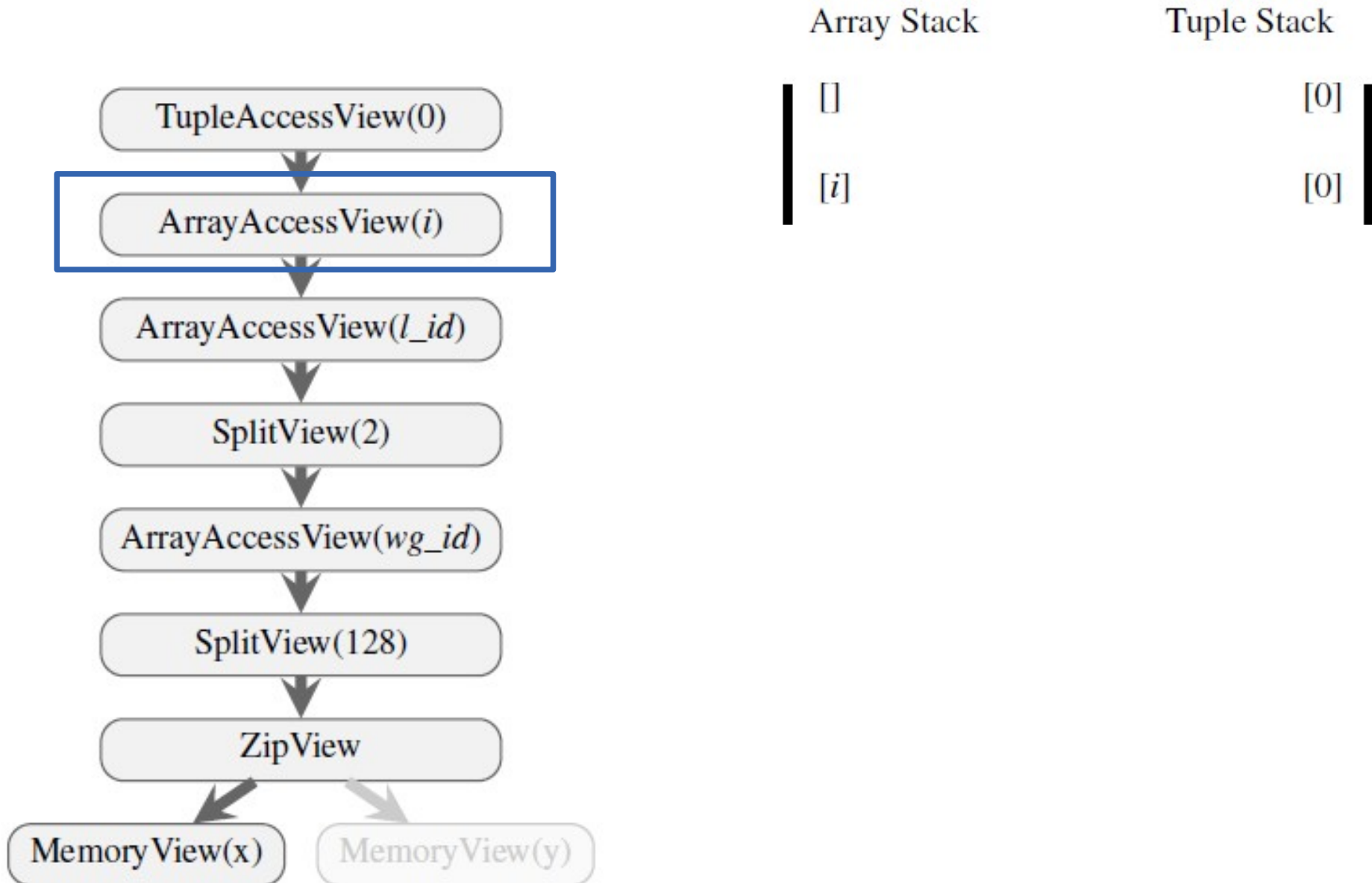
Tuple Stack



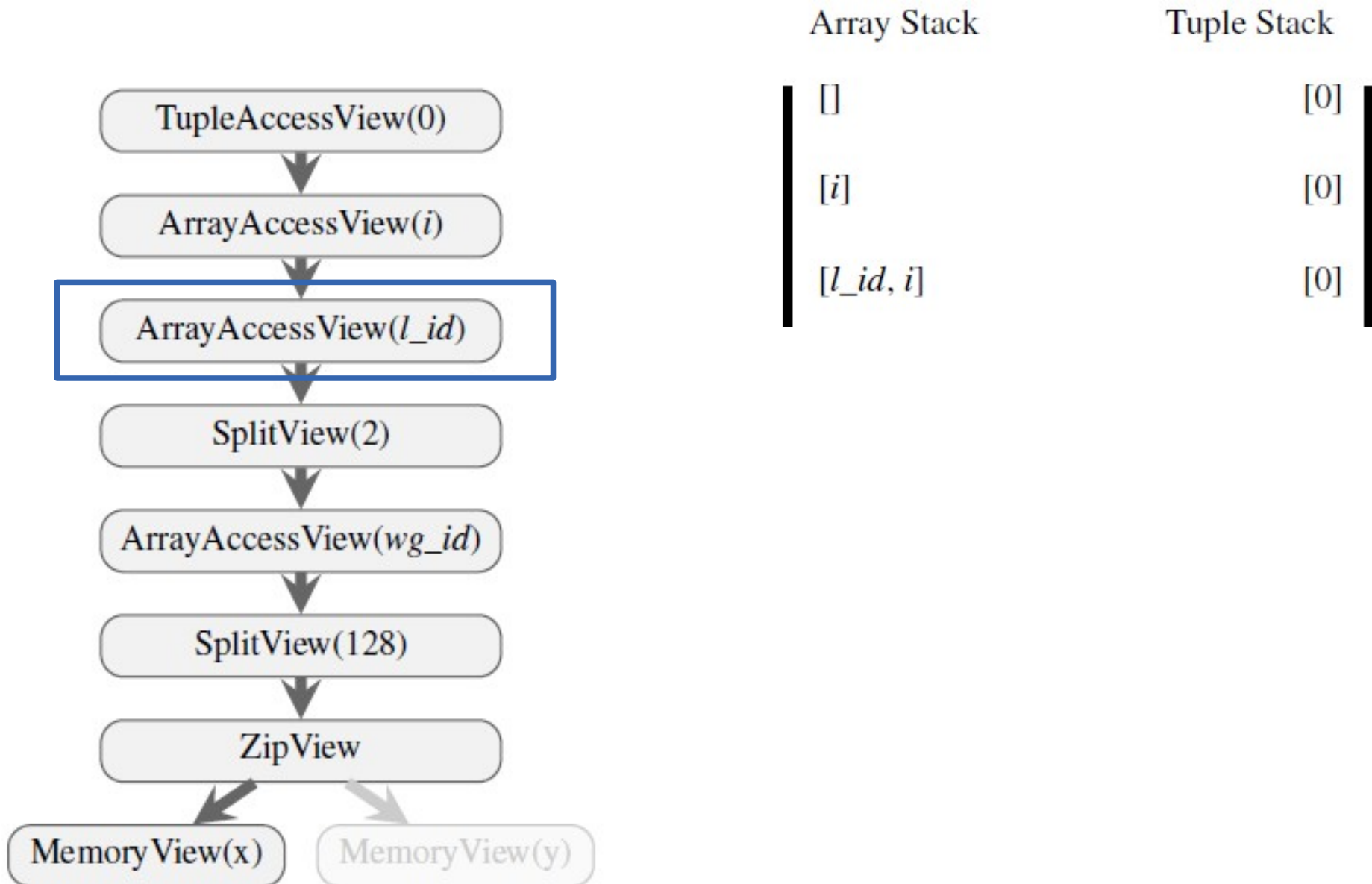
View Consumption



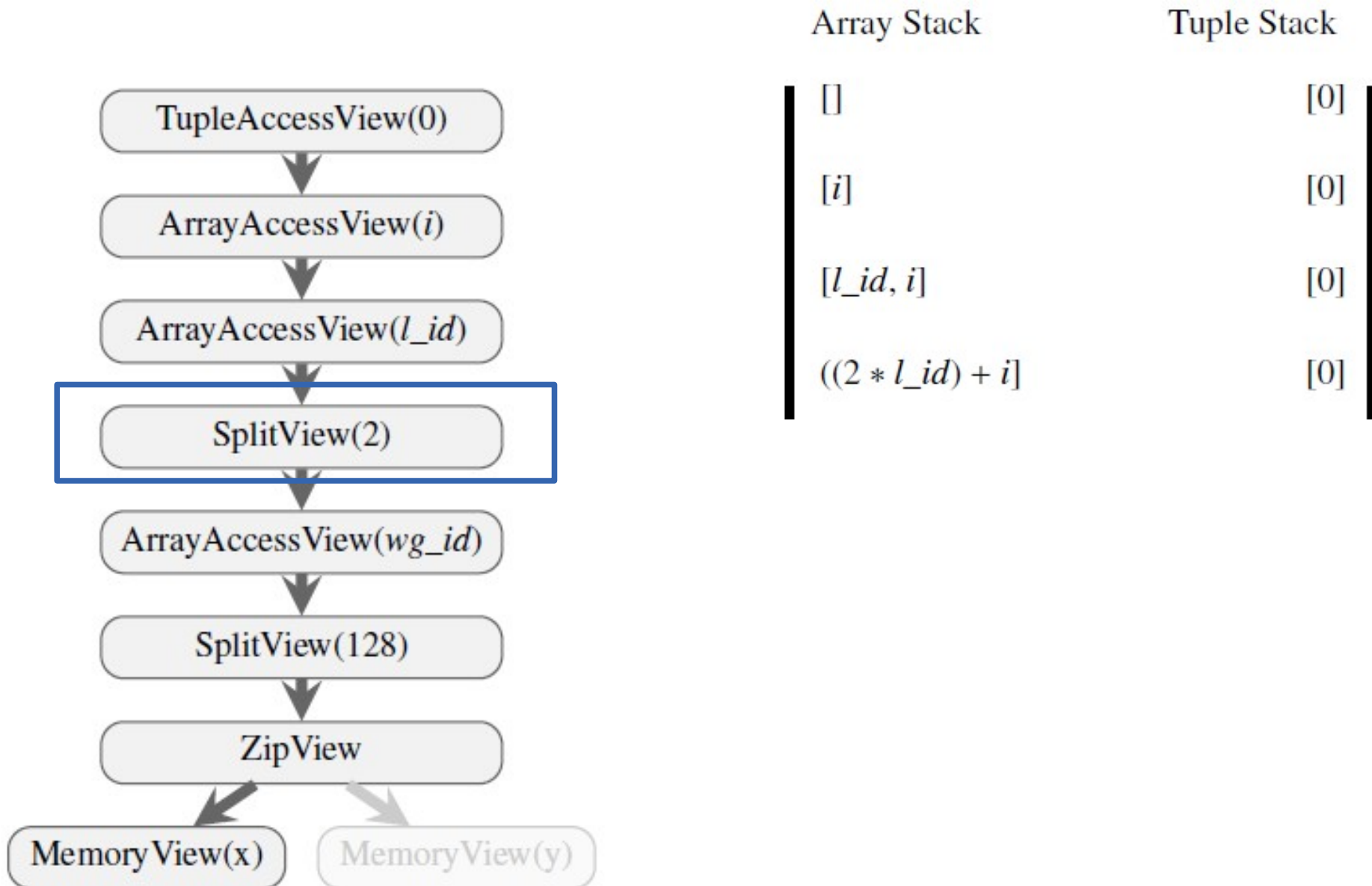
View Consumption



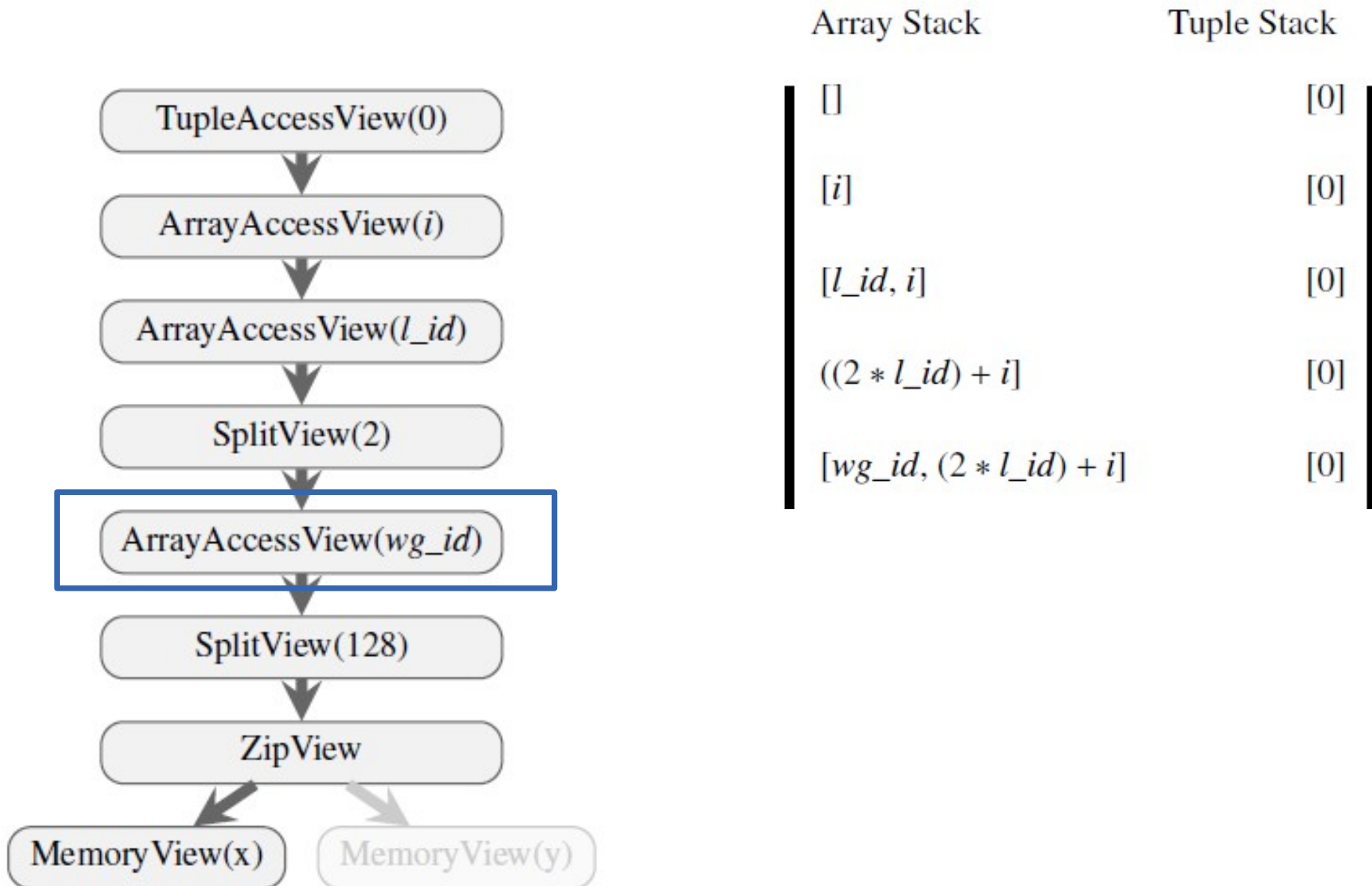
View Consumption



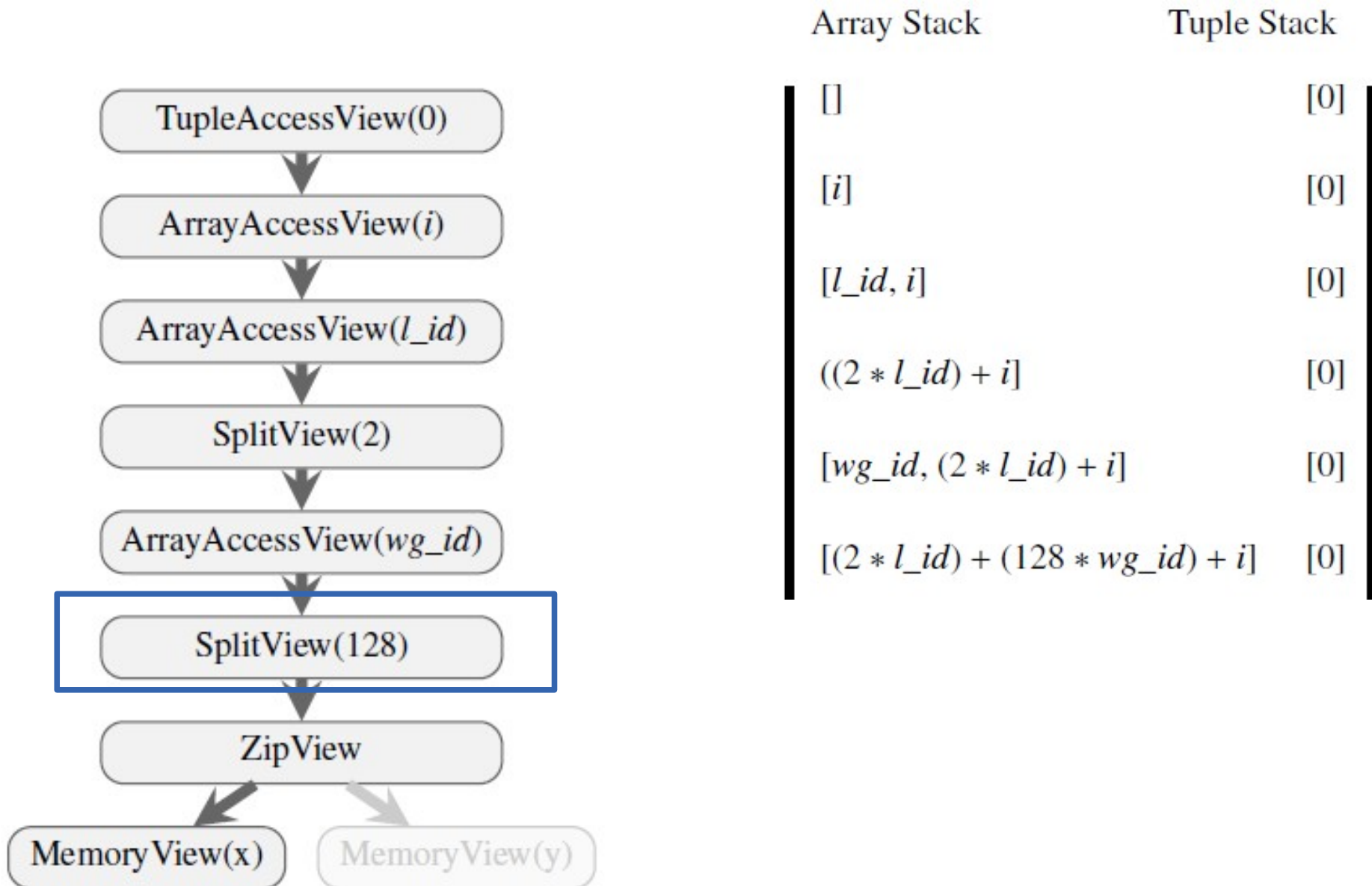
View Consumption



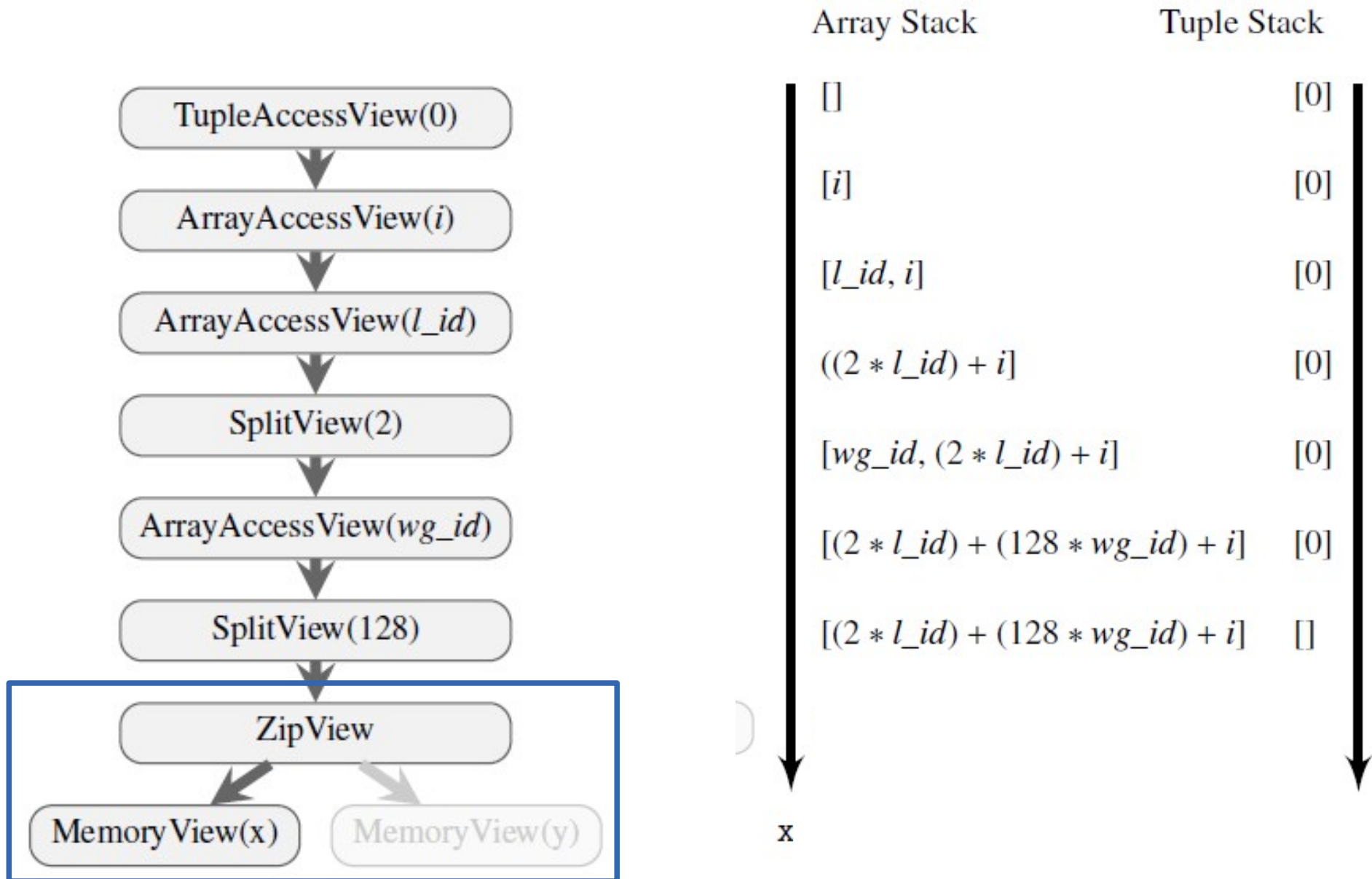
View Consumption



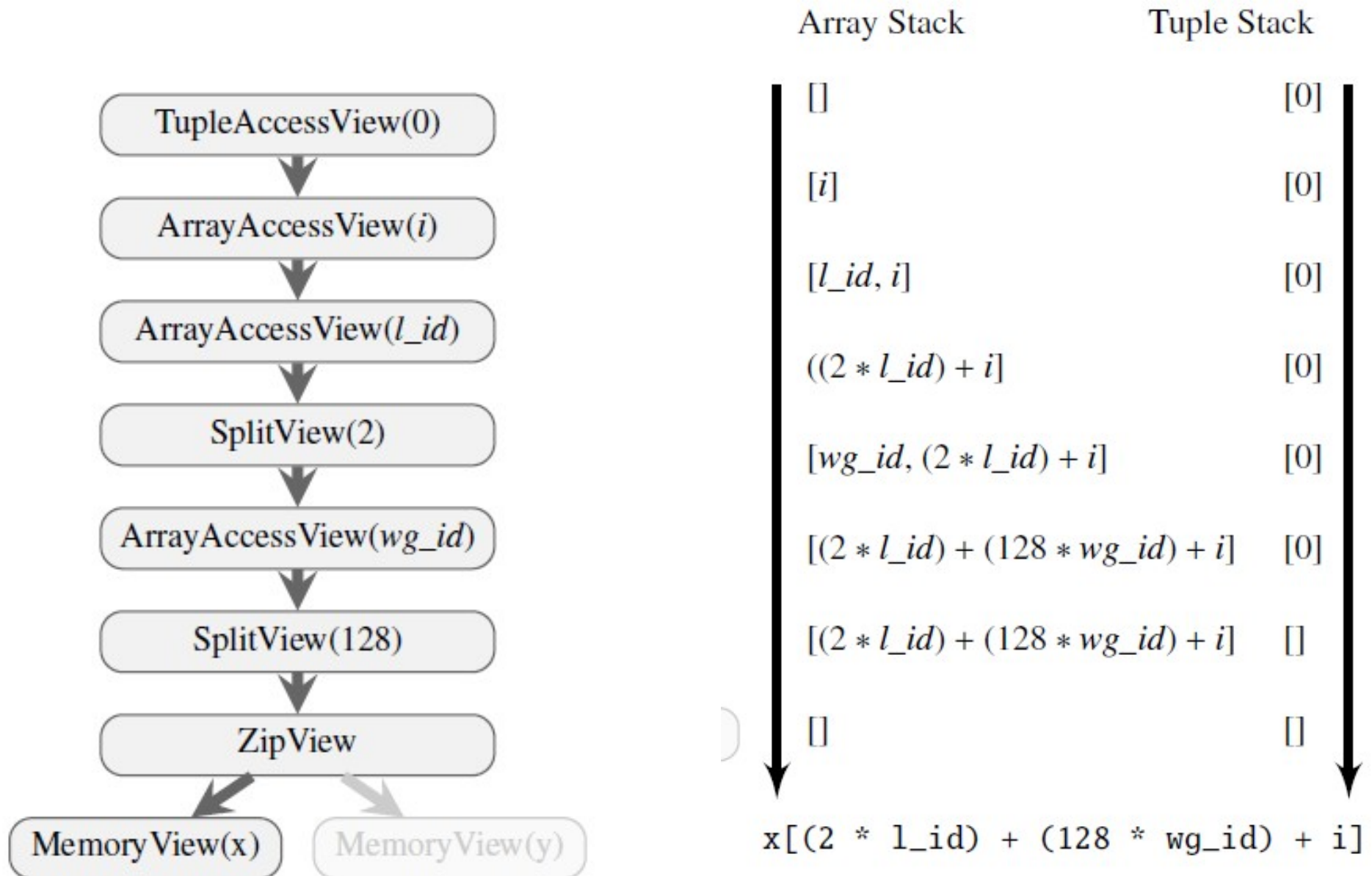
View Consumption



View Consumption



View Consumption



dot-product (a tiny bit of it):

```
(join ◦ mapWrg0( ...  
  join ◦ mapLcl0( ...  
    reduceSeq(λ(a, xy) ↦ a + (xy0 * xy1), 0)) ◦ split2  
  ) ◦ split128)( zip(x, y) )
```

produced OpenCL output:

```
int wg_id = get_group_id(0) ;  
int l_id = get_local_id(0) ;  
...  
acc = 0.0f;  
for (int i = 0 ; i < 2 ; I += 1) {  
  acc = acc +  
    x [ 2 * l_id + 128 * wg_id + I ] *  
    y [ 2 * l_id + 128 * wg_id + I ] ;  
}
```

Array indices not so simple

$$\left(\left(\left(\text{wg_id} * M + l_id \right) / M \right) + \left(\left(\text{wg_id} * M + l_id \right) \% M \right) * N \right) / N * N + \left(\left(\text{wg_id} * M + l_id \right) / M \right) + \left(\left(\text{wg_id} * M + l_id \right) \% M \right) * N$$

Array indices not so simple

$((((wg_id * M + l_id) / M) + (((wg_id * M + l_id) \% M) * N)) / N) * N + (((wg_id * M + l_id) / M) + (((wg_id * M + l_id) \% M) * N)) \% N$

Can use well-known algebraic rules

$$x / y = 0, \quad \text{if } x < y \text{ and } y \neq 0 \quad (1)$$

$$(x \times y + z) / y = x + z / y, \quad \text{if } y \neq 0 \quad (2)$$

$$x \% y = x, \quad \text{if } x < y \text{ and } y \neq 0 \quad (3)$$

$$(x + y) \% z = (x \% z + y \% z) \% z, \quad \text{if } z \neq 0 \quad (4)$$

$$(x / y) * y + x \% y = x, \quad \text{if } y \neq 0 \quad (5)$$

$$(x * y) \% y = 0, \quad \text{if } y \neq 0 \quad (6)$$

Array indices not so simple

$((((wg_id * M + l_id) / M) + (((wg_id * M + l_id) \% M) * N)) / N) * N + (((wg_id * M + l_id) / M) + (((wg_id * M + l_id) \% M) * N)) \% N$

Can use well-known algebraic rules

$$x / y = 0, \quad \text{if } x < y \text{ and } y \neq 0 \quad (1)$$

$$(x \times y + z) / y = x + z / y, \quad \text{if } y \neq 0 \quad (2)$$

$$x \% y = x, \quad \text{if } x < y \text{ and } y \neq 0 \quad (3)$$

$$(x + y) \% z = (x \% z + y \% z) \% z, \quad \text{if } z \neq 0 \quad (4)$$

$$(x / y) * y + x \% y = x, \quad \text{if } y \neq 0 \quad (5)$$

$$(x * y) \% y = 0, \quad \text{if } y \neq 0 \quad (6)$$

and simplify array index expression

$((((wg_id * M + l_id) / M) + (((wg_id * M + l_id) \% M) * N)) / N) * N + (((wg_id * M + l_id) / M) + (((wg_id * M + l_id) \% M) * N)) \% N$



l_id

$*N +$

wg_id

Putting it all together

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

rewriting



```
(join ◦ mapWrg0(
  join ◦ toGlobal(mapLcl0(mapSeq(id))) ◦ split1 ◦
  iterate6( join ◦
    mapLcl0( toLocal(mapSeq(id)) ◦
      reduceSeq(add, 0) ) ◦
    split2 ) ◦
  join ◦ mapLcl0( toLocal(mapSeq(id)) ◦
    reduceSeq(multAndSumUp, 0) ) ◦ split2
) ◦ split128)( zip(x, y) )
```



code generation

```
kernel void KERNEL(const global float *restrict x,
                  const global float *restrict y,
                  global float *z, int N) {
  local float tmp1[64]; local float tmp2[64];
  local float tmp3[32];
  float acc1; float acc2;
  for (int wg_id = get_group_id(0); wg_id < N/128;
       wg_id += get_num_groups(0)) {
    { int l_id = get_local_id(0);
      acc1 = 0.0f;
      for (int i = 0; i < 2; i += 1) {
        acc1 = multAndSumUp(acc1,
                           x[2 * l_id + 128 * wg_id + i],
                           y[2 * l_id + 128 * wg_id + i]); }
      tmp1[l_id] = id(acc1); }
    barrier(CLK_LOCAL_MEM_FENCE);
    int size = 64;
    local float *in = tmp1; local float *out = tmp2;
    for (int iter = 0; iter < 6; iter += 1) {
      if (get_local_id(0) < size / 2) {
        acc2 = 0.0f;
        for (int i = 0; i < 2; i += 1) {
          acc2 = add(acc2, in[2 * l_id + i]); }
        out[l_id] = id(acc2); }
      barrier(CLK_LOCAL_MEM_FENCE);
      size = size / 2;
      in = (out == tmp1) ? tmp1 : tmp3;
      out = (out == tmp1) ? tmp3 : tmp1;
      barrier(CLK_LOCAL_MEM_FENCE); }
    if (get_local_id(0) < 1) {
      z[wg_id] = id(tmp3[l_id]); }
    barrier(CLK_GLOBAL_MEM_FENCE); } }
```

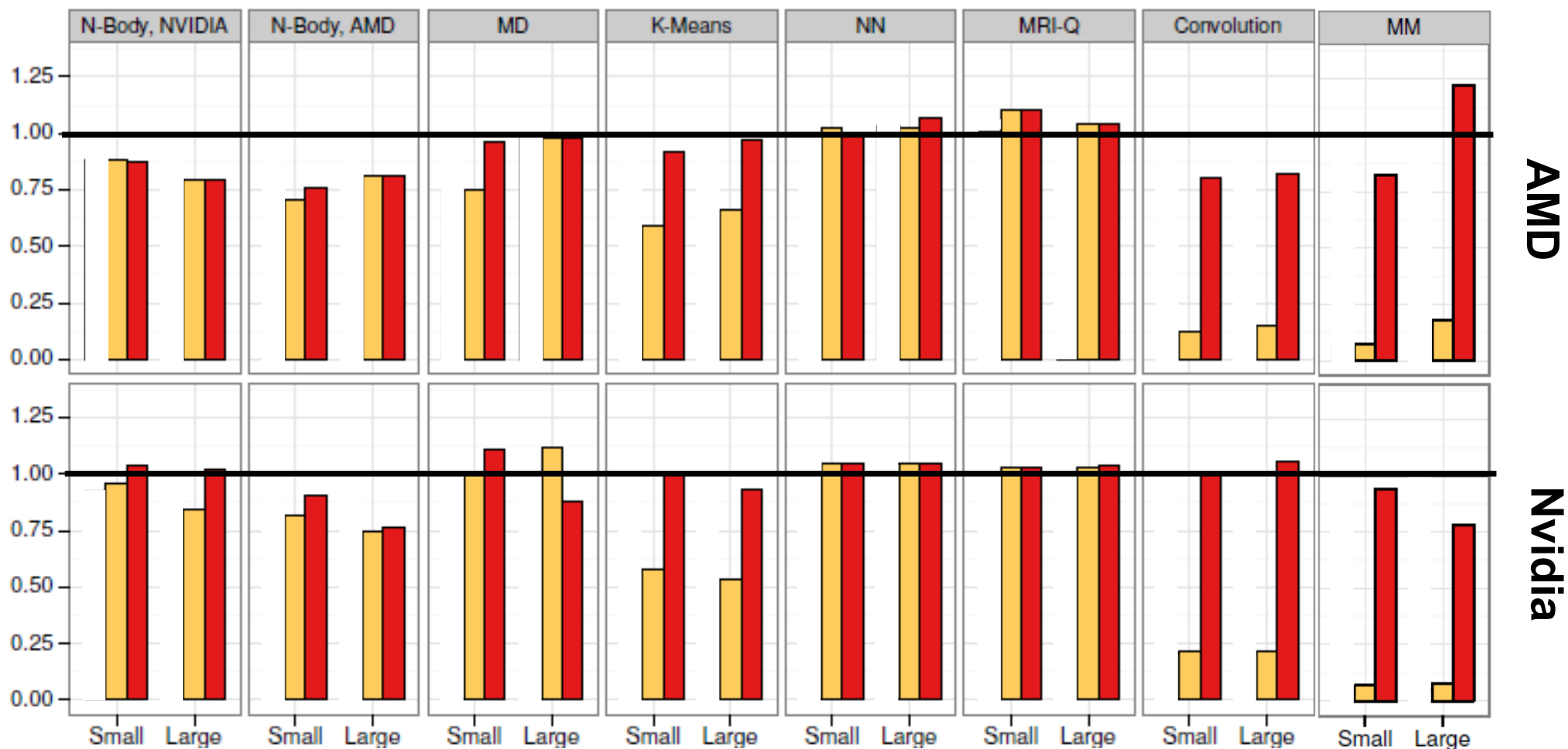
Code Generator vs Human



no simplification



arithmetic simplification



The Lift Approach: Summary

- ▶ Hardware-agnostic data-parallel functional intermediate language
 - hides hardware complexity, can be targeted by DSLs
- ▶ **Low-level functional language for each hardware type**
 - OpenCL language (this talk), in the future: MPI, OpenMP
- ▶ **Rewrite rules based optimisation**
 - maps the hardware-agnostic language to the hardware-specific language
 - extensible
- ▶ Possible to achieve high performance through exploration
 - opportunity for smarter technique (e.g. predictive modelling)

if you want to know more: www.lift-project.org