

**Naums Mogers**

Aaron Smith, Dimitrios Vytiniotis

Michel Steuwer, Christophe Dubach

Ryota Tomioka

# Towards Mapping Lift To Deep Neural Networks

---



THE UNIVERSITY of EDINBURGH  
**informatics**

EPSRC Centre for Doctoral Training in  
**Pervasive Parallelism**

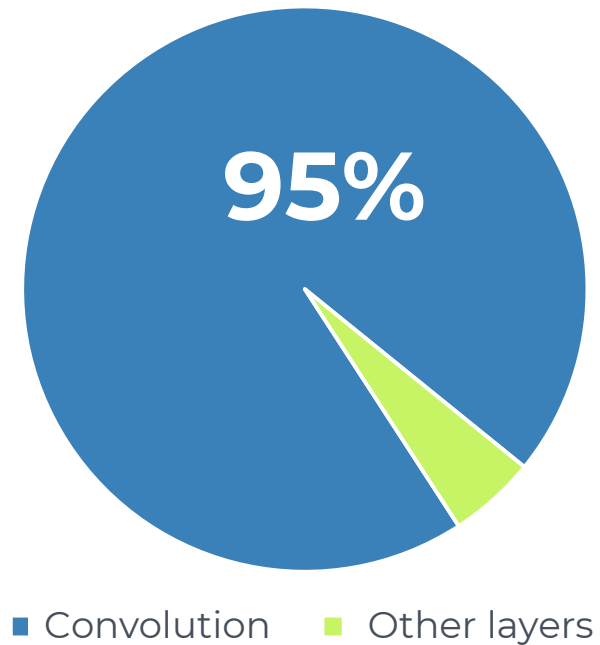
Microsoft  
**Research**

# Convolution

■ Many popular NN architectures depend on **convolution**

- AlexNet (2012)
- ZFNet(2013)
- GoogleNet(2014)
- VGGNet(2014)
- ResNet(2015)

Runtime of convolution and other layers  
(VGG measured on Mali GPU with cIBLAS)

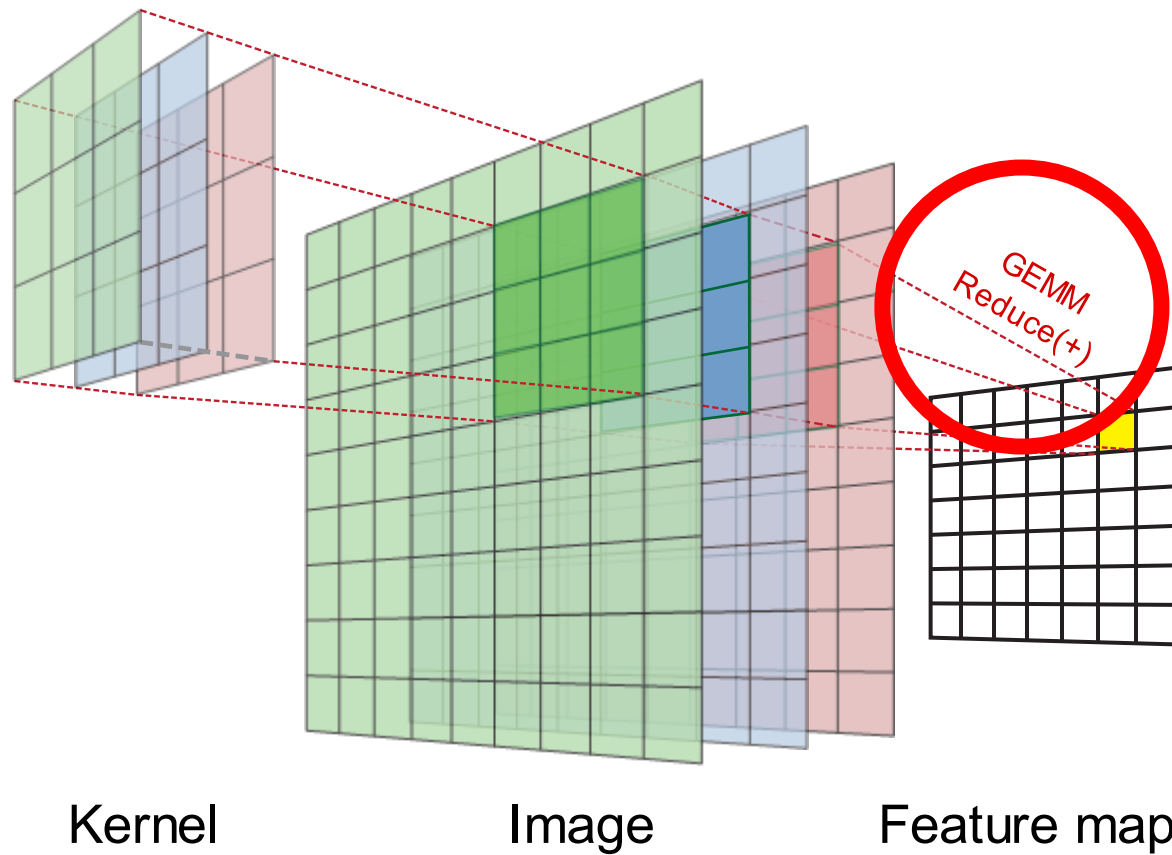


# Convolution: algorithms

---

- Stencils
- *Im2col*
- *Fast Fourier transform*
- *Winograd*

# Convolution: **stencils**



# Convolution: im2col

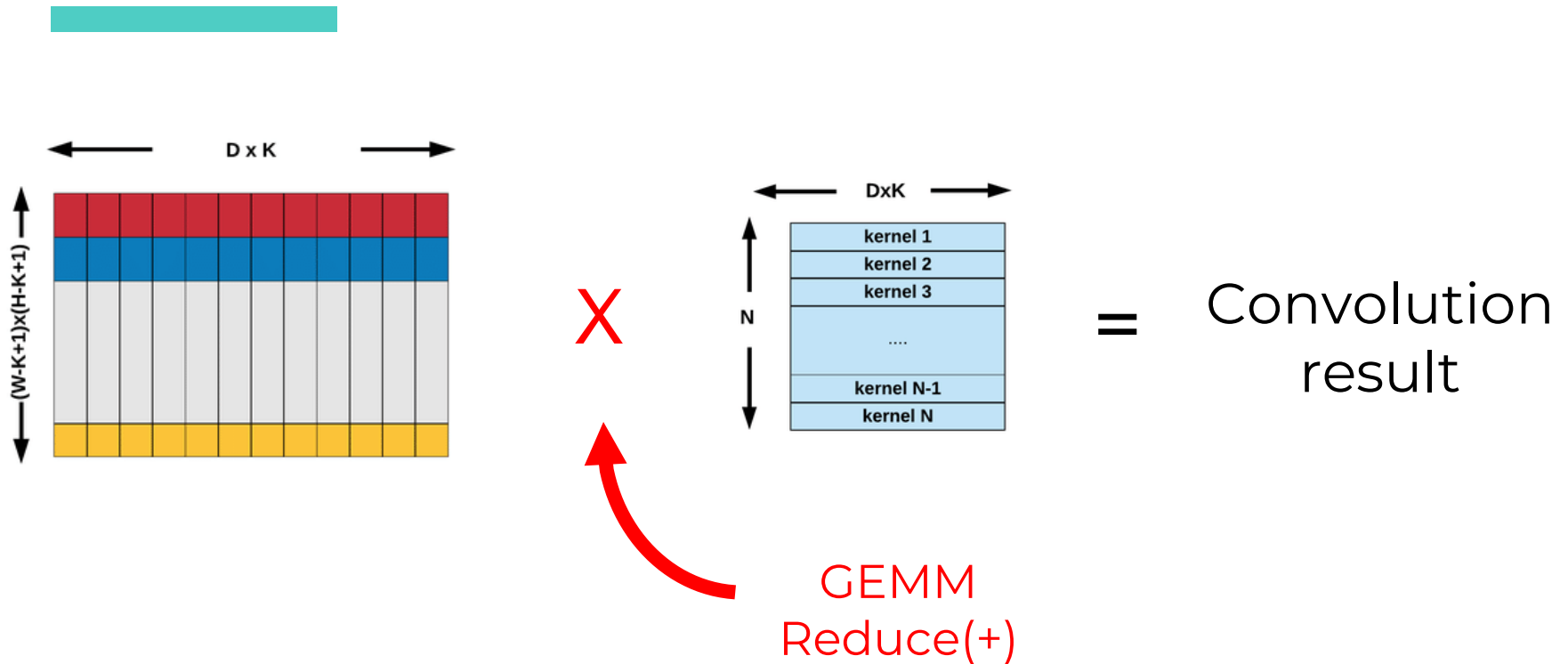
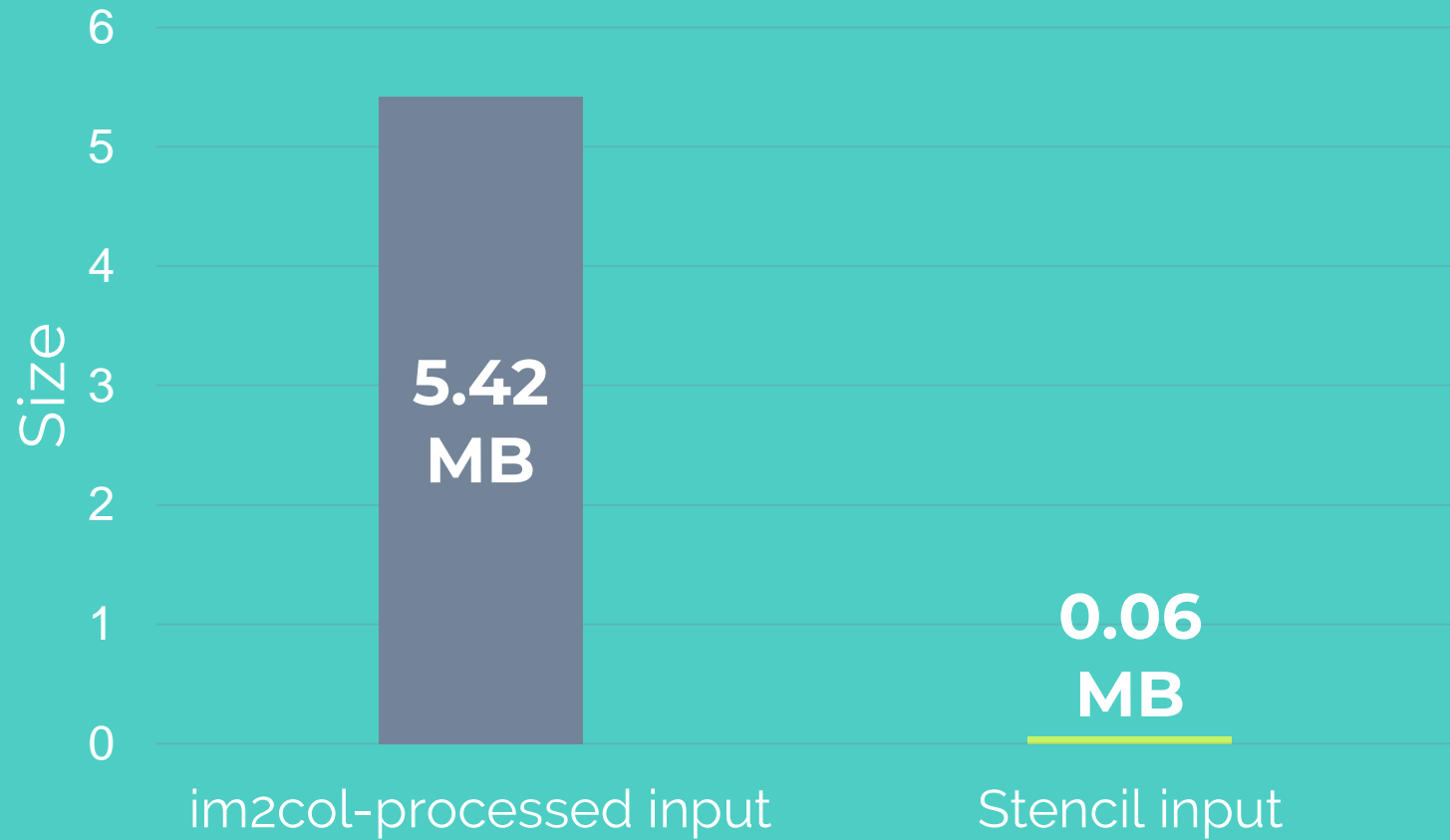


Image source: Loukidakis et al (2018)

## Input image size in the two methods (1<sup>st</sup> layer of VGG)



# GEMM in neural networks

---

- Both stencil and im2col methods for convolution rely heavily on GEMM
- GEMM is also the basis of fully connected layers
  - MLPs make for a large portion of server workloads

# Hardware accelerators

## GPUs



Image source:  
<https://tomshardware.co.uk>

## FPGAs



Image source:  
<https://eetimes.com>

## ASICs

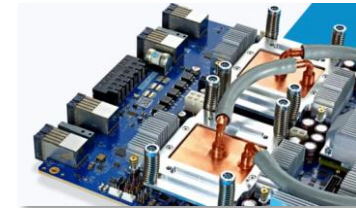


Image source:  
<https://blogs.microsoft.com>

**Examples:** TPU, BrainWave, DianNao, Huawei Da Vinci, Movidius Myriad

**All accelerators feature:**

Large memory

High bandwidth

Multidimensional  
computational units



# Compiling for hardware accelerators

---

- Multidimensional computational units are **exposed** in instruction sets:
  - VVAdd32, VVAdd64
  - VVMul32, VVMul64, VVMul128
  - MVMul32, MVMul64, MVMul128
- Finding opportunities to use these primitives is **challenging**

# Compiling for hardware accelerators

---

□ Using built-in primitives is made **harder** by:

□ **Other optimisations**

- Parallelisation
- Tiling (*for limited shared memory*)
- Memory coalescing
- Prefetching
- (*etc*)

□ **Individual device characteristics**

# Compiling for hardware accelerators

## GEMM

```
for (int i = 0; i < M; i++) {
  for (int j = 0; j < N; j++) {
    for (int k = 0; k < K; k++) {
      temp[k + K*N*i + K*j] =
        A[k + K*i] * B[k + K*j];
    }
    for (int k = 0; k < K; k++) {
      C[j + N*i] +=
        temp[k + K*N*i + K*j];
    }
  }
}
```

# Compiling for hardware accelerators

## GEMM

```
for (int i = 0; i < M; i++) {
  for (int j = 0; j < N; j++) {
    for (int k = 0; k < K; k++) {
      temp[k + K*N*i + K*j] =
        A[k + K*i] * B[k + K*j];
    }
    for (int k = 0; k < K; k++) {
      C[j + N*i] +=
        temp[k + K*N*i + K*j];
    }
  }
}
```



## GEMM<sup>til</sup>

```
for (int i = 0; i < M; i++) {
  for (int j = 0; j < N; j++) {
    for (int k = 0; k < K/4; k++) {
      for (int l = 0; l < 4; l++) {
        temp[l + k*4 + K*N*i + K*j] =
          A[l + k*4 + K*i] * B[l + k*4 + K*j];
      }
    }
    for (int k = 0; k < K; k++) {
      C[j + N*i] +=
        temp[k + K*N*i + K*j];
    }
  }
}
```

## Questions:

Data contiguity

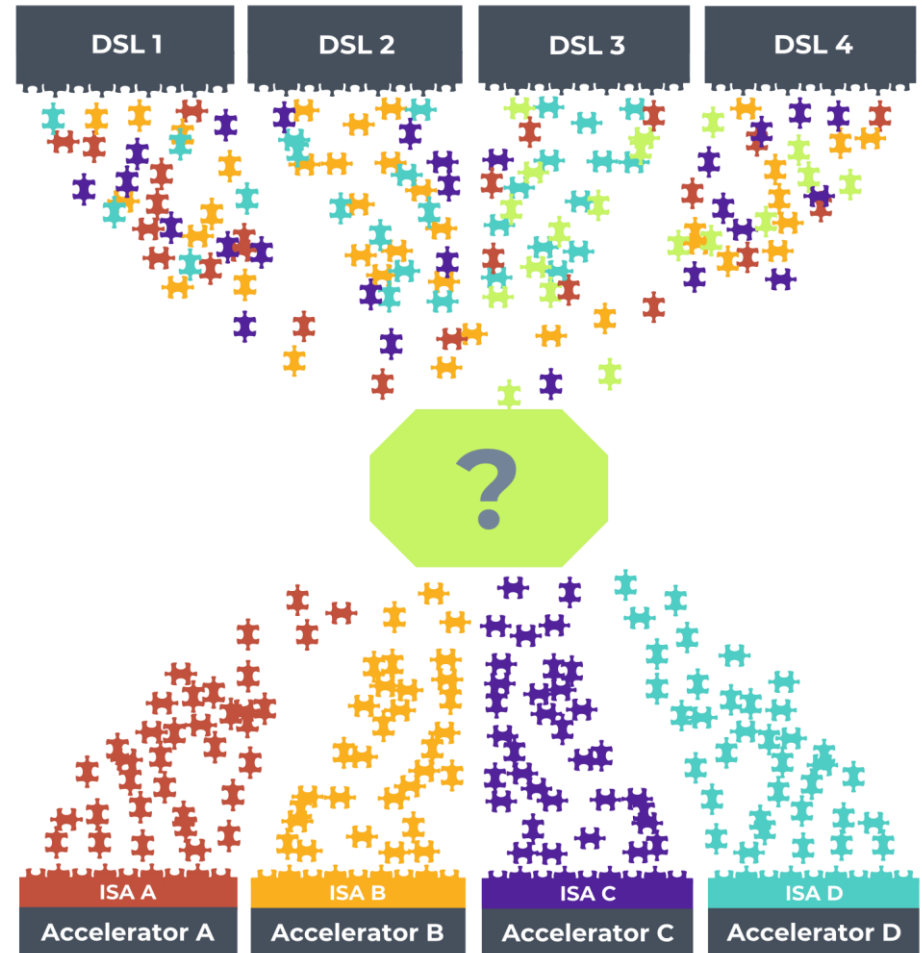
Data size

Efficient memory accesses

Data type

# The problem

- Combine device-specific operators optimally
- Design a performance portable approach
- Automate and abstract the process from the user



# Lift: the approach

1. Separate algorithm (**WHAT**) from implementation (**HOW**)

**GEMM**<sup>imperative</sup>

```
for (int i = 0; i < M; i++) {
  for (int j = 0; j < N; j++) {
    for (int k = 0; k < K; k++) {
      temp[k + K*N*i + K*j] =
        A[k + K*i] * B[k + K*j];
    }
    for (int k = 0; k < K; k++) {
      C[j + N*i] +=
        temp[k + K*N*i + K*j];
    }
  }
}
```



**GEMM**<sup>functional</sup>

```
A >> Map(Arow =>
  B >> Map(Bcol =>
    Zip(Arow, Brow) >>
    Map(ScalarMul) >>
    Reduce(0, Add)
  )
)
```

# Lift: the approach

---

## 2. Detect and rewrite patterns

**GEMM**<sup>functional</sup>

```
A >> Map(Arow =>
  B >> Map(Bcol =>
    Zip(Arow, Brow) >>
    Map(ScalarMul) >>
    Reduce(0, Add)
  )
)
```

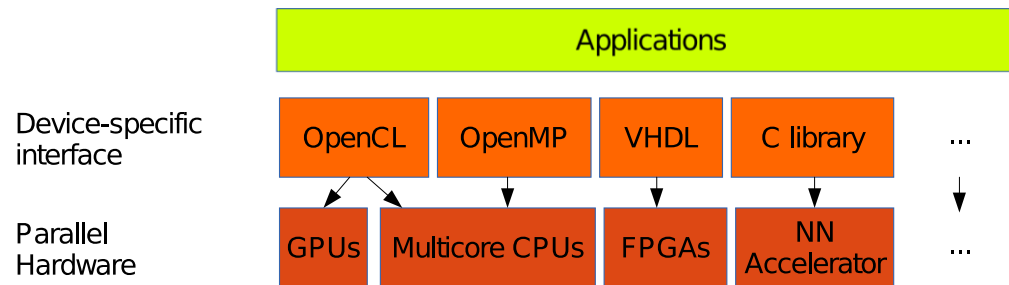


**GEMM**<sup>optimised</sup>

```
MMul(A, B)
```

# Lift in context

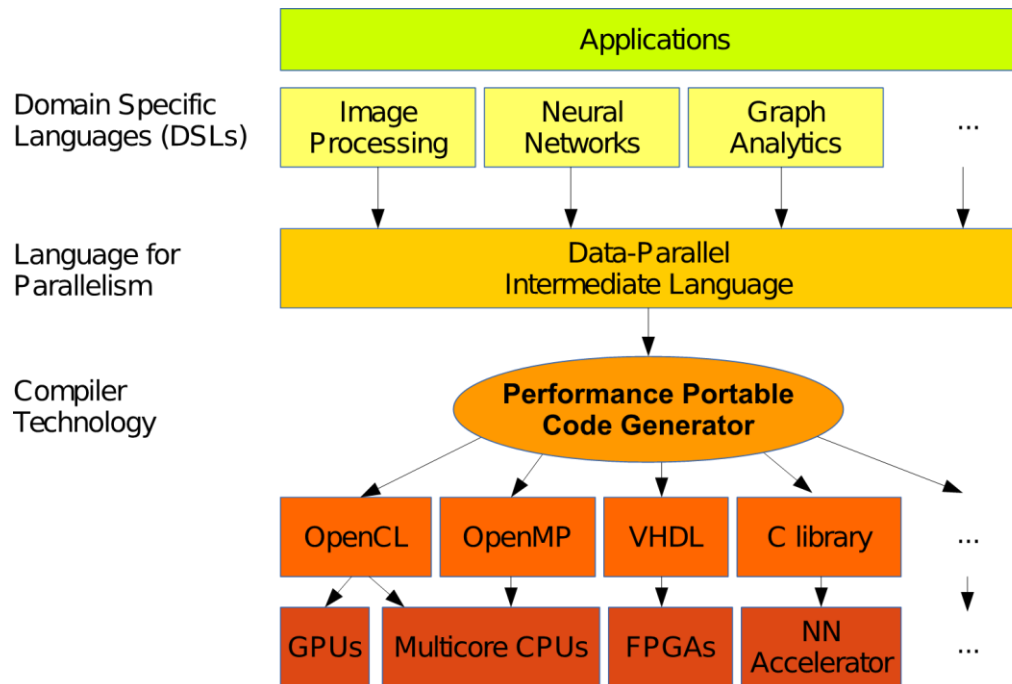
## Current landscape





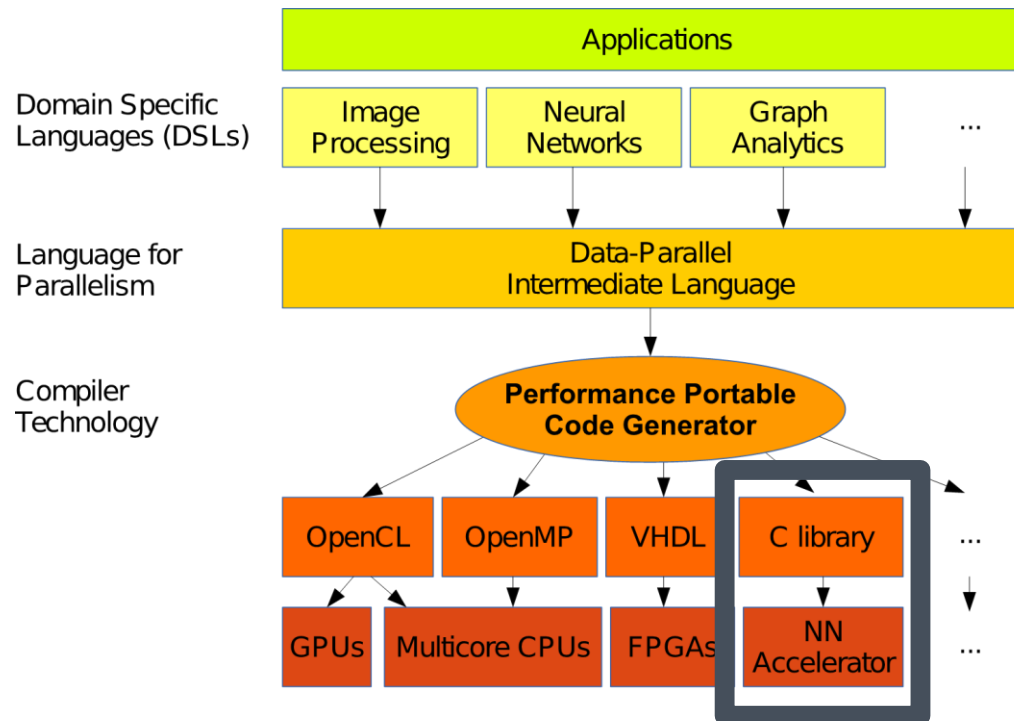
# Lift in context

## The goal



# Lift in context

## The focus of this talk



# Lift: IR

---

## □ Data types

- `Int`, `Float8` / `Float16` / `Float32`, `Arrays`

## □ Algorithmic patterns

- `Map`, `Slide`, `Reduce`, `Zip`, `Join`, `Split`

## □ Address space operators

- `toChip`, `toDram`, `toOutput`

## □ Arithmetic operators

- `ScalarAdd`, `VVAdd`, `MVAdd`, `MMAdd`
- `ScalarMul`, `VVMul`, `MVMul`, `MMMul`
- `VVRelu`, `VVTanh`

# Lift: example rewrite rules

## Split-join rule

```
Map(f)
  ↓ ↓ ↓
Split(n) >>
Map(Map(f)) >>
Join
```

## Map fusion rule

```
Map(f) >> Map(g)
  ↓ ↓ ↓
Map(f >> g)
```

## GEMV rule

```
matrix >> Map(row =>
  VVMul(row, vector) >>
  Reduce(ScalarAdd, 0))
  ↓ ↓ ↓
MVMul(matrix, vector)
```

## Rewrite rule system

- Generic and customisable
- 3 levels: DSL, algorithmic, hardware
- Extensible

# Lift: rewriting

A fully connected layer

```
1 for (int i = 0; i < n_neurons; i++) {
2   for (int j = 0; j < x_length; j++) {
3     temp[j + x_length*i] =
4       X[j] * W[j + x_length*i];
5   }
6
7   OUT[i] = B[i];
8   for (int j = 0; j < x_length; j++) {
9     OUT[j] += temp[j + x_length*i];
10  }
11 }
```



```
1 temp = MVMul(W, X);
2 OUT = VVAdd(temp, B);
```

# Lift: rewriting

```
1 Zip(  
2   W,  
3   B) >>  
4 Map((neuronW, neuronB) =>  
5   VVMul(neuronW, toChip(X)) >>  
6   Reduce(ScalarAdd, neuronB)) >>  
7 VVRelu() >> toOutput()
```



```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6  
7   OUT[i] = B[i];  
8   for (int j = 0; j < x_length; j++) {  
9     OUT[j] += temp[j + x_length*i];  
10  }  
11 }
```

# Lift: rewriting

```
1 Zip(  
2   W,  
3   B) >>  
4 Map((neuronW, neuronB) =>  
5   VVMul(neuronW, toChip(X)) >>  
6   Reduce(ScalarAdd, neuronB)) >>  
7 VVRelu() >> toOutput()
```

```
// GEMV rewrite  
matrix >> Map(row =>  
  VVMul(row, vector) >>  
  Reduce(ScalarAdd, 0))  
          ↓ ↓ ↓  
MVMul(matrix, vector)
```

```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6  
7   OUT[i] = B[i];  
8   for (int j = 0; j < x_length; j++) {  
9     OUT[j] += temp[j + x_length*i];  
10  }  
11 }
```

# Lift: rewriting

```
1 Zip(  
2   W,  
3   B) >>  
4 Map((neuronW, neuronB) =>  
5   VVMul(neuronW, toChip(X)) >>  
6   Reduce(ScalarAdd, neuronB)) >>  
7 VVRelu() >> toOutput()
```

<Extract Initializer From Reduce>

```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6  
7   OUT[i] = B[i];  
8   for (int j = 0; j < x_length; j++) {  
9     OUT[j] += temp[j + x_length*i];  
10  }  
11 }
```



# Lift: rewriting

```
1 Zip(  
2   W,  
3   B) >>  
4 Map((neuronW, neuronB) =>  
5   VVMul(neuronW, toChip(X)) >>  
6   Reduce(ScalarAdd, neuronB)) >>  
7 VVRelu() >> toOutput()
```



```
1 Zip(  
2   W,  
3   B) >>  
4 Map((neuronW, neuronB) =>  
5   VVMul(neuronW, toChip(X)) >>  
6   Reduce(ScalarAdd, 0) >>  
7   ScalarAdd(neuronB)) >>  
8 VVRelu() >> toOutput()
```

```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6  
7   OUT[i] = B[i];  
8   for (int j = 0; j < x_length; j++) {  
9     OUT[j] += temp[j + x_length*i];  
10  }  
11 }
```

```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6  
7   OUT[i] = 0;  
8   for (int j = 0; j < x_length; j++) {  
9     OUT[j] += temp[j + x_length*i];  
10  }  
11   OUT[i] += B[i];  
12 }
```

# Lift: rewriting

```
1 Zip(  
2   W,  
3   B) >>  
4 Map((neuronW, neuronB) =>  
5   VVMul(neuronW, toChip(X)) >>  
6   Reduce(ScalarAdd, 0) >>  
7   ScalarAdd(neuronB)) >>  
8 VVRelu() >> toOutput()
```

```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6  
7   OUT[i] = 0;  
8   for (int j = 0; j < x_length; j++) {  
9     OUT[j] += temp[j + x_length*i];  
10  }  
11  OUT[i] += B[i];  
12 }
```

# Lift: rewriting

```
1 Zip(  
2   W,  
3   B) >>  
4 Map((neuronW, neuronB) =>  
5   VVMul(neuronW, toChip(X)) >>  
6   Reduce(ScalarAdd, 0) >>  
7   ScalarAdd(neuronB)) >>  
8 VVRelu() >> toOutput()
```

<Map fission>

```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6  
7   OUT[i] = 0;  
8   for (int j = 0; j < x_length; j++) {  
9     OUT[j] += temp[j + x_length*i];  
10  }  
11  OUT[i] += B[i];  
12 }
```

# Lift: rewriting

```
1 Zip(  
2   W,  
3   B) >>  
4 Map((neuronW, neuronB) =>  
5   VVMul(neuronW, toChip(X)) >>  
6   Reduce(ScalarAdd, 0) >>  
7   ScalarAdd(neuronB)) >>  
8 VVRelu() >> toOutput()
```



```
1 Zip(  
2   W >> Map(neuronW =>  
3     VVMul(neuronW, toChip(X))),  
4   B) >>  
5 Map((neuronWX, neuronB) =>  
6   neuronWX >>  
7   Reduce(ScalarAdd, 0) >>  
8   ScalarAdd(neuronB)) >>  
9 VVRelu() >> toOutput()
```

```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6  
7   OUT[i] = 0;  
8   for (int j = 0; j < x_length; j++) {  
9     OUT[j] += temp[j + x_length*i];  
10  }  
11  OUT[i] += B[i];  
12 }
```

```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6 }  
7 for (int i = 0; i < n_neurons; i++) {  
8   OUT[i] = 0;  
9   for (int j = 0; j < x_length; j++) {  
10    OUT[j] += temp[j + x_length*i];  
11  }  
12  OUT[i] += B[i];  
13 }
```

# Lift: rewriting

```
1 Zip(  
2   W >> Map(neuronW =>  
3     VVMul(neuronW, toChip(X))),  
4   B) >>  
5 Map((neuronWX, neuronB) =>  
6   neuronWX >>  
7   Reduce(ScalarAdd, 0) >>  
8   ScalarAdd(neuronB)) >>  
9 VVRelu() >> toOutput()
```

```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6 }  
7 for (int i = 0; i < n_neurons; i++) {  
8   OUT[i] = 0;  
9   for (int j = 0; j < x_length; j++) {  
10    OUT[j] += temp[j + x_length*i];  
11  }  
12  OUT[i] += B[i];  
13 }
```

# Lift: rewriting

```
1 Zip(  
2   W >> Map(neuronW =>  
3     VVMul(neuronW, toChip(X))),  
4   B) >>  
5 Map((neuronWX, neuronB) =>  
6   neuronWX >>  
7   Reduce(ScalarAdd, 0) >>  
8   ScalarAdd(neuronB)) >>  
9 VVRelu() >> toOutput()
```

<A couple rewrites later...>

```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6 }  
7 for (int i = 0; i < n_neurons; i++) {  
8   OUT[i] = 0;  
9   for (int j = 0; j < x_length; j++) {  
10    OUT[j] += temp[j + x_length*i];  
11  }  
12  OUT[i] += B[i];  
13 }
```

# Lift: rewriting

```
1 Zip(  
2   W >> Map(neuronW =>  
3     VVMul(neuronW, toChip(X))),  
4   B) >>  
5 Map((neuronWX, neuronB) =>  
6   neuronWX >>  
7   Reduce(ScalarAdd, 0) >>  
8   ScalarAdd(neuronB)) >>  
9 VVRelu() >> toOutput()
```



```
1 Zip(  
2   W >> Map(neuronW =>  
3     VVMul(neuronW, toChip(X)) >>  
4     Reduce(ScalarAdd, 0)),  
5   B) >>  
6 Map((neuronWXreduced, neuronB) =>  
7   neuronWXreduced >>  
8   ScalarAdd(neuronB)) >>  
9 VVRelu() >> toOutput()
```

```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6 }  
7 for (int i = 0; i < n_neurons; i++) {  
8   OUT[i] = 0;  
9   for (int j = 0; j < x_length; j++) {  
10    OUT[j] += temp[j + x_length*i];  
11  }  
12  OUT[i] += B[i];  
13 }
```

```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6   OUT[i] = 0;  
7   for (int j = 0; j < x_length; j++) {  
8     OUT[j] += temp[j + x_length*i];  
9   }  
10 }  
11 for (int i = 0; i < n_neurons; i++) {  
12   OUT[i] += B[i];  
13 }
```

# Lift: rewriting

```
1 Zip(  
2   W >> Map(neuronW =>  
3     VVMul(neuronW, toChip(X)) >>  
4     Reduce(ScalarAdd, 0)),  
5   B) >>  
6 Map((neuronWXreduced, neuronB) =>  
7   neuronWXreduced >>  
8   ScalarAdd(neuronB)) >>  
9 VVRelu() >> toOutput()
```

```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6   OUT[i] = 0;  
7   for (int j = 0; j < x_length; j++) {  
8     OUT[j] += temp[j + x_length*i];  
9   }  
10 }  
11 for (int i = 0; i < n_neurons; i++) {  
12   OUT[i] += B[i];  
13 }
```



# Lift: rewriting

```
1 Zip(  
2   W >> Map(neuronW =>  
3     VVMul(neuronW, toChip(X)) >>  
4     Reduce(ScalarAdd, 0)),  
5   B) >>  
6 Map((neuronWXreduced, neuronB) =>  
7   neuronWXreduced >>  
8   ScalarAdd(neuronB)) >>  
9 VVRelu() >> toOutput()
```

```
// GEMV rewrite  
matrix >> Map(row =>  
  VVMul(row, vector) >>  
  Reduce(ScalarAdd, 0))  
      ↓ ↓ ↓  
MVMul(matrix, vector)
```

```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6   OUT[i] = 0;  
7   for (int j = 0; j < x_length; j++) {  
8     OUT[j] += temp[j + x_length*i];  
9   }  
10 }  
11 for (int i = 0; i < n_neurons; i++) {  
12   OUT[i] += B[i];  
13 }
```

# Lift: rewriting

```
1 Zip(  
2   W >> Map(neuronW =>  
3     VVMul(neuronW, toChip(X)) >>  
4     Reduce(ScalarAdd, 0)),  
5   B) >>  
6 Map((neuronWXreduced, neuronB) =>  
7   neuronWXreduced >>  
8   ScalarAdd(neuronB)) >>  
9 VVRelu() >> toOutput()
```



```
1 Zip(  
2   MVMul(W, X >> toChip),  
3   B) >>  
4 Map((neuronWXreduced, neuronB) =>  
5   neuronWXreduced >>  
6   ScalarAdd(neuronB)) >>  
7 VVRelu() >> toOutput()
```

```
1 for (int i = 0; i < n_neurons; i++) {  
2   for (int j = 0; j < x_length; j++) {  
3     temp[j + x_length*i] =  
4       X[j] * W[j + x_length*i];  
5   }  
6   OUT[i] = 0;  
7   for (int j = 0; j < x_length; j++) {  
8     OUT[j] += temp[j + x_length*i];  
9   }  
10 }  
11 for (int i = 0; i < n_neurons; i++) {  
12   OUT[i] += B[i];  
13 }
```

```
1 temp = MVMul(W, X);  
2 for (int i = 0; i < n_neurons; i++) {  
3   OUT[i] += B[i];  
4 }
```

# Lift: rewriting

```
1 Zip(  
2   MVMul(W, X >> toChip)  
3   B) >>  
4 Map((neuronWXreduced, neuronB) =>  
5   neuronWXreduced >>  
6   ScalarAdd(neuronB)) >>  
7 VVRelu() >> toOutput()
```

```
1 temp = MVMul(W, X);  
2 for (int i = 0; i < n_neurons; i++) {  
3   OUT[i] += B[i];  
4 }
```

# Lift: rewriting

```
1 Zip(  
2   MVMul(W, X >> toChip)  
3   B) >>  
4 Map((neuronWXreduced, neuronB) =>  
5   neuronWXreduced >>  
6   ScalarAdd(neuronB)) >>  
7 VVRelu() >> toOutput()
```

< Vectorise Map-Zip >

```
1 temp = MVMul(W, X);  
2 for (int i = 0; i < n_neurons; i++) {  
3   OUT[i] += B[i];  
4 }
```

# Lift: rewriting

```
1 Zip(  
2   MVMul(W, X >> toChip)  
3   B) >>  
4 Map((neuronWXreduced, neuronB) =>  
5   neuronWXreduced >>  
6   ScalarAdd(neuronB)) >>  
7 VVRelu() >> toOutput()
```



```
1 VVAdd(  
2   MVMul(W, X >> toChip)  
3   B) >>  
4 VVRelu() >> toOutput()
```

```
1 temp = MVMul(W, X);  
2 for (int i = 0; i < n_neurons; i++) {  
3   OUT[i] += B[i];  
4 }
```

```
1 temp = MVMul(W, X);  
2 OUT = VVAdd(temp, B);
```

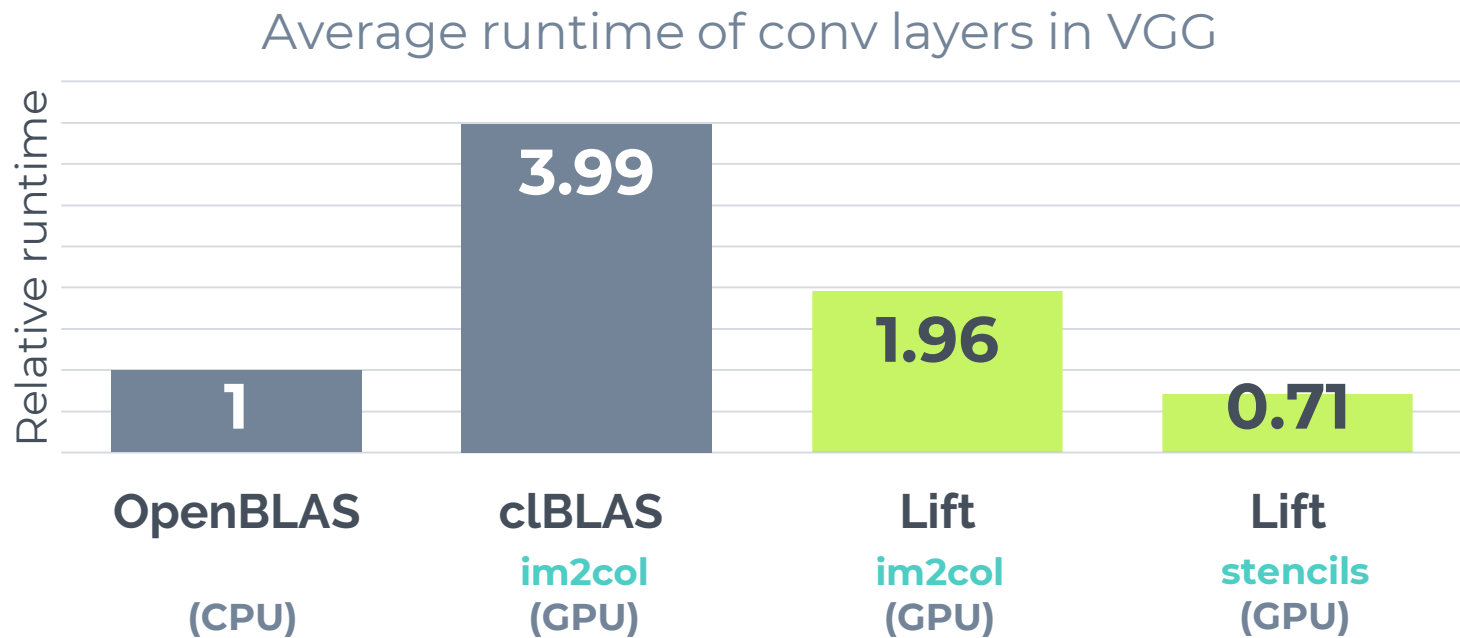
# Lift: rewriting

```
1 VVAdd(  
2   MVMul(W, X >> toChip)  
3   B) >>  
4 VVRelu() >> toOutput()
```

```
1 temp = MVMul(W, X);  
2 OUT = VVAdd(temp, B);
```

# Preliminary results

- Functional correctness on the BrainWave accelerator
- Performance measurements on Mali GPU



# Future work

---

- Generation of both OpenCL kernel and host runtime
- Performance evaluation on popular DNN architectures
- Support for more hardware accelerators