THE UNIVERSITY of EDINBURGH
**informatics**

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

www.lift-project.org

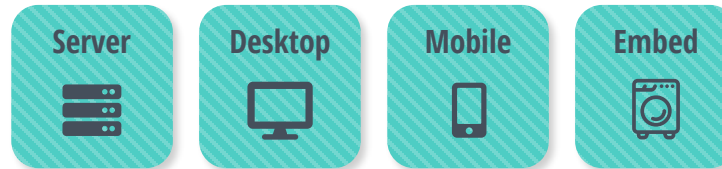# Functional Interface for Performance Portability on Parallel Accelerators
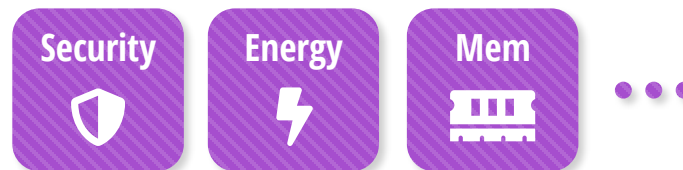
**Naums Mogers**

Christophe Dubach

ARM Research Summit, September 2019
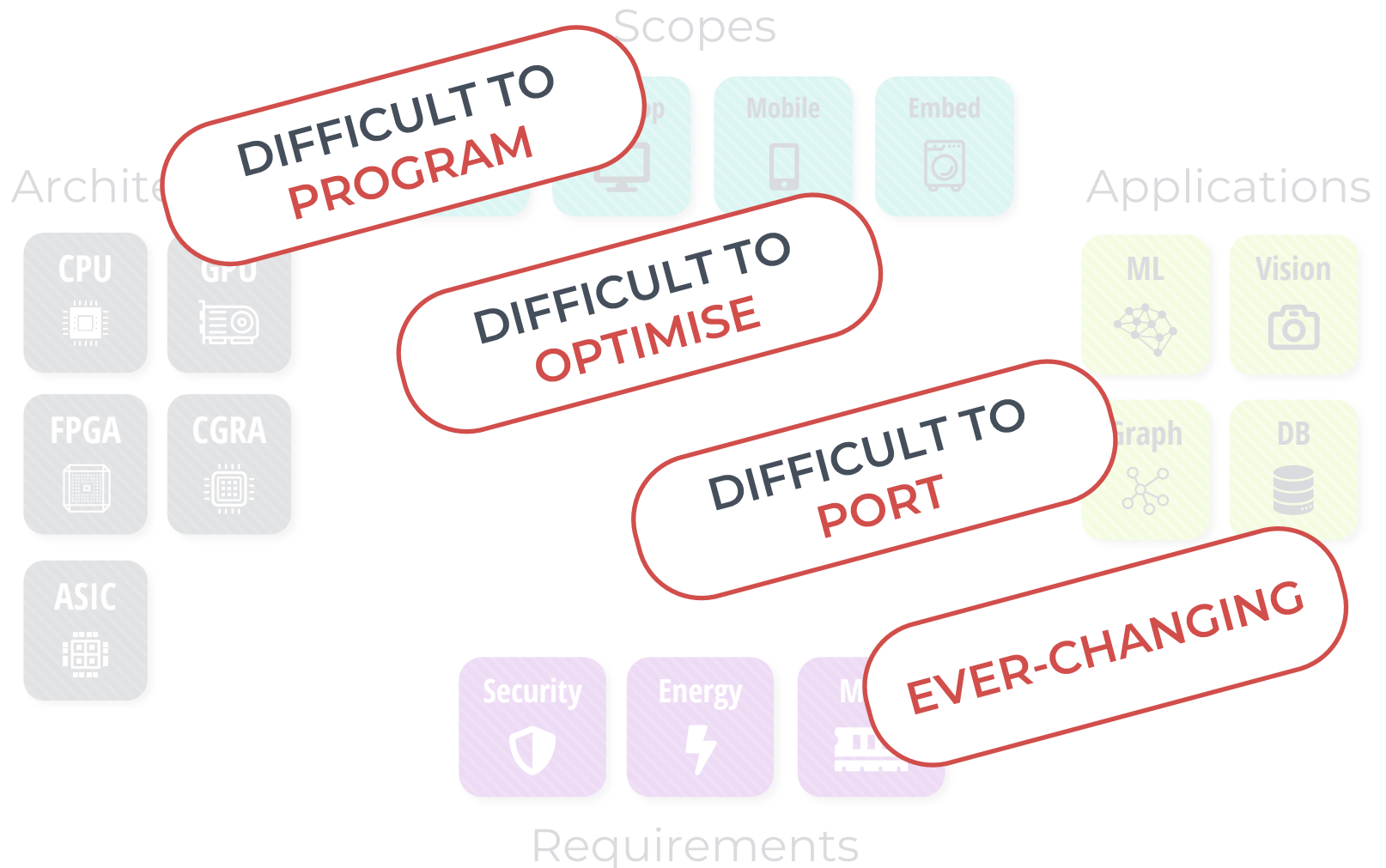
# Hardware accelerators

## Scopes

Server

Desktop

Mobile

Embed

## Architectures

CPU

GPU

FPGA

CGRA

ASIC

## Applications

ML

Vision

Graph

DB

## Requirements

Security

Energy

Mem

2

# Hardware accelerators

Scopes

Architectures

Applications

DIFFICULT TO PROGRAM

DIFFICULT TO OPTIMISE

DIFFICULT TO PORT

EVER-CHANGING

CPU

GPU

FPGA

CGRA

ASIC

Mobile

Embed
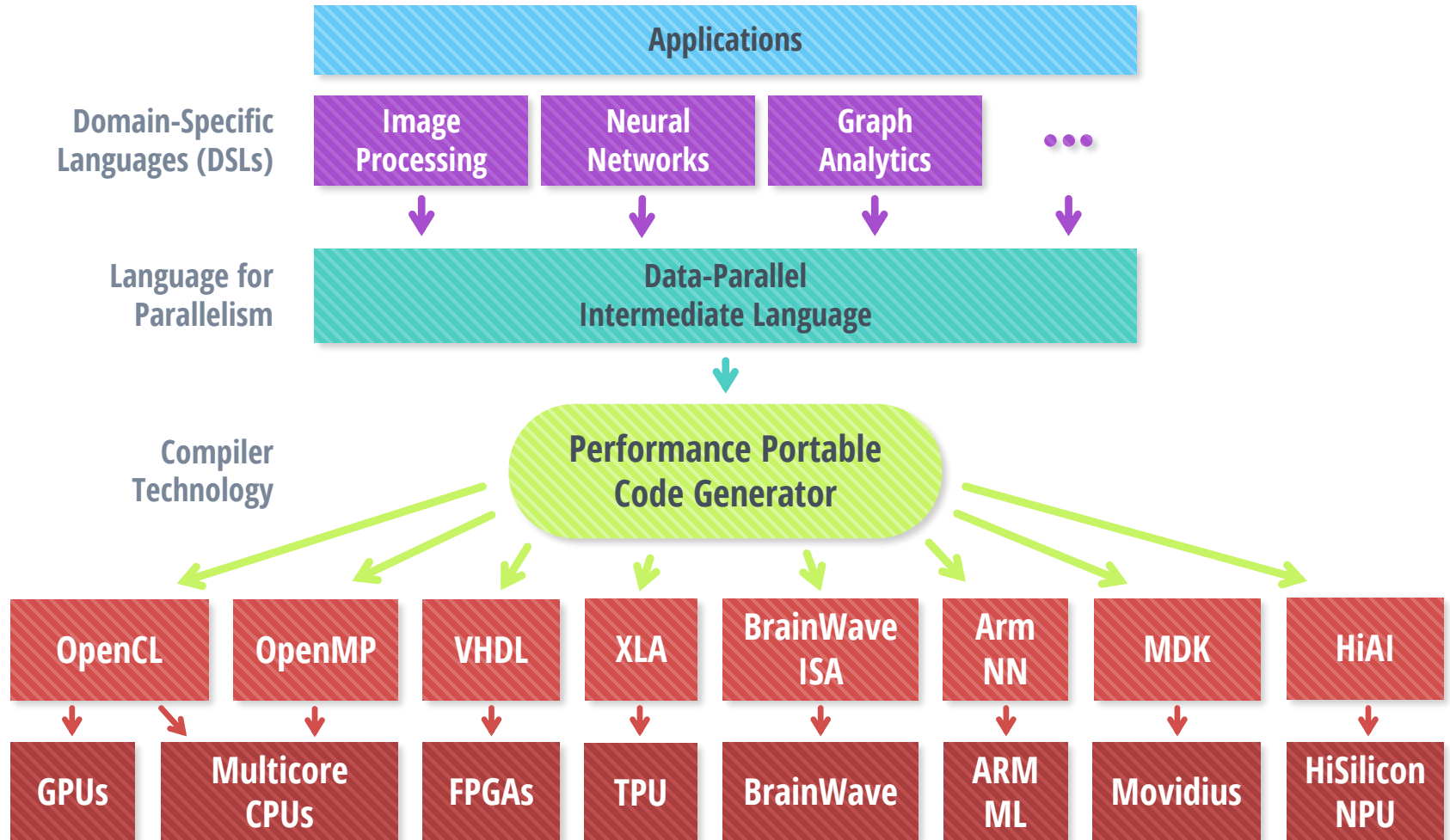
ML

Vision

Graph

DB

Security

Energy

Requirements

# Current landscape

| Applications | | | | | | | |
|---|---|---|---|---|---|---|---|
| OpenCL | OpenMP | VHDL | XLA | BrainWave ISA | Arm NN | MDK | HiAI |
| GPUs | Multicore CPUs | FPGAs | TPU | BrainWave | ARM ML | Movidius | HiSilicon NPU |

# 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 | XLA | BrainWave ISA | Arm NN | MDK | HiAI |

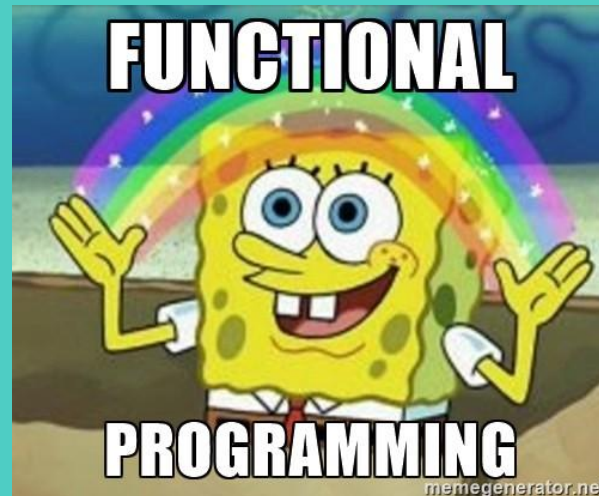| GPUs | Multicore CPUs | FPGAs | TPU | BrainWave | ARM ML | Movidius | HiSilicon NPU |

# What is the right interface for HW accelerators?

# What is the right interface for HW accelerators?

# Functional approach

## Abstract

- Expresses algorithm (WHAT), not implementation (HOW)

## High-level

- Captures plenty of algorithmic meta-info for analysis

## Pure

- Easy to transform

## Safe

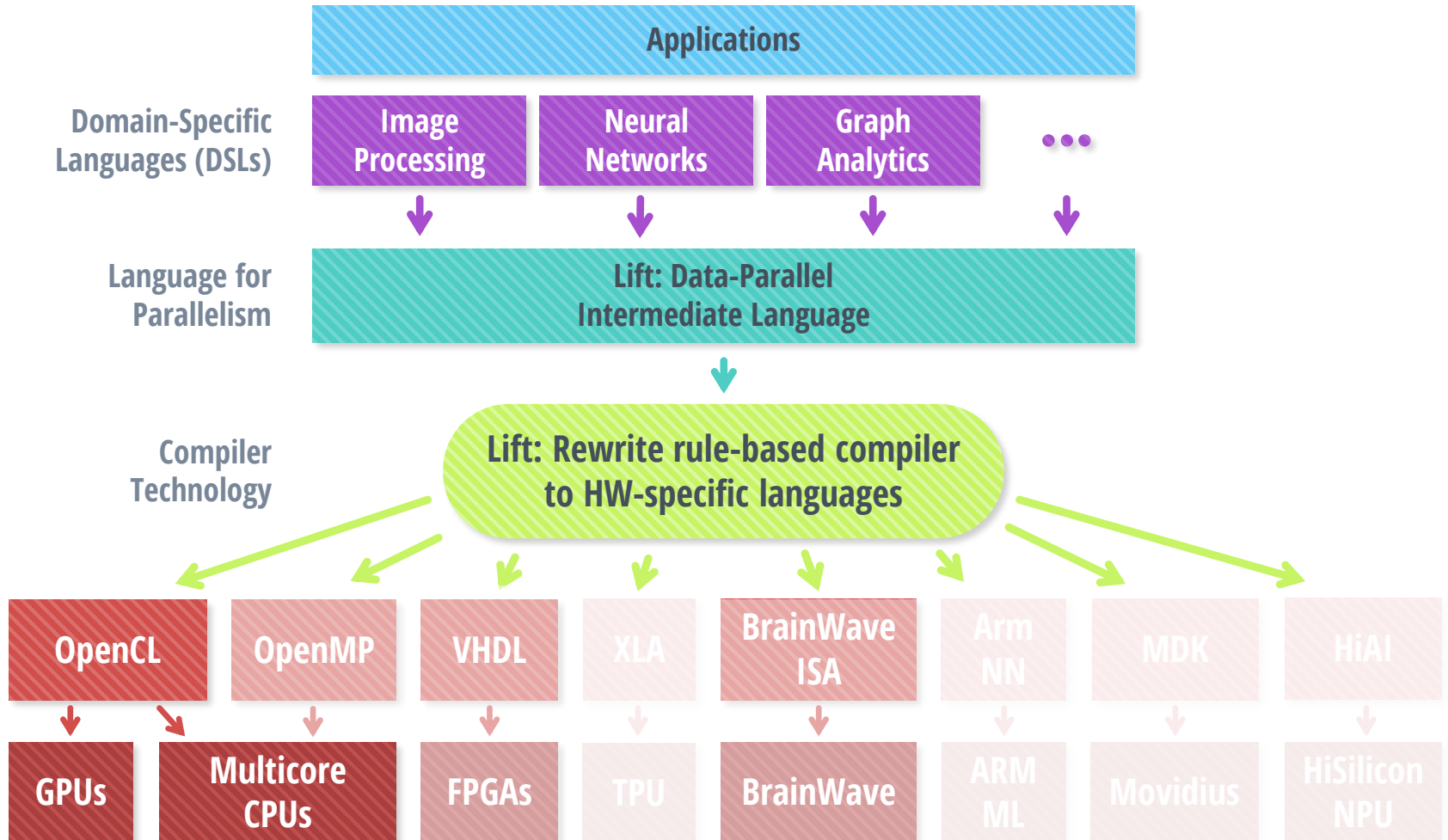- Easier to use and parallelise

## Expressive

- Control flow
- Memory management

## Composable

- Easier to maintain, code-reuse

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

# Lift

# Lift

🔳 Functional data-parallel language and compiler

```
1 A >> map(λ(rowA) =>
2   B >> map(λ(colB) =>
3     zip(rowA, colB) >>
4     map(*) >> reduce(0, +)))
```

➡️

```
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*j + K*N*i] =
5                  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*j + K*N*i];
10         }
11     }
12 }
```

# Lift IR

| Algorithmic patterns | Data types | Address space operators | Casters | Data operators |
|---|---|---|---|---|
| Map, Reduce | Int, Float | toGlobal | toVector | add, mul |
| Zip, Split | Vector | toLocal | toScalar | dot |
| Scatter, Gather | Array | toPrivate | | tanh |
| Slide | | | | |

# Lift IR: Views

```
Reorder(stride(s)) >> Map(f)
```

**NO WRITES TO MEMORY**     **TRANSFORMED READS**

▣ Virtual composable data layout transformations
  ▫ Reorder, Transpose, Slide, Slice, etc
▣ Expressed with Views
▣ Help avoid extra memory writes

12

# Lift IR

| Algorithmic patterns | Data types | Address space operators | Casters | Data operators |
|---|---|---|---|---|
| Map, Reduce | Int, Float | toGlobal | toVector | add, mul |
| Zip, Split | Vector | toLocal | toScalar | dot |
| Scatter, Gather | Array | toPrivate | | tanh |
| Slide | | | | |

# Lift IR

| IR level | Algorithmic patterns | Data types | Address space operators | Casters | Data operators |
|---|---|---|---|---|---|
| **DSL** | **conv, lstm** **blur, sharpen** **select** | Vector Matrix Tensor | | | |
| **Generic** | Map, Reduce Zip, Split Scatter, Gather Slide | Int, Float Vector Array | toGlobal toLocal toPrivate | toVector toScalar | add, mul dot tanh |
| **Platform-specific** | MapGlobal MapLocal ReduceSeq | Int8 Float16 Float32 | toDRAM toSRAM toRegistor | toInt8 toFloat16 | VVMul MVMul VTanh |

# How do we achieve performance portability?

# Lift: 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)
```

- ▣ Express algorithmic implementation choices
- ▣ Preserve semantic correctness
- ▣ Leverage algorithmic info

- ▣ Decouples optimisation from code generation

# Lift: Rewrite Rules

| IR level | Rewrite rules | |
|---|---|---|
| **DSL** | ▣ Algorithm choices for high-level primitives<br>▣ Precision level<br>▣ ... | ```conv(..)  ↦  stencilConv(..)conv(..)  ↦  gemmConv(..)conv(..)  ↦  winogradConv(..)``` |
| **Generic** | ▣ Split-join rule<br>▣ Map fusion rule<br>▣ Reduce rules<br>▣ ... | ```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)``` |
| **Platform-specific** | ▣ Using built-ins<br>▣ Lowering to the platform programming model<br>▣ ... | ```map(*) >> reduce(0, +)  ↦  dot()map(map(f))               ↦  mapWrg(mapLcl(f))map(f)                    ↦  asVector() >>                          map(vectorise(f)) >>                          asScalar()``` |

# Lift: rewriting

**HOW TO OPTIMISE?**

```
1  for (int n = 0; n < N; n++) {
2      out[n] = B[n];
3      for (int d = 0; d < D; d++) {
4          out[n] += X[d] * W[d + D*n];
5      }
6  }
```

# Lift: rewriting

```
1 layer(W: float[N][D], B: float[N],
2        X: float[D]):    float[N] =
3   zip(W, B) >> map(λ(Wn, Bn) =>
4     zip(Wn, X) >> map(*) >> reduce(Bn, +))
```

**HARD STARTING POINT**

```
1 for (int n = 0; n < N; n++) {
2     out[n] = B[n];
3     for (int d = 0; d < D; d++) {
4         out[n] += X[d] * W[d + D*n];
5     }
6 }
```

# Lift: rewriting

```
1 layer(W: float[N][D], B: float[N],
2       X: float[D]):    float[N] =
3   zip(W, B) >> map(λ(Wn, Bn) =>
4     zip(Wn, X) >> map(*) >> reduce(Bn, +))
```

```
1  layer(W: float[N][D], B: float[N],
2         X: float[D]):    float[N] =
3    zip(W, B) >> map(λ(Wn, Bn) =>
4      zip(Wn, X) >>
5      concat(
6          slice(0, (D/64)*64) >> split(64) >>
7          map(map(*) >> reduce(0, +)) >>
8          reduce(0, +),
9
10         slice((D/64)*64, D) >>
11         map(*) >> reduce(0, +)) >>
12
13     reduce(Bn, +))
```

```
1 for (int n = 0; n < N; n++) {
2     out[n] = B[n];
3     for (int d = 0; d < D; d++) {
4         out[n] += X[d] * W[d + D*n];
5     }
6 }
```

# Lift: rewriting

```
1 layer(W: float[N][D], B: float[N],
2       X: float[D]):    float[N] =
3  zip(W, B) >> map(λ(Wn, Bn) =>
4    zip(Wn, X) >> map(*) >> reduce(Bn, +))
```

**Built-in primitive**

```
map(*) >> reduce(0, +)
      ↧    ↧    ↧
  dot_product_accel()
```

```
1 layer(W: float[N][D], B: float[N],
2       X: float[D]):    float[N] =
3  zip(W, B) >> map(λ(Wn, Bn) =>
4    zip(Wn, X) >>
5    concat(
6       slice(0, (D/64)*64) >> split(64) >>
7       map(map(*) >> reduce(0, +)) >>
8       reduce(0, +),
9
10      slice((D/64)*64, D) >>
11      map(*) >> reduce(0, +)) >>
12
13   reduce(Bn, +))
```

**Exploitation**

```
1 for (int n = 0; n < N; n++) {
2    out[n] = B[n];
3    for (int d = 0; d < D; d++) {
4       out[n] += X[d] * W[d + D*n];
5    }
6 }
```

```
1 layer(W: float[N][D], B: float[N],
2       X: float[D]):    float[N] =
3  zip(W, B) >> mapWrg(λ(Wn, Bn) =>
4    zip(Wn, X) >>
5    concat(
6       slice(0, (D/64)*64) >> split(64) >>
7       mapLcl(dot_product_accel()) >>
8       reduce(0, +),
9
10      slice((D/64)*64, D) >>
11      mapSeq(*) >> reduce(0, +)) >>
12
13   reduce(Bn, +))
```

# Lift: rewriting

```
1 layer(W: float[N][D], B: float[N],
2       X: float[D]):    float[N] =
3  zip(W, B) >> map(λ(Wn, Bn) =>
4     zip(Wn, X) >> map(*) >> reduce(Bn, +))
```

```
1 layer(W: float[N][D], B: float[N],
2       X: float[D]):    float[N] =
3  zip(W, B) >> map(λ(Wn, Bn) =>
4     zip(Wn, X) >>
5     concat(
6         slice(0, (D/64)*64) >> split(64) >>
7         map(map(*) >> reduce(0, +)) >>
8         reduce(0, +),
9
10         slice((D/64)*64, D) >>
11         map(*) >> reduce(0, +)) >>
12
13     reduce(Bn, +))
```

**Built-in primitive**

```
map(*) >> reduce(0, +)
     ↧     ↧      ↧
  dot_product_accel()
```

**Parallelisation choice**

```
map(.. >> map(..) >> ..)
    ↧        ↧       ↧
mapWrg(.. >> mapLcl(..) >> ..)
```

```
1 for (int n = 0; n < N; n++) {
2     out[n] = B[n];
3     for (int d = 0; d < D; d++) {
4         out[n] += X[d] * W[d + D*n];
5     }
6 }
```

**Exploitation**

```
1 layer(W: float[N][D], B: float[N],
2       X: float[D]):    float[N] =
3  zip(W, B) >> mapWrg(λ(Wn, Bn) =>
4     zip(Wn, X) >>
5     concat(
6         slice(0, (D/64)*64) >> split(64) >>
7         mapLcl(dot_product_accel()) >>
8         reduce(0, +),
9
10         slice((D/64)*64, D) >>
11         mapSeq(*) >> reduce(0, +)) >>
12
13     reduce(Bn, +))
```

# Lift: rewriting

Search

```
1  layer(W: float[N][D], B: float[N],
2        X: float[D]):    float[N] =
3   zip(W, B) >> map(λ(Wn, Bn) =>
4      zip(Wn, X) >> map(*) >> reduce(Bn, +))
```

```
1  layer(W: float[N][D], B: float[N],
2        X: float[D]):    float[N] =
3   zip(W, B) >> map(λ(Wn, Bn) =>
4      zip(Wn, X) >>
5      concat(
6          slice(0, (D/64)*64 >> split(64) >>
7          map(map(*) >> reduce(0, +)) >>
8          reduce(0, +),
9
10         slice((D/64)*64, D) >>
11         map(*) >> reduce(0, +)) >>
12
13         reduce(Bn, +))
```

**Built-in primitive**

```
map(*) >> reduce(0, +)
  ↨      ↨      ↨
dot_product_accel()
```

**Parallelisation choice**

```
map(.. >> map(..) >> ..)
   ↨       ↨       ↨
mapWrg(.. >> mapLcl(..) >> ..)
```

Exploitation

```
1  layer(W: float[N][D], B: float[N],
2        X: float[D]):    float[N] =
3   zip(W, B) >> mapWrg(λ(Wn, Bn) =>
4      zip(Wn, X) >>
5      concat(
6          slice(0, (D/64)*64) >> split(64) >>
7          mapLcl(dot_product_accel()) >>
8          reduce(0, +),
9
10         slice((D/64)*64, D) >>
11         mapSeq(*) >> reduce(0, +)) >>
12
13         reduce(Bn, +))
```

```
1  for (int n = get_group_id(0); n < N; n += get_num_groups(0)) {
2      out[n] = B[n];
3      for (int t = get_local_id(0); t < (D/64); t += get_local_size(0)) {
4          out[n] += dot_product_accel(W + (t*64 + n*D), X + (t*64));
5      }
6      barrier(CLK_LOCAL_MEM_FENCE);
7
8      if (get_local_id(0) < 1) {
9          for (int d = D/64*64; d < D; d += 1) {
10             out[n] += W[d + n*D] * X[d];
11         }
12     }
13 }
```

Code generation

23

# Lift: Rewrite Rules

- Domain-specific and generic
- Reusable
- Provably correct
- Self-contained, extensible

# Lift: Constraint Inference

- Required for valid search space generation when using tuning parameters
- Leverages algorithmic meta-info
- Can express heuristics and HW restrictions

$$\text{Split}(s) \ \$ \ [T]_N \qquad => \qquad N \ \% \ s == 0$$

$$\text{asVector}(v) \ \$ \ [T]_N \qquad => \qquad N \ \% \ v == 0$$

$$\text{Slide}(len, \ step) \ o \ [T]_N \qquad => \qquad N >= len$$

$$\text{SlideStrict}(len, \ step) \ o \ [T]_N \qquad => \qquad N >= len \ \&\& $$
$$N \ \% \ ((N - (len - step)) \ / \ step) == 0$$

# Lift: Search Space Exploration

- ▣ Uniform random sampling
- ▣ Predictor models
- ▣ Genetic algorithms
- ▣ ...

# Lift: Research Directions

- Linear algebra
- Sparse data parallelism
- Optimising reductions
- Stencil computations
- 3D wave modelling
- High-level synthesis for FPGAs
- Machine Learning

# Lift for Machine Learning

- ▣ Machine Learning
  - ▢ Convolution inference optimisation
  - ▢ Platforms: Mali GPUs, BrainWave

```
1  def convLayer(kernelsWeights : [[[[float]_inputChannels]kernelWidth]kernelHeight]numKernel,
2                kernelsBiases  : [float]_numKernel,
3                inputData      : [[[float]_inputChannels]inputWidth]inputHeight,
4                padSize        : (int,int,int,int),           kernelStride   : (int,int))
5                : [[[float]_outWidth]outHeight]numKernel = {
6
7    val paddedInput = pad2D(padSize, value = 0 ,inputData)
8    val slidingWindows = slide2D(kernelHeight, kernelWidth, kernelStride._1, kernelStride._2, paddedInput)
9    map2D(slidingWindow ->
10     map((singleKernelWeights, singleKernelBias) ->
11       reduce(init = singleKernelBias, f = (acc, (x, w)) -> {acc + x * w},
12         zip(join(join(slidingWindow)), join(join(singleKernelWeights)))),
13       zip(kernelsWeights, kernelsBiases)),
14     slidingWindows)}
```

# Lift for Machine Learning

▣ Machine Learning
- □ Convolution inference optimisation
- □ Platforms: Mali GPUs, BrainWave

# Lift for Machine Learning

- ▣ Machine Learning
  - ▫ Convolution inference optimisation
  - ▫ Platforms: Mali GPUs, BrainWave



■ Lift vs ARM-C Direct   ■ ARM-C GEMM vs ARM-C Direct

**VGG layers on Mali G72**

# Lift for Machine Learning

Naums Mogers, PhD student, Edinburgh
*How to best exploit HW accelerators?*

Christof Schlaak, PhD student, Edinburgh
*How to generate accelerator architectures?*

Lu Li, Postdoctoral Researcher, Edinburgh
*How to optimise the host code?*
*How to drive the rewriting process?*

Christophe Dubach, Reader, Edinburgh
*All of the above*

# Lift source code is published

https://github.com/lift-project/lift

http://www.lift-project.org

References

*(icons) Noun Project,* https://thenounproject.com
*(icons) Font Awesome,* https://fontawesome.com