# Matrix Multiplication Beyond Auto-Tuning:
# Rewrite Based GPU Code Generation

Michel Steuwer, Toomas Remmelg, Christophe Dubach

{michel.steuwer, toomas.remmelg, christophe.dubach}@ed.ac.uk

ICSA | Institute for Computing Systems Architecture

---

**Lift** is a novel approach to achieve high performance on parallel accelerators

Matrix multiplication → Image processing → Graph analytics → ...

**Lift** data parallel language

**Lift** rewrite-based system

OpenCL

NVIDIA GPU — AMD GPU — ARM GPU

lift-project.org

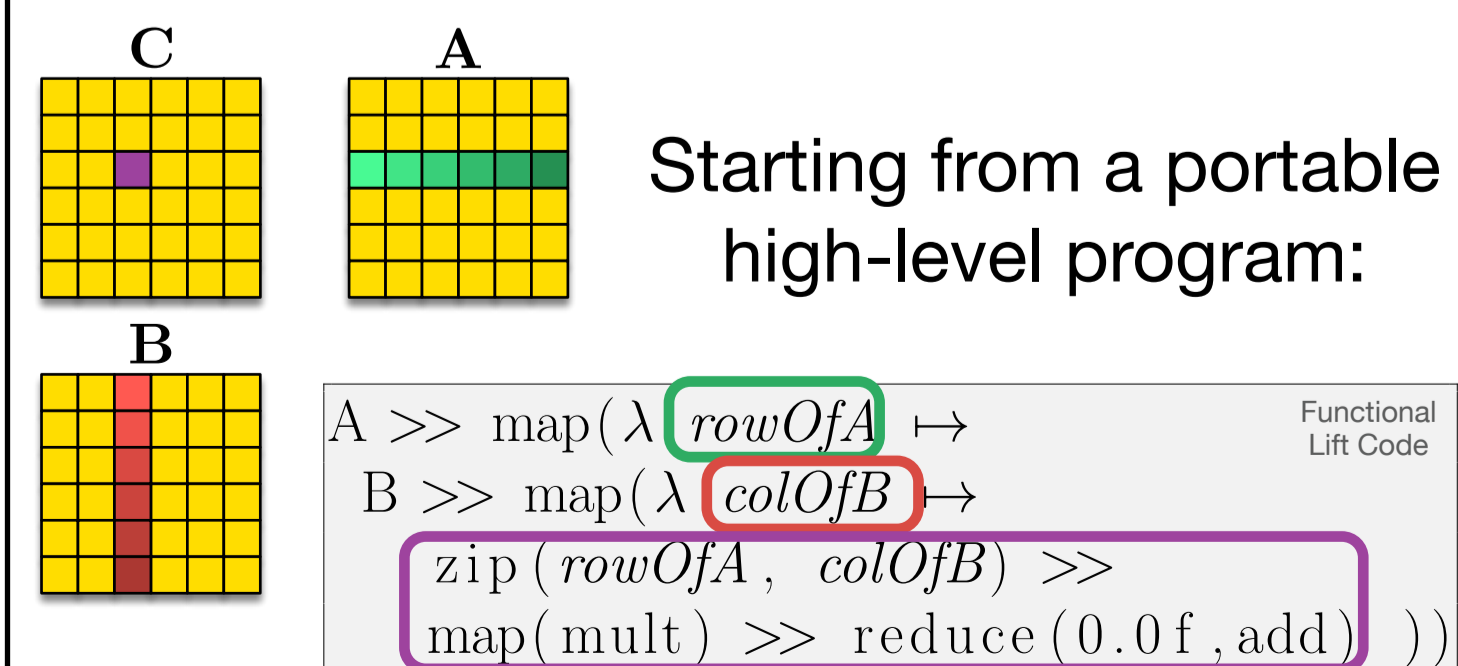Starting from a *high-level language*, rewrite rules are used to derive *optimised implementations*.

Goal: Achieving *Performance Portability*, i.e. high performance across different GPUs

---

## Inherent Limitation of *Auto-Tuning*

Chooses between a fixed number of optimisations and tuning parameters

$\Longrightarrow$ Falls short on new architectures! Performance portability cannot be achieved by only using auto-tuning

---

## Matrix Multiplication in the **Lift** Data Parallel Language

C    A
B

Starting from a portable high-level program:

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

---

## **Lift** Rewrite-Based System:
### Compilation and Optimisation Using Provably Correct Rewrite Rules

*Algorithmic rewrite rules* express optimisation choices

**Split-join rule**

$$\mathrm{map}(f) \implies \mathrm{split}(k) \gg \mathrm{map}(\mathrm{map}(f)) \gg \mathrm{join}$$

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

```
                                       Imperative
                                       OpenCL Code
1  for (int i = 0; i<M/2; i++) {
2    for (int l = 0; l<2; 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-map interchange rule**

$$X \gg \mathrm{map}(\lambda\ x \mapsto Y \gg \mathrm{map}(\lambda\ y \mapsto f))$$
$$\implies$$
$$Y \gg \mathrm{map}(\lambda\ y \mapsto X \gg \mathrm{map}(\lambda\ x \mapsto f)) \gg \mathrm{transpose}$$

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

```
                                       Imperative
                                       OpenCL Code
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 }
```

**More algorithmic rules**

...

*OpenCL specific rules* express *mapping choices*

**OpenCL specific rules**

$$\mathrm{map}(f) \implies \mathrm{mapGlb}_{\{0,1,2\}}(f)$$
$$\mathrm{map}(f) \implies \mathrm{mapSeq}(f)$$
$$f \implies \mathrm{toGlobal}(f)$$
$$f \implies \mathrm{toPrivate}(f)$$
$$\mathrm{map}(f) \implies \mathrm{asVector}(n,\ b) \gg \mathrm{map}(\mathrm{vectorize}(n,\ f)) \gg \mathrm{asScalar}$$

```
                                                   Functional
                                                   Lift Code
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, 0.0f) >>
            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(dot) >> reduceSeq(acc, add)
  ) >> join
  )
) >> toGlobal(mapSeq(mapSeq(mapSeq(id))))
>> transpose() >>
map(transpose) >> transpose
) >> join >> transpose
) >> join
```
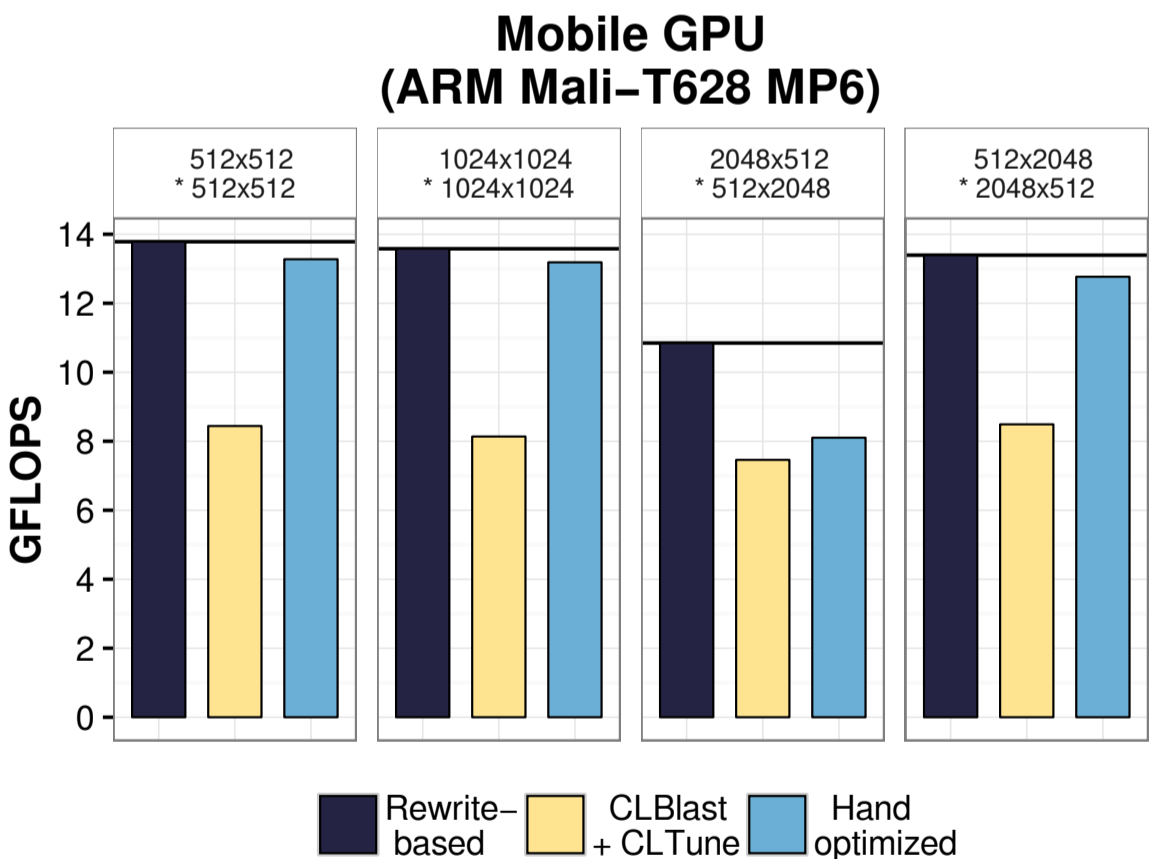
---

## Highly-Optimised Code for the Embedded ARM Mali GPU

```
                                       Imperative
                                       OpenCL Code
1  int i = get_global_id(0);
2  int j = get_global_id(1);
3
4  float temp_0; float temp_1;
5  float temp_2; float 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 = dot(vload4(k + K*i/2, A),
12     vload4(k + K*j/2, B));
13   acc_0 += temp_0;
14
15   temp_1 = dot(vload4(k + K*i/2, A),
16     vload4(k + K + 2*K*j/4, B));
17   acc_1 += temp_1;
18
19   temp_2 = dot(vload4(k + K + 2*K*i/4, A),
20     vload4(k + K*j/2, B));
21   acc_2 += temp_2;
22
23   temp_3 = dot(vload4(k + K + 2*K*i/4, A),
24     vload4(k + K + 2*K*j/4, B));
25   acc_3 += temp_3;
26 }
27 C[2*N*i + 2*j] = id(acc_0);
28 C[1 + 2*N*i + 2*j] = id(acc_1);
29 C[N + 2*N*i + 2*j] = id(acc_2);
30 C[1 + N + 2*N*i + 2*j] = id(acc_3);
```

---

## High Performance on Mali

**Mobile GPU (ARM Mali-T628 MP6)**

512x512 * 512x512 | 1024x1024 * 1024x1024 | 2048x512 * 512x2048 | 512x2048 * 2048x512

GFLOPS

Legend: Rewrite-based | CLBlast + CLTune | Hand optimized

---

## Performance Portability
### Automatic Exploration Using the **Lift** Rewrite-Based System

**Desktop GPU (Nvidia GeForce GTX Titan Black)**
**Desktop GPU (AMD Radeon HD 7970)**

GFLOPS

Legend: Rewrite-based | CLBlast + CLTune

---