

Performance Portable GPU Code Generation for Matrix Multiplication



Toomas Remmelg



Thibaut Lutz



Michel Steuwer



Christophe Dubach



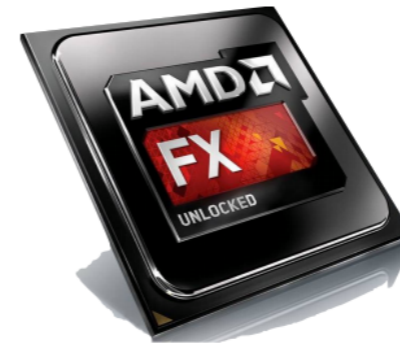
THE UNIVERSITY
of EDINBURGH

Supported by:

Google Oracle Labs

The Problem

- Parallel processors everywhere
- Many different types: CPUs, GPUs, ...
- Parallel programming is hard
- Optimising even harder
- **Problem:**
No portability of performance!



CPU



GPU

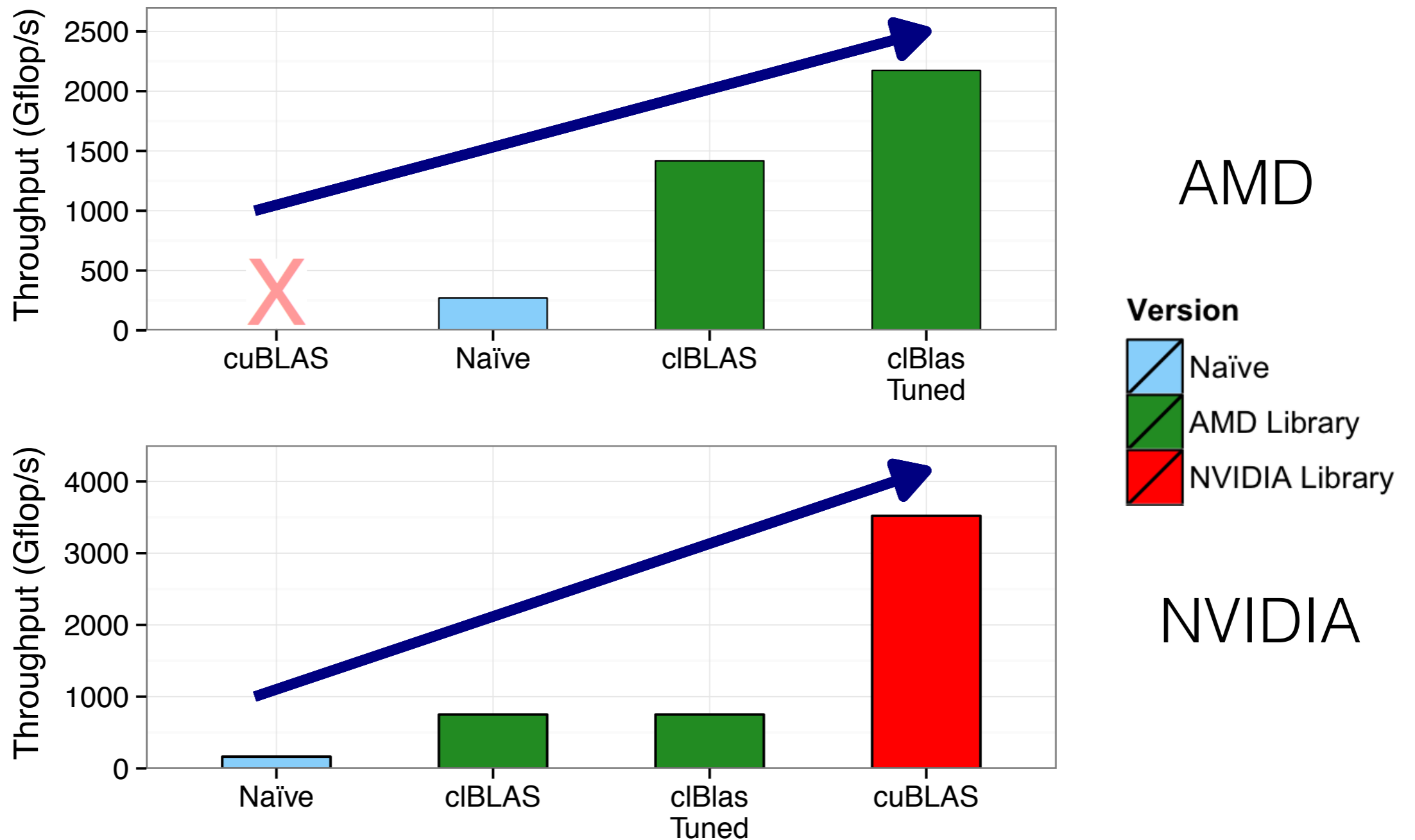


Accelerator



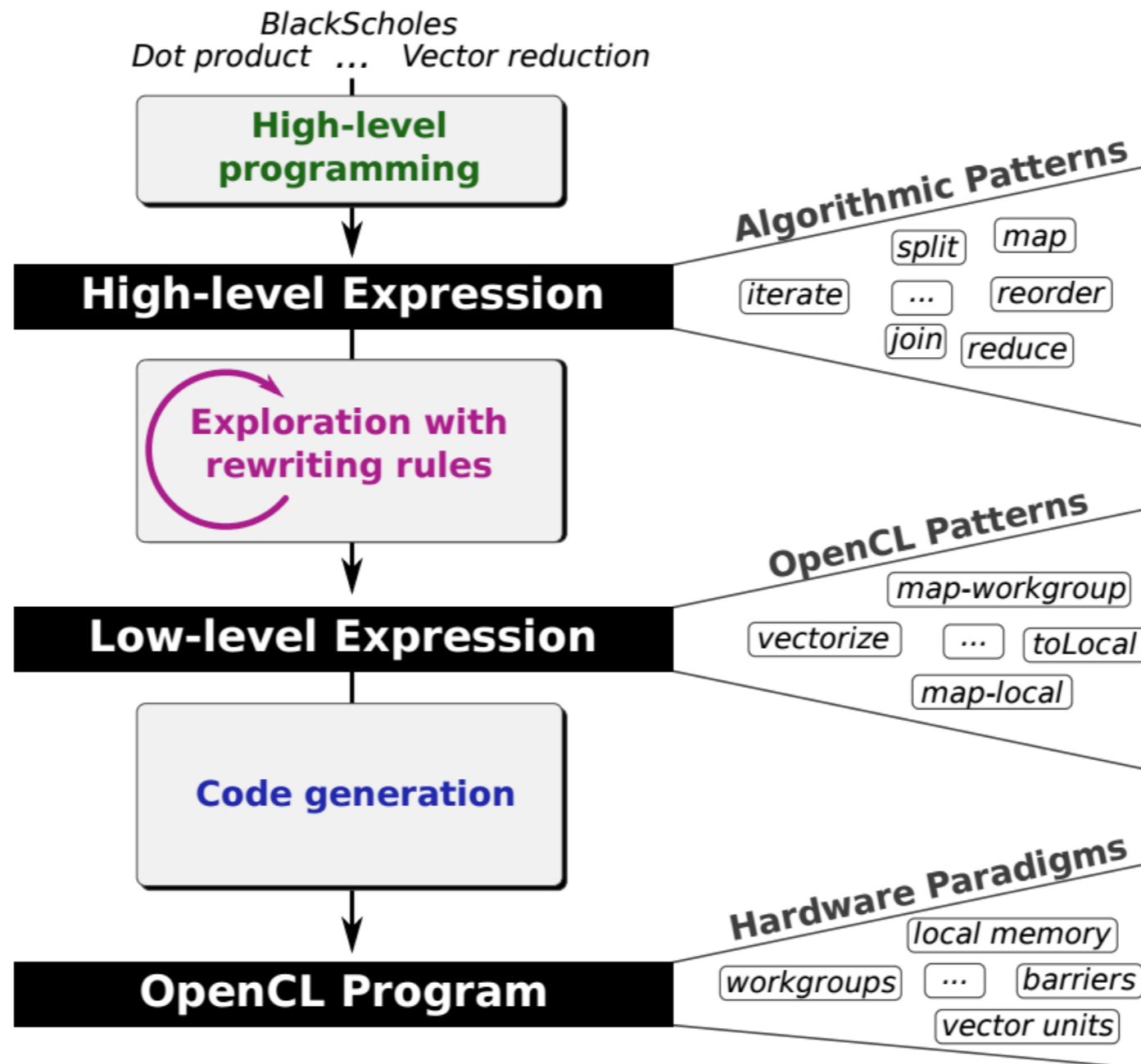
FPGA

Performance Portability of Matrix Multiplication

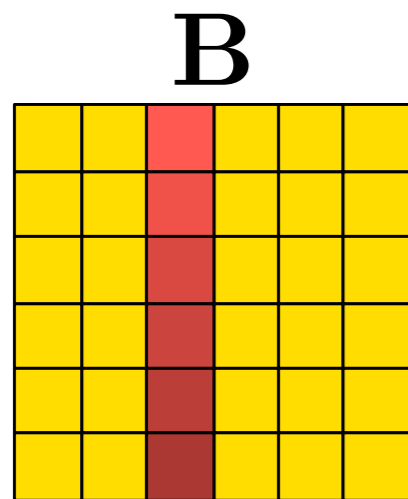
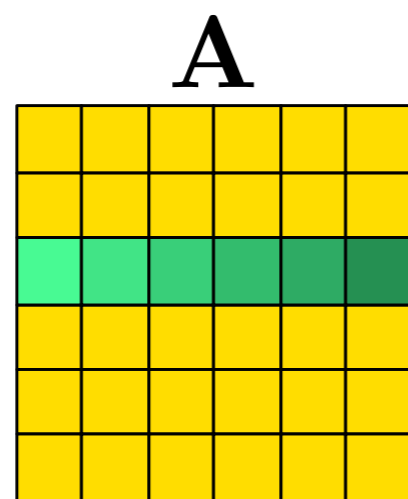
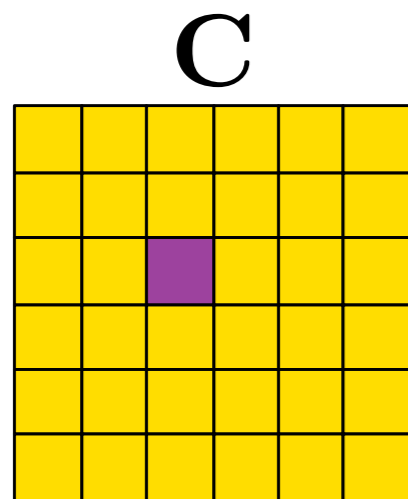


How to achieve performance portability?

Rewrite Rules



Matrix Multiplication Expressed Functionally



Functional Representation

$$\mathbf{A} * \mathbf{B} =$$

$$Map(\overrightarrow{row A} \mapsto$$

$$Map(\overrightarrow{col B} \mapsto$$

$$DotProduct(\overrightarrow{row A}, \overrightarrow{col B})$$

$$) \circ Transpose() \$ \mathbf{B}$$

$$) \$ \mathbf{A}$$

OpenCL

```

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 }

```

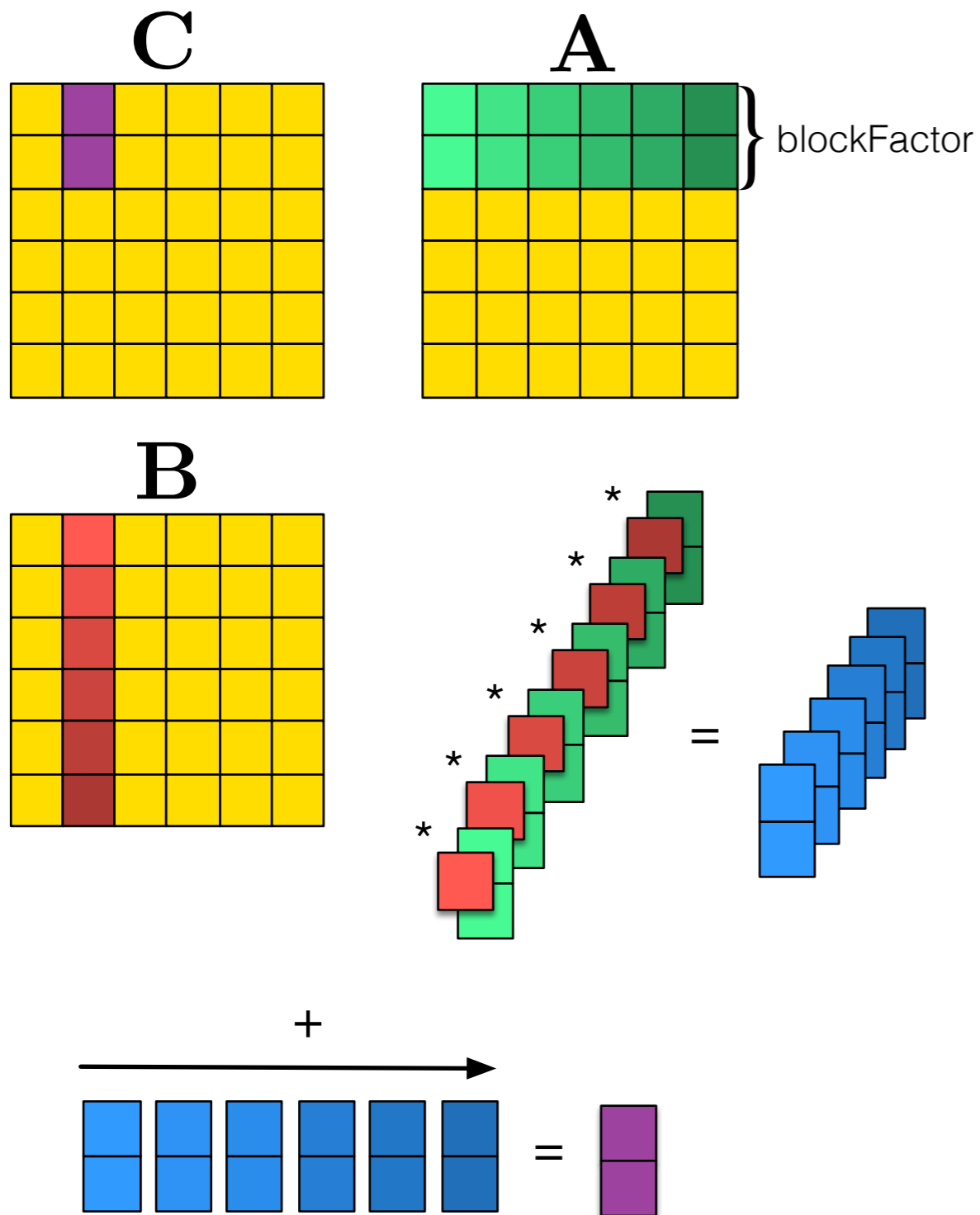
Traditional Optimisations

- **Register Blocking** Loading elements into registers and reusing them.
- **Tiling** Solving the problem by dividing matrices into smaller tiles.
- **Vectorisation** Using wider vector units if available.

Why can't this be automated by traditional compilers?

- **Complex analysis** Proving the optimisations are legal.
- **Conservative** Must always be correct.
- **No obvious defaults for parameters** Good tile and block sizes depend on hardware capabilities.

Register Blocking

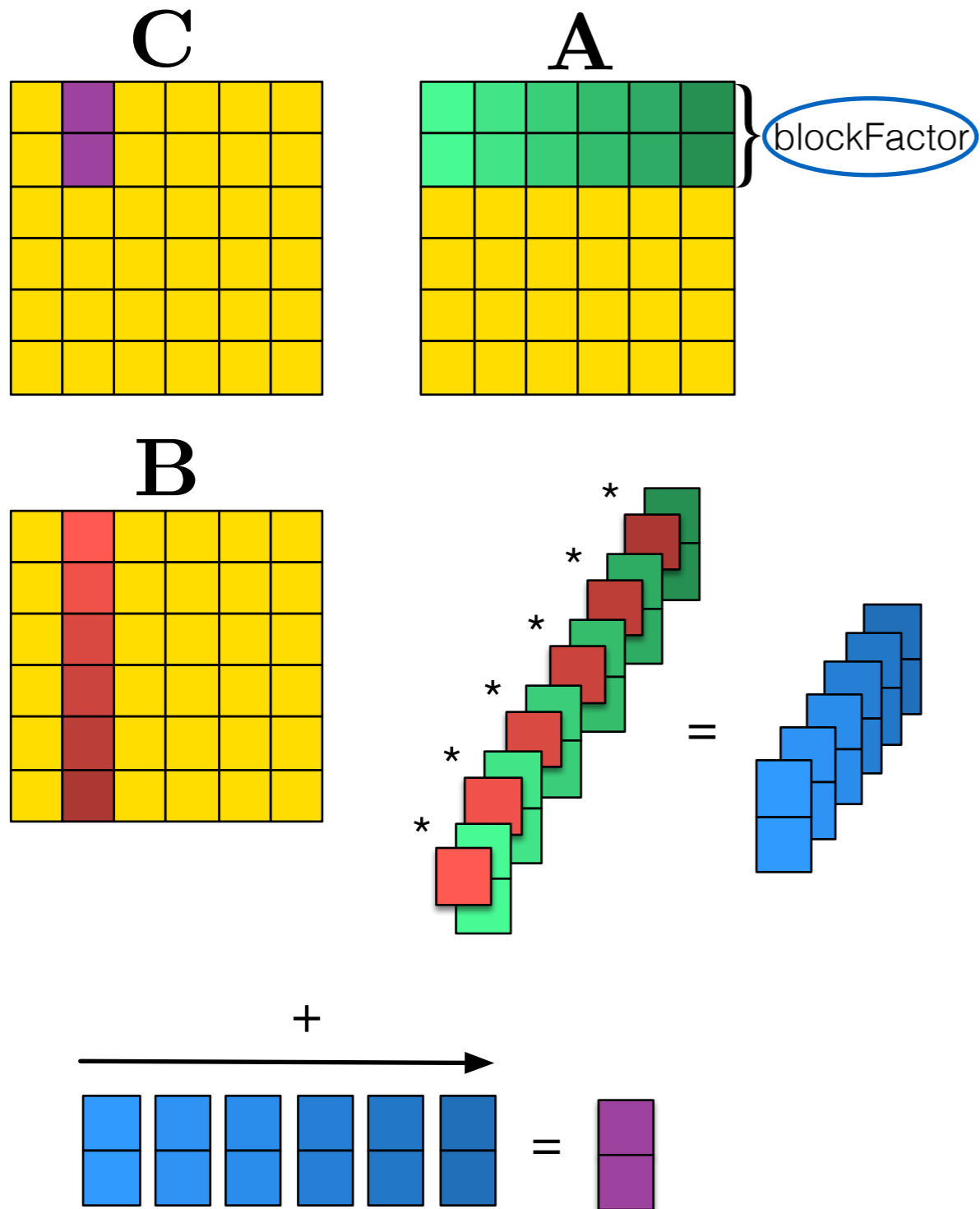


```

1 kernel void KERNEL(
2   const global float* restrict A,
3   const global float* restrict B,
4   global float* C, int K, int M, int N)
5 {
6   float acc[blockFactor];
7
8   for (int glb_id_1 = get_global_id(1);
9        glb_id_1 < M / blockFactor;
10        glb_id_1 += get_global_size(1)) {
11     for (int glb_id_0 = get_global_id(0); glb_id_0 < N;
12          glb_id_0 += get_global_size(0)) {
13
14         for (int i = 0; i < K; i += 1)
15             float temp = B[i * N + glb_id_0];
16         for (int j = 0; j < blockFactor; j += 1)
17             acc[j] +=
18                 A[blockFactor * glb_id_1 * K + j * K + i]
19                 * temp;
20
21         for (int j = 0; j < blockFactor; j += 1)
22             C[blockFactor * glb_id_1 * N + j * N + glb_id_0]
23               = acc[j];
24     }
25 }
26

```

Register Blocking

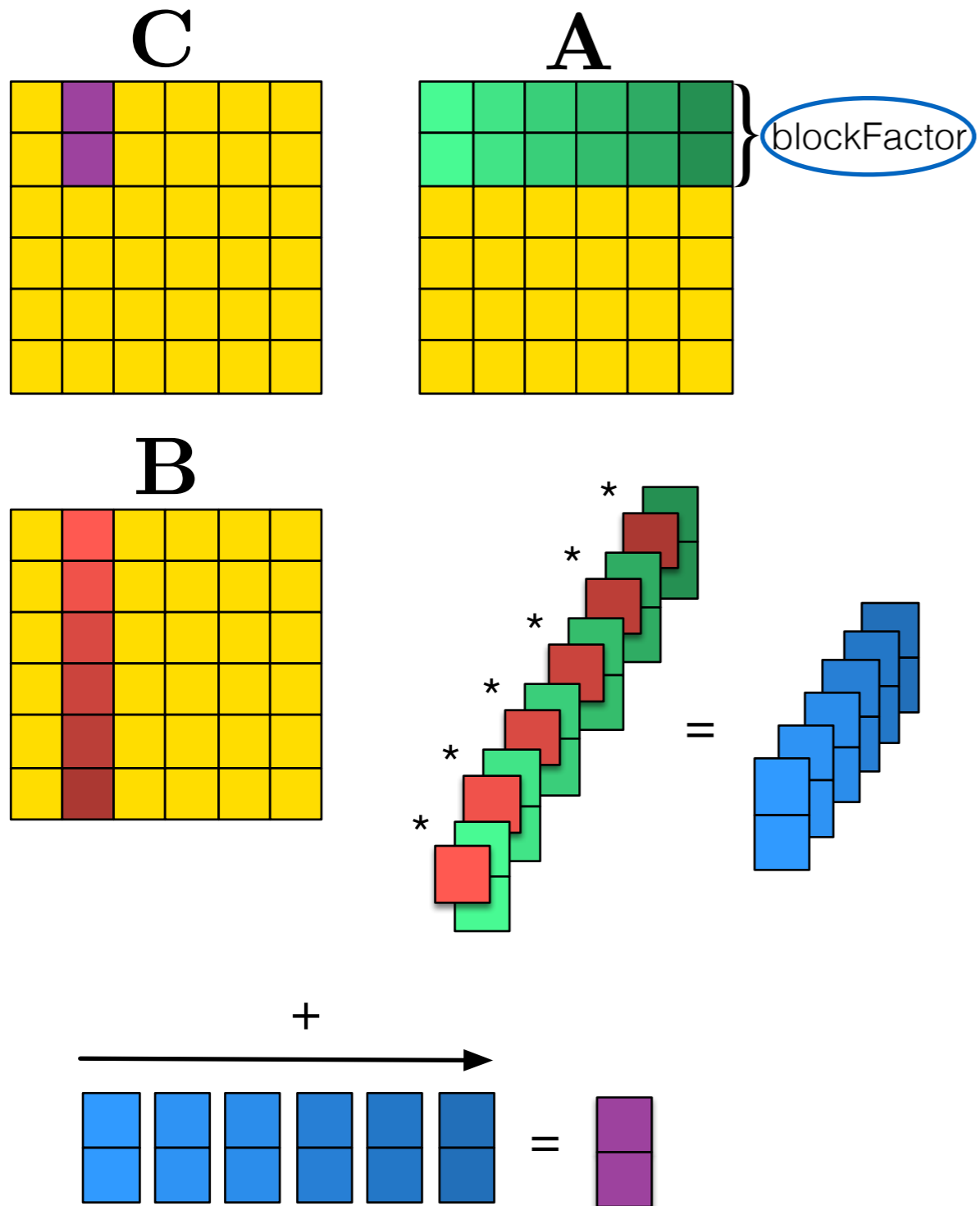


```

1 kernel void KERNEL(
2   const global float* restrict A,
3   const global float* restrict B,
4   global float* C, int K, int M, int N)
5 {
6   float acc[blockFactor];
7
8   for (int glb_id_1 = get_global_id(1);
9        glb_id_1 < M / blockFactor;
10        glb_id_1 += get_global_size(1)) {
11     for (int glb_id_0 = get_global_id(0); glb_id_0 < N;
12          glb_id_0 += get_global_size(0)) {
13
14         for (int i = 0; i < K; i += 1)
15             float temp = B[i * N + glb_id_0];
16         for (int j = 0; j < blockFactor; j += 1)
17             acc[j] +=
18                 A[blockFactor * glb_id_1 * K + j * K + i]
19                 * temp;
20
21         for (int j = 0; j < blockFactor; j += 1)
22             C[blockFactor * glb_id_1 * N + j * N + glb_id_0]
23               = acc[j];
24     }
25 }
26

```


Register Blocking



```

1 kernel void KERNEL(
2   const global float* restrict A,
3   const global float* restrict B,
4   global float* C, int K, int M, int N)
5 {
6   float acc[blockFactor];
7
8   for (int glb_id_1 = get_global_id(1);
9        glb_id_1 < M / blockFactor;
10        glb_id_1 += get_global_size(1)) {
11     for (int glb_id_0 = get_global_id(0); glb_id_0 < N;
12          glb_id_0 += get_global_size(0)) {
13
14         for (int i = 0; i < K; i += 1)
15             float temp = B[i * N + glb_id_0];
16         for (int j = 0; j < blockFactor; j += 1)
17             acc[j] +=
18                 A[blockFactor * glb_id_1 * K + j * K + i]
19                 * temp;
20
21         for (int j = 0; j < blockFactor; j += 1)
22             C[blockFactor * glb_id_1 * N + j * N + glb_id_0]
23               = acc[j];
24     }
25 }
26 }

```

Register Blocking

$$\begin{aligned} & \text{Map}(\overrightarrow{\text{row } A} \mapsto \\ & \quad \text{Map}(\overrightarrow{\text{col } B} \mapsto \\ & \quad \quad \text{Reduce}(+) \circ \text{Map}(*)) \\ & \quad \quad \$ \text{Zip}(\overrightarrow{\text{row } A}, \overrightarrow{\text{col } B}) \\ & \quad) \circ \text{Transpose}() \$ \mathbf{B} \\ &) \$ \mathbf{A} \end{aligned}$$
$$\text{Map}(f) \Rightarrow \text{Join}() \circ \text{Map}(\text{Map}(f)) \circ \text{Split}(k)$$

Register Blocking

$Map(\overrightarrow{rowA} \mapsto$
 $Map(\overrightarrow{colB} \mapsto$
 $Reduce(+) \circ Map(*)$
 $\$ Zip(\overrightarrow{rowA}, \overrightarrow{colB})$
 $) \circ Transpose() \$ B$
 $) \$ A$



$Join() \circ Map(rowsA \mapsto$
 $Map(\overrightarrow{rowA} \mapsto$
 $Map(\overrightarrow{colB} \mapsto$
 $Reduce(+) \circ Map(*)$
 $\$ Zip(\overrightarrow{rowA}, \overrightarrow{colB})$
 $) \circ Transpose() \$ B$
 $) \$ rowsA$
 $) \circ Split(blockFactor) \$ A$

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

Register Blocking

$Join() \circ Map(rowsA \mapsto$
 $Map(\overrightarrow{rowA} \mapsto$
 $Map(\overrightarrow{colB} \mapsto$
 $Reduce(+) \circ Map(*)$
 $\$ Zip(\overrightarrow{rowA}, \overrightarrow{colB})$
 $) \circ Transpose() \$ \mathbf{B}$
 $) \$ rowsA$
 $) \circ Split(blockFactor) \$ \mathbf{A}$

$Map(a \mapsto Map(b \mapsto f(a, b))) \Rightarrow$
 $Transpose() \circ Map(b \mapsto Map(a \mapsto f(a, b)))$

Register Blocking

$Join() \circ Map(rowsA \mapsto$

$Map(\overrightarrow{rowA} \mapsto$

$Map(\overrightarrow{colB} \mapsto$

$Reduce(+) \circ Map(*)$

$\$ Zip(\overrightarrow{rowA}, \overrightarrow{colB})$

$) \circ Transpose() \$ B$

$) \$ rowsA$

$) \circ Split(blockFactor) \$ A$



$Join() \circ Map(rowsA \mapsto$

$Transpose() \circ Map(\overrightarrow{colB} \mapsto$

$Map(\overrightarrow{rowA} \mapsto$

$Reduce(+) \circ Map(*)$

$\$ Zip(\overrightarrow{rowA}, \overrightarrow{colB})$

$) \$ rowsA$

$) \circ Transpose() \$ B$

$) \circ Split(blockFactor) \$ A$

$Map(a \mapsto Map(b \mapsto f(a, b))) \Rightarrow$

$Transpose() \circ Map(b \mapsto Map(a \mapsto f(a, b)))$

Register Blocking

$Join() \circ Map(rowsA \mapsto$
 $Transpose() \circ Map(\overrightarrow{colB} \mapsto$
 $Map(\overrightarrow{rowA} \mapsto$
 $Reduce(+) \circ Map(*$
 $\$ Zip(\overrightarrow{rowA}, \overrightarrow{colB})$
 $) \$ rowsA$
 $) \circ Transpose() \$ \mathbf{B}$
 $) \circ Split(blockFactor) \$ \mathbf{A}$

$$Map(f \circ g) \Rightarrow Map(f) \circ Map(g)$$

Register Blocking

$Join() \circ Map(rowsA \mapsto$
 $Transpose() \circ Map(\overrightarrow{colB} \mapsto$
 $Map(\overrightarrow{rowA} \mapsto$
 $Reduce(+) \circ Map(*)$
 $\$ Zip(\overrightarrow{rowA}, \overrightarrow{colB})$
 $) \$ rowsA$
 $) \circ Transpose() \$ \mathbf{B}$
 $) \circ Split(blockFactor) \$ \mathbf{A}$



$Join() \circ Map(rowsA \mapsto$
 $Transpose() \circ Map(\overrightarrow{colB} \mapsto$
 $Map($
 $Reduce(+)$
 $) \circ Map(\overrightarrow{rowA} \mapsto$
 $Map(*) \$ Zip(\overrightarrow{rowA}, \overrightarrow{colB})$
 $) \$ rowsA$
 $) \circ Transpose() \$ \mathbf{B}$
 $) \circ Split(blockFactor) \$ \mathbf{A}$

$$Map(f \circ g) \Rightarrow Map(f) \circ Map(g)$$

Register Blocking

$Join() \circ Map(rowsA \mapsto$
 $Transpose() \circ Map(\overrightarrow{colB} \mapsto$
 $Map($
 $Reduce(+)$
 $) \circ Map(\overrightarrow{rowA} \mapsto$
 $Map(*) \$ Zip(\overrightarrow{rowA}, \overrightarrow{colB})$
 $) \$ rowsA$
 $) \circ Transpose() \$ \mathbf{B}$
 $) \circ Split(blockFactor) \$ \mathbf{A}$

$Map(Reduce(f)) \Rightarrow$
 $Transpose() \circ Reduce(Map(f) \circ Zip())$

Register Blocking

$Join() \circ Map(rowsA \mapsto$
 $Transpose() \circ Map(\overrightarrow{colB} \mapsto$
 $Map($
 $Reduce(+)$
 $) \circ Map(\overrightarrow{rowA} \mapsto$
 $Map(*) \$ Zip(\overrightarrow{rowA}, \overrightarrow{colB})$
 $) \$ rowsA$
 $) \circ Transpose() \$ \mathbf{B}$
 $) \circ Split(blockFactor) \$ \mathbf{A}$



$Join() \circ Map(rowsA \mapsto$
 $Transpose() \circ Map(\overrightarrow{colB} \mapsto$
 $Transpose() \circ Reduce((\overrightarrow{acc}, \overrightarrow{next}) \mapsto$
 $Map(+)$ Zip(\overrightarrow{acc}, \overrightarrow{next})$
 $) \circ Transpose() \circ Map(\overrightarrow{rowA} \mapsto$
 $Map(*) \$ Zip(\overrightarrow{rowA}, \overrightarrow{colB})$
 $) \$ rowsA$
 $) \circ Transpose() \$ \mathbf{B}$
 $) \circ Split(blockFactor) \$ \mathbf{A}$

$Map(Reduce(f)) \Rightarrow$
 $Transpose() \circ Reduce(Map(f) \circ Zip())$

Register Blocking

$Join() \circ Map(rowsA \mapsto$

$Transpose() \circ Map(\overrightarrow{colB} \mapsto$

$Transpose() \circ Reduce((\overrightarrow{acc}, \overrightarrow{next}) \mapsto$

$Map(+) \$ Zip(\overrightarrow{acc}, \overrightarrow{next})$

$) \circ Transpose() \circ Map(\overrightarrow{rowA} \mapsto$

$Map(*) \$ Zip(\overrightarrow{rowA}, \overrightarrow{colB})$

$) \$ rowsA$

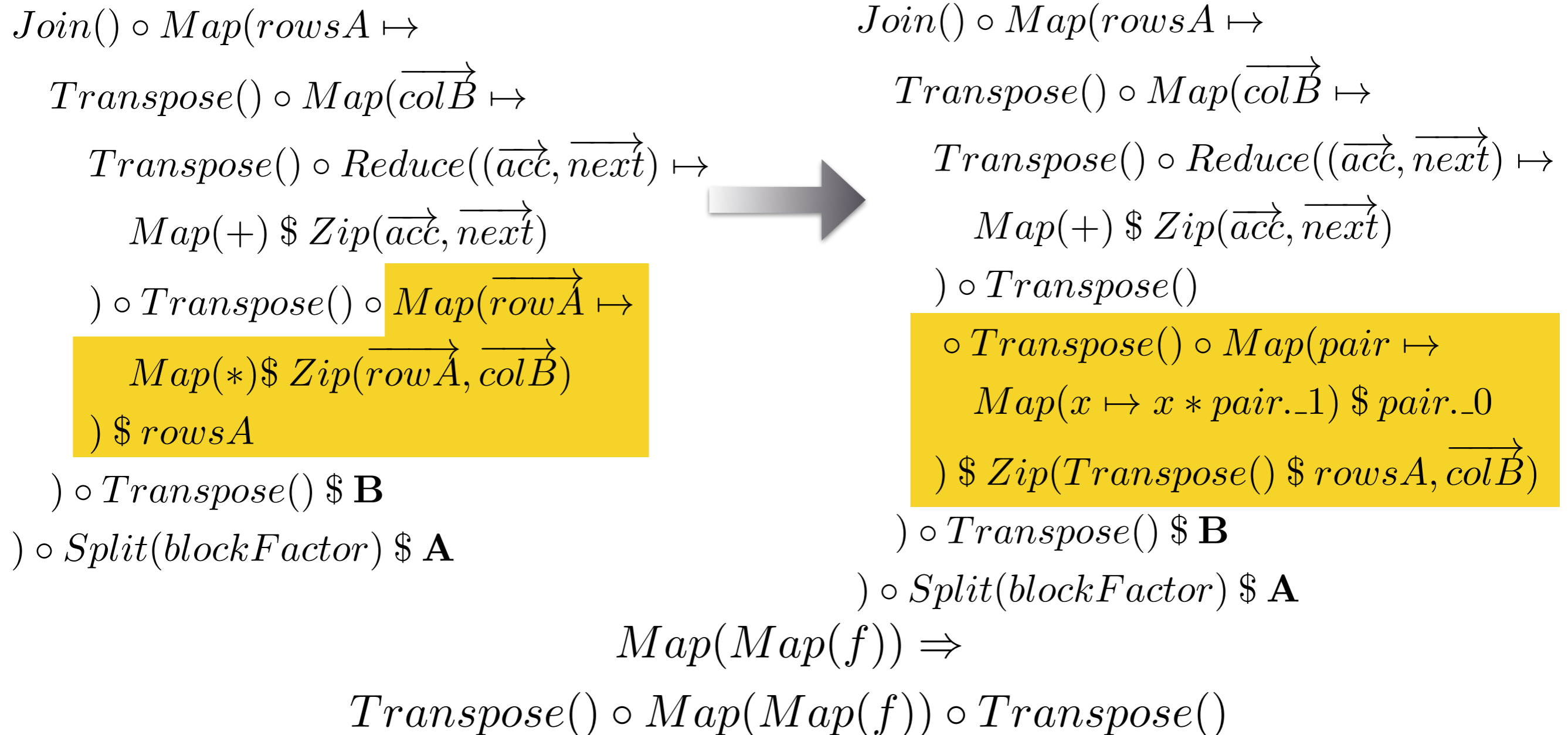
$) \circ Transpose() \$ \mathbf{B}$

$) \circ Split(blockFactor) \$ \mathbf{A}$

$Map(Map(f)) \Rightarrow$

$Transpose() \circ Map(Map(f)) \circ Transpose()$

Register Blocking



Register Blocking

$Join() \circ Map(rowsA \mapsto$
 $Transpose() \circ Map(\overrightarrow{colB} \mapsto$
 $Transpose() \circ Reduce((\overrightarrow{acc}, \overrightarrow{next}) \mapsto$
 $Map(+) \$ Zip(\overrightarrow{acc}, \overrightarrow{next})$
 $) \circ Transpose()$
 $\circ Transpose() \circ Map(pair \mapsto$
 $Map(x \mapsto x * pair._1) \$ pair._0$
 $) \$ Zip(Transpose() \$ rowsA, \overrightarrow{colB})$
 $) \circ Transpose() \$ \mathbf{B}$
 $) \circ Split(blockFactor) \$ \mathbf{A}$

$Transpose() \circ Transpose() \Rightarrow id$

Register Blocking

$$\begin{aligned}
 &Join() \circ Map(rowsA \mapsto \\
 &Transpose() \circ Map(\overrightarrow{colB} \mapsto \\
 &Transpose() \circ Reduce((\overrightarrow{acc}, \overrightarrow{next}) \mapsto \\
 &Map(+) \$ Zip(\overrightarrow{acc}, \overrightarrow{next}) \\
 &)\circ Transpose() \\
 &\circ Transpose() \circ Map(pair \mapsto \\
 &Map(x \mapsto x * pair._1) \$ pair._0 \\
 &)\$ Zip(Transpose() \$ rowsA, \overrightarrow{colB}) \\
 &)\circ Transpose() \$ \mathbf{B} \\
 &)\circ Split(blockFactor) \$ \mathbf{A}
 \end{aligned}$$


$$\begin{aligned}
 &Join() \circ Map(rowsA \mapsto \\
 &Transpose() \circ Map(\overrightarrow{colB} \mapsto \\
 &Transpose() \circ Reduce((\overrightarrow{acc}, \overrightarrow{next}) \mapsto \\
 &Map(+) \$ Zip(\overrightarrow{acc}, \overrightarrow{next}) \\
 &)\circ Map(pair \mapsto \\
 &Map(x \mapsto x * pair._1) \$ pair._0 \\
 &)\$ Zip(Transpose() \$ rowsA, \overrightarrow{colB}) \\
 &)\circ Transpose() \$ \mathbf{B} \\
 &)\circ Split(blockFactor) \$ \mathbf{A}
 \end{aligned}$$

$$Transpose() \circ Transpose() \Rightarrow id$$

Register Blocking

$$\begin{aligned} & \text{Join}() \circ \text{Map}(\text{rows}A \mapsto \\ & \quad \text{Transpose}() \circ \text{Map}(\overrightarrow{\text{col}B} \mapsto \\ & \quad \quad \text{Transpose}() \circ \text{Reduce}((\overrightarrow{\text{acc}}, \overrightarrow{\text{next}}) \mapsto \\ & \quad \quad \quad \text{Map}(+) \$ \text{Zip}(\overrightarrow{\text{acc}}, \overrightarrow{\text{next}}) \\ & \quad \quad \quad) \circ \text{Map}(\text{pair} \mapsto \\ & \quad \quad \quad \quad \text{Map}(x \mapsto x * \text{pair}._1) \$ \text{pair}._0 \\ & \quad \quad \quad \quad) \$ \text{Zip}(\text{Transpose}() \$ \text{rows}A, \overrightarrow{\text{col}B}) \\ & \quad \quad \quad) \circ \text{Transpose}() \$ \mathbf{B} \\ & \quad \quad) \circ \text{Split}(\text{blockFactor}) \$ \mathbf{A} \end{aligned}$$
$$\begin{aligned} & \text{Reduce}(f) \circ \text{Map}(g) \Rightarrow \\ & \text{Reduce}((\text{acc}, x) \mapsto f(\text{acc}, g(x))) \end{aligned}$$

Register Blocking

$$\begin{aligned}
 &Join() \circ Map(rowsA \mapsto \\
 &Transpose() \circ Map(\overrightarrow{colB} \mapsto \\
 &Transpose() \circ \text{Reduce}((\overrightarrow{acc}, \overrightarrow{next}) \mapsto \\
 &Map(+) \$ Zip(\overrightarrow{acc}, \overrightarrow{next}) \\
 &)\circ Map(pair \mapsto \\
 &Map(x \mapsto x * pair._1) \$ pair._0 \\
 &)\$ Zip(Transpose() \$ rowsA, \overrightarrow{colB}) \\
 &)\circ Transpose() \$ \mathbf{B} \\
 &)\circ Split(blockFactor) \$ \mathbf{A}
 \end{aligned}$$


$$\begin{aligned}
 &Join() \circ Map(rowsA \mapsto \\
 &Transpose() \circ Map(\overrightarrow{colB} \mapsto \\
 &Transpose() \circ Reduce((\overrightarrow{acc}, \overrightarrow{pair}) \mapsto \\
 &Map(+) \$ Zip(\overrightarrow{acc}, \\
 &Map(x \mapsto x * pair._1) \$ pair._0) \\
 &)\$ Zip(Transpose() \$ rowsA, \overrightarrow{colB}) \\
 &)\circ Transpose() \$ \mathbf{B} \\
 &)\circ Split(blockFactor) \$ \mathbf{A}
 \end{aligned}$$

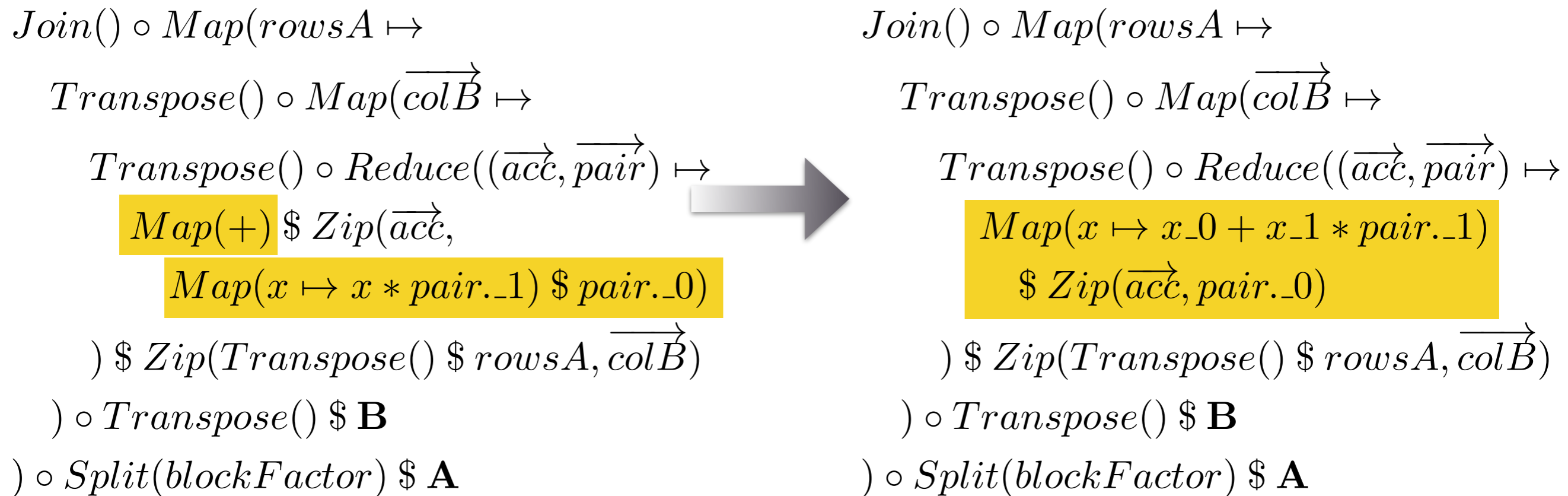
$$\begin{aligned}
 &Reduce(f) \circ Map(g) \Rightarrow \\
 &Reduce((acc, x) \mapsto f(acc, g(x)))
 \end{aligned}$$

Register Blocking

$Join() \circ Map(rowsA \mapsto$
 $Transpose() \circ Map(\overrightarrow{colB} \mapsto$
 $Transpose() \circ Reduce((\overrightarrow{acc}, \overrightarrow{pair}) \mapsto$
 $Map(+) \$ Zip(\overrightarrow{acc},$
 $Map(x \mapsto x * pair._1) \$ pair._0)$
 $) \$ Zip(Transpose() \$ rowsA, \overrightarrow{colB})$
 $) \circ Transpose() \$ \mathbf{B}$
 $) \circ Split(blockFactor) \$ \mathbf{A}$

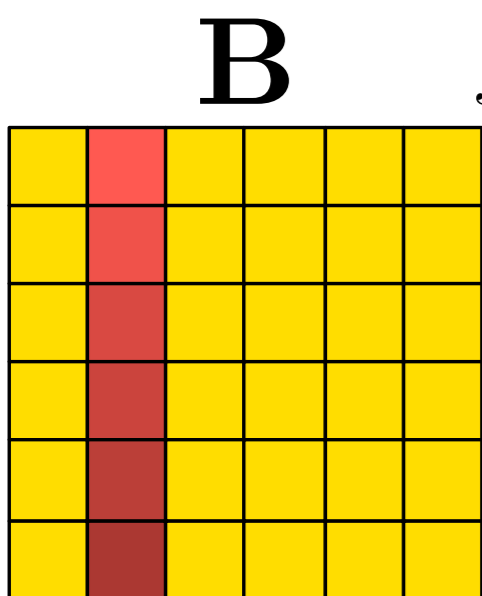
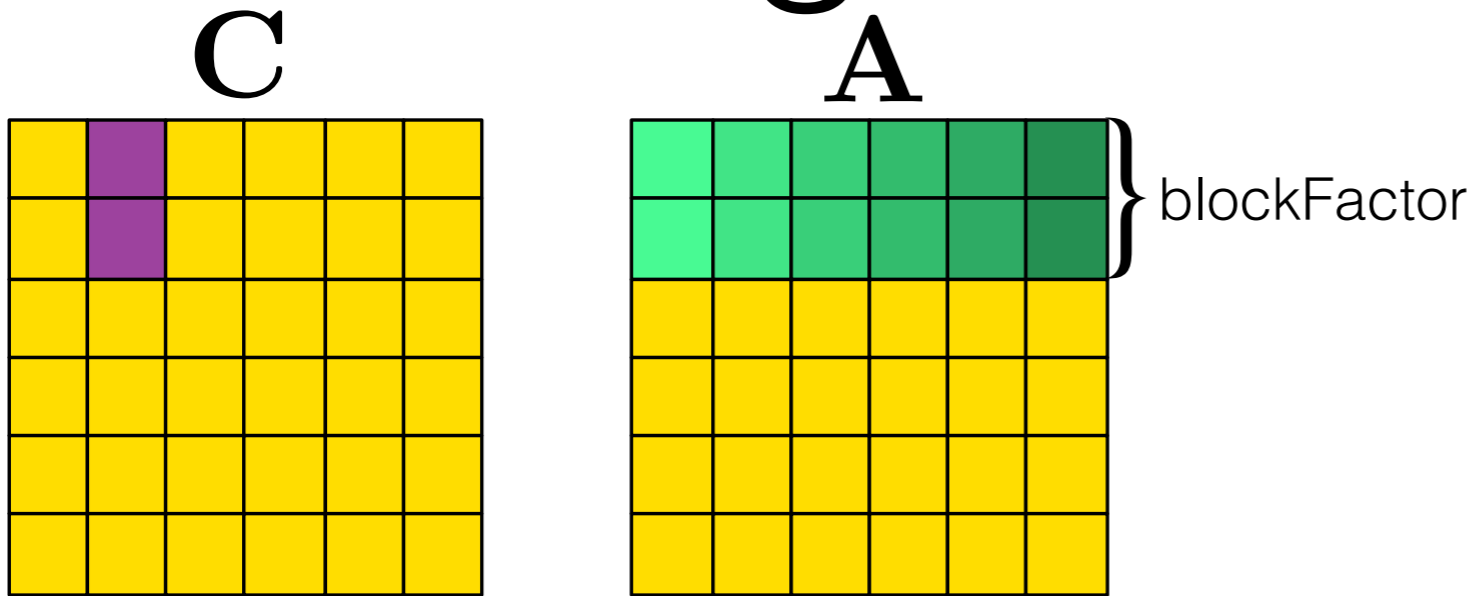
$$Map(f) \circ Map(g) \Rightarrow Map(f \circ g)$$

Register Blocking



$$Map(f) \circ Map(g) \Rightarrow Map(f \circ g)$$

Register Blocking



Functional Representation
 $Join() \circ Map(rowsA \mapsto$

$Transpose() \circ Map(\overrightarrow{colB} \mapsto$

$Transpose() \circ Reduce((\overrightarrow{acc}, \overrightarrow{pair}) \mapsto$

$Map(x \mapsto x_0 + x_1 * pair._1)$

$\$ Zip(\overrightarrow{acc}, pair._0)$

$) \$ Zip(Transpose() \$ rowsA, \overrightarrow{colB})$

$) \circ Transpose() \$ B$

$) \circ Split(blockFactor) \$ A$

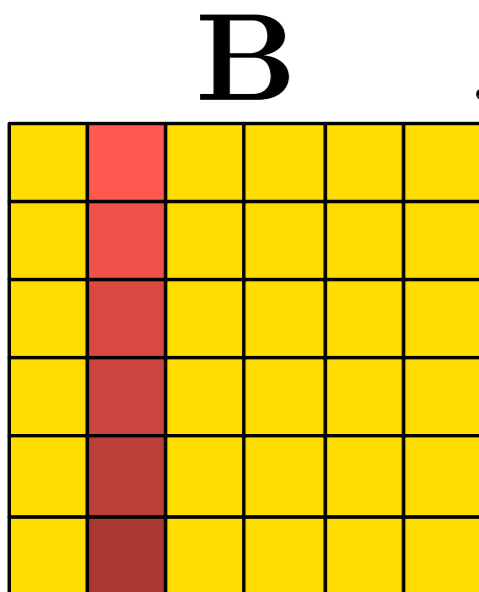
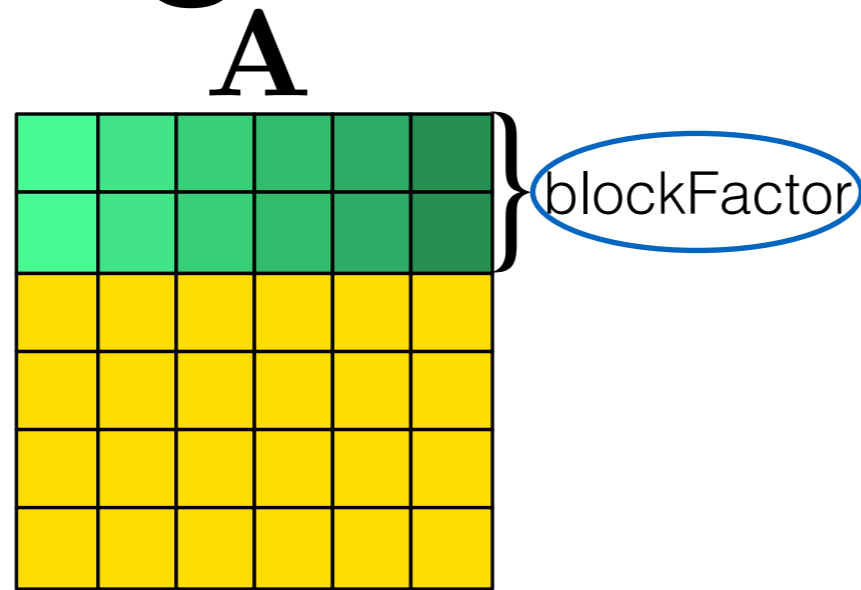
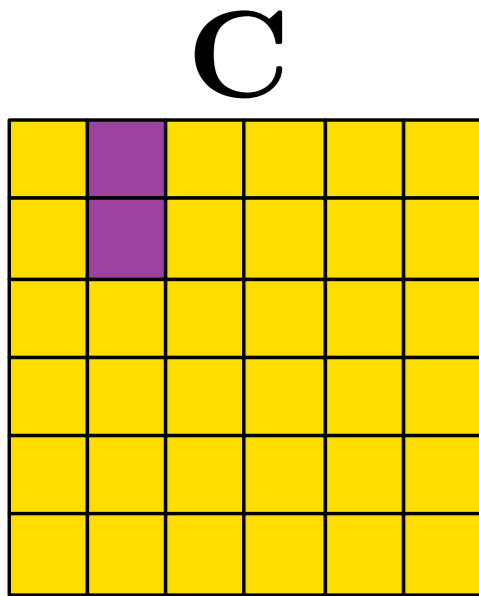
OpenCL

```

1 kernel void KERNEL(
2   const global float* restrict A,
3   const global float* restrict B,
4   global float* C, int K, int M, int N)
5 {
6   float acc[blockFactor];
7
8   for (int glb_id_1 = get_global_id(1);
9         glb_id_1 < M / blockFactor;
10        glb_id_1 += get_global_size(1)) {
11     for (int glb_id_0 = get_global_id(0); glb_id_0 < N;
12          glb_id_0 += get_global_size(0)) {
13
14         for (int i = 0; i < K; i += 1)
15             float temp = B[i * N + glb_id_0];
16         for (int j = 0; j < blockFactor; j += 1)
17             acc[j] +=
18                 A[blockFactor * glb_id_1 * K + j * K + i]
19                 * temp;
20
21         for (int j = 0; j < blockFactor; j += 1)
22             C[blockFactor * glb_id_1 * N + j * N + glb_id_0]
23               = acc[j];
24     }
25 }
26 }

```

Register Blocking



Functional Representation

$Join() \circ Map(rowsA \mapsto$

$Transpose() \circ Map(\overrightarrow{colB} \mapsto$

$Transpose() \circ Reduce((\overrightarrow{acc}, \overrightarrow{pair}) \mapsto$

$Map(x \mapsto x_0 + x_1 * pair._1)$

$\$ Zip(\overrightarrow{acc}, pair._0)$

$) \$ Zip(Transpose() \$ rowsA, \overrightarrow{colB})$

$) \circ Transpose() \$ B$

$) \circ Split(blockFactor) \$ A$

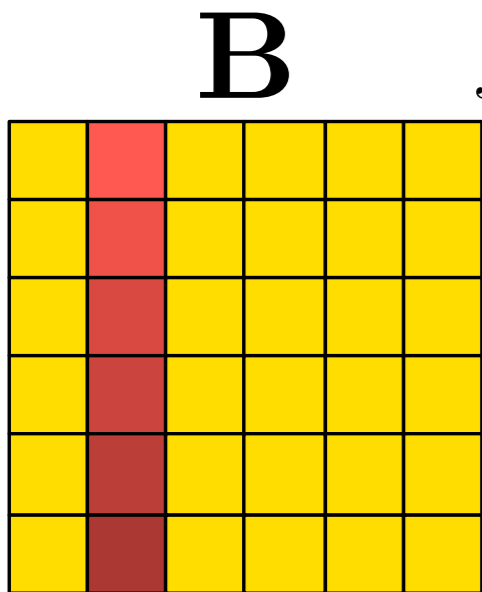
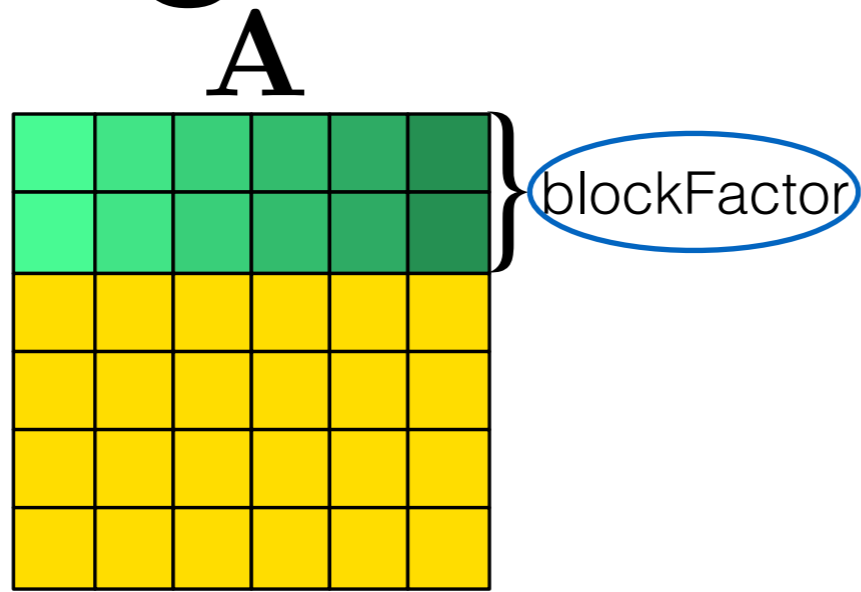
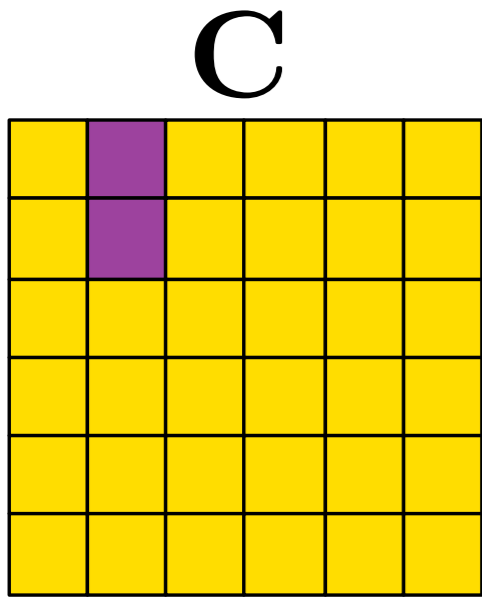
OpenCL

```

1 kernel void KERNEL(
2   const global float* restrict A,
3   const global float* restrict B,
4   global float* C, int K, int M, int N)
5 {
6   float acc[blockFactor];
7
8   for (int glb_id_1 = get_global_id(1);
9         glb_id_1 < M / blockFactor;
10        glb_id_1 += get_global_size(1)) {
11     for (int glb_id_0 = get_global_id(0); glb_id_0 < N;
12          glb_id_0 += get_global_size(0)) {
13
14         for (int i = 0; i < K; i += 1)
15             float temp = B[i * N + glb_id_0];
16         for (int j = 0; j < blockFactor; j += 1)
17             acc[j] +=
18                 A[blockFactor * glb_id_1 * K + j * K + i]
19                 * temp;
20
21         for (int j = 0; j < blockFactor; j += 1)
22             C[blockFactor * glb_id_1 * N + j * N + glb_id_0]
23               = acc[j];
24     }
25 }
26 }

```

Register Blocking



Functional Representation

$Join() \circ Map(rowsA \mapsto$

$Transpose() \circ Map(\overrightarrow{colB} \mapsto$

$Transpose() \circ Reduce((\overrightarrow{acc}, \overrightarrow{pair}) \mapsto$

$Map(x \mapsto x_0 + x_1 * pair._1)$

$\$ Zip(\overrightarrow{acc}, pair._0)$

$) \$ Zip(Transpose() \$ rowsA, \overrightarrow{colB})$

$) \circ Transpose() \$ B$

$) \circ Split(blockFactor) \$ A$

OpenCL

```

1 kernel void KERNEL(
2   const global float* restrict A,
3   const global float* restrict B,
4   global float* C, int K, int M, int N)
5 {
6   float acc[blockFactor];
7
8   for (int glb_id_1 = get_global_id(1);
9         glb_id_1 < M / blockFactor;
10        glb_id_1 += get_global_size(1)) {
11     for (int glb_id_0 = get_global_id(0); glb_id_0 < N;
12          glb_id_0 += get_global_size(0)) {
13
14         for (int i = 0; i < K; i += 1)
15             float temp = B[i * N + glb_id_0];
16             for (int j = 0; j < blockFactor; j += 1)
17                 acc[j] +=
18                     A[blockFactor * glb_id_1 * K + j * K + i]
19                     * temp;
20
21         for (int j = 0; j < blockFactor; j += 1)
22             C[blockFactor * glb_id_1 * N + j * N + glb_id_0]
23               = acc[j];
24     }
25 }
26 }
    
```

Combining Optimisations

$\mathbf{A} * \mathbf{B} =$
 $Map(\overrightarrow{rowA} \mapsto$
 $Map(\overrightarrow{colB} \mapsto$
 $DotProduct(\overrightarrow{rowA}, \overrightarrow{colB}))$
 $) \circ Transpose() \$ \mathbf{B}$
 $) \$ \mathbf{A}$


 80 rewrites

$(p239, p36 \mapsto$
 $Join() \circ Map((p179 \mapsto$
 $Transpose() \circ Join() \circ Map((p70 \mapsto$
 $Transpose() \circ Join() \circ Map((p20 \mapsto$
 $Transpose() \circ Map((p65 \mapsto$
 $Transpose()(p65)$
 $)) \circ Transpose()(p20)$
 $)) \circ Transpose() \circ Reduce((p75, p0 \mapsto$
 $Map((p164 \mapsto$
 $Join() \circ Map((p81 \mapsto$
 $Reduce((p136, p90 \mapsto$
 $Map((p163 \mapsto$
 $Get(0)(p163) + Get(1)(p163) * Get(1)(p90)$
 $)) \circ Zip(2)(p136, Get(0)(p90))$
 $))(Get(0)(p81), Zip(2)(Transpose() \circ Get(1)(p164), Get(1)(p81)))$
 $)) \circ Zip(2)(Get(0)(p164), Get(1)(p0))$
 $)) \circ Zip(2)(p75, Split(blockFactor) \circ Transpose() \circ Get(0)(p0))$
 $))(Zip(2)(Split(sizeK) \circ Transpose()(p179), p70))$
 $)) \circ Transpose() \circ Map((p4 \mapsto$
 $Split(sizeN) \circ Transpose()(p4)$
 $)) \circ Split(sizeK)(p36)$
 $)) \circ Split(sizeM)(p239)$
 $)$

How do we apply
these optimisations?

Exploration Strategy

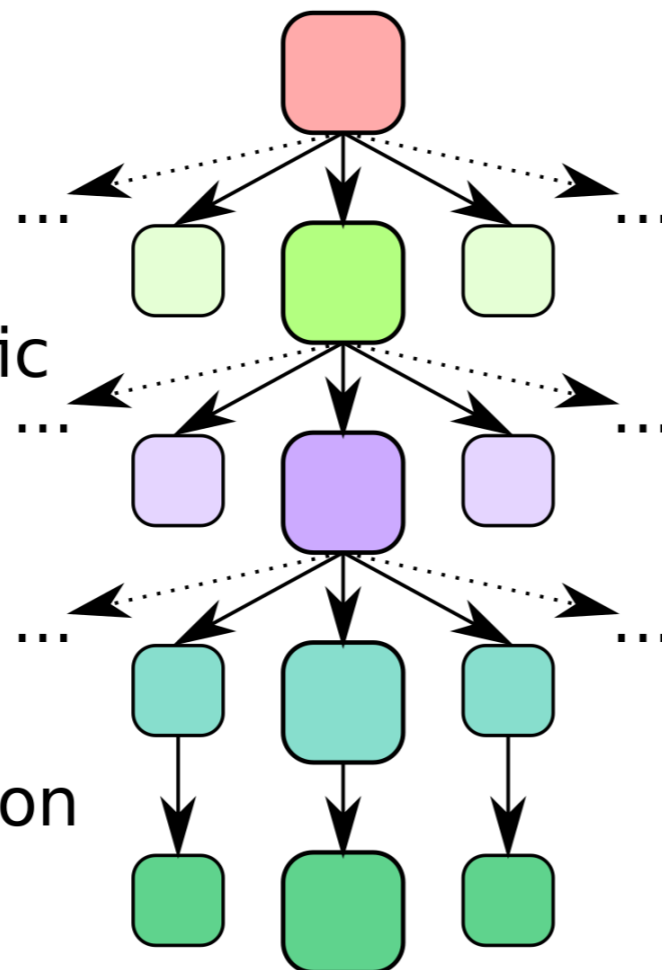
Phases:

Algorithmic
Exploration

OpenCL specific
Exploration

Parameter
Exploration

Code Generation



Program Variants:

High-Level Program 1

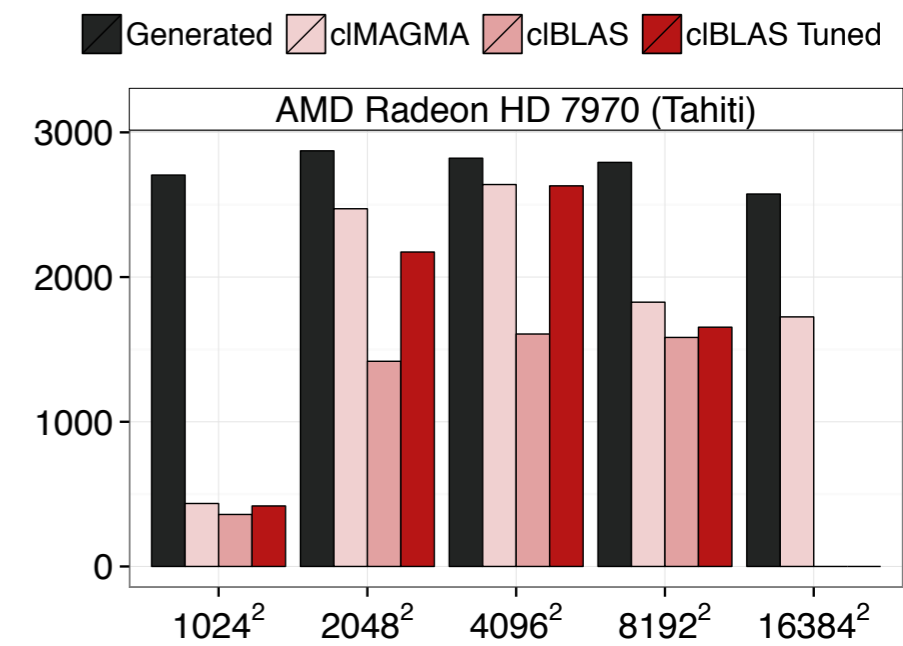
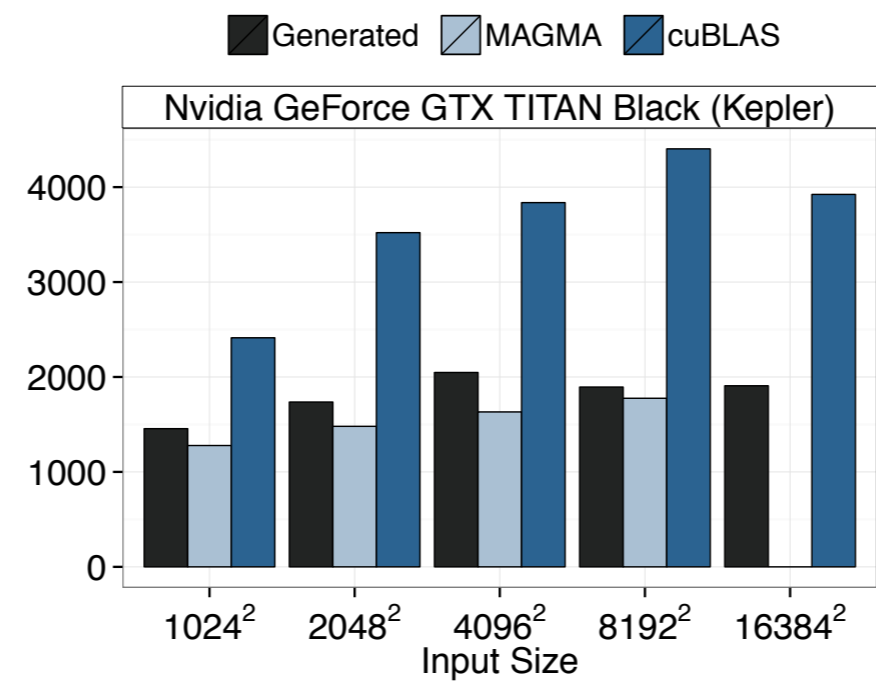
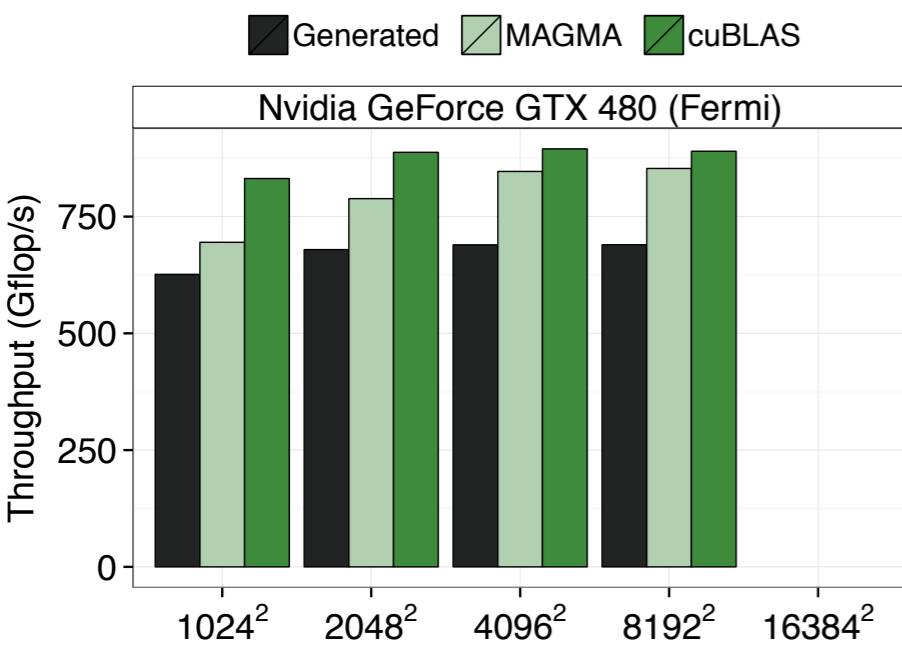
Algorithmic
Rewritten Program 8

OpenCL Specific
Program 760

Fully Specialized
Program 46,000

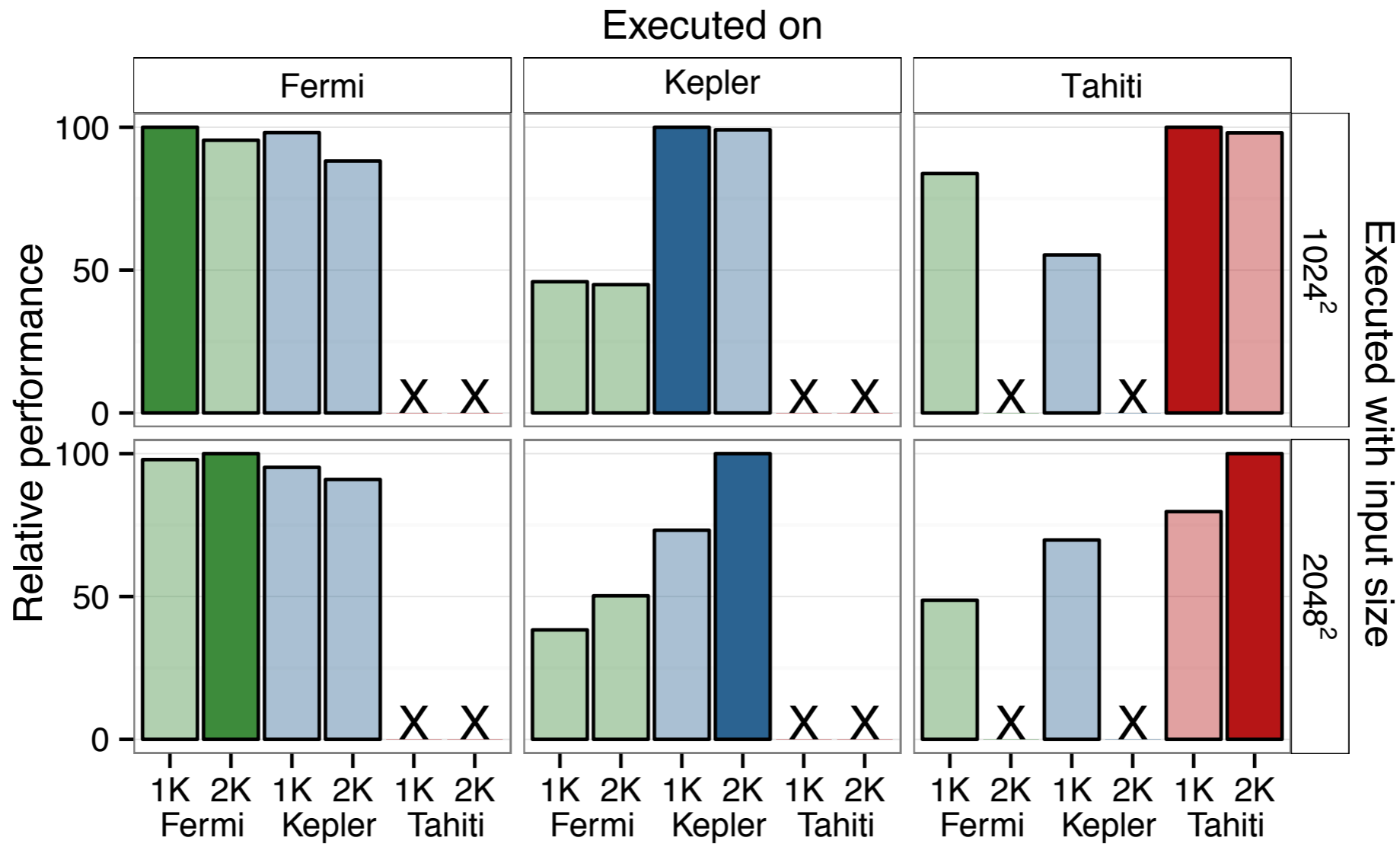
OpenCL Code 46,000

Performance Results



Performance close or better than hand-tuned MAGMA library

Performance Portability



The six specialized OpenCL kernels

Performance is not portable
across architectures and input sizes

Conclusion

- OpenCL code is not performance portable
- Using a functional approach along with rewrite rules we can achieve performance portability
- Performance of matrix multiplication on par with tuned OpenCL code

Toomas Remmelg - toomas.remmelg@ed.ac.uk

Supported by:

  Labs

