# A Functional Approach to Performance Portable GPU Code Generation

## A Case Study on Matrix Multiplication
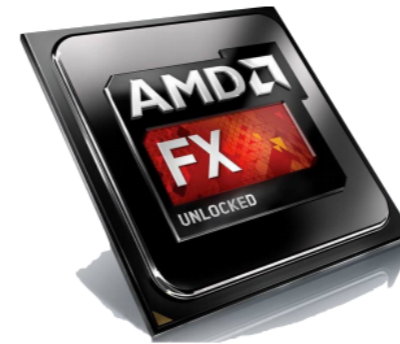
**Toomas Remmelg**, Thibaut Lutz, Michel Steuwer, Christophe Dubach

THE UNIVERSITY *of* EDINBURGH

# The Problem

- Parallel processors everywhere

- Many different types:
  CPUs, GPUs, …

- Parallel programming is hard

- Optimising even harder

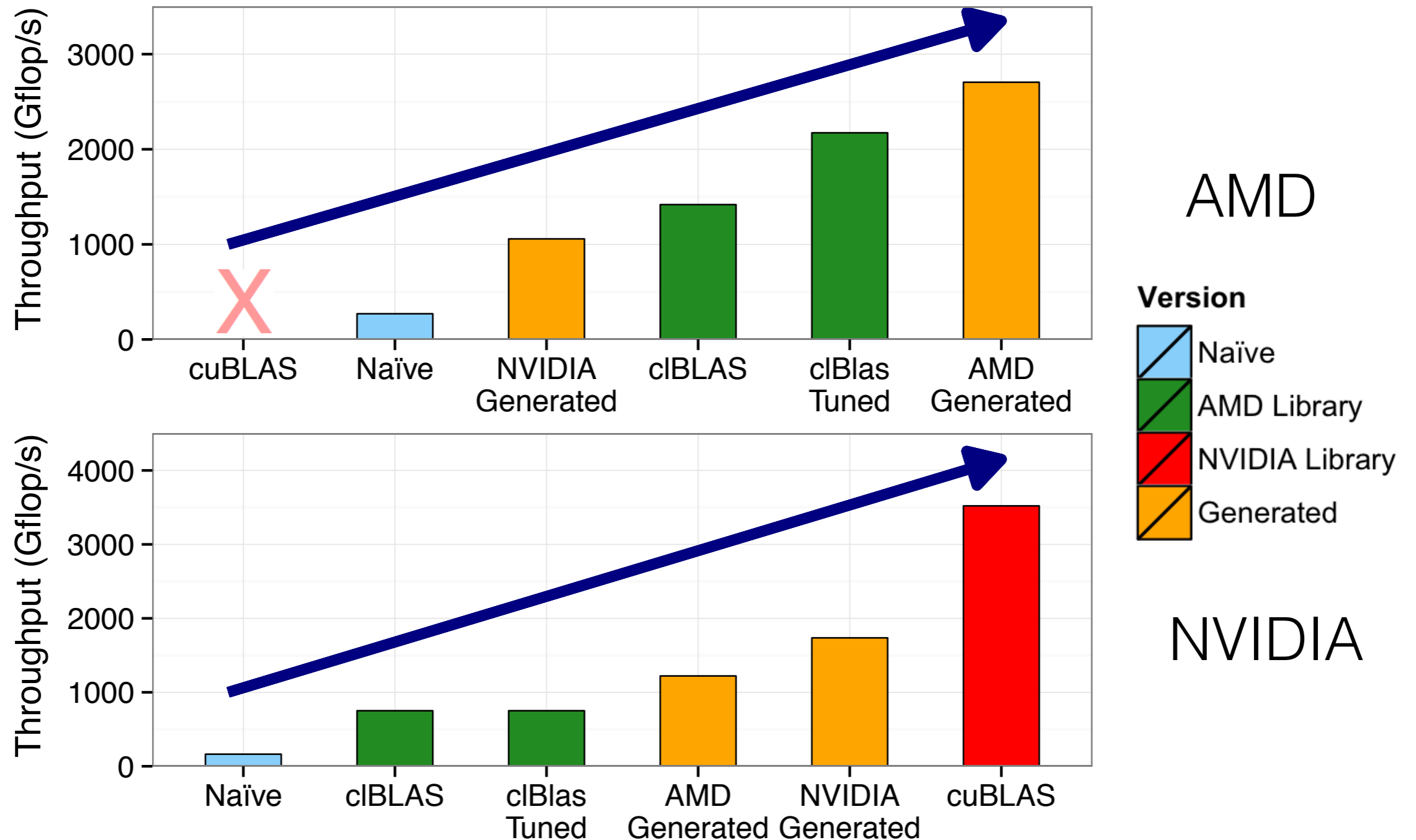- **Problem:**
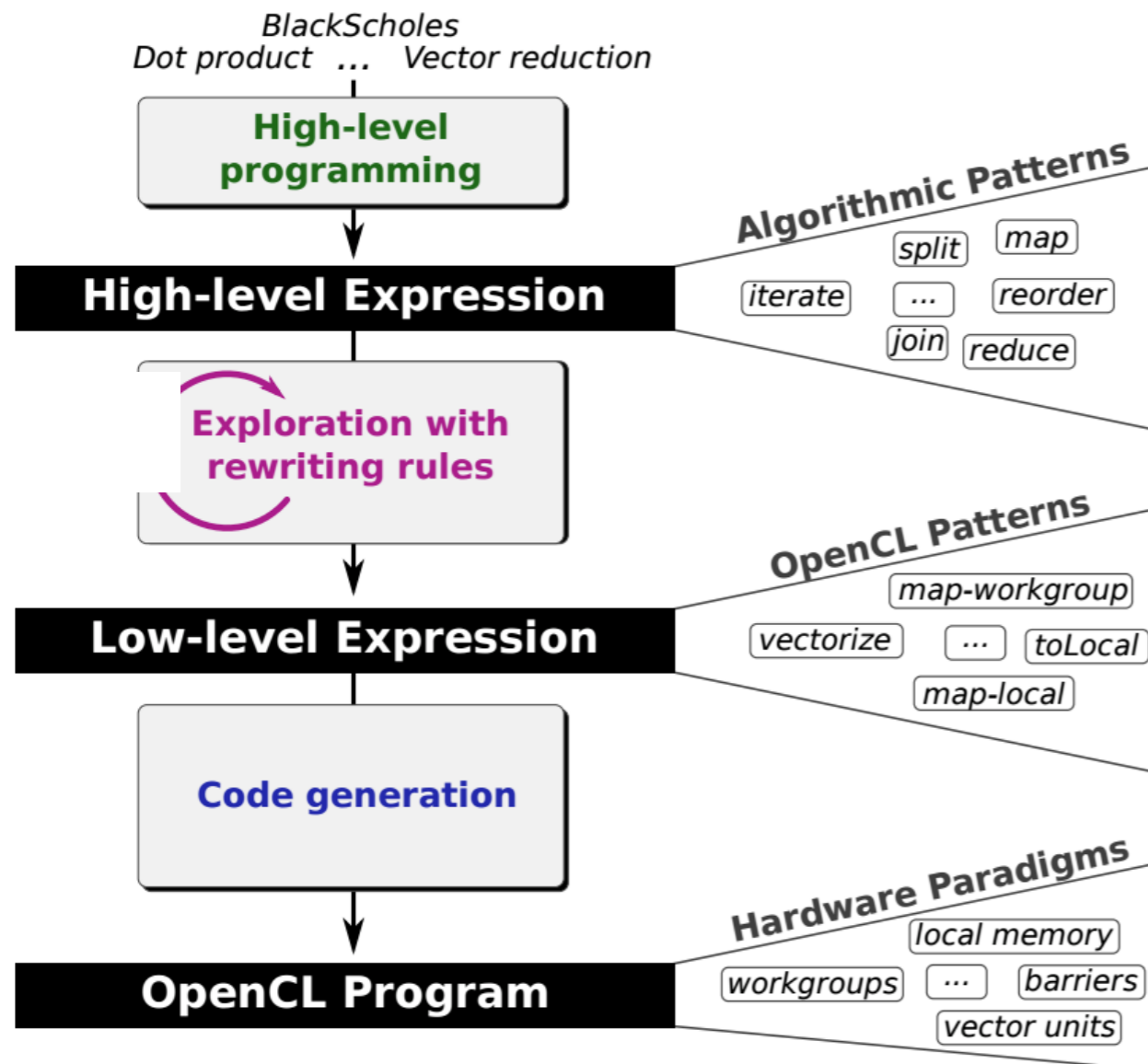  No portability of performance!



CPU

GPU

Accelerator

FPGA

THE UNIVERSITY *of* EDINBURGH
**informatics**

# Performance Portability of Matrix Multiplication

# Generating Performance Portable Code Using Rewrite Rules



BlackScholes
Dot product ... Vector reduction

**High-level programming**

**High-level Expression**

Algorithmic Patterns: split, map, iterate, ..., reorder, join, reduce

**Exploration with rewriting rules**

**Low-level Expression**

OpenCL Patterns: map-workgroup, vectorize, ..., toLocal, map-local

**Code generation**

Hardware Paradigms: local memory, workgroups, ..., barriers, vector units

**OpenCL Program**

$$\lambda\, xs \,.\, map\,(\lambda\, x \,.\, x * 3)\, xs$$

**(a) High-level expression** written by the programmer.

rewrite rules

$$\lambda\, xs \,.\, (join \circ mapWorkgroup\,(joinVec \circ$$
$$mapLocal\,(mapVec\,(\lambda\, x \,.\, x * 3))$$
$$\circ splitVec\, 4) \circ split\, 1024)\, xs$$

**(b) Low-level expression** derived using rewrite rules and search.

code generator

```
1   int4 mul3(int4 x) { return x * 3; }
2   kernel vectorScal(global int* in,out, int len){
3    for (int i=get_group_id; i < len/1024;
4        i+=get_num_groups) {
5     global int* grp_in  = in+(i*1024);
6     global int* grp_out = out+(i*1024);
7     for (int j=get_local_id; j < 1024/4;
8         j+=get_local_size) {
9      global int4* in_vec4 =(int4*)grp_in+(j*4);
10     global int4* out_vec4=(int4*)grp_out+(j*4);
11     *out_vec4 = mul3(*in_vec4);
12   } } }
```

**(c) OpenCL program** produced by our code generator.

THE UNIVERSITY of EDINBURGH
**informatics**

4

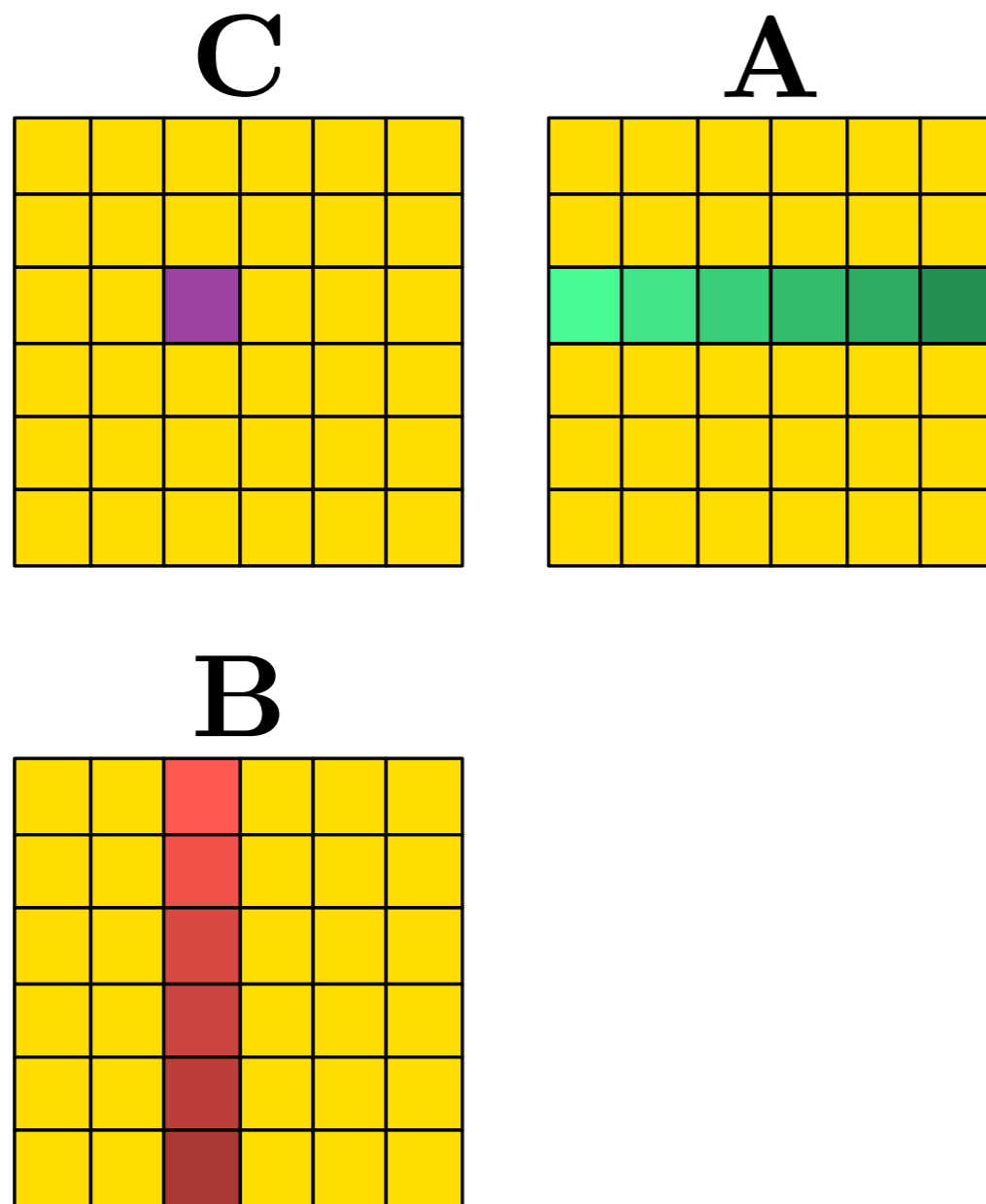# Matrix Multiplication in OpenCL

```
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 }
```

Naïve

```
1  kernel mm_amd_opt(global float * A, B, C,
2              int K, M, N) {
3  local float tileA [512]; tileB [512];
4
5  private float acc_0;        ...; acc_31;
6  private float blockOfB_0; ...; blockOfB_3;
7  private float blockOfA_0; ...; blockOfA_7;
8
9  int lid0 = local_id(0); lid1 = local_id(1);
10 int wid0 = group_id(0); wid1 = group_id(1);
11
12 for (int w1=wid1; w1<M/64; w1+=num_grps(1)) {
13   for (int w0=wid0; w0<N/64; w0+=num_grps(0)) {
14
15     acc_0 = 0.0f; ...; acc_31 = 0.0f;
16     for (int i=0; i<K/8; i++) {
17      vstore4(vload4(lid1*M/4+2*i*M+16*w1+lid0,A), 16*lid1+lid0, tileA);
18      vstore4(vload4(lid1*N/4+2*i*N+16*w0+lid0,B), 16*lid1+lid0, tileB);
19      barrier (...) ;
20
21      for (int j = 0; j<8; j++) {
22       blockOfA_0 = tileA[0+64*j+lid1*8];    ...; blockOfA_7 = tileA[7+64*j+lid1*8];
23       blockOfB_0 = tileB[0 +64*j+lid0];      ...; blockOfB_3 = tileB[48+64*j+lid0];
24
25       acc_0  += blockOfA_0 * blockOfB_0; ...; acc_28 += blockOfA_7 * blockOfB_0;
26       acc_1  += blockOfA_0 * blockOfB_1; ...; acc_29 += blockOfA_7 * blockOfB_1;
27       acc_2  += blockOfA_0 * blockOfB_2; ...; acc_30 += blockOfA_7 * blockOfB_2;
28       acc_3  += blockOfA_0 * blockOfB_3; ...; acc_31 += blockOfA_7 * blockOfB_3;
29      }
30      barrier (...) ;
31     }
32
33     C[ 0+8*lid1*N+64*w0+64*w1*N+0*N+lid0]=acc_0; ...; C[ 0+8*lid1*N+64*w0+64*w1*N+7*N+lid0]=acc_28;
34     C[16+8*lid1*N+64*w0+64*w1*N+0*N+lid0]=acc_1; ...; C[16+8*lid1*N+64*w0+64*w1*N+7*N+lid0]=acc_29;
35     C[32+8*lid1*N+64*w0+64*w1*N+0*N+lid0]=acc_2; ...; C[32+8*lid1*N+64*w0+64*w1*N+7*N+lid0]=acc_30;
36     C[48+8*lid1*N+64*w0+64*w1*N+0*N+lid0]=acc_3; ...; C[48+8*lid1*N+64*w0+64*w1*N+7*N+lid0]=acc_31;
37 } } }
```

Generated for AMD

THE UNIVERSITY of EDINBURGH
**informatics**

# Functional Definition of Matrix Multiplication

**C**       **A**

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

$$Map(\overrightarrow{rowA} \mapsto$$

$$Map(\overrightarrow{colB} \mapsto$$

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

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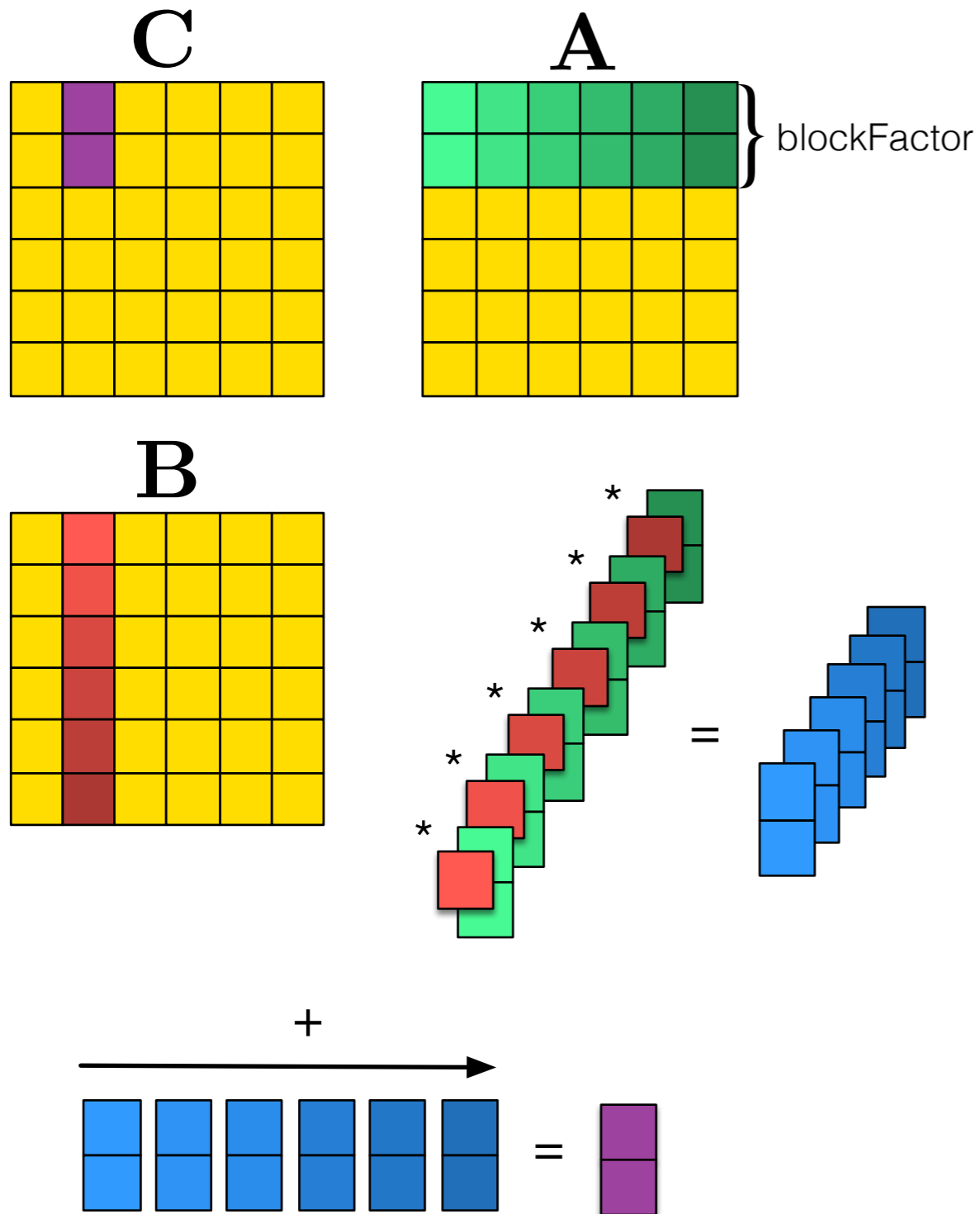
$$) \, \$ \, \mathbf{A}$$

**B**

# Traditional Optimisations

- **Register Blocking** Loading elements into registers and reusing them.

- **Tiling** Solving the problem by diving 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.

THE UNIVERSITY *of* EDINBURGH
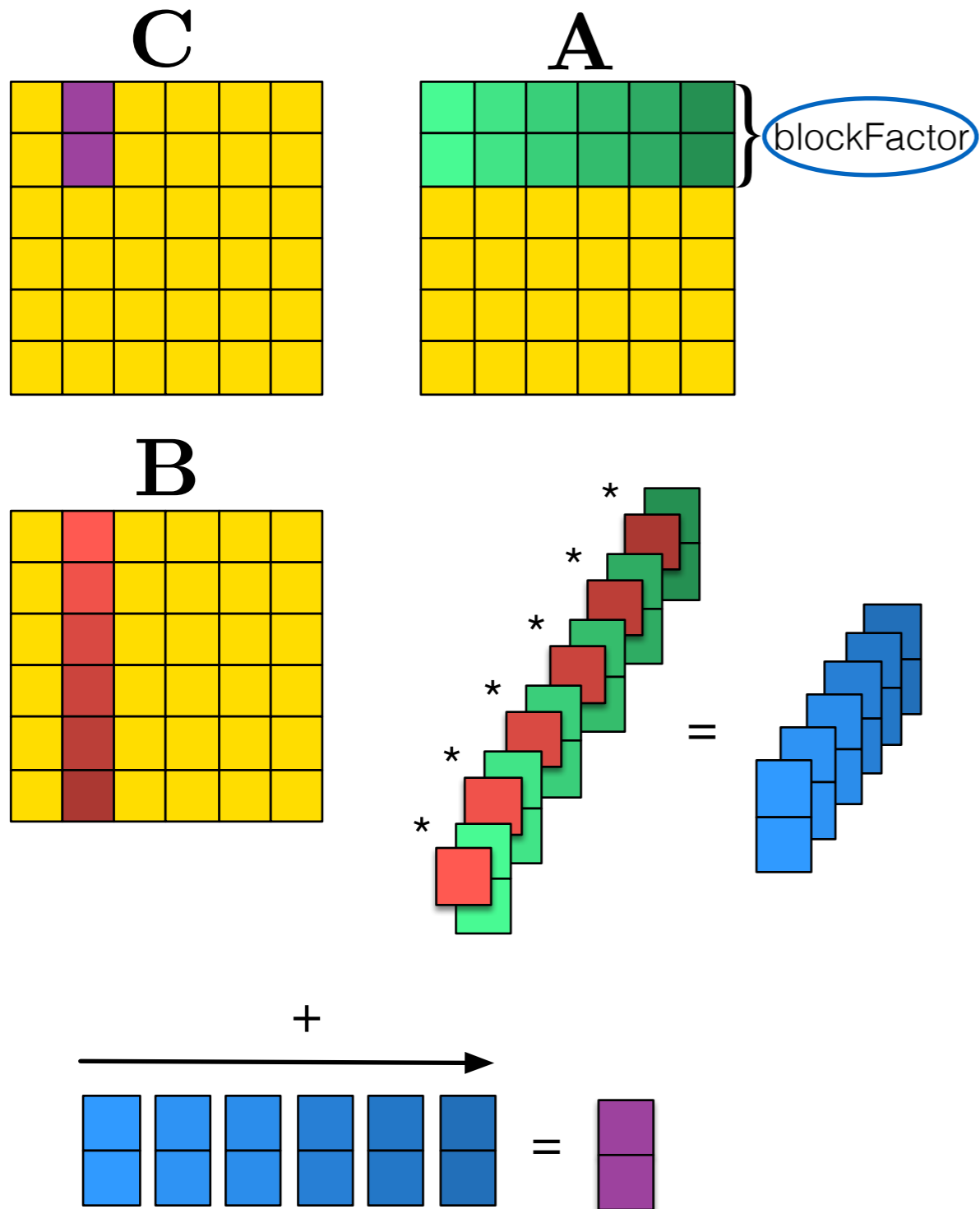**informatics**
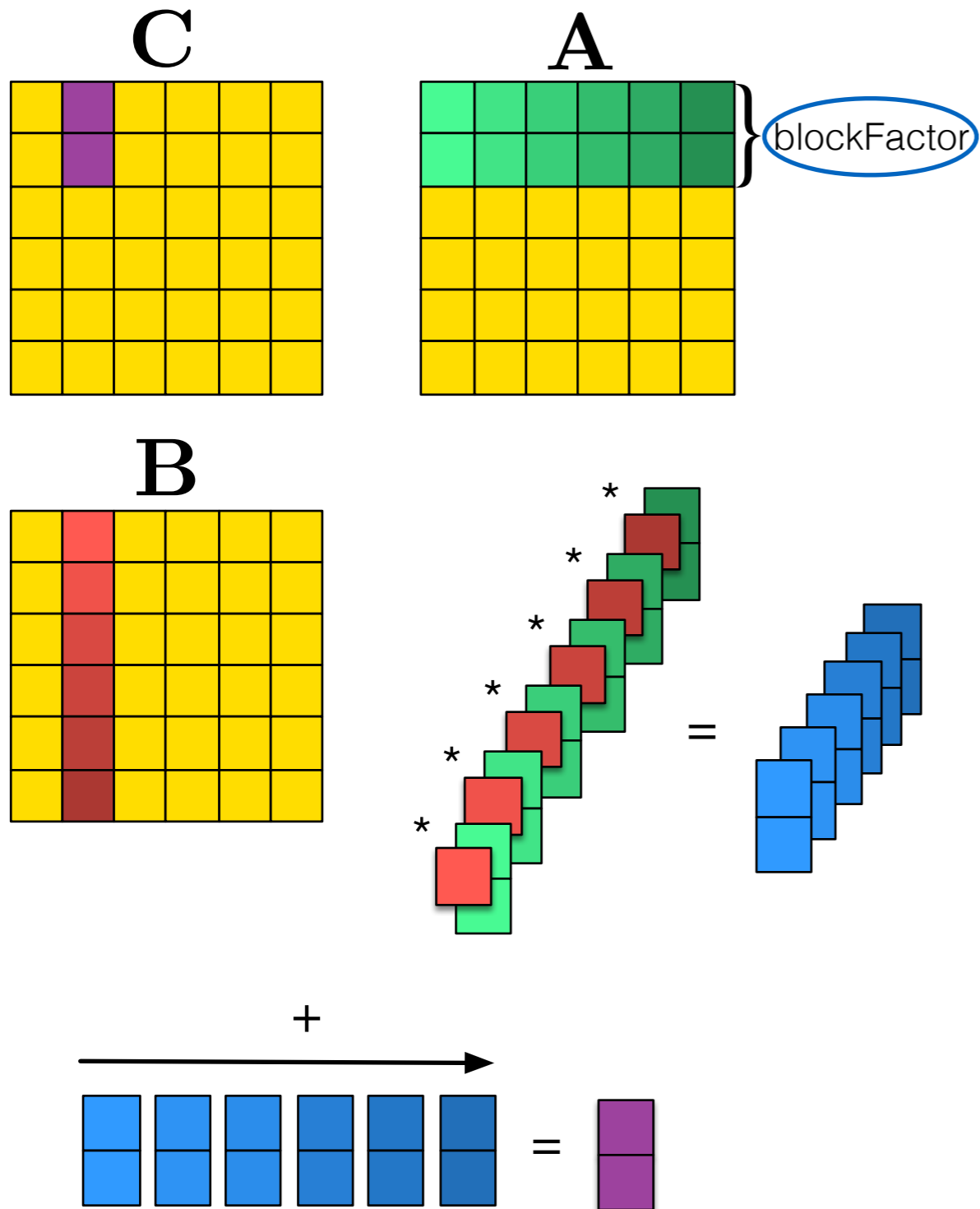
# Register Blocking



```
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  }
```

8

# Register Blocking
# as a Series of Rewrites

$registerBlocking =$

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

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

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

$\quad Map(Reduce(f)) \Rightarrow Transpose() \circ Reduce((acc, x) \mapsto Map(f) \circ Zip(acc, x))$

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

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

$\quad Reduce(f) \circ Map(g) \Rightarrow Reduce((acc, x) \mapsto f(acc, g(x)))$

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

# Register Blocking

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

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

# Register Blocking

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

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

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

# 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

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

$Map(\overrightarrow{rowA} \mapsto$

$Map(\overrightarrow{colB} \mapsto$

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

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

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

$) \, \$ \, rowsA$

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

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

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

$Map(\overrightarrow{rowA} \mapsto$

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

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

$) \, \$ \, rowsA$

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

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

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

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

# Register Blocking

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

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

$$Map(\overrightarrow{rowA} \mapsto$$

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

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

$$)\ \$\ rowsA$$

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

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

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

# Register Blocking

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

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

$\quad\quad \boxed{Map(\overrightarrow{rowA} \mapsto}$

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

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

$\quad\quad ) \, \$ \, rowsA$

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

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

$\Rightarrow$

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

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

$\quad\quad \boxed{Map(}$

$\quad\quad\quad Reduce(+)$

$\quad\quad \boxed{) \circ Map(\overrightarrow{rowA} \mapsto}$

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

$\quad\quad ) \, \$ \, rowsA$

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

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

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

# 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

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

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

$\quad\quad Map($

$\quad\quad\quad Reduce(+)$

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

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

$\quad\quad ) \, \$ \, rowsA$

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

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

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

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

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

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

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

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

$\quad\quad ) \, \$ \, rowsA$

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

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

$$Map(Reduce(f)) \Rightarrow$$

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

# 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

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

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

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

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

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

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

    $) \$ rowsA$

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

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

$\Longrightarrow$

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

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

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

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

    $) \circ Transpose()$

    $\circ Transpose() \circ Map(pair \mapsto$

      $Map(x \mapsto x * pair.\_1) \$ pair.\_0$

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

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

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

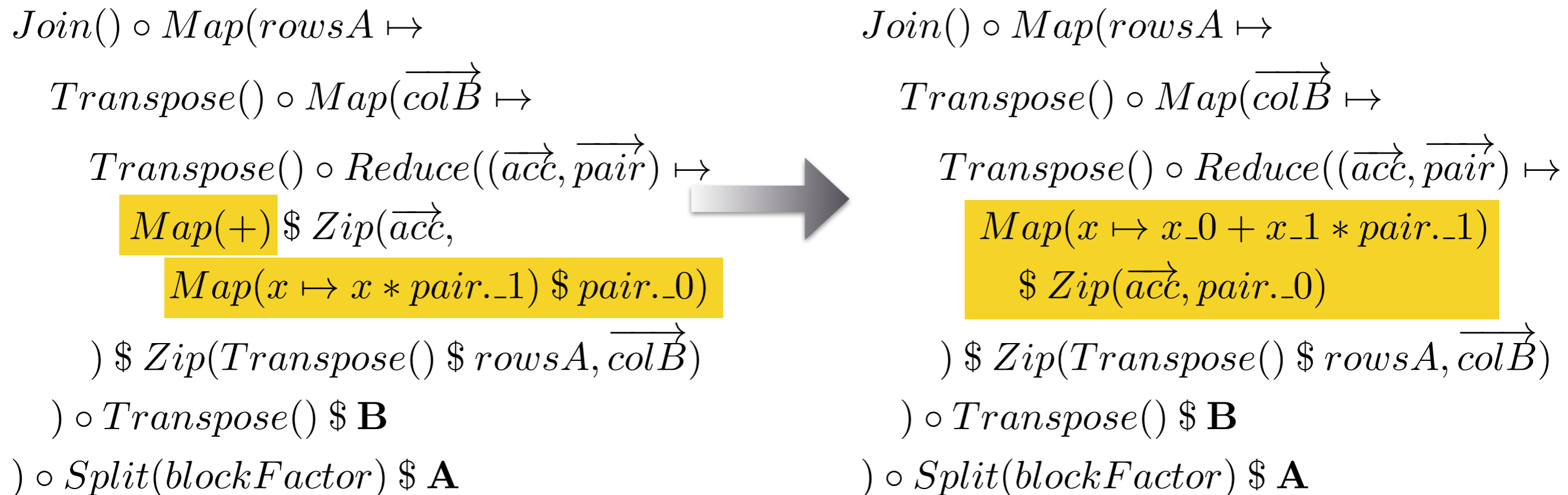$Map(Map(f)) \Rightarrow$

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

14

# Register Blocking

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

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

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

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

$$) \circ Transpose()$$

$$\circ Transpose() \circ Map(pair \mapsto$$

$$Map(x \mapsto x * pair._1) \$ pair._0$$

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

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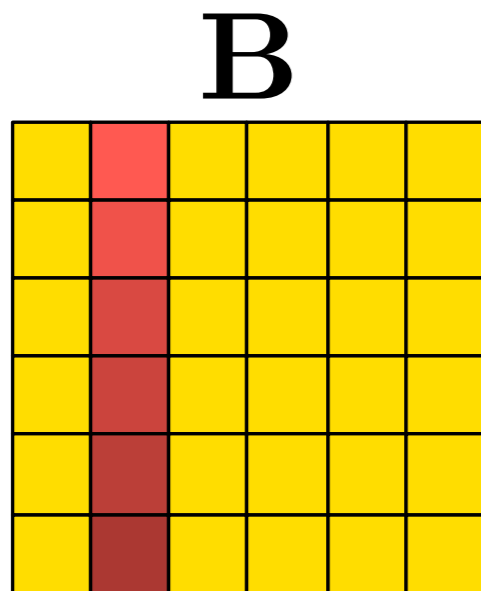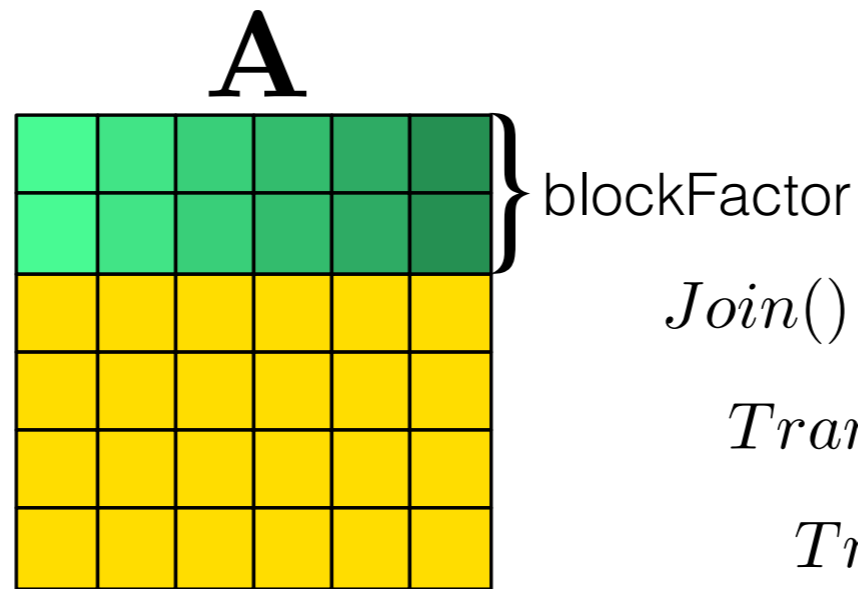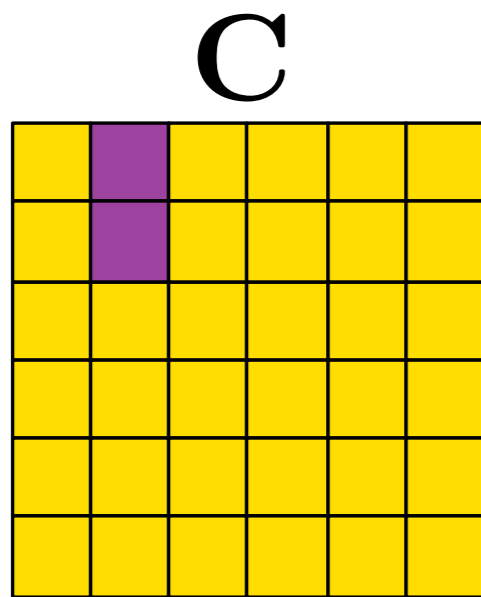
$$) \circ Split(blockFactor) \$ \mathbf{A}$$

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

# Register Blocking

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

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

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

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

$\quad\quad \boxed{) \circ Transpose()}$

$\quad\quad \boxed{\circ Transpose()} \circ Map(pair \mapsto$

$\quad\quad\quad Map(x \mapsto x * pair.\_1) \$ pair.\_0$

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

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

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

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

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

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

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

$\quad\quad ) \boxed{\circ} Map(pair \mapsto$

$\quad\quad\quad Map(x \mapsto x * pair.\_1) \$ pair.\_0$

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

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

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

$$Transpose() \circ Transpose() \Rightarrow id$$
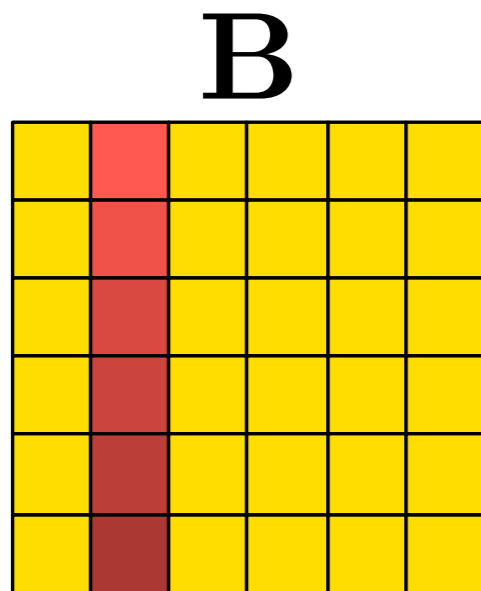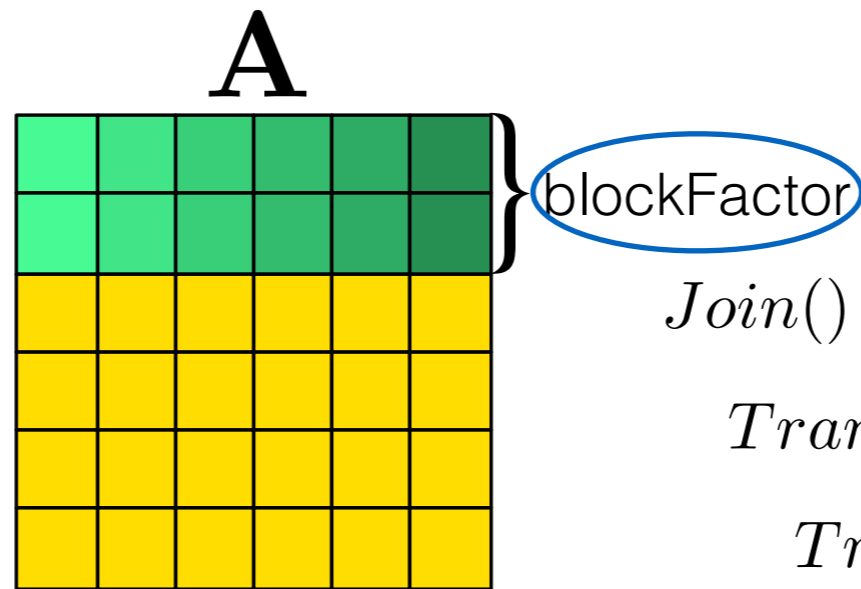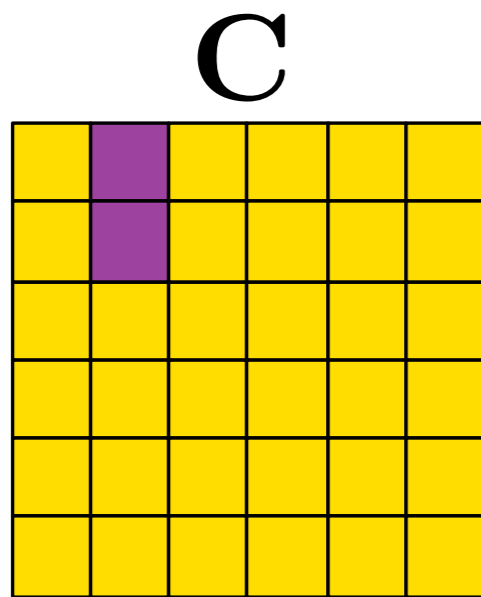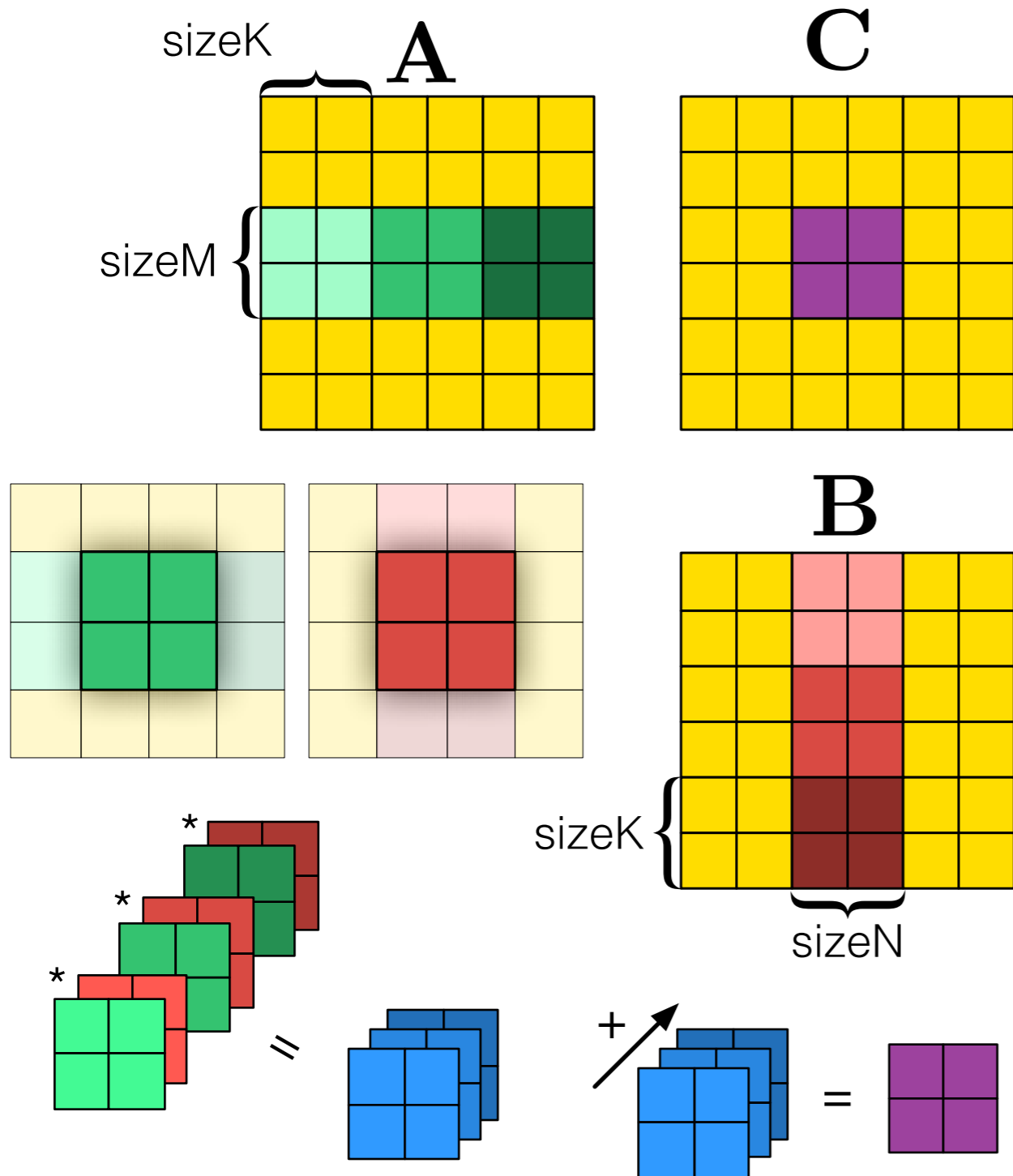
# Register Blocking

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

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

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

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

$$) \circ Map(pair \mapsto$$

$$Map(x \mapsto x * pair._1) \,\$\, pair._0$$

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

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

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

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

$$Reduce((acc, x) \mapsto f(acc, g(x))$$

# Register Blocking

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

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

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

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

$\quad\quad )$ ==$\circ Map$==$(pair \mapsto$

$\quad\quad\quad Map(x \mapsto x * pair.\_1) \ \$ \ pair.\_0$

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

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

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

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

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

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

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

$\quad\quad\quad\quad$ ==$Map(x \mapsto x * pair.\_1) \ \$ \ pair.\_0)$==

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

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

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

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

$$Reduce((acc, x) \mapsto f(acc, g(x))$$

# Register Blocking

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

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

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

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

$$Map(x \mapsto x * pair.\_1) \mathbin{\$} pair.\_0)$$

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

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

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

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

# Register Blocking

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

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

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

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

        $Map(x \mapsto x * pair.\_1) \ \$ \ pair.\_0)$

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

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

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

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

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

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

      $Map(x \mapsto x\_0 + x\_1 * pair.\_1)$

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

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

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

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

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

# Register Blocking



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

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

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

$$Map(x \mapsto x\_0 + x\_1 * pair.\_1)$$

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

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

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

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

18

# Register Blocking



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

# Register Blocking

**C**

**A**

blockFactor

**B**

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

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

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

$$Map(x \mapsto x\_0 + x\_1 * pair.\_1)$$

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

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

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

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

# Tiling



```
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      local float  A_lcl[sizeM*sizeK];
8      local float  B_lcl[sizeK*sizeN];
9      float acc = 0.0f;
10
11     for (int i = 0; i < K / sizeK; i += 1) {
12       A_lcl[l_idA( lcl_id_0 ,  lcl_id_1 )] =
13         A[idA(i,  lcl_id_0 ,  lcl_id_1 ,  grp_id_0)];
14       B_lcl[l_idB( lcl_id_0 ,  lcl_id_1 )] =
15         B[idB(grp_id_1,  lcl_id_0 ,  lcl_id_1 ,  i)];
16
17       barrier(CLK_LOCAL_MEM_FENCE);
18
19       for (int j = 0; j < sizeK; j += 1)
20         acc += A_lcl[l_idA(j,  lcl_id_1 )]
21              * B_lcl[l_idB( lcl_id_0 ,  j)];
22
23       barrier(CLK_LOCAL_MEM_FENCE);
24     }
25     C[idC(grp_id_0,  lcl_id_1 ,  grp_id_1,  lcl_id_0 )] = acc;
26   }
```

19

# Tiling



```
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     local float  A_lcl[sizeM*sizeK];
8     local float  B_lcl[sizeK*sizeN];
9     float acc = 0.0f;
10
11    for (int i = 0; i < K / sizeK; i += 1) {
12      A_lcl[l_idA( lcl_id_0 ,  lcl_id_1 )] =
13        A[idA(i,  lcl_id_0 ,  lcl_id_1 ,  grp_id_0)];
14      B_lcl[l_idB( lcl_id_0 ,  lcl_id_1 )] =
15        B[idB(grp_id_1,  lcl_id_0 ,  lcl_id_1 ,  i)];
16
17      barrier(CLK_LOCAL_MEM_FENCE);
18
19      for (int j = 0; j < sizeK; j += 1)
20        acc += A_lcl[l_idA(j,  lcl_id_1 )]
21             * B_lcl[l_idB( lcl_id_0 ,  j)];
22
23      barrier(CLK_LOCAL_MEM_FENCE);
24    }
25    C[idC(grp_id_0,  lcl_id_1 ,  grp_id_1,  lcl_id_0 )] = acc;
26  }
```

19

# Tiling



```
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     local float A_lcl[sizeM*sizeK];
8     local float B_lcl[sizeK*sizeN];
9     float acc = 0.0f;
10
11    for (int i = 0; i < K / sizeK; i += 1) {
12      A_lcl[l_idA( lcl_id_0 ,  lcl_id_1 )] =
13        A[idA(i,  lcl_id_0 ,  lcl_id_1 ,  grp_id_0)];
14      B_lcl[l_idB( lcl_id_0 ,  lcl_id_1 )] =
15        B[idB(grp_id_1,  lcl_id_0 ,  lcl_id_1 ,  i)];
16
17      barrier(CLK_LOCAL_MEM_FENCE);
18
19      for (int j = 0; j < sizeK; j += 1)
20        acc += A_lcl[l_idA(j,  lcl_id_1 )]
21             * B_lcl[l_idB( lcl_id_0 ,  j)];
22
23      barrier(CLK_LOCAL_MEM_FENCE);
24    }
25    C[idC(grp_id_0,  lcl_id_1 ,  grp_id_1,  lcl_id_0 )] = acc;
26  }
```

19

# Tiling

# Tiling

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

Rewriting

$$TiledMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Untile() \circ$$
$$Map(\overrightarrow{aRows} \mapsto$$
$$Map(\overrightarrow{bCols} \mapsto$$
$$Reduce((\mathbf{acc}, pairOfTiles) \mapsto$$
$$\mathbf{acc} + pairOfTiles.\_0 *$$
$$Transpose() \circ pairOfTiles.\_1$$
$$) \$ Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$$
$$) \circ Transpose() \circ Tile(sizeN, sizeK) \$ \mathbf{B}$$
$$) \circ Tile(sizeM, sizeK) \$ \mathbf{A}$$

# Tiling



$$TiledMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Untile() \circ$$
$$Map(\overrightarrow{aRows} \mapsto$$
$$Map(\overrightarrow{bCols} \mapsto$$
$$Reduce((\mathbf{acc}, pairOfTiles) \mapsto$$
$$\mathbf{acc} + pairOfTiles.\_0 *$$
$$Transpose() \circ pairOfTiles.\_1$$
$$) \$ \, Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$$
$$) \circ Transpose() \circ Tile(sizeN, sizeK) \$ \, \mathbf{B}$$
$$) \circ Tile(sizeM, sizeK) \$ \, \mathbf{A}$$

# Tiling

**A**



sizeM

sizeK

**C**



**B**



sizeK

sizeN

$TiledMultiply(\mathbf{A}, \mathbf{B}) =$
$\quad Untile() \circ$
$\quad Map(\overrightarrow{aRows} \mapsto$
$\quad\quad Map(\overrightarrow{bCols} \mapsto$
$\quad\quad\quad Reduce((\mathbf{acc}, pairOfTiles) \mapsto$
$\quad\quad\quad\quad \mathbf{acc} + pairOfTiles.\_0 *$
$\quad\quad\quad\quad\quad Transpose() \circ pairOfTiles.\_1$
$\quad\quad\quad) \$ Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$
$\quad\quad) \circ Transpose() \circ Tile(sizeN, sizeK) \$ \mathbf{B}$
$\quad) \circ Tile(sizeM, sizeK) \$ \mathbf{A}$

# Tiling



$$TiledMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Untile() \circ$$
$$Map(\overrightarrow{aRows} \mapsto$$
$$Map(\overrightarrow{bCols} \mapsto$$
$$Reduce((\mathbf{acc}, pairOfTiles) \mapsto$$
$$\mathbf{acc} + pairOfTiles.\_0 *$$
$$Transpose() \circ pairOfTiles.\_1$$
$$) \ \$ \ Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$$
$$) \circ Transpose() \circ Tile(sizeN, sizeK) \ \$ \ \mathbf{B}$$
$$) \circ Tile(sizeM, sizeK) \ \$ \ \mathbf{A}$$

THE UNIVERSITY of EDINBURGH
**informatics**

# Tiling



$$TiledMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Untile() \circ$$
$$Map(\overrightarrow{aRows} \mapsto$$
$$Map(\overrightarrow{bCols} \mapsto$$
$$Reduce((\mathbf{acc}, pairOfTiles) \mapsto$$
$$\mathbf{acc} + pairOfTiles.\_0 *$$
$$Transpose() \circ pairOfTiles.\_1$$
$$) \$ Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$$
$$) \circ Transpose() \circ Tile(sizeN, sizeK) \$ \mathbf{B}$$
$$) \circ Tile(sizeM, sizeK) \$ \mathbf{A}$$

# Combining Optimisations

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

80 rewrites

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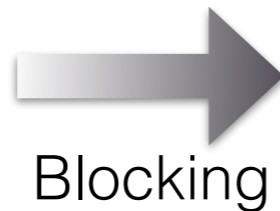$$)$$

22

# Combining Optimisations

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

Tiling →

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

23

# Combining Optimisations

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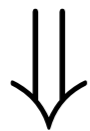$$) \circ Transpose() \ \$ \ \mathbf{B}$$
$$) \ \$ \ \mathbf{A}$$

Blocking

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

24

# Combining Optimisations

$A * B =$

$Map(\overrightarrow{rowA} \mapsto$

$\quad Map(\overrightarrow{colB} \mapsto$

$\quad\quad DotProduct(\overrightarrow{rowA}, \overrightarrow{colB})$

$\quad) \circ Transpose() \,\$\, B$

$) \,\$\, A$

80 rewrites

$(p239, p36 \mapsto$

$Join() \circ Map((p179 \mapsto$

$Transpose() \circ Join() \circ Map((p70 \mapsto$

$Transpose() \circ Join() \circ Map((p20 \mapsto$

$Transpose() \circ Map((p65 \mapsto$

$Transpose()(p65)$

$)) \circ Transpose()(p20)$

$)) \circ Transpose() \circ Reduce((p75, p0 \mapsto$

$Map((p164 \mapsto$

$Join() \circ Map((p81 \mapsto$

$Reduce((p136, p90 \mapsto$

$Map((p163 \mapsto$

$Get(0)(p163) + Get(1)(p163) * Get(1)(p90)$

$)) \circ Zip(2)(p136, Get(0)(p90))$

$))(Get(0)(p81), Zip(2)(Transpose() \circ Get(1)(p164), Get(1)(p81)))$

$)) \circ Zip(2)(Get(0)(p164), Get(1)(p0))$

$)) \circ Zip(2)(p75, Split(blockFactor) \circ Transpose() \circ Get(0)(p0))$

$))(Zip(2)(Split(sizeK) \circ Transpose()(p179), p70))$

$)) \circ Transpose() \circ Map((p4 \mapsto$

$Split(sizeN) \circ Transpose()(p4)$

$)) \circ Split(sizeK)(p36)$

$)) \circ Split(sizeM)(p239)$

$)$

25

# Macro Rules

- A *macro rule* is a rewrite rule that has a particular goal and can apply different rewrite rules to achieve it.

- Examples:

  - 1D Register Blocking, 2D Register Blocking

  - Tiling

  - Map-Map Interchange

# Map-Map Interchange Macro Rule

$Map(a \mapsto$
$\quad Map(b \mapsto$
$\quad\quad f(a,b)) \$ B$
$) \$ A$

$\Downarrow$

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

$Map(a \mapsto$
$\quad Map(b \mapsto$
$\quad\quad f(b)) \$ a$
$) \$ A$

$\Downarrow$

$Transpose() \circ$
$Map(b \mapsto$
$\quad Map(a \mapsto$
$\quad\quad f(a)) \$ b$
$) \circ Transpose() \$ A$

$Map(a \mapsto$
$\quad Map(b \mapsto$
$\quad\quad f(b.\_0, b.\_1)$
$\quad) \$ Zip(a,c)$
$) \$ A$

$\Downarrow$

$Transpose() \circ$
$Map(b \mapsto$
$\quad Map(a \mapsto$
$\quad\quad f(a, b.\_1)$
$\quad) \$ b.\_0$
$) \$ Zip($
$Transpose()\$A,$
$c)$

$Map(a \mapsto$
$\quad Map(b \mapsto$
$\quad\quad f(b, a.\_1)$
$\quad) \$ a.\_0$
$) \$ Zip(A,c)$

$\Downarrow$

$Transpose() \circ$
$Map(b \mapsto$
$\quad Map(a \mapsto$
$\quad\quad f(a.\_0, a.\_1)$
$\quad) \$ Zip(b.\_0, c)$
$) \circ Transpose() \$ A$

# Exploration Strategy

High-Level Expression

Macro Rules

Rewritten Expression

1

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

# Exploration Strategy

High-Level Expression

Macro Rules

Rewritten Expression

**1**

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

**1.1**

$$TiledMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Untile() \circ$$
$$Map(\overrightarrow{aRows} \mapsto$$
$$\quad Map(\overrightarrow{bCols} \mapsto$$
$$\quad\quad Reduce((\mathbf{acc}, pairOfTiles) \mapsto$$
$$\quad\quad\quad \mathbf{acc} + pairOfTiles.\_0 * pairOfTiles.\_1$$
$$\quad ) \,\$\, Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$$
$$\quad ) \circ Transpose() \circ Tile(sizeN, sizeK) \,\$\, \mathbf{B}$$
$$) \circ Tile(sizeM, sizeK) \,\$\, \mathbf{A}$$

**1.2**

$$BlockedMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Join() \circ Map(Transpose()) \circ$$
$$Map(\overrightarrow{rowsA} \mapsto$$
$$\quad Map(\overrightarrow{colB} \mapsto$$
$$\quad\quad Transpose() \circ$$
$$\quad\quad Reduce(((\overline{acc}, rowElemPair) \mapsto$$
$$\quad\quad\quad Map(p \mapsto p.\_0 + p.\_1 * rowElemPair.\_1) \,\$\,$$
$$\quad\quad\quad Zip(\overline{acc}, rowElemPair.\_0)$$
$$\quad\quad ) \,\$\, Zip(Transpose() \,\$\, \overrightarrow{rowsA}, \overrightarrow{colB})$$
$$\quad ) \circ Transpose() \,\$\, \mathbf{B}$$
$$) \circ Split(blockFactor) \,\$\, \mathbf{A}$$

**1.3**

$$TiledMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Untile() \circ$$
$$Map(\overrightarrow{aRows} \mapsto$$
$$\quad Map(\overrightarrow{bCols} \mapsto$$
$$\quad\quad Reduce((\mathbf{acc}, pairOfTiles) \mapsto$$
$$\quad\quad\quad \mathbf{acc} + pairOfTiles.\_0 * pairOfTiles.\_1$$
$$\quad ) \,\$\, Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$$
$$\quad ) \circ Transpose() \circ Tile(sizeN, sizeK) \,\$\, \mathbf{B}$$
$$) \circ Tile(sizeM, sizeK) \,\$\, \mathbf{A}$$

**1.4**

$$BlockedMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Join() \circ Map(Transpose()) \circ$$
$$Map(\overrightarrow{rowsA} \mapsto$$
$$\quad Map(\overrightarrow{colB} \mapsto$$
$$\quad\quad Transpose() \circ$$
$$\quad\quad Reduce(((\overline{acc}, rowElemPair) \mapsto$$
$$\quad\quad\quad Map(p \mapsto p.\_0 + p.\_1 * rowElemPair.\_1) \,\$\,$$
$$\quad\quad\quad Zip(\overline{acc}, rowElemPair.\_0)$$
$$\quad\quad ) \,\$\, Zip(Transpose() \,\$\, \overrightarrow{rowsA}, \overrightarrow{colB})$$
$$\quad ) \circ Transpose() \,\$\, \mathbf{B}$$
$$) \circ Split(blockFactor) \,\$\, \mathbf{A}$$

# Exploration Strategy

High-Level Expression

Rewritten Expression

Macro Rules

1.3

$$TiledMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Untile() \circ$$
$$Map(\overrightarrow{aRows} \mapsto$$
$$Map(\overrightarrow{bCols} \mapsto$$
$$Reduce((\mathbf{acc}, pairOfTiles) \mapsto$$
$$\mathbf{acc} + pairOfTiles.\_0 * pairOfTiles.\_1$$
$$) \$ Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$$
$$) \circ Transpose() \circ Tile(sizeN, sizeK) \$ \mathbf{B}$$
$$) \circ Tile(sizeM, sizeK) \$ \mathbf{A}$$

# Exploration Strategy

High-Level Expression

Macro Rules

Rewritten Expression

$1.3$

$$TiledMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Untile() \circ$$
$$Map(\overrightarrow{aRows} \mapsto$$
$$Map(\overrightarrow{bCols} \mapsto$$
$$Reduce((\mathbf{acc}, pairOfTiles) \mapsto$$
$$\mathbf{acc} + pairOfTiles._0 * pairOfTiles._1$$
$$) \$ Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$$
$$) \circ Transpose() \circ Tile(sizeN, sizeK) \$ \mathbf{B}$$
$$) \circ Tile(sizeM, sizeK) \$ \mathbf{A}$$

# Exploration Strategy



1.3

$$
\begin{aligned}
&TiledMultiply(\mathbf{A}, \mathbf{B}) = \\
&\quad Untile() \circ \\
&\quad Map(\overrightarrow{aRows} \mapsto \\
&\qquad Map(\overrightarrow{bCols} \mapsto \\
&\qquad\quad Reduce((\mathbf{acc}, pairOfTiles) \mapsto \\
&\qquad\qquad \mathbf{acc} + pairOfTiles.\_0 * pairOfTiles.\_1 \\
&\qquad\quad) \ \$ \ Zip(\overrightarrow{aRows}, \overrightarrow{bCols}) \\
&\qquad) \circ Transpose() \circ Tile(sizeN, sizeK) \ \$ \ \mathbf{B} \\
&\quad) \circ Tile(sizeM, sizeK) \ \$ \ \mathbf{A}
\end{aligned}
$$

# Exploration Strategy



High-Level Expression

Macro Rules

Rewritten Expression

Map to OpenCL

Lowered Expression

**1.3**

$TiledMultiply(\mathbf{A}, \mathbf{B}) =$
$\quad Untile() \circ$
$\quad Map(\overrightarrow{aRows} \mapsto$
$\quad\quad Map(\overrightarrow{bCols} \mapsto$
$\quad\quad\quad Reduce((\mathbf{acc}, pairOfTiles) \mapsto$
$\quad\quad\quad\quad \mathbf{acc} + pairOfTiles.\_0 * pairOfTiles.\_1$
$\quad\quad\quad ) \$ Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$
$\quad\quad ) \circ Transpose() \circ Tile(sizeN, sizeK) \$ \mathbf{B}$
$\quad ) \circ Tile(sizeM, sizeK) \$ \mathbf{A}$

**1.3.1**

$TiledMultiply(\mathbf{A}, \mathbf{B}) =$
$\quad Untile() \circ$
$\quad MapWrg(1)(\overrightarrow{aRows} \mapsto$
$\quad\quad MapWrg(0)(\overrightarrow{bCols} \mapsto$
$\quad\quad\quad ReduceSeq((\mathbf{acc}, pairOfTiles) \mapsto$
$\quad\quad\quad\quad \mathbf{acc} + toLocal(pairOfTiles.\_0)$
$\quad\quad\quad\quad\quad * toLocal(pairOfTiles.\_1)$
$\quad\quad\quad ) \$ Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$
$\quad\quad ) \circ Transpose() \circ Tile(sizeN, sizeK) \$ \mathbf{B}$
$\quad ) \circ Tile(sizeM, sizeK) \$ \mathbf{A}$

**1.3.2**

$TiledMultiply(\mathbf{A}, \mathbf{B}) =$
$\quad Untile() \circ$
$\quad MapWrg(1)(\overrightarrow{aRows} \mapsto$
$\quad\quad MapWrg(0)(\overrightarrow{bCols} \mapsto$
$\quad\quad\quad ReduceSeq((\mathbf{acc}, pairOfTiles) \mapsto$
$\quad\quad\quad\quad \mathbf{acc} + toLocal(pairOfTiles.\_0)$
$\quad\quad\quad\quad\quad * toLocal(pairOfTiles.\_1)$
$\quad\quad\quad ) \$ Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$
$\quad\quad ) \circ Transpose() \circ Tile(sizeN, sizeK) \$ \mathbf{B}$
$\quad ) \circ Tile(sizeM, sizeK) \$ \mathbf{A}$

**1.3.3**

$TiledMultiply(\mathbf{A}, \mathbf{B}) =$
$\quad Untile() \circ$
$\quad MapWrg(1)(\overrightarrow{aRows} \mapsto$
$\quad\quad MapWrg(0)(\overrightarrow{bCols} \mapsto$
$\quad\quad\quad ReduceSeq((\mathbf{acc}, pairOfTiles) \mapsto$
$\quad\quad\quad\quad \mathbf{acc} + toLocal(pairOfTiles.\_0)$
$\quad\quad\quad\quad\quad * toLocal(pairOfTiles.\_1)$
$\quad\quad\quad ) \$ Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$
$\quad\quad ) \circ Transpose() \circ Tile(sizeN, sizeK) \$ \mathbf{B}$
$\quad ) \circ Tile(sizeM, sizeK) \$ \mathbf{A}$

# Exploration Strategy

High-Level Expression

Macro Rules

Rewritten Expression

Map to OpenCL

Lowered Expression

1.3.2

$$TiledMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Untile() \circ$$
$$MapWrg(1)(\overrightarrow{aRows} \mapsto$$
$$MapWrg(0)(\overrightarrow{bCols} \mapsto$$
$$ReduceSeq((\mathbf{acc}, pairOfTiles) \mapsto$$
$$\mathbf{acc} + toLocal(pairOfTiles.\_0)$$
$$* toLocal(pairOfTiles.\_1)$$
$$) \$ Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$$
$$) \circ Transpose() \circ Tile(sizeN, sizeK) \$ \mathbf{B}$$
$$) \circ Tile(sizeM, sizeK) \$ \mathbf{A}$$

# Exploration Strategy



High-Level Expression

Macro Rules

Rewritten Expression

Map to OpenCL

Lowered Expression

1.3.2

$$TiledMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Untile() \circ$$
$$MapWrg(1)(\overrightarrow{aRows} \mapsto$$
$$MapWrg(0)(\overrightarrow{bCols} \mapsto$$
$$ReduceSeq((\mathbf{acc}, pairOfTiles) \mapsto$$
$$\mathbf{acc} + toLocal(pairOfTiles.\_0)$$
$$* toLocal(pairOfTiles.\_1)$$
$$) \$ Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$$
$$) \circ Transpose() \circ Tile(sizeN, sizeK) \$ \mathbf{B}$$
$$) \circ Tile(sizeM, sizeK) \$ \mathbf{A}$$

# Exploration Strategy



Macro Rules

Map to OpenCL

Parameter Mapping

1.3.2

$$TiledMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Untile() \circ$$
$$MapWrg(1)(\overrightarrow{aRows} \mapsto$$
$$MapWrg(0)(\overrightarrow{bCols} \mapsto$$
$$ReduceSeq((\mathbf{acc}, pairOfTiles) \mapsto$$
$$\mathbf{acc} + toLocal(pairOfTiles.\_0)$$
$$* toLocal(pairOfTiles.\_1)$$
$$) \; \$ \; Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$$
$$) \circ Transpose() \circ Tile(sizeN, sizeK) \; \$ \; \mathbf{B}$$
$$) \circ Tile(sizeM, sizeK) \; \$ \; \mathbf{A}$$

# Exploration Strategy

# Exploration Strategy



1.3.2.5

$$TiledMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Untile() \circ$$
$$MapWrg(1)(\overrightarrow{aRows} \mapsto$$
$$MapWrg(0)(\overrightarrow{bCols} \mapsto$$
$$ReduceSeq((\mathbf{acc}, pairOfTiles) \mapsto$$
$$\mathbf{acc} + toLocal(pairOfTiles._0)$$
$$* toLocal(pairOfTiles._1)$$
$$) \, \$ \, Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$$
$$) \circ Transpose() \circ Tile(128, 16) \, \$ \, \mathbf{B}$$
$$) \circ Tile(128, 16) \, \$ \, \mathbf{A}$$

Macro Rules

Map to OpenCL

Parameter Mapping

High-Level Expression

Rewritten Expression

Lowered Expression

Specialised Expression

# Exploration Strategy



High-Level Expression

Macro Rules

Rewritten Expression

Map to OpenCL

Lowered Expression

Parameter Mapping

Specialised Expression

1.3.2.5

$$TiledMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Untile() \circ$$
$$MapWrg(1)(\overrightarrow{aRows} \mapsto$$
$$MapWrg(0)(\overrightarrow{bCols} \mapsto$$
$$ReduceSeq((\mathbf{acc}, pairOfTiles) \mapsto$$
$$\mathbf{acc} + toLocal(pairOfTiles._0)$$
$$* toLocal(pairOfTiles._1)$$
$$) \$ Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$$
$$) \circ Transpose() \circ Tile(128, 16) \$ \mathbf{B}$$
$$) \circ Tile(128, 16) \$ \mathbf{A}$$

# Exploration Strategy



**High-Level Expression** → Macro Rules → **Rewritten Expression** → Map to OpenCL → **Lowered Expression** → Parameter Mapping → **Specialised Expression** → Code Generation → **OpenCL Code**

1.3.2.5

$$TiledMultiply(\mathbf{A}, \mathbf{B}) =$$
$$Untile() \circ$$
$$MapWrg(1)(\overrightarrow{aRows} \mapsto$$
$$MapWrg(0)(\overrightarrow{bCols} \mapsto$$
$$ReduceSeq((\mathbf{acc}, pairOfTiles) \mapsto$$
$$\mathbf{acc} + toLocal(pairOfTiles.\_0)$$
$$* toLocal(pairOfTiles.\_1)$$
$$) \$ Zip(\overrightarrow{aRows}, \overrightarrow{bCols})$$
$$) \circ Transpose() \circ Tile(128, 16) \$ \mathbf{B}$$
$$) \circ Tile(128, 16) \$ \mathbf{A}$$

# Exploration Strategy



High-Level Expression

→ Macro Rules

Rewritten Expression

→ Map to OpenCL

Lowered Expression

→ Parameter Mapping

Specialised Expression

→ Code Generation

OpenCL Code

```
1   kernel mm_amd_opt(global float * A, B, C,
2                      int K, M, N) {
3   local float  tileA [512];  tileB [512];
4
5   private float acc_0;       ...;  acc_31;
6   private float blockOfB_0; ...;  blockOfB_3;
7   private float blockOfA_0; ...;  blockOfA_7;
8
9   int lid0 = local_id (0);  lid1 = local_id (1);
10  int wid0 = group_id(0); wid1 = group_id(1);
11
12  for (int w1=wid1; w1<M/64; w1+=num_grps(1)) {
13   for (int w0=wid0; w0<N/64; w0+=num_grps(0)) {
14
15    acc_0 = 0.0f;  ...;  acc_31 = 0.0f;
16    for (int i=0; i<K/8; i++) {
17     vstore4(vload4(lid1*M/4+2*i*M+16*w1+lid0,A), 16*lid1+lid0, tileA);
18     vstore4(vload4(lid1*N/4+2*i*N+16*w0+lid0,B), 16*lid1+lid0, tileB);
19     barrier (...) ;
20
21     for (int j = 0; j<8; j++) {
22      blockOfA_0 = tileA[0+64*j+lid1*8];   ...;  blockOfA_7 = tileA[7+64*j+lid1*8];
23      blockOfB_0 = tileB[0 +64*j+lid0];      ...;  blockOfB_3 = tileB[48+64*j+lid0];
24
25      acc_0  += blockOfA_0 * blockOfB_0;  ...;  acc_28 += blockOfA_7 * blockOfB_0;
26      acc_1  += blockOfA_0 * blockOfB_1;  ...;  acc_29 += blockOfA_7 * blockOfB_1;
27      acc_2  += blockOfA_0 * blockOfB_2;  ...;  acc_30 += blockOfA_7 * blockOfB_2;
28      acc_3  += blockOfA_0 * blockOfB_3;  ...;  acc_31 += blockOfA_7 * blockOfB_3;
29     }
30     barrier (...) ;
31    }
32
33    C[ 0+8*lid1*N+64*w0+64*w1*N+0*N+lid0]=acc_0; ...; C[ 0+8*lid1*N+64*w0+64*w1*N+7*N+lid0]=acc_28;
34    C[16+8*lid1*N+64*w0+64*w1*N+0*N+lid0]=acc_1; ...; C[16+8*lid1*N+64*w0+64*w1*N+7*N+lid0]=acc_29;
35    C[32+8*lid1*N+64*w0+64*w1*N+0*N+lid0]=acc_2; ...; C[32+8*lid1*N+64*w0+64*w1*N+7*N+lid0]=acc_30;
36    C[48+8*lid1*N+64*w0+64*w1*N+0*N+lid0]=acc_3; ...; C[48+8*lid1*N+64*w0+64*w1*N+7*N+lid0]=acc_31;
37  } } }
```

# Heuristics

**For Macro Rules:**

- Nesting depth

- Distance of addition and multiplication
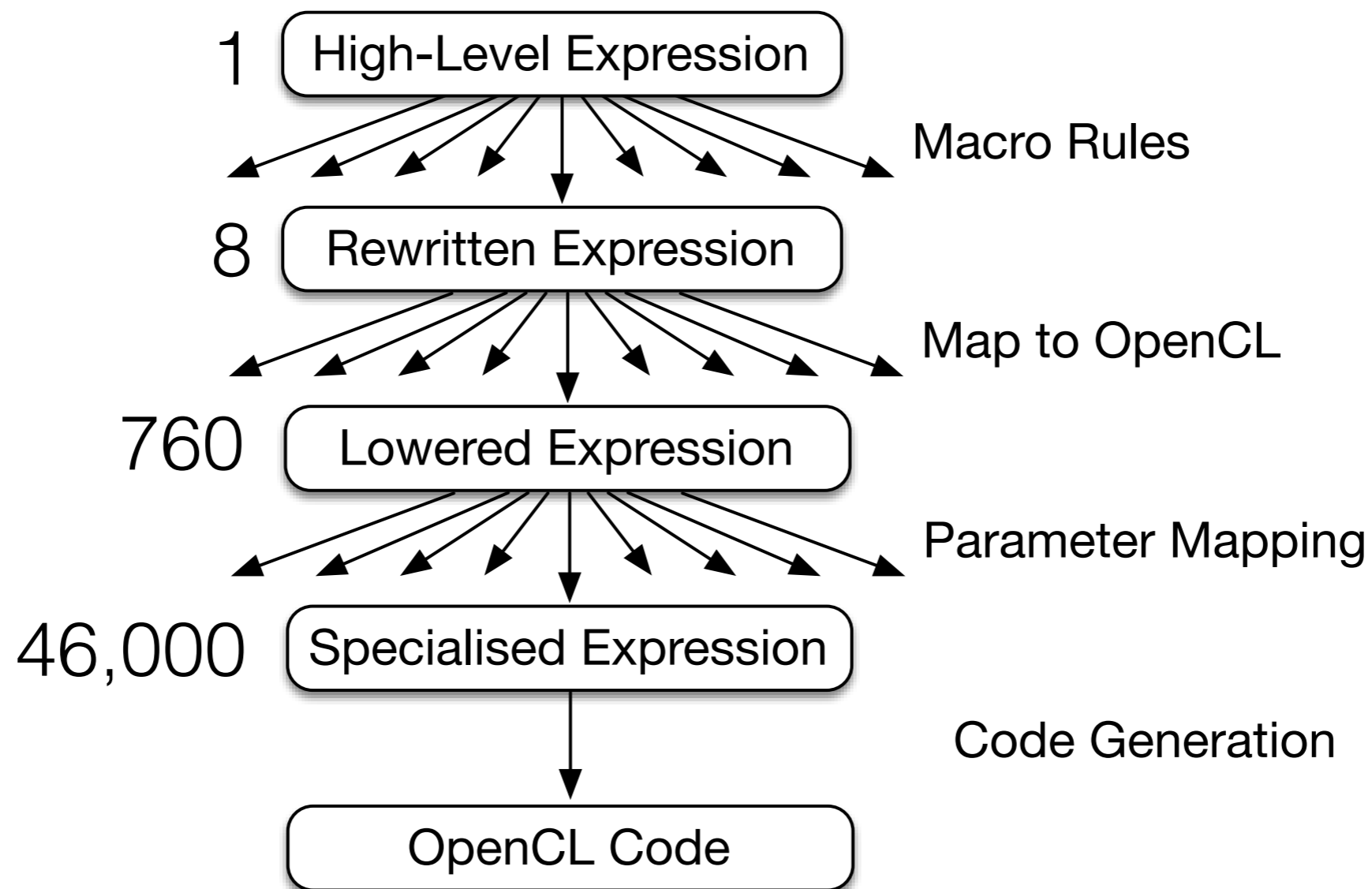
- Number of times rules are applied

**For Map to OpenCL:**

- Fixed parallelism mapping

- Limited choices for mapping to local and global memory
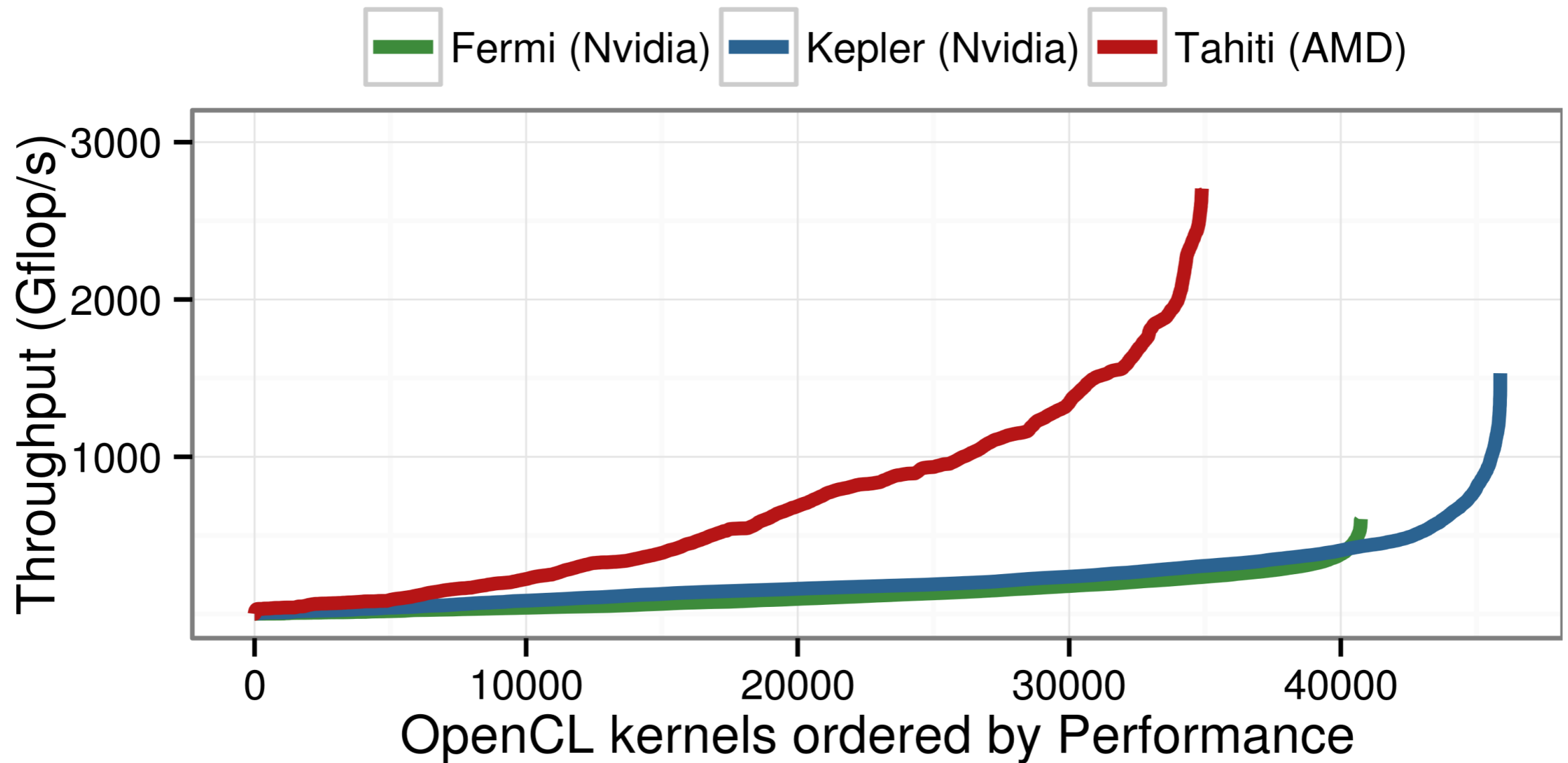
- Follows best practice

**For Parameter Mapping:**

- Amount of memory used

  - Global

  - Local

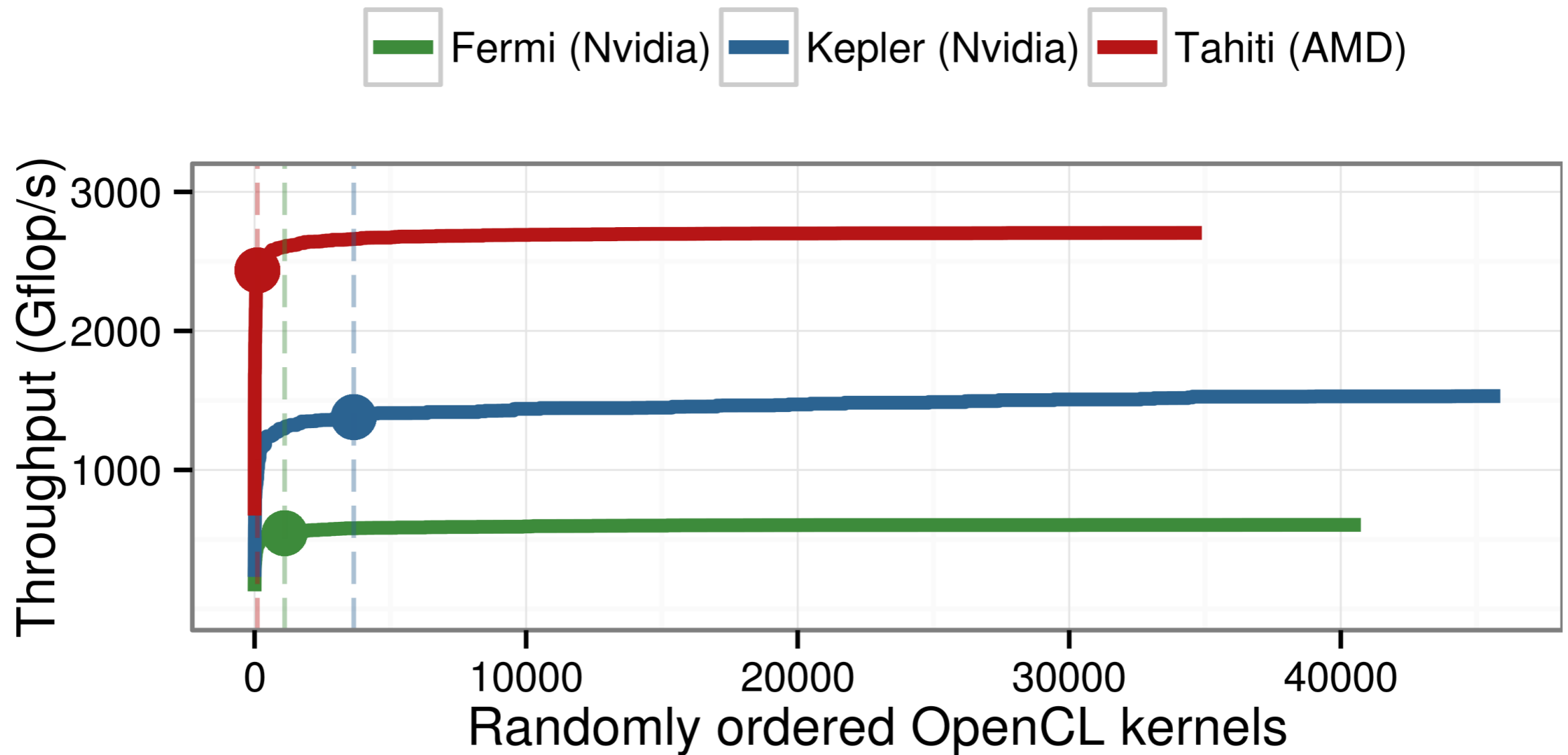  - Registers

- Amount of parallelism

  - Work-items

  - Workgroup

THE UNIVERSITY *of* EDINBURGH
**informatics**

# Exploration in Numbers

1 — High-Level Expression

Macro Rules

8 — Rewritten Expression

Map to OpenCL

760 — Lowered Expression

Parameter Mapping

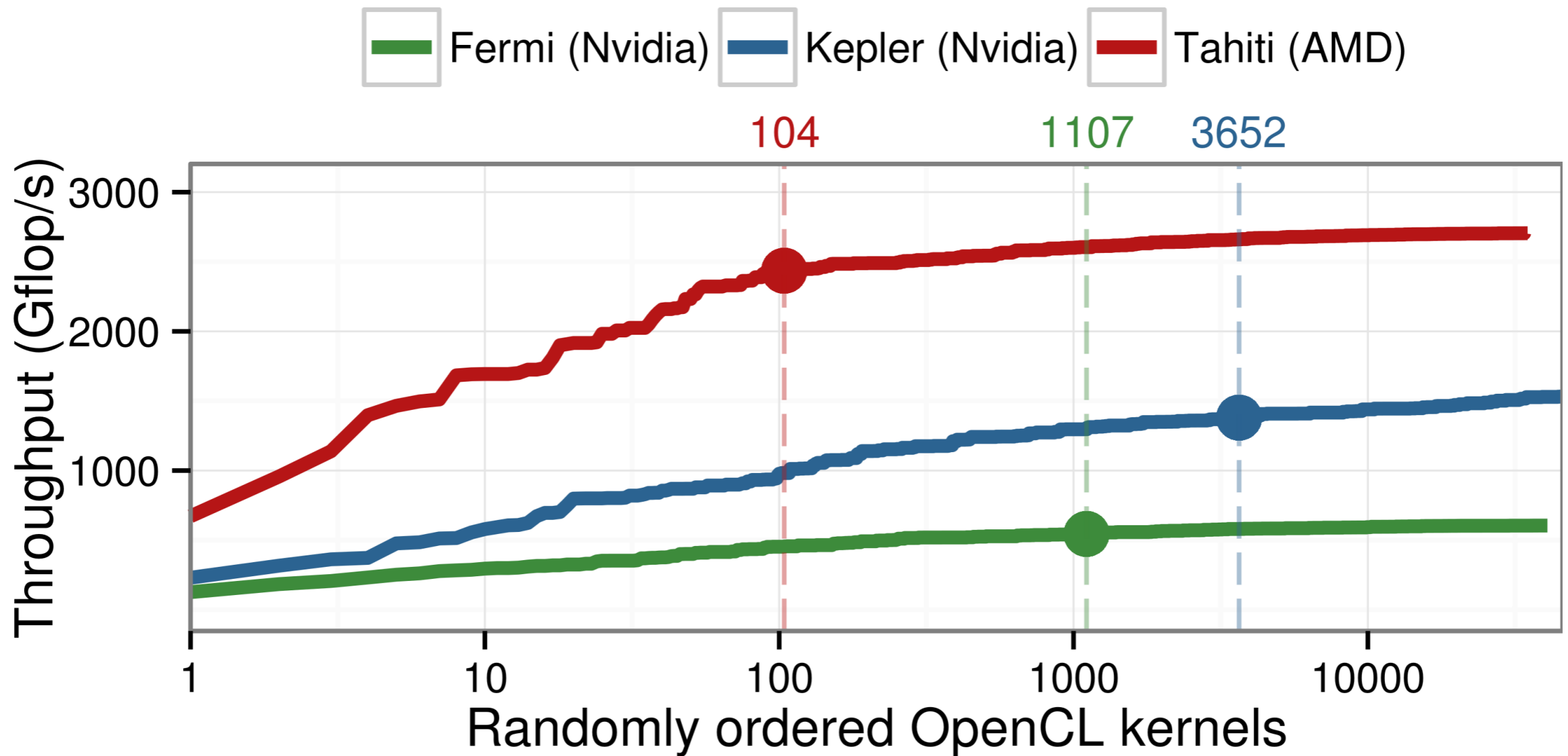46,000 — Specialised Expression
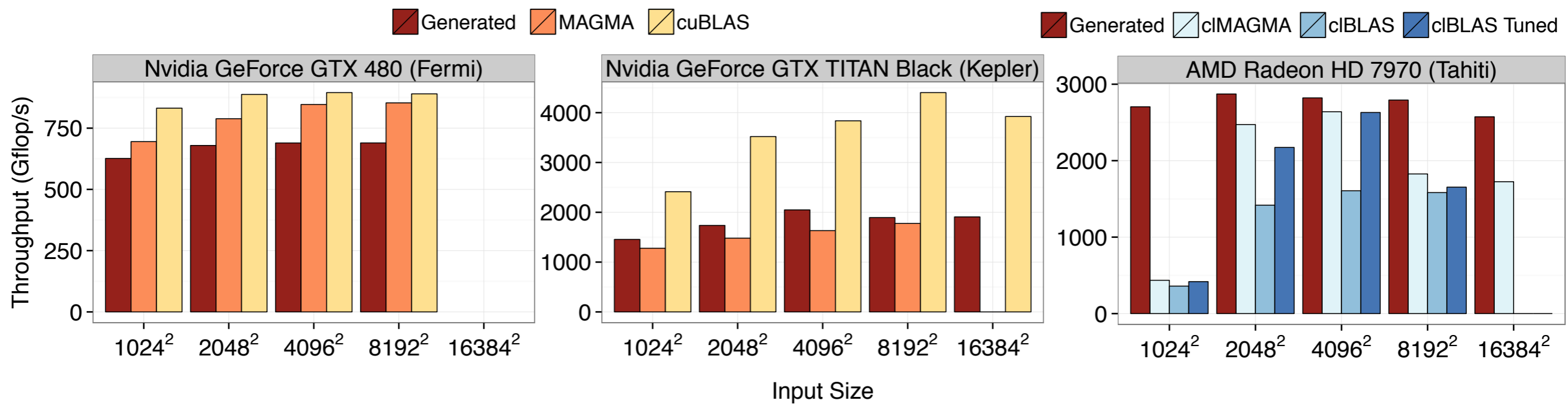
Code Generation
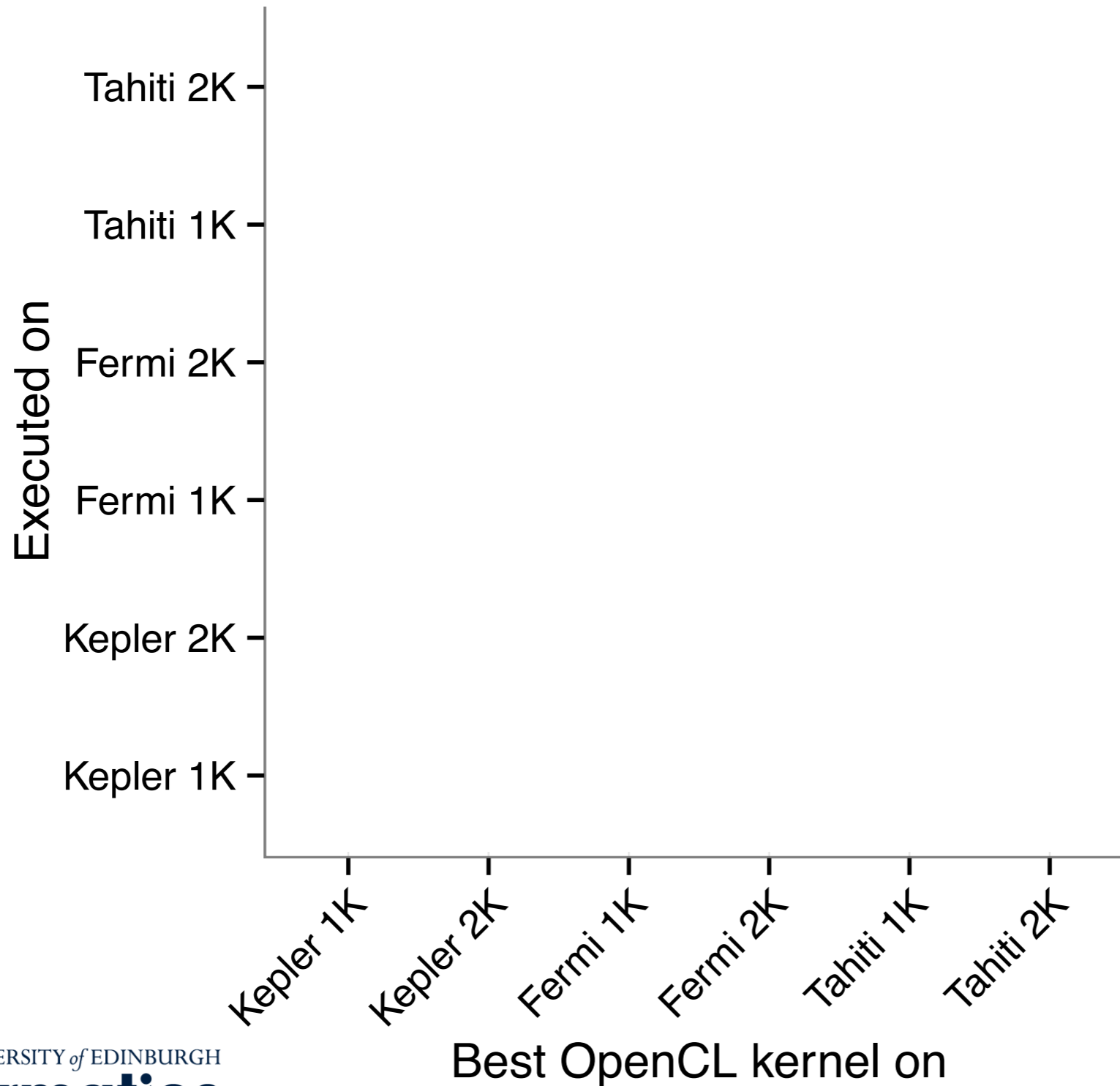
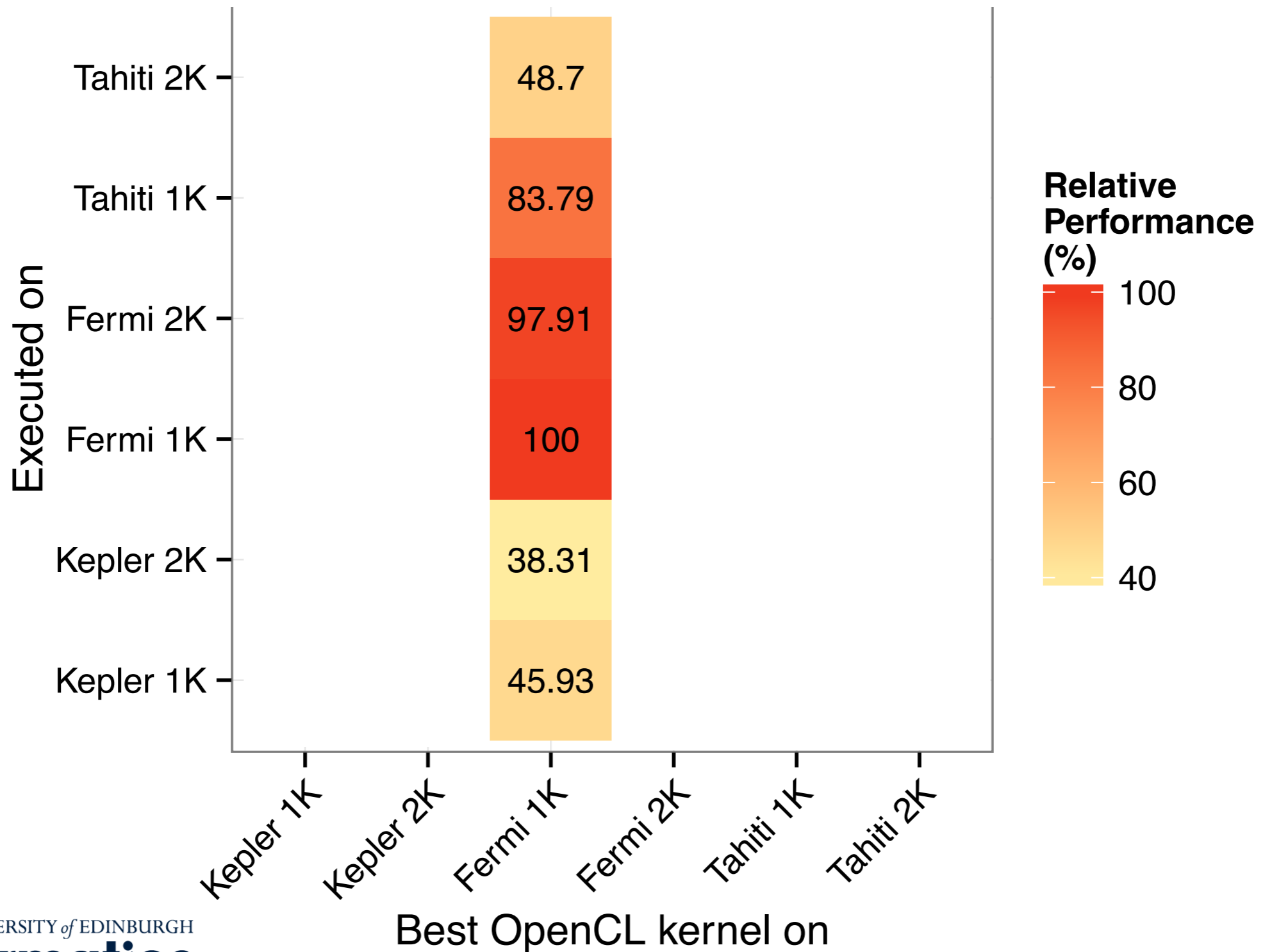OpenCL Code

# Exploration Space

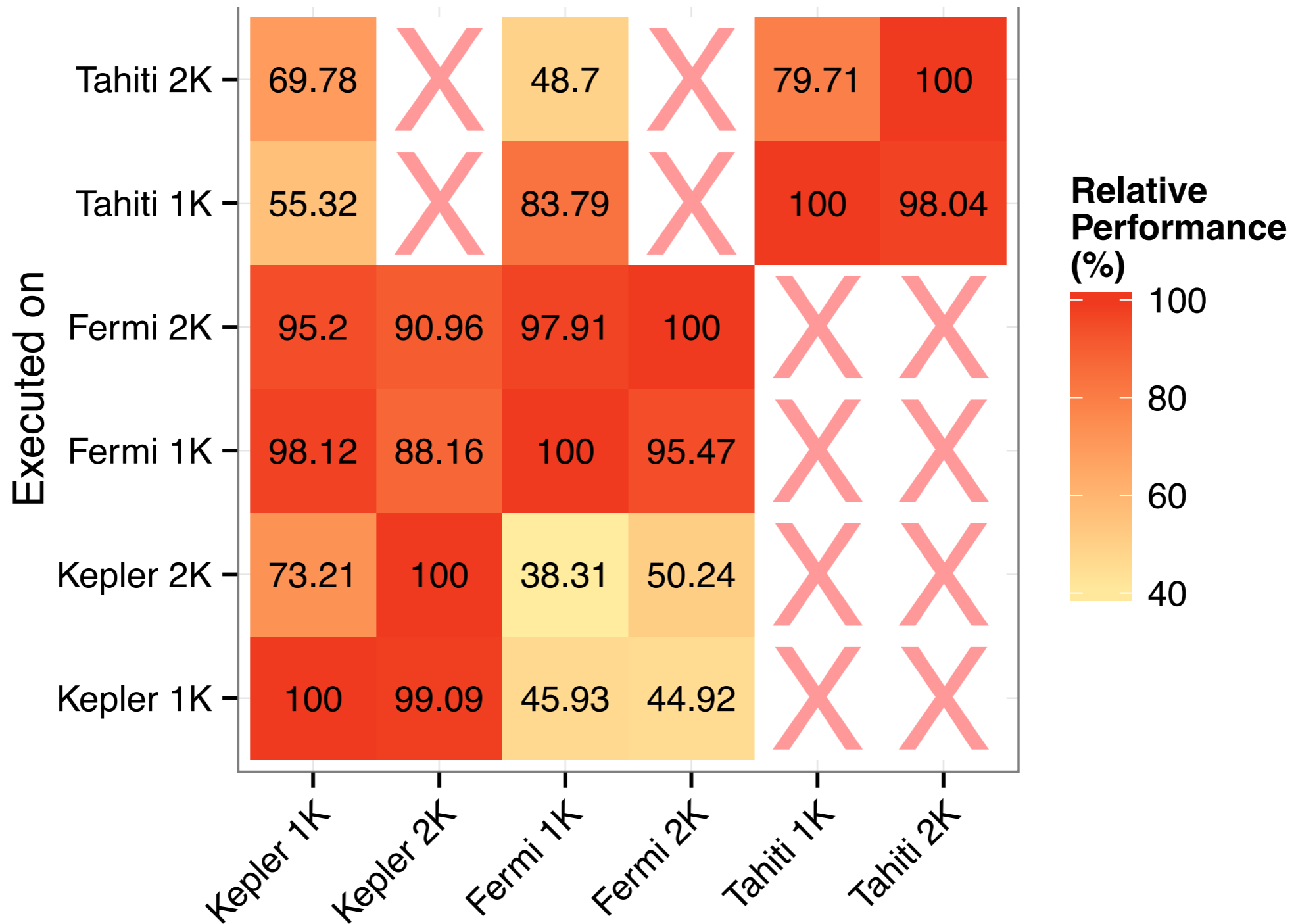# Performance Evolution

# Performance Evolution

# Performance Results

# Performance Portability

# Performance Portability

# Performance Portability

# Conclusion

- OpenCL code is not performance portable

- Using a functional approach along with rewrite rules we can generate performance portable code

- Performance of matrix multiplication on par with tuned OpenCL code

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

Supported by:

41