



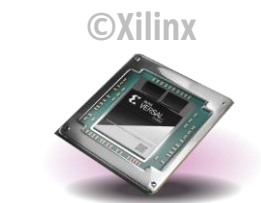
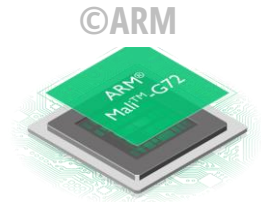
# **Guided Rewriting and Constraint Satisfaction for Parallel GPU Code Generation**

---

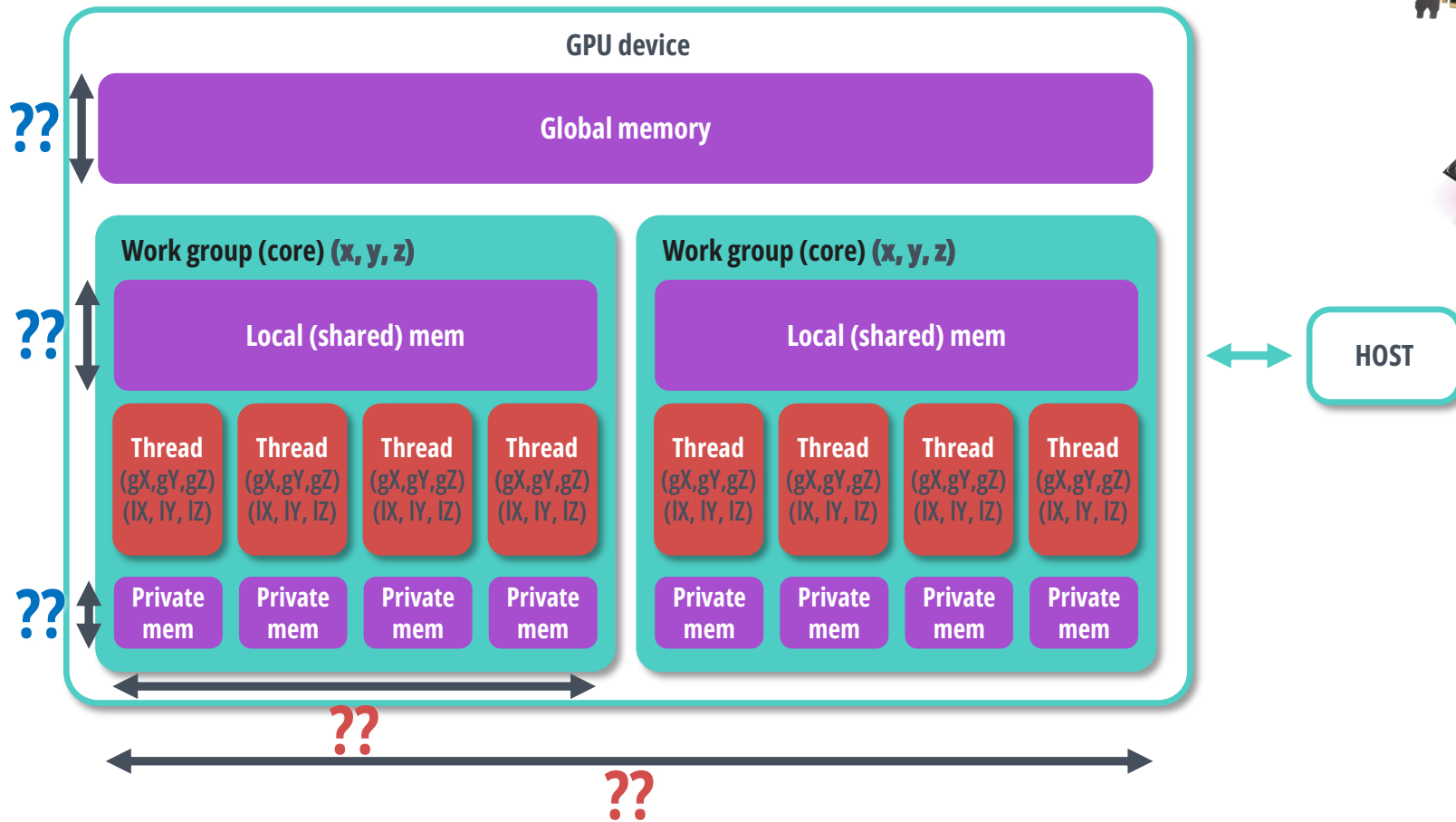
**Naums Mogers**

May 9<sup>th</sup>, 2023

# Parallel Architectures

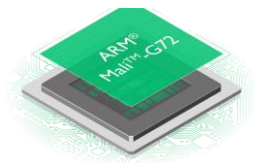


Parallel architectures are hard to optimize for



# Parallel Architectures

©ARM



©NVIDIA

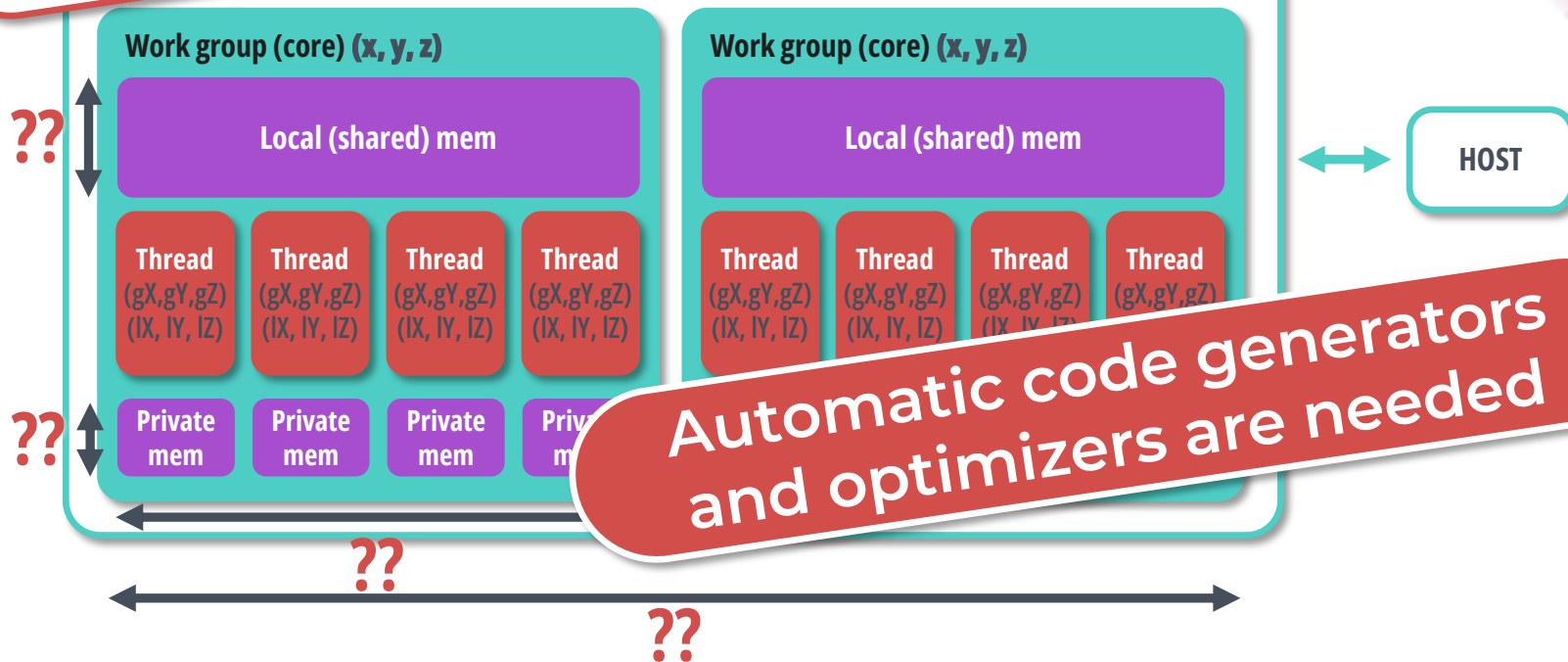


©Xilinx



Parallel architectures are hard to optimize for

Large implementation space



Automatic code generators and optimizers are needed

# The Programmability Challenge

---

- Hierarchical execution and memory models
- Diverse and heterogeneous accelerator architectures
- Manual optimisation is too costly
- Heuristic optimisation strategies are over-constrained
- Automatic optimisation suffers from the search space explosion

# Current Approaches

---

- Kernel libraries
  - Costly to maintain and extend
- User-provided schedules (*Halide*) and design choices (*PetaBricks*, *Tangram*)
  - Burden on the user
- Polyhedral compilation (*Tensor Comprehensions*)
  - Limited to affine loops
- Functional rewriting (*Futhark*, *Accelerate*, *Lift*, *RISE*)
  - Dependence on heuristics

Low-level expression  
close to OpenCL

```
1 join o
2 transposeW o
3 mapWrg(0)(
4   mapLcl(0)(f) o
5   asVector(4)) o
6 transpose o
7 split(64)
```



High-level  
expression

```
1 map(f)
```

### Low-level expression close to OpenCL

```
1 join o
2 transposeW o
3 mapWrg(0)(
4   mapLcl(0)(f) o
5   asVector(4)) o
6 transpose o
7 split(64)
```

### Parallelised expression

```
1 join o
2 transposeW o
3 mapWrg(0)(
4   mapLcl(0)(f) o
5   asVector(v)) o
6 transpose o
7 split(s)
```

### High-level expression

```
1 map(f)
```

Chapter 4

*Tuning*

# Chapter 4:

## Functional IR for Auto-Tuning

- Shows that a functional IR can represent low-level optimisations in the context of convolution

```
1 def partialConv(kernelsWeights : [[[[float]]inputChannels]kernelWidth]kernelHeight]numKernels ,
2                 paddedInput    : [[[[float]]inputChannels]paddedInputWidth]paddedInputHeight ,
3                 kernelStride    : (int, int))
4   : [[[[[[float]]windowSize/ω]σ]nWindowsInTile/σ]κ]numKernels/κ]nTilesInInput = {
5   val tiledInput4D = join(slide2D(0, tilingStride, paddedInput))
6   val tiledSlidedInput5D = map(join(slide2D((kernelHeight, kernelWidth), kernelStride)), tiledInput4D)
7   val windowSize = inputChannels * kernelWidth * kernelHeight
8   def coalesceChunkVectorizeWindow(window : [[[[float]]inputChannels]kernelWidth]kernelHeight])
9     : [[floatv]ω]windowSize/ω = {
10     val flatWindow1D = join(join(window))
11     val flatCoalescedWindow1D = reorder(striddenIndex(windowSize/ω), flatWindow1D)
12     val flatCoalescedChunkedWindow1D = split(ω, flatCoalescedWindow1D)
13     asVector(v, flatCoalescedChunkedWindow1D) }
14   val tiledSlidedCoalescedChunkedVectorizedInput4D = map(tile4D -> split(σ, map
15     coalesceChunkVectorizeWindow(window3D), tile4D)), tiledSlidedInput5D)
16   val groupedCoalescedChunkedVectorizedKernelsWeights4D = split(κ, map(singleKernelWeights ->
17     coalesceChunkVectorizeWindow(singleKernelWeights), kernelsWeights))
18   mapWrg(1, inputTile3D ->
19     mapWrg(0, kernelsGroupWeights3D -> transpose(
20       mapLcl(1, inputWindows2D -> transpose(
21         mapLcl(0, (inputWindowsChunk1D, kernelsGroupChunk2D) ->
22           mapSeq(singleKernelReducedChunk -> toGlobal(singleKernelReducedChunk),
23             join(
24               reduceSeq(
25                 init = mapSeq(toPrivate(id(Value(0, [float]κ)))),
26                 f = (acc, (inputsValue, kernelsGroupValue1D)) ->
27                   let(inputsValuePrivate ->
28                     mapSeq((accValue, singleKernelValue) ->
29                       mapSeq((inputValuePrivate) ->
30                         accValue + vectorize(v, dot(inputValuePrivate, singleKernelValue)),
31                         inputsValuePrivate,
32                         zip(acc, kernelsGroupValue1D),
33                         mapSeq(toPrivate(vectorize(v, id(inputValue))))),
34                         zip(transpose(inputWindowsChunk1D, transpose(kernelsGroupChunk2D))),
35                         zip(inputWindows2D, transpose(kernelsGroupWeights3D))),
36                         inputTile3D)),
37                   groupedCoalescedChunkedVectorizedKernelsWeights4D),
38                   tiledSlidedCoalescedChunkedVectorizedInput4D)
```

Tiling

Coalescing

Vectorization

Loop unrolling

Kernel grouping

Thread coarsening

Memory optimization



# Chapter 4:

## Functional IR for Auto-Tuning

---

- Shows that a **functional IR** can in the context of convolution
- Describes **an auto-tuning approach** which leverages strongly typed functional patterns

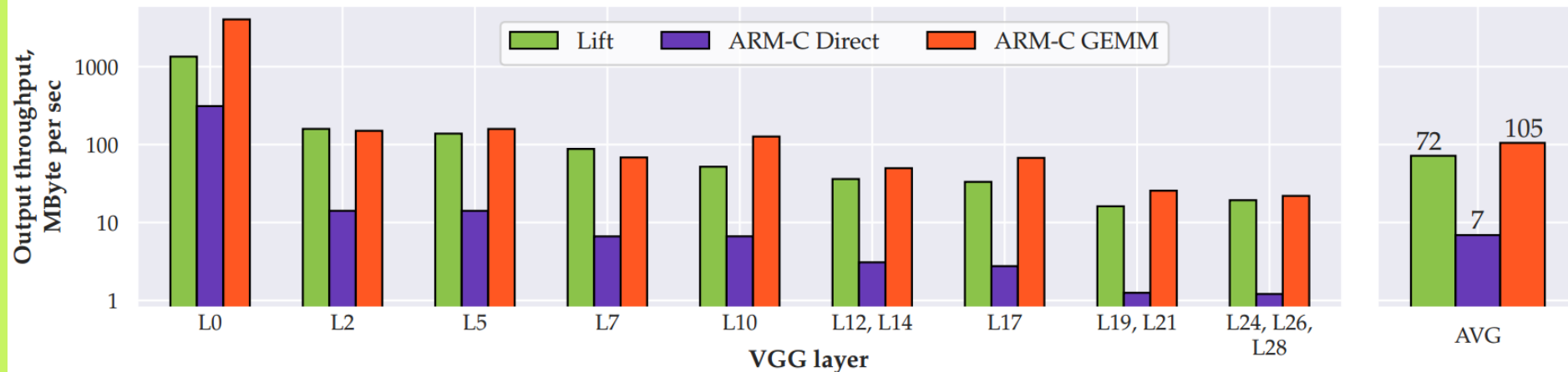
`Split(s) $ [T]N`  $\Rightarrow$  `N % s == 0`

`asVector(v) $ [T]N`  $\Rightarrow$  `N % v == 0`

`Slide(len, step) o [T]N`  $\Rightarrow$  `N >= len`

`SlideStrict(len, step) o [T]N`  $\Rightarrow$  `N >= len &&`  
`N % ((N - (len - step)) / step) == 0`

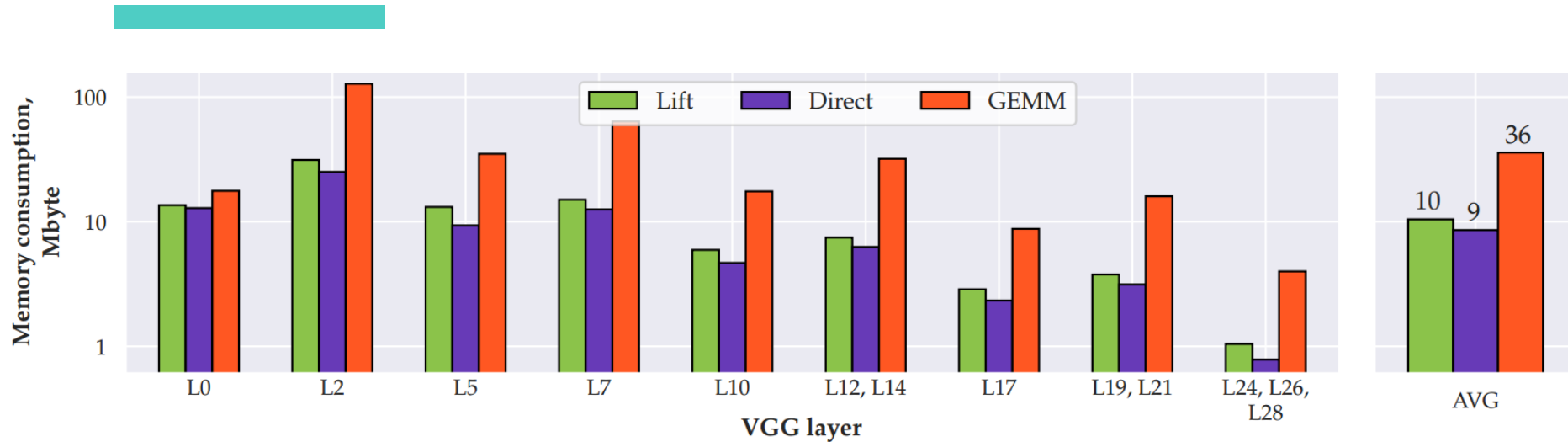
# Results: Throughput



□ Lift kernels are:

- Always faster than **direct convolution** in ARM-C (**x10 on average**)
- In some cases, on par or better than **GEMM** in ARM-C (**x0.7 on average**)

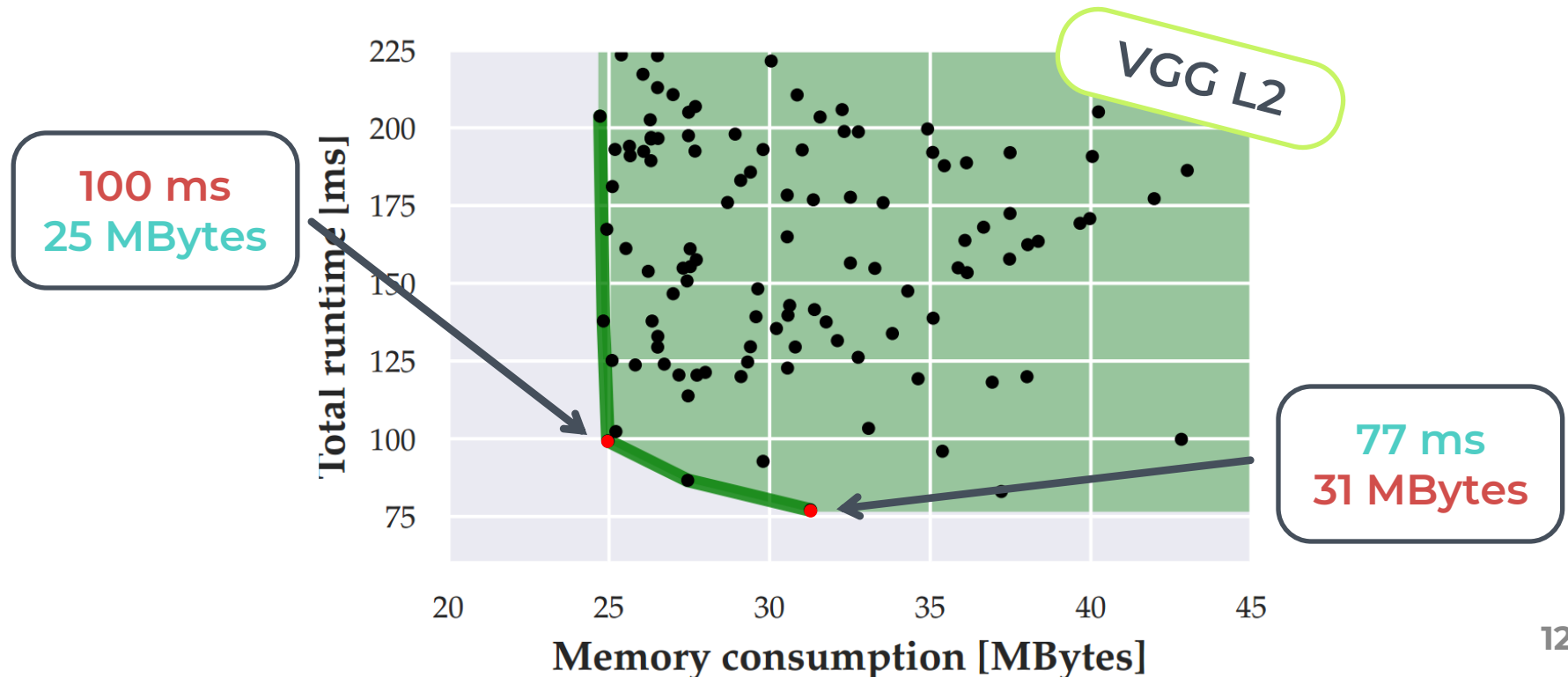
# Results: Memory Consumption



- Memory consumption of Lift kernels is:
  - On par with **direct convolution** in ARM-C (**x1.1 on average**)
  - Always better than **GEMM** in ARM-C (**x3.6 on average**)

# Results: Multi-Objective Optimisation

- Shifting priorities:
  - Low memory footprint vs throughput/latency
- Search space exploration for multi-objective optimization



## Low-level expression close to OpenCL

```
1 join o
2 transposeW o
3 mapWrg(0)(
4   mapLcl(0)(f) o
5   asVector(4)) o
6 transpose o
7 split(64)
```

## Parallelised expression

```
1 join o
2 transposeW o
3 mapWrg(0)(
4   mapLcl(0)(f) o
5   asVector(v)) o
6 transpose o
7 split(s)
```

## Structurally optimised expression

```
1 join o
2 transposeW o
3 map(
4   map(f)) o
5 transpose o
6 split(s)
```

## High-level expression

```
1 map(f)
```

Chapter 4

*Tuning*

Chapter 5

*Parallelization  
Vectorization*

# Chapter 5:

## Parallelism Mapping Through Constraint Satisfaction

---

*I found the approach proposed by the author extremely convincing and fairly natural (to the point I'm almost surprised it wasn't proposed earlier).*

CC'22 reviewer

# Chapter 5:

## Parallelism Mapping Through Constraint Satisfaction

---

### ■ Functional patterns

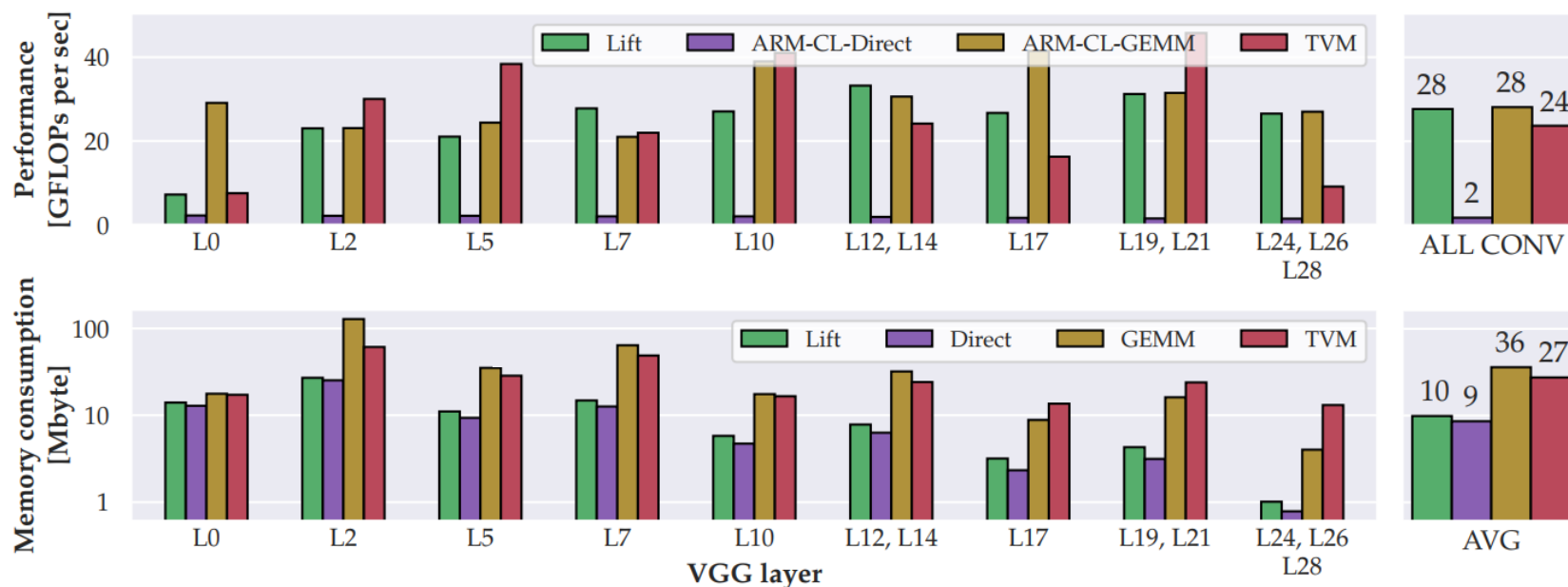
- ...expose parallelism
- ...express parallel restrictions succinctly
- ...aid dependency analysis for synchronization

### ■ Chapter 5:

- Expresses the GPU parallel programming model as arithmetic constraints on functional patterns
- Uses a constraint solver to explore parallel mappings
- Describes a functional IR-based barrier insertion method

# Results: Performance & Memory

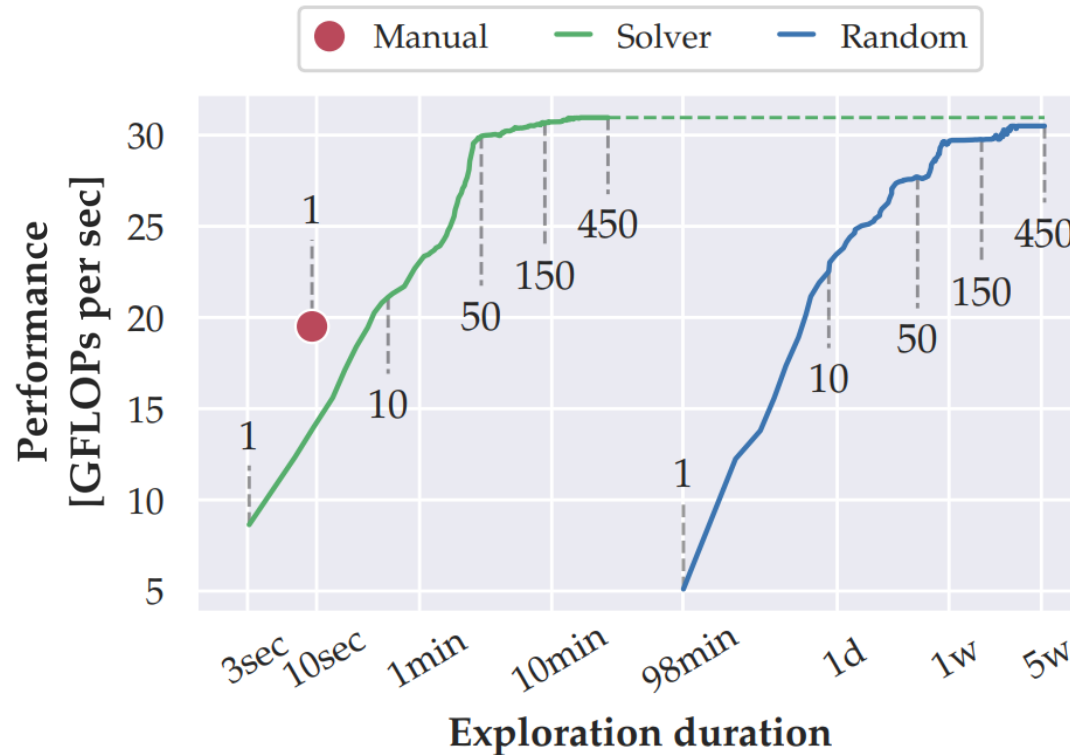
VGG-16 on Mali-G72 GPU



- Stencil performance **on par** with ARM-CL GEMM and **0.86x** of TVM's
- **3.6x** less memory on average than ARM-CL GEMM
- **2.7x** less memory on average than TVM



# Results: Exploration Efficiency



- Peak performance after 95 minutes
- Peaks before the random approach produces even 1 result (a bad one)

Low-level expression  
close to OpenCL

```
1 join o
2 transposeW o
3 mapWrg(0)(
4   mapLcl(0)(f) o
5   asVector(4)) o
6 transpose o
7 split(64)
```

Parallelised  
expression

```
1 join o
2 transposeW o
3 mapWrg(0)(
4   mapLcl(0)(f) o
5   asVector(v)) o
6 transpose o
7 split(s)
```

Structurally  
optimised  
expression

```
1 join o
2 transposeW o
3 map(
4   map(f)) o
5 transpose o
6 split(s)
```

High-level  
expression

```
1 map(f)
```

Chapter 4

*Tuning*

Chapter 5

*Parallelization  
Vectorization*

Chapter 6

*Algorithmic optimisations  
HW-specific optimisations*

# Chapter 6: Towards Guided Rewriting

---

- Expressing design choices directly in a functional IR:
  - **Decouples** optimisation from code generation
  - **Truncates** the search to valid implementations
  - Helps the user **drive rewriting loosely** through composable optimisations
- Chapter 6:
  - Defines **eleven rewrite points** expressing a range of low-level optimisations
  - Shows how **two convolution algorithms** are produced from **one expression**

```

1 def conv( inputData      : [[float]inChs]inW]inH r,
2     kernelsWeights : [[float]inChs]kerW]kerH]outChs r,
3     padSize       : (int,int,int,int),
4     kernelStride   : (int,int) ) : [[float]outChs]outW]outH =
5
6 toHost o
7
8 padExtra2D( kernelFission(
9     joinSplitND( padExtra2D( kernelFission( slideWindows' =>
10         kernelFission( kernelsWeights' =>
11             oclKernel( (slideWindows'', kernelsWeights'') =>
12                 abTile( abTile(
13                     tileNestedMapReduce( reduceFun = + ) {
14                         tileNestedMapReduce( reduceFun = + ) {
15                             abTile(
16                                 mapND2( slideWin: [T]inChs + kerW + kerH ) =>
17                                     map( singleK: [T]inChs + kerW + kerH ) =>
18                                         privateAccumulator(
19                                             reduce( 0, + ) o map( * ) << zip( slideWin, singleK )
20                                         ) << kernelsWeights''
21                                     ) << slideWindows'' ) ) ) )
22                             ) << ( slideWindows', kernelsWeights' )
23                         ) o map( joinND2 ) o toGPU << kernelsWeights
24                     ) ) o mapND2( joinND2 ) o
25                         slideND2( kerH, kerW, kernelStride_1, kernelStride_2 )
26                     ) o padND2( padSize, value = 0 ) o toGPU << inputData

```

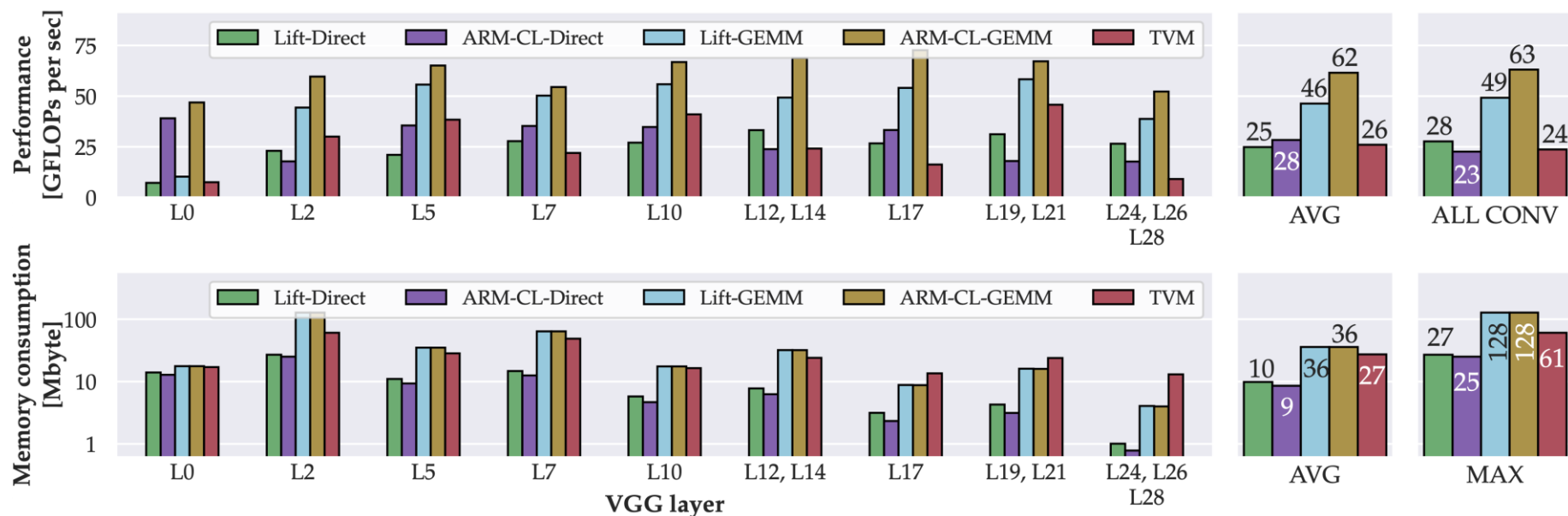
# High-performance **direct** convolution

# High-performance **im2col+GEMM** #1

## High-performance **im2col+GEMM** #2

# Results: Performance & Memory

VGG-16 on Mali-G72 GPU



- Using three Lift-generated kernels:
  - Stencil performance **on par** with ARM-CL and TVM
  - Lift-GEMM **outperforms** TVM and **0.77x** of ARM-CL-GEMM
  - Lift-GEMM uses 33% more memory than TVM on average

Low-level expression  
close to OpenCL

```
1 join o
2 transposeW o
3 mapWrg(0)(
4   mapLcl(0)(f) o
5   asVector(4)) o
6 transpose o
7 split(64)
```

Parallelised  
expression

```
1 join o
2 transposeW o
3 mapWrg(0)(
4   mapLcl(0)(f) o
5   asVector(v)) o
6 transpose o
7 split(s)
```

Structurally  
optimised  
expression

```
1 join o
2 transposeW o
3 map(
4   map(f)) o
5 transpose o
6 split(s)
```

High-level  
expression

```
1 map(f)
```

Chapter 4

***Tuning***

Chapter 5

***Parallelization  
Vectorization***

Chapter 6

***Algorithmic optimisations  
HW-specific optimisations***

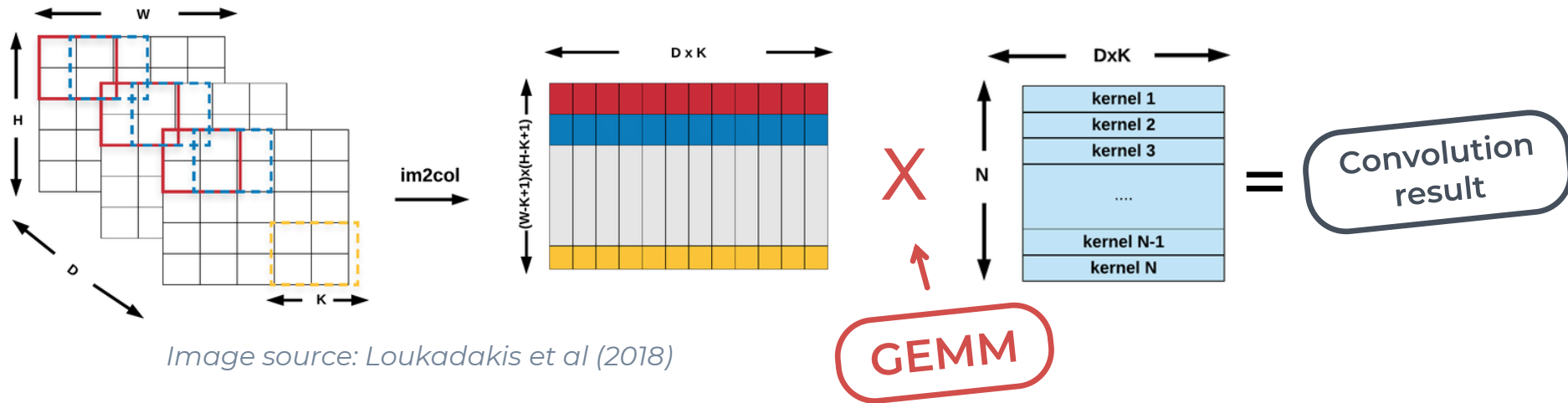


# Auto-Tuning



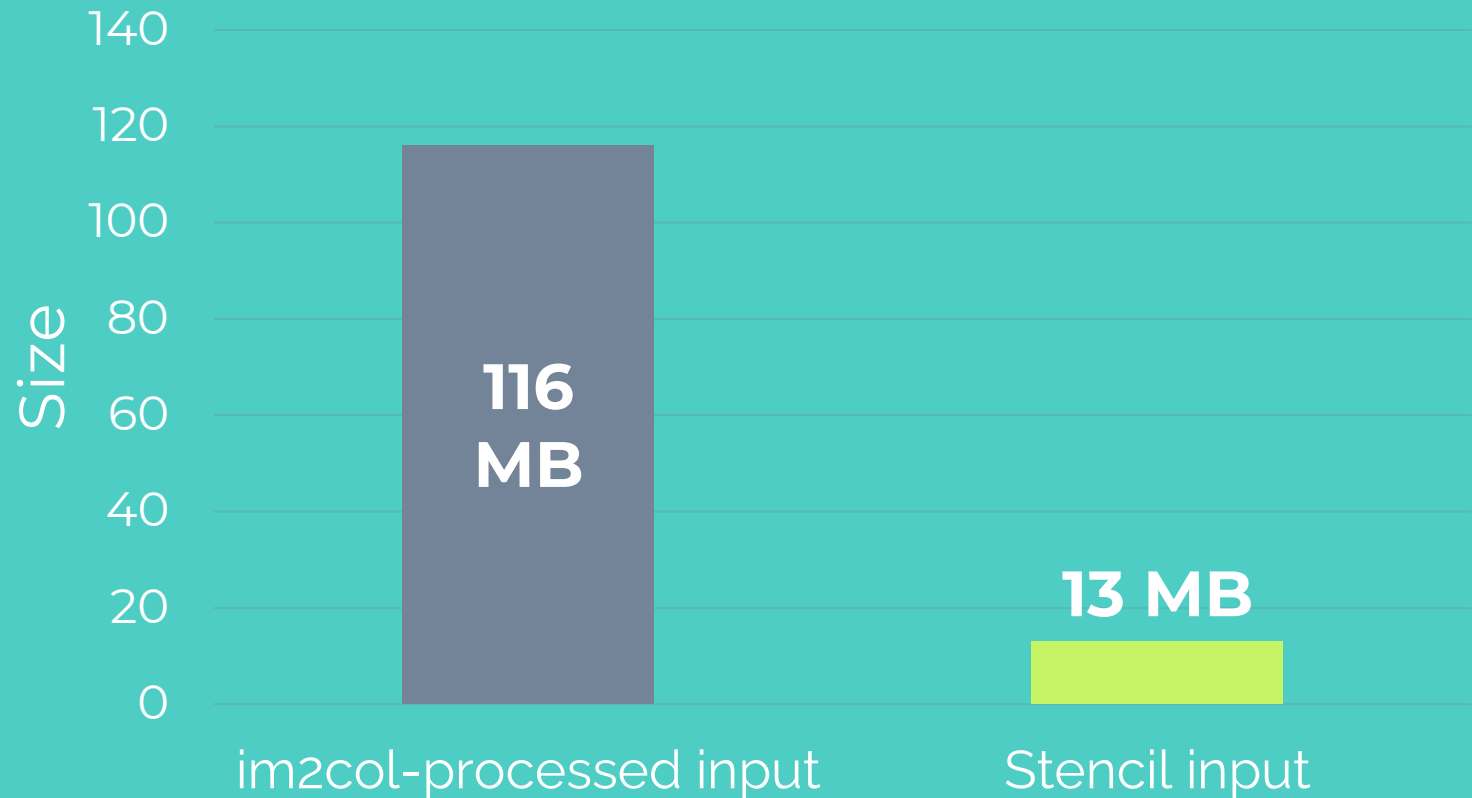
# Convolution: state of the art

- Express as General Matrix Multiplication (GEMM) using **im2col**
- High-performance libraries provide fast solutions on most devices



- im2col increases memory consumption**
  - This is a problem for resource-constrained platforms

## Input image size in the two methods (2<sup>nd</sup> layer of VGG)



# Constraints In Convolution

---

- Tuning parameter constraints **depend** on the applied rewrites
- Too many possible constraints to enumerate
  - Tile size **has to be a factor** of spatial input size
  - Vector length **has to be a factor** of a window size
  - Buffer size **must not exceed** device limit
  - ...etc

# Experimental Setup


---

- Benchmark:
  - **VGG-16**
- OpenCL kernel generation:
  - **1000 random candidates per layer**
- Comparison:
  - **Autotuned ARM Compute Library kernels**
- Platform:
  - **ARM Mali-G72 (12 cores) mobile GPU with Kirin 970 SoC**

# Parallelism Mapping

# Parallelization Is Hard

```
1  for i in 1 to I do
2    for j in 1 to J do
3      for k in 1 to K do
4        shared_buf[j][k] = f( input[i][j][k] );
5
6      for k in 1 to K do
7        for j in 1 to J do
8          output[i][j][k] = g( shared_buf[j][k] );
```



Transposed!

# Parallelization Is Hard

```
1 for i in global_id(0) to I do
2   for j in 1 to J do
3     for k in 1 to K do
4       shared_buf[j][k] = f( input[i][j][k] );
5
6   for k in 1 to K do
7     for j in 1 to J do
8       output[i][j][k] = g( shared_buf[j][k] );
```

Under-saturated cores

# Parallelization Is Hard

```
1  for i in global_id(1) to I do
2    for j in global_id(0) to J do
3      for k in 1 to K do
4        shared_buf[j][k] = f( input[i][j][k] );
5
6    for k in global_id(0) to K do
7      for j in 1 to J do
8        output[i][j][k] = g( shared_buf[j][k] );
```

Different threads

Cannot synchronize



# Parallelization Is Hard

```
1 for i in 1 to I do
2   for j in group_id(0) to J do
3     for k in local_id(0) to K do
4       shared_buf[j][k] = f( input[i][j][k] );
5
6   for k in group_id(0) to K do
7     for j in local_id(0) to J do
8       output[i][j][k] = g( shared_buf[j][k] );
```

Different thread blocks

Out-of-scope reads

# Parallelization Is Hard

```
1  for i in group_id(0) to I do
2    for j in local_id(0) to J do
3      for k in 1 to K do
4        shared_buf[j][k] = f( input[i][j][k] );
5    barrier();
6    for k in local_id(0) to K do
7      for j in 1 to J do
8        output[i][j][k] = g( shared_buf[j][k] );
```

Synchronization overhead

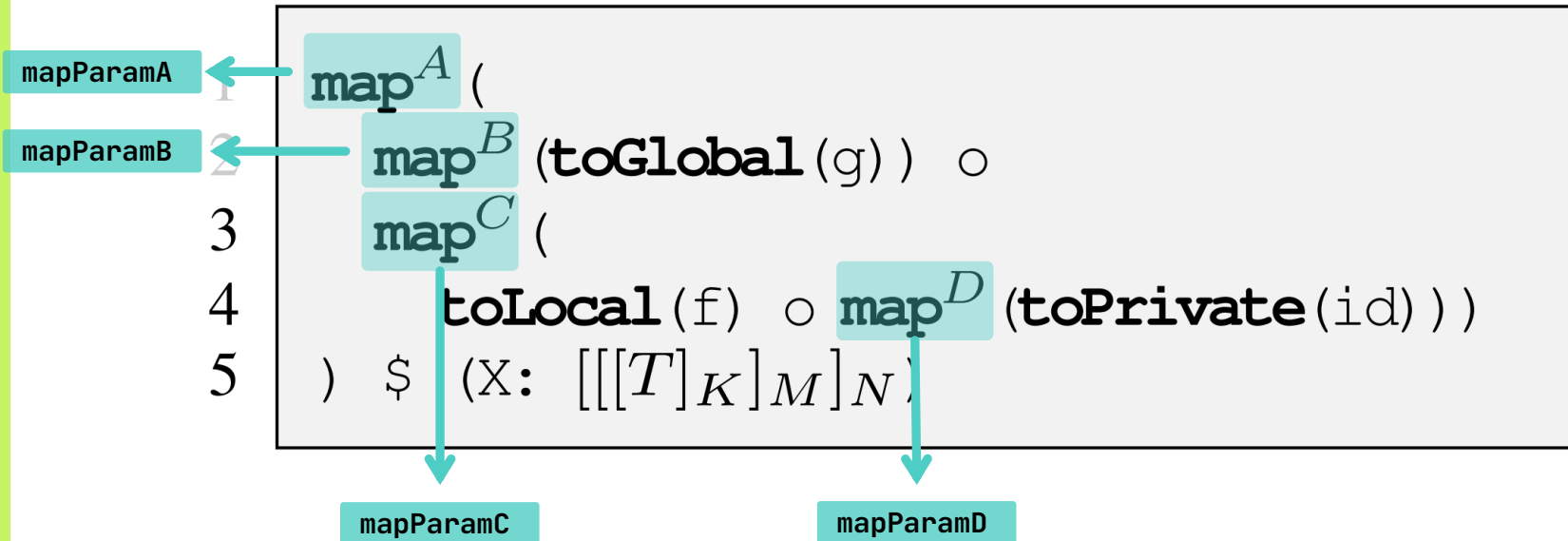
# Parallelization Is Hard

```
1  for i in group_id(0) to I do
2    for j in local_id(0) to J do
3      for k in 1 to K do
4        shared_buf[j][k] = f( input[i][j][k] );
5
6      for k in 1 to K do
7        for j in local_id(0) to J do
8          output[i][j][k] = g( shared_buf[j][k] );
```

Same indices

Might be optimal

# Parallelization: Scheduling Parameters



- Associate a **parameter** with each map
- **Encode** scheduling choices as integers
- Model parallelization restrictions as **integer constraints**

# Parallelization:

## Encoding of Choices



Code value	Map transformation
0	mapSeq
1	Fused with the outer map
10, 11, 12	mapLcl in dimension 0, 1 or 2 respectively
20, 21, 22	mapWrg in dimension 0, 1 or 2 respectively
30, 31, 32	mapGlb in dimension 0, 1 or 2 respectively

# Constraint:

## Private Memory Scope

```
1  mapA (  
2    mapB (toGlobal (g)) ○  
3    mapC (  
4      toLocal (f) ○ mapD (toPrivate (id)) )  
5  ) $ (X: [[[T]K]M]N)
```

IF: accesses private memory  
THEN: cannot be parallel



*Private memory access constraint.* maps that consume or produce private memory cannot be parallelized because private memory is restricted in scope to a single thread.

# Constraint:

## Private Memory Scope

```
1  mapA (  
2    mapB (toGlobal (g)) ○  
3    mapC (  
4      toLocal (f) ○ mapD (toPrivate (id)) )  
5  ) $ (X: [[[T]K]M]N)
```

IF: accesses private memory  
THEN: cannot be parallel



*Private memory access constraint.* **maps** that consume or produce private memory cannot be parallelized because private memory is restricted in scope to a single thread.

$\forall m : m.\text{usesPrivateMemory}$

**GEN CONSTRAINT:**  $\text{mapEncoding}(m)/10 < 1$

Code value	Map transformation
0	mapSeq
1	Fused with the outer map
10, 11, 12	mapLcl in dimension 0, 1 or 2
20, 21, 22	mapWrg in dimension 0, 1 or 2
30, 31, 32	mapGlb in dimension 0, 1 or 2

# Constraint:

## Private Memory Scope

```
1  mapA (  
2    mapB (toGlobal (g)) ○  
3    mapC (  
4      toLocal (f) ○ mapD (toPrivate (id)) )  
5  ) $ (X: [[[T]K]M]N)
```

(mapEncoding (mapParamD) / 10 < 1)



# Constraint:

## Shared Memory Scope

```
1  mapA (  
2    mapB (toGlobal(g)) ∘  
3    mapC (  
4      toLocal(f) ∘ mapD (toPrivate(id)) )  
5  ) $ (X: [[[T]K]M]N)
```

(mapEncoding (mapParamD) / 10 < 1)

map<sup>C</sup> is local or sequential

∧

map<sup>A</sup> is a mapWrg

# Constraint:

## Shared Memory Scope

```
1  mapA (  
2    mapB (toGlobal(g)) ∘  
3    mapC (  
4      toLocal(f) ∘ mapD (toPrivate(id)) )  
5  ) $ (X: [[[T]K]M]N)
```

$(\text{mapEncoding}(\text{mapParamD})/10 < 1)$

$(\text{mapEncoding}(\text{mapParamC})/10 < 2) \wedge \text{mapEncoding}(\text{mapParamA}) = 20 + w$

# Constraint:

## Shared Memory Scope

```
1  mapA (  
2    mapB (toGlobal(g)) ○  
3    mapC (  
4      toLocal(f) ○ mapD (toPrivate(id)) )  
5  ) $ (X: [[[T]K]M]N)
```

$(\text{mapEncoding}(\text{mapParamD})/10 < 1)$

$(\text{mapEncoding}(\text{mapParamC})/10 < 2) \wedge \text{mapEncoding}(\text{mapParamA}) = 20 + w$

$(\text{mapEncoding}(\text{mapParamB})/10 < 2) \wedge \text{mapEncoding}(\text{mapParamA}) = 20 + w$

# Constraint:

## Hierarchical Parallelism

```
1  mapA (  
2    mapB (toGlobal(g)) ∘  
3    mapC (  
4      toLocal(f) ∘ mapD (toPrivate(id)) )  
5  ) $ (X: [[[T]K]M]N)
```

$(\text{mapEncoding}(\text{mapParamD}) / 10 < 1)$

$(\text{mapEncoding}(\text{mapParamC}) / 10 < 2) \wedge \text{mapEncoding}(\text{mapParamA}) = 20 + w$

$(\text{mapEncoding}(\text{mapParamB}) / 10 < 2) \wedge \text{mapEncoding}(\text{mapParamA}) = 20 + w$

**map<sup>A</sup>** is not parallel

∨

**map<sup>A</sup>** and **map<sup>B</sup>** are not parallelized in the same way

# Constraint:

## Hierarchical Parallelism

```
1  mapA (  
2    mapB (toGlobal(g)) ∘  
3    mapC (  
4      toLocal(f) ∘ mapD (toPrivate(id)) )  
5  ) $ (X: [[[T]K]M]N)
```

$(\text{mapEncoding}(\text{mapParamD})/10 < 1)$

$(\text{mapEncoding}(\text{mapParamC})/10 < 2) \wedge \text{mapEncoding}(\text{mapParamA}) = 20 + w$

$(\text{mapEncoding}(\text{mapParamB})/10 < 2) \wedge \text{mapEncoding}(\text{mapParamA}) = 20 + w$

$(\text{mapEncoding}(\text{mapParamA})/10 < 1) \vee (\text{mapEncoding}(\text{mapParamA}) \neq \text{mapEncoding}(\text{mapParamB}))$

$(\text{mapEncoding}(\text{mapParamA})/10 < 1) \vee (\text{mapEncoding}(\text{mapParamA}) \neq \text{mapEncoding}(\text{mapParamC}))$

$\neg((\text{mapEncoding}(\text{mapParamA})/10 = 1) \wedge (\text{mapEncoding}(\text{mapParamB})/10 = 2) \wedge$

$(\text{mapEncoding}(\text{mapParamB}) \% 10 = \text{mapEncoding}(\text{mapParamA}) \% 10))$

**+10 more hierarchical parallelism constraints**

# Constraint Satisfaction

```

1  mapA (
2    mapB (toGlobal (g)) ∘
3    mapC (
4      toLocal (f) ∘ mapD (toPrivate (id))
5    ) $ (X: [[[T]K]M]N)

```

```

mapParamA = 20
mapParamB = 10
mapParamC = 10
mapParamD = 0

```

(mapEncoding<sub>0</sub>(mapParamD) / 10 < 1)

(mapEncoding<sub>10</sub>(mapParamC) / 10 < 2) ∧ mapEncoding<sub>20</sub>(mapParamA) = 20 + w

(mapEncoding<sub>10</sub>(mapParamB) / 10 < 2) ∧ mapEncoding<sub>20</sub>(mapParamA) = 20 + w

(mapEncoding<sub>20</sub>(mapParamA) / 10 < 1) ∨ (mapEncoding<sub>20</sub>(mapParamA) ≠ mapEncoding<sub>10</sub>(mapParamB))

(mapEncoding<sub>20</sub>(mapParamA) / 10 < 1) ∨ (mapEncoding<sub>20</sub>(mapParamA) ≠ mapEncoding<sub>10</sub>(mapParamC))

¬((mapEncoding<sub>20</sub>(mapParamA) / 10 = 1) ∧ (mapEncoding<sub>10</sub>(mapParamB) / 10 = 2) ∧

(mapEncoding<sub>10</sub>(mapParamB) % 10 = mapEncoding<sub>20</sub>(mapParamA) % 10))

+10 more hierarchical parallelism constraints

# Constraint Satisfaction

```

1  mapWrg(0) (
2    mapLcl(0) (toGlobal(g)) ○
3    mapLcl(0) (
4      toLocal(f) ○ mapSeq(toPrivate(i
5    ) $ (X: [[[T]K]M]N)
  
```

```

mapParamA = 20
mapParamB = 10
mapParamC = 10
mapParamD = 0
  
```

```

(mapEncoding0(mapParamD)/10 < 1)
(mapEncoding10(mapParamC)/10 < 2) ∧ mapEncoding20(mapParamA) = 20 + w
(mapEncoding10(mapParamB)/10 < 2) ∧ mapEncoding20(mapParamA) = 20 + w
(mapEncoding20(mapParamA)/10 < 1) ∨ (mapEncoding20(mapParamA) ≠ mapEncoding10(mapParamB))
(mapEncoding20(mapParamA)/10 < 1) ∨ (mapEncoding20(mapParamA) ≠ mapEncoding10(mapParamC))
¬((mapEncoding20(mapParamA)/10 = 1) ∧ (mapEncoding10(mapParamB)/10 = 2) ∧
  (mapEncoding10(mapParamB)%10 = mapEncoding20(mapParamA)%10))
  
```

+10 more hierarchical parallelism constraints

# Heuristics

## Sequential Map Fusion Heuristic

Perfectly nested sequential maps can always be fused to reduce search space

$\forall \text{Chain} \in \text{MapNestingChains}, \forall m1 \in \text{Chain}, \forall m2 \in \text{Chain},$   
 $m2.\text{perfectlyNestedIn}(m1)$

**GEN CONSTRAINT:**  $\neg((\text{mapEncoding}(m1) = 0) \wedge$   
 $(\text{mapEncoding}(m2) = 0))$

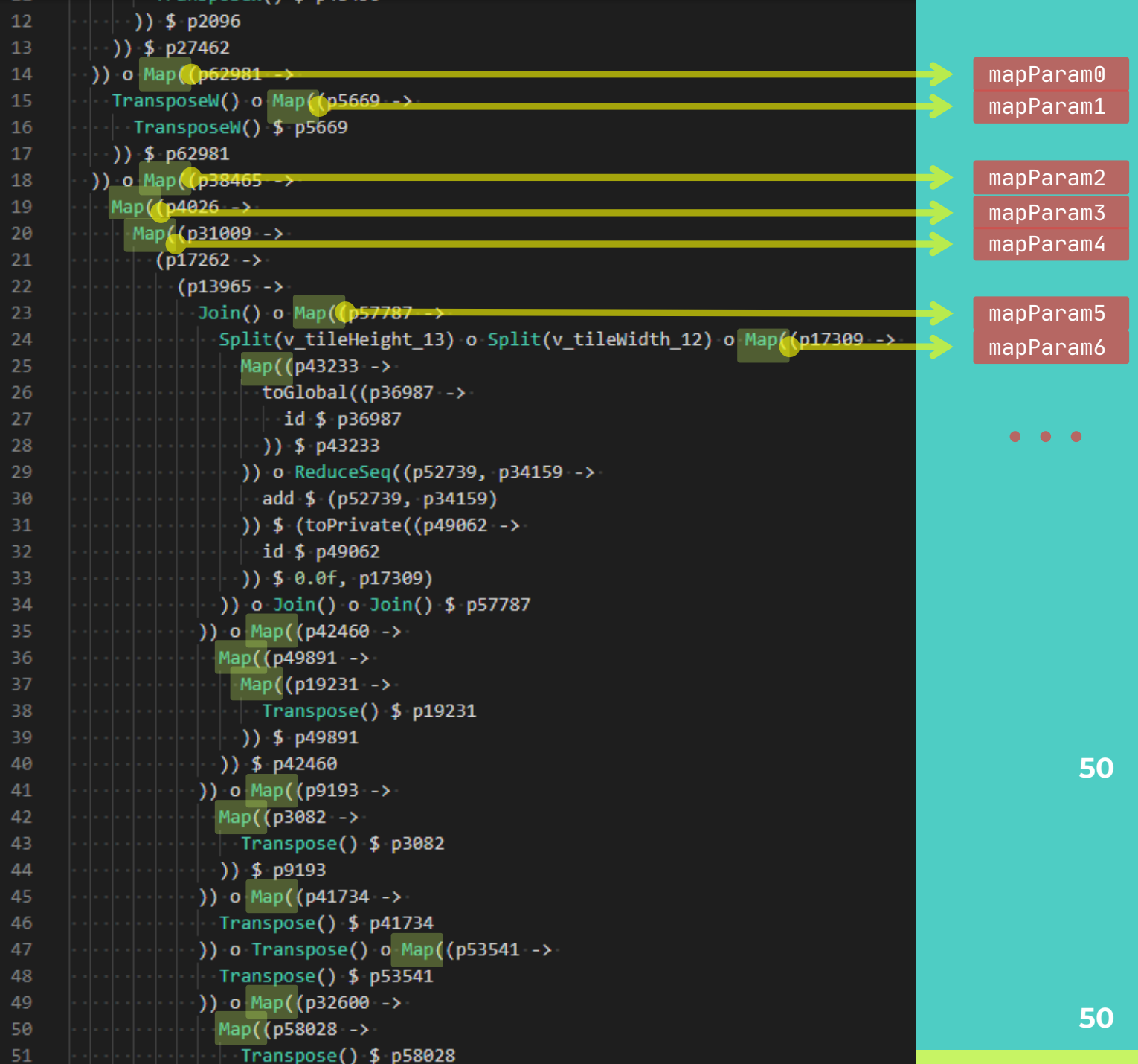
Code value	Map transformation
0	mapSeq
1	Fused with the outer map
10, 11, 12	mapLcl in dimension 0, 1 or 2 respectively
20, 21, 22	mapWrg in dimension 0, 1 or 2 respectively
30, 31, 32	mapGlb in dimension 0, 1 or 2 respectively



```

12      ...)) $ p2096
13      ...)) $ p27462
14      ...)) o Map((p62981 ->
15      ... TransposeW() o Map((p5669 ->
16      ... TransposeW() $ p5669
17      ...)) $ p62981
18      ...)) o Map((p38465 ->
19      ... Map((p4026 ->
20      ... Map((p31009 ->
21      ... (p17262 ->
22      ... (p13965 ->
23      ... Join() o Map((p57787 ->
24      ... Split(v_tileHeight_13) o Split(v_tileWidth_12) o Map((p17309 ->
25      ... Map((p43233 ->
26      ... toGlobal((p36987 ->
27      ... id $ p36987
28      ...)) $ p43233
29      ...)) o ReduceSeq((p52739, p34159 ->
30      ... add $ (p52739, p34159)
31      ...)) $ (toPrivate((p49062 ->
32      ... id $ p49062
33      ...)) $ 0.0f, p17309)
34      ...)) o Join() o Join() $ p57787
35      ...)) o Map((p42460 ->
36      ... Map((p49891 ->
37      ... Map((p19231 ->
38      ... Transpose() $ p19231
39      ...)) $ p49891
40      ...)) $ p42460
41      ...)) o Map((p9193 ->
42      ... Map((p3082 ->
43      ... Transpose() $ p3082
44      ...)) $ p9193
45      ...)) o Map((p41734 ->
46      ... Transpose() $ p41734
47      ...)) o Transpose() o Map((p53541 ->
48      ... Transpose() $ p53541
49      ...)) o Map((p32600 ->
50      ... Map((p58028 ->
51      ... Transpose() $ p58028

```



```

11      Transpose() $ p18133
12      )) $ p2096
13      )) $ p27462
14      )) o Map((p62981 ->
15      TransposeW() o Map((p5669 ->
16      TransposeW() $ p5669
17      )) $ p62981
18      )) o MapWrg(2, (p38465 ->
19      MapWrg(0, (p4026 ->
20      MapSeq((p31009 ->
21      (p17262 ->
22      (p13965 ->
23      Join() o MapLcl(2, (p57787 ->
24      Split(v_tileHeight_13) o Split(v_tilewidth_12) o MapLcl(0, (p17309
25      MapSeq((p43233 ->
26      toGlobal((p36987 ->
27      id $ p36987
28      )) $ p43233
29      )) o ReduceSeq((p52739, p34159 ->
30      add $ (p52739, p34159)
31      )) $ (toPrivate((p49062 ->
32      id $ p49062
33      )) $ 0.0f, p17309)
34      )) o Join() o Join() $ p57787
35      )) o Map((p42460 ->
36      Map((p49891 ->
37      Map((p19231 ->
38      Transpose() $ p19231
39      )) $ p49891
40      )) $ p42460
41      )) o Map((p9193 ->
42      Map((p3082 ->
43      Transpose() $ p3082
44      )) $ p9193
45      )) o Map((p41734 ->
46      Transpose() $ p41734
47      )) o Transpose() o Map((p53541 ->
48      Transpose() $ p53541
49      )) o Map((p32600 ->
50      Map((p58028 ->
51      Transpose() $ p58028

```

# Guided Rewriting

```

9  Map((p23724 ->
10  TransposeW() $ p23724
11  )) o TransposeW() o Map((p53064 ->
12  Map((p59387 ->
13  Join() $ p59387
14  )) $ p53064
15  )) o Join() o Map((p22048 ->
16  Map((p3291 ->
17  Map((p36687 ->
18  Join() $ p36687
19  )) o Join() o Map((p38111 ->
20  TransposeW() $ p38111
21  )) $ p3291
22  )) $ p22048
23  )) o Map((p14415 ->
24  TransposeW() o Map((p3478 ->
25  TransposeW() $ p3478
26  )) $ p14415
27  )) o Map((p8637 ->
28  Map((p56646 ->
29  Map((p60274 ->
30  (p44931 ->
31  (p16296 ->
32  Join() o Map((p61748 ->
33  Map((p3056 ->
34  Map((p10788 ->
35  Join() $ p10788
36  )) o Join() o Map((p2536 ->
37  TransposeW() $ p2536
38  )) $ p3056
39  )) $ p61748
40  )) o Map((p44436 ->
41  TransposeW() o Map((p27595 ->
42  TransposeW() $ p27595
43  )) $ p44436
44  )) o Map((p59960 ->
45  Map((p8566 ->
46  Map((p32269 ->
47  (p63116 ->
48  (p62040 ->
49  (p26213 ->
50  (p23403 ->
51  Join() o TransposeW() o Map((p24164 ->
52  TransposeW() $ p24164
53  )) o Map((p11838 ->
54  Map((p62211 ->
55  TransposeW() $ p62211
56  )) $ p11838
57  )) o Map((p38162 ->
58  Map((p37670 ->
59  Map((p34694 ->
60  Map((p11272 ->
61  Map((p63724 ->
62  toGlobal((p53305 ->
63
64
65  )) $ p11272
66  )) o ReduceSeq((p50867, p35848 ->
67  Map((p13698 ->
68  add $ (Get(0) $ p13698, Get(1) $ p13698)
69  )) o Zip() $ (p50867, p35848)
70  )) $ (Map((p29056 ->
71  toPrivate((p64395 ->
72  id $ p64395
73  )) $ p29056
74  )) $ 0.0f, p34694)
75  )) $ p37670
76  )) $ p38162
77  )) o Map((p11356 ->
78  Map((p61280 ->
79  Transpose() $ p61280
80  )) $ p11356
81  )) o Map((p27341 ->
82  Transpose() $ p27341
83  )) o Transpose() o Map((p3491 ->
84  p3491
85  )) o Join() o Map((p23870 ->
86  Map((p12189 ->
87  (p18636 ->
88  Join() o Map((p18649 ->
89  Map((p26353 ->
90  Map((p8803 ->
91  Map((p19946 ->
92  Map((p18234 ->
93  toLocal((p55370 ->
94  id $ p55370
95  )) $ p18234
96  )) $ p19946
97  )) $ p8803
98  )) $ p26353
99  )) $ p18649
100  )) o ReduceSeq((p38106, p46524 ->
101  Map((p25181 ->
102  Map((p31204 ->
103  Map((p3344 ->
104  Map((p40326 ->
105  add $ (Get(0) $ p40326, Get(1) $ p40326)
106  )) $ p3344
107  )) $ p31204
108  )) $ p25181
109  )) o Split(v_seqWindowsPerThreadY_21) o
110  Split(v_seqWindowsPerThreadX_20) o Split(1) o
111  Zip() $ (Join() o Join() o Join() $ p38106,
112  Join() o Join() o Join() o Map((p54992 ->
113  Map((p7921 ->
114  Map((p750 ->
115  Join() $ p750
116  )) o Join() o Map((p23842 ->
117  TransposeW() $ p23842
118  )) $ p7921
119  )) $ p54992
120  )) o Map((p20319 ->
121  TransposeW() o Map((p8345 ->
122
123  )) o Map((p28717 ->
124  Map((p2203 ->
125  Map((p257 ->
126  (p38712 ->
127  (p53378 ->
128  Map((p52451 ->
129  Map((p26784 ->
130  Map((p58124 ->
131  Map((p59290, p5770 ->
132  add $ (p59290, mult $ (
133  Get(0) $ p5770, Get(1) $ p5770))
134  )) $ (toPrivate((p9979 ->
135  id $ p9979
136  )) $ 0.0f, Zip() $ (Join() o Join() $
137  p58124, Join() o Join() $ p52451))
138  )) $ p26784
139  )) $ p53378
140  )) $ p38712
141  )) o Map((p7268 ->
142  Map((p15357 ->
143  Transpose() o Map((p5313 ->
144  Transpose() $ p5313
145  )) $ p15357
146  )) $ p7268
147  )) o Map((p16531 ->
148  Transpose() o Map((p10152 ->
149  Transpose() $ p10152
150  )) $ p16531
151  )) o Map((p42971 ->
152  Map((p48512 ->
153  Map((p26122 ->
154  Map((p52770 ->
155  Map((p42776 ->
156  toPrivate((p60560 ->
157  id $ p60560
158  )) $ p42776
159  )) $ p52770
160  )) $ p26122
161  )) $ p48512
162  )) $ p42971
163  )) o Map((p54387 ->
164  Map((p36304 ->
165  Transpose() o Map((p30379 ->
166  Transpose() $ p30379
167  )) $ p36304
168  )) $ p54387
169  )) o Map((p37721 ->
170  Transpose() o Map((p36420 ->
171  Transpose() $ p36420
172  )) $ p37721
173  )) $ p257
174  )) o Map((p57134 ->
175  Transpose() o Map((p10654 ->
176  Transpose() $ p10654
177  )) $ p57134
178  )) o Map((p58070 ->

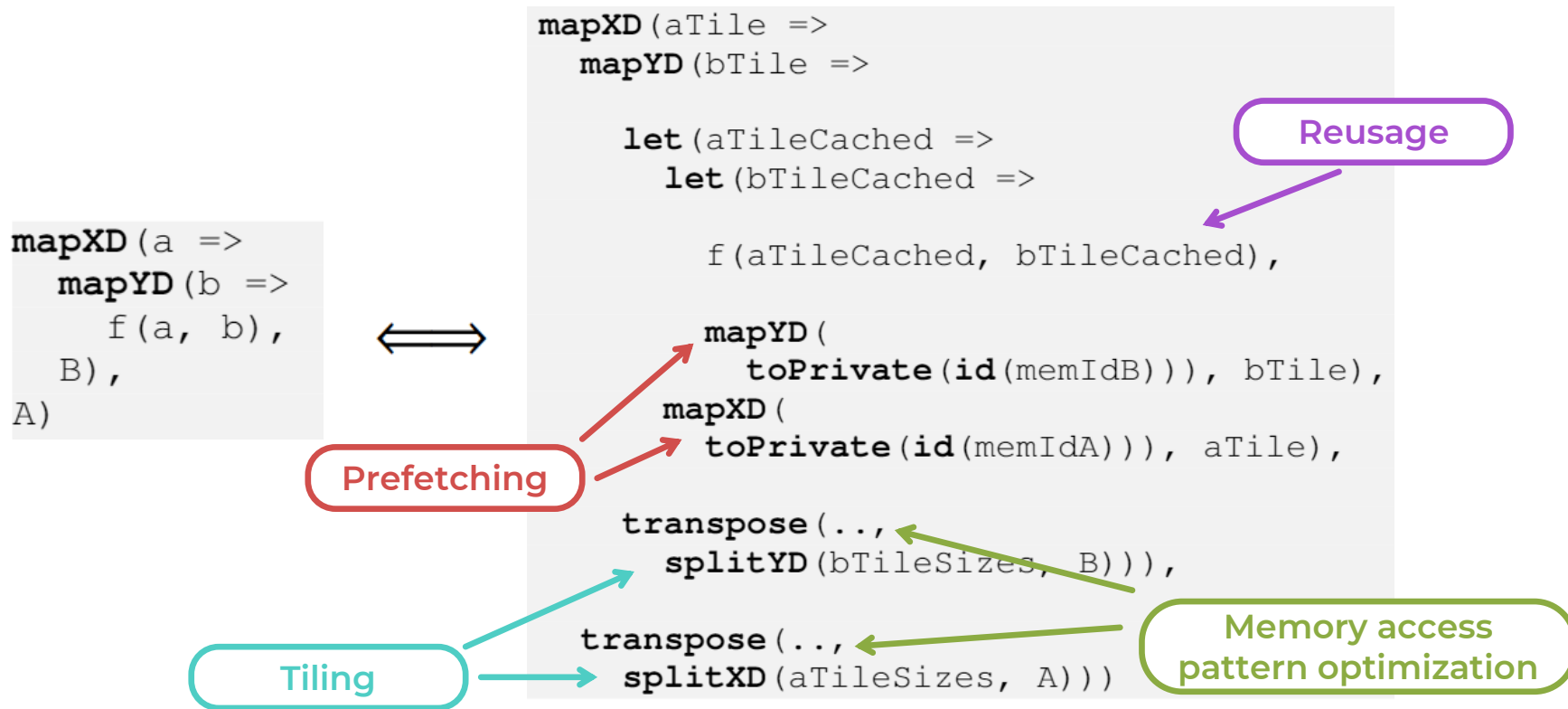
```

# Rewrite points expanded

x3

# Guided Rewriting: Tiling

- When two nested ND-maps iterate over different buffers, tiles can be prefetched and reused



# Guided Rewriting: Tiling

- When two nested ND-maps iterate over different buffers, tiles can be prefetched and reused

