

EPSRC Centre for Doctoral Training in Pervasive Parallelism

Mapping Parallelism in a Functional IR through Constraint Satisfaction

A Case Study on Convolution for Mobile GPUs

Artifact: <u>https://gitlab.com/naummo/liftpar-cc-2022-artifact</u>

Naums Mogers, Lu Li,

Valentin Radu, Christophe Dubach

April 6th, 2022

Parallel Architectures

Parallel architectures are hard to optimize for



©NVIDIA

©Xilinx



Parallel Architectures

Parallel architectures are hard to optimize for

GPU device Global memory HOST Work group (core) (x, y, z) Work group (core) (x, y, z) Local (shared) mem Local (shared) mem 3 Thread Thread Thread Thread Thread Thread Thread Thread Private **Private Private** Private **Private Private Private** Private mem mem mem mem mem mem mem mem

CARM

©NVIDIA

©Xilinx















Parallelization: Approaches

- User-provided schedules
- Polyhedral compilation + heuristics
- Functional IR (Accelerate, Futhark) + heuristics
- Proposed approach:
 - Functional IR
 - Generate parallelization constraints
 - Explore the design space automatically

Avoid evaluating invalid parallel mappings













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



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



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.

	Code value	Map transformation
∀m: m.usesPrivateMemory	0	mapSeq
Ş	1	Fused with the outer map
GEN CONSTRAINT: mapEncoding(m)/10 < 1	10, 11, 12	<pre>mapLcl in dimension 0, 1 or 2</pre>
	20, 21, 22	mapWrg in dimension 0, 1 or 2
	30, 31, 32	mapGlb in dimension 0, 1 or 2

Constraint: Private Memory Scope



(mapEncoding(mapParamD)/10 < 1)

Constraint: Shared Memory Scope





Constraint: Shared Memory Scope



(mapEncoding(mapParamD)/10 < 1) (mapEncoding(mapParamC)/10 < 2) \land mapEncoding(mapParamA) = 20 + w

Constraint: Shared Memory Scope



(mapEncoding(mapParamD)/10 < 1) (mapEncoding(mapParamC)/10 < 2) \land mapEncoding(mapParamA) = 20 + w(mapEncoding(mapParamB)/10 < 2) \land mapEncoding(mapParamA) = 20 + w

Constraint: Hierarchical Parallelism





Constraint: Hierarchical Parallelism

1
$$\operatorname{map}^{A}($$

2 $\operatorname{map}^{B}(\operatorname{toGlobal}(g)) \circ$
3 $\operatorname{map}^{C}($
4 $\operatorname{toLocal}(f) \circ \operatorname{map}^{D}(\operatorname{toPrivate}(\operatorname{id})))$
5) \$ (X: $[[[T]_{K}]_{M}]_{N})$

 $\begin{array}{l} \left(\texttt{mapEncoding(mapParamD)} / 10 < 1 \right) \\ \left(\texttt{mapEncoding(mapParamC)} / 10 < 2 \right) \land \texttt{mapEncoding(mapParamA)} = 20 + \texttt{w} \\ \left(\texttt{mapEncoding(mapParamB)} / 10 < 2 \right) \land \texttt{mapEncoding(mapParamA)} = 20 + \texttt{w} \\ \left(\texttt{mapEncoding(mapParamA)} / 10 < 1 \right) \lor \left(\texttt{mapEncoding(mapParamA)} \neq \texttt{mapEncoding(mapParamB)} \right) \\ \left(\texttt{mapEncoding(mapParamA)} / 10 < 1 \right) \lor \left(\texttt{mapEncoding(mapParamA)} \neq \texttt{mapEncoding(mapParamC)} \right) \\ \neg \left((\texttt{mapEncoding(mapParamA)} / 10 = 1 \right) \land \left(\texttt{mapEncoding(mapParamB)} / 10 = 2 \right) \land \\ \left(\texttt{mapEncoding(mapParamA)} / 10 = \texttt{mapEncoding(mapParamA)} \right) \\ \end{array}$

+10 more hierarchical parallelism constraints

Constraint Satisfaction



+10 more hierarchical parallelism constraints

Constraint Satisfaction



(mapEncodin@(mapParamD)/10 < 1) (mapEncodin10mapParamC)/10 < 2) ∧ mapEncodin20mapParamA) = 20 + w (mapEncodin10mapParamB)/10 < 2) ∧ mapEncodin20mapParamA) = 20 + w (mapEncodir20mapParamA)/10 < 1) ∨ (mapEncodin20mapParamA) ≠ mapEncodin10mapParamB)) (mapEncodir20mapParamA)/10 < 1) ∨ (mapEncodin20mapParamA) ≠ mapEncodin10mapParamB)) ¬((mapEncodin20mapParamA)/10 = 1) ∧ (mapEncodin10mapParamB)/10 = 2) ∧ (mapEncodin10mapParamB)%10 = mapEncodin20mapParamA)%10))

+10 more hierarchical parallelism constraints

Heuristics

Sequential Map Fusion Heuristic

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

∀Chain ∈ MapNestingChains, ∀m1 ∈ Chain, ∀m2 ∈ Chain, m2.perfectlyNestedIn(m1) GEN CONSTRAINT: ¬((mapEncoding(m1) = 0) ∧ (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

Convolution in Lift

1	<pre>def conv(in: [[[T]_{inChs}]_{inW}]_{inH},</pre>
2	ks: $[[[T]_{inChs}]_{kerW}]_{kerH}]_{outChs}$,
3	kerStepX: int, kerStepY: int
4) : $[[T]_{outChs}]_{outW}]_{outH} =$
5	$mapND_2(slideWin: [[[T]_{inChs}]_{kerW}]_{kerH} \Rightarrow$
6	$map(singleK: [[T]_{inChs}]_{kerW}]_{kerH} \Rightarrow$
7	<pre>reduce(+, 0, map(*,</pre>
8	<pre>joinND₂(zipND₃(slideWin, singleK)))),</pre>
9	ks), slideND 2(kerH, kerStepY, kerW, kerStepX, in))

	Map((p23724 ->
	TransposeW() \$ p23724
)) o TransposeW() o Map((p53064 ->
2	Map((p59387 ->
	Join()-\$-p59387
	··)) \$ n53064
)) $\alpha = \ln(1) \alpha = Man((n22048) - 2)$
6	- Man((n3201
7	Map((p36687 ->
	$hap((poool) \neq p$
	(p^{2})
)) \$ h3291
)) 0 Map((p14415 ->)
4	Transposew() 0 Map((p34/8 ->
	ransposeW() \$ p34/8
6)) \$ p14415
)) o Map((p8637 ->
	Map((p56646 ->
	Map((p60274 ->
	·····(p44931 · -> ·
	······································
2	Join() o Map((p61748 ->
	······································
	••••••••••••••••••••••••••••••••••••••
5	Join() \$ p10788
	<pre>>></pre>
	······································
))\$_n3056
g	() () () () () () () () () () () () () (
0 0	(0.44436)
	(p27505 -)
2)) \$ p44426
)) \$ \p44430
4 =)) 0 Map((p39900 - 2
	Map((p8000 ->
D 	(
	(pb311b ->>
	(p62040 ->
	(p26213 ->
	(p23403 ->
	Join() o TransposeW() o Map((p24164
2	······································
	······································
	••••••••••••••••••••••••••••••••••••••
	••••••••••••••••••••••••••••••••••••••
)) \$ p11838
	(p38162 ->
	Map((p34694 ->
0	Man((n11272 ->
	Map((p112)2 ->
	toGlobal((p52205>
2	id \$ p53205
)) \$ p3300
)) \$ p05/24

• • • • • • • •)) \$ p11272
)) o ReduceSeq((p50867, p35848 ->
• • • • • • • • Map((p13698 -> •
add \$ (Get(0) \$ p13698, Get(1) \$ p13698)
)) o Zip() \$ (p50867, p35848)
····)) \$ (Map((p29056 -> ·
••••••toPrivate((p64395•->•
·····id \$ p64395
····)) \$ p29056
····)) \$ 0.0f, p34694)
····)) \$ p37670
)) \$ p38162
)) o Map((p11356 ->
Map((p61280 ->
· · · · Transpose() · \$ · p61280
)) \$ p11356
)) o Map((p27341 ->
Transpose() \$ p27341
)) o Transpose() o Map((p3491 ->
p3491
)) o Join() o Map((p23870 ->
Map((p12189 ->
(p18636 ->
Join() o Map((p18649 ->
Map((D8803 -> -
Map((p19946>-
tol ocal ((n55370 ->
id \$-055370
)) \$ n18234
))\$ n10046
)) \$ pp340
)) \$ p0003
)) \$ 20000
)) \$ P18049
)) 0 ReduceSeq((p38100, p40524 ->
Map((p25181 ->)
Map((p31204 - 2))
$\operatorname{Map}((p3344) \rightarrow p)$
Map((p4032b->
add \$*(Get(0) \$*p40326, Get(1) \$*p40326)
······································
···························)) · \$-p31204
(
)) o Split(v_seqWindowsPerThreadY_21) o
Split(v_seqWindowsPerThreadX_20) o Split(1) o
Zip() \$ (Join() o Join() o Join() \$ p38106,
Join() o Join() o Join() o Join() o Map((p54992 -
Map((p7921 ->
Map((p750->
<pre>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>></pre>
<pre></pre>
·····)) \$ p7921
·····)) \$ p54992
)) o Map((p20319 ->
<pre></pre>
TransposeW() \$ p9345
()) \$ p20319
(n/28717 -)

)) o Map((p28717 ->
Map((p2203 ->
Map((p257 ->
(p38712 ->
(p53378 ->
• • • • • • • • • • • • • Map((p52451 • -> •
······································
······································
<pre></pre>
add \$ (p59290, mult \$ (
<pre> • • • • • • • • • • • • • • • • • • •</pre>
<pre></pre>
••••••••••••••••••••••••••••••••••••••
<pre>''''''''''''''''''''''''''''''''''''</pre>
<pre>p58124, Join() o Join() \$ p52451))</pre>
·····················))·\$ p26784
···············))·\$·p53378
·····))·\$·p38712
·····)·o·Map((p7268·->·
••••••••••••••••••••••••••••••••••••••
<pre> ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' '</pre>
······································
·····················))·\$·p15357
· · · · · · · · · ·)) · \$ · p7268
····)) o Map((p16531 ->
Transpose() o Map((p10152 ->
· · · · · · · · · · · · · Transpose() · \$ · p10152
· · · · · · · · · · ·)) · \$ · p16531
·····)) o Map((p42971 ->
Map((p48512 ->
Map((p26122 ->
Map((p52770>-
Map((p42//6·->·
toPrivate((p60560 ->
)) \$ P427/b
)) \$ µ52770
)) \$ µ20122
)) \$ 040512
)) p 429/1
)) 0 Piap((p34367 - 7
Transpose() \$ n30379
)) \$-n36304
)) \$ n54387
)) o Man((n37721 ->
Transnose() = Man((n36420 - >
Transpose() \$ n36420
······)) \$ p37721
()) \$ p257
) o Map((p57134 ->
Transpose() o Map((p10654 ->
········Transpose()·\$·p10654
)) \$ p57134
····)) · o · Map((p58070 · -> ·
Map((p34527 ->
••••••••••••••••••••••••••••••••••••••

- \$

12)) \$ p2096
13)) <u>\$ p27462</u>
14)) o Map((p62981 ->
15	TransposeW() o Map((p5669 ->
16	TransposeW() \$-p5669
17)) \$ p62981
18)) o Map((p38465 ->
19	Map((p4026 ->
20	Мар((р31009 ->
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	
27	id \$ p36987
28	••••••••••••••••••••••••••••••••••••••
29	••••••••••••••••••••••••••••••••••••••
30	•••••••••••••••••add \$•(p52739, p34159)
31	••••••••••••••)) \$ (toPrivate((p49062>•
32	••••••••••••••••id \$ p49062
33	•••••••••••••••)) \$ 0.0f, p17309)
34	••••••••••••••••••••••••••••••••••••••
35))_o_Map((p42460 ->
36	Map((p49891 ->
37	·····································
38	Transpose() \$ p19231
39	• • • • • • • • • • • • • • • • • • •
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	Iranspose() \$ p53541
49)) 0 Map((p32600 ->
50	$\operatorname{Map}((p58028 -)$
51	Transpose() \$ p58028

->

12)) \$ p2096		
13)) <u>\$ p27462</u>		
14)) o Map ((p62981 ->		mapParam0
15	TransposeW() o Map((p5669 ->		mapParam1
16	TransposeW() \$-p5669		
17	····)) \$ p62981		
18)) o Map((p38465 ->		mapParam2
19	Map((p4026 ->		mapParam3
20	Map <mark>((</mark> p31009 ->	LS I	manParam4
21	(p17262 ->		
22	(p13965 ->		
23	Join() o Map((p57787 ->		mapParam5
24	Split(v_tileHeight_13) o Split(v_tileWidth_12) o Map((p17309 ->_		mapParam6
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	• • • • • • • • • • • • • • • • • • •		
35)) o Map((p42460 ->		
36	Map((p49891 ->		
37	Map((p19231 ->		
38	Transpose() \$ p19231		
39)) \$ p49891		
40)) \$ -p42460		32
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 ->		72
50	Map((p58028 ->		52
51	Transpose() \$ p58028		

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,(p1730)
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	
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

Results: Exploration Efficiency



Results: Exploration Efficiency



Peak performance after 95 minutes

Results: Exploration Efficiency



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

Results: Performance



Results: Performance



Results: Performance



Results: Performance & Memory



Results: Performance & Memory

VGG-16 on Mali-G72 GPU



3.6x less memory on average than ARM-CL GEMM

Results: Performance & Memory



- 3.6x less memory on average than ARM-CL GEMM
- **2.7x** less memory on average than TVM

Conclusions



- Automatic parallelization is too complex to navigate using hand-coded heuristics
 - Deeply nested scheduling model on GPUs
- Constraints help prune the search space

We are Open Source!

Project: www.lift-project.org

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

Artifact: https://gitlab.com/naummo/liftpar-cc-2022-artifact

Naums Mogers University of Edinburgh, UK naums.mogers@ed.ac.uk

Valentin Radu University of Sheffield, UK v.radu@sheffield.ac.uk Lu Li University of Edinburgh, UK lu.li@ed.ac.uk

Christophe Dubach McGill University, Canada christophe.dubach@mcgill.ca



EPSRC Centre for Doctoral Training in Pervasive Parallelism