# Locality-Aware Automatic Parallelization for GPGPU with OpenHMPP Directives

José M. Andión, Manuel Arenaz, François Bodin, Gabriel Rodríguez and Juan Touriño

7th International Symposium on High-Level Parallel Programming and Applications (HLPP 2014) July 3-4, 2014 — Amsterdam, Netherlands



- Motivation: General Purpose Computation with GPUs
- GPGPU with CUDA & OpenHMPP
- The KIR: an IR for the Detection of Parallelism
- Locality-Aware Generation of Efficient GPGPU Code

CORUÑA UNIVERSIDADE DA CORUÑA

- Case Studies: CONV3D & SGEMM
- Performance Evaluation
- Conclusions & Future Work

- Motivation: General Purpose Computation with GPUs
- GPGPU with CUDA & OpenHMPP
- The KIR: an IR for the Detection of Parallelism
- Locality-Aware Generation of Efficient GPGPU Code

UNIVERSIDADE DA CORUÑA

- Case Studies: CONV3D & SGEMM
- Performance Evaluation
- Conclusions & Future Work



#### The Parallel Challenge

David A. Patterson and John L. Hennessy. Computer Organization and Design: The Hardware/Software Interface. Elsevier, 2014.

IITECTURE GROUP

UNIVERSIDADE DA CORUÑA

#### **ACCELERATORS/CO-PROCESSORS**



- Motivation: General Purpose Computation with GPUs
- GPGPU with CUDA & OpenHMPP
- The KIR: an IR for the Detection of Parallelism
- Locality-Aware Generation of Efficient GPGPU Code

UNIVERSIDADE DA CORUÑA

- Case Studies: CONV3D & SGEMM
- Performance Evaluation
- Conclusions & Future Work

#### GPGPU with CUDA

• First GPGPU programs look like graphics applications





< UNIVERSIDADE DA CORUÑA

CUDA enables the use of C

CUDA kernel: specifies the operation of a single GPU thread

• Main ideas:

1.Lightweight parallel threads in hierarchy: grid, block

2.Shared memory

3.Barriers

#### GPU Programming Features in CUDA

1	Threadification					
2	Thread grouping: warps					
3	Minimization of CPU-GPU data transfers					
4	Coalescing					
5	Maximization of the usage of registers and shared memory					
6	Divergency					
7	Occupancy					
8	Threads per block					

e universidade da coruña

#### GPGPU with OpenHMPP

- Directive-based approaches provide several advantages:
  - More readable codes
  - Only one file
  - Independent from the hardware platform
  - Reasonable performance









#### GPGPU with OpenHMPP

- Directive-based approaches provide several advantages:
  - More readable codes
  - Only one file
  - Independent from the hardware platform
  - Reasonable performance







#### GPGPU with OpenHMPP

- Directive-based approaches provide several advantages:
  - More readable codes
  - Only one file
  - Independent from the hardware platform
  - Reasonable performance





🗲 UNIVERSIDADE DA CORUÑA

GAC COMPUTER ARCHITECTURE GROUP UNIVERSITY OF A CORUNA

#### GPU Programming Features with OpenHMPP

1	Threadification						
2	Thread grouping						
3	Minimization of CPU-GPU data transfers						
4	Coalescing						
5	Maximization of the usage of registers and shared memory						
6	Divergency gridify advancedLoad						
7	OccupancydelegatedStore permute unroll fuse						
8	Threads per block tile shared						

e universidade da coruña

- Motivation: General Purpose Computation with GPUs
- GPGPU with CUDA & OpenHMPP
- The KIR: an IR for the Detection of Parallelism
- Locality-Aware Generation of Efficient GPGPU Code

UNIVERSIDADE DA CORUÑA

- Case Studies: CONV3D & SGEMM
- Performance Evaluation
- Conclusions & Future Work

#### diKernel: Domain-Independent Computational Kernel



- Characterizes the computations carried out in a program without being affected by how they are coded
- •Exposes multiple levels of parallelism

M. Arenaz et al. XARK: An Extensible Framework for Automatic Recognition of Computational Kernels. ACM Transactions on Programming Languages and Systems, 30(6), 2008.

COMPUTER ARCHITECTURE GROUP UNIVERSITY OF A CORUNA GAC.UDC.ES

🗲 UNIVERSIDADE DA CORUÑA

### Building the KIR

• Non-statement based, high-level, hierarchical IR

1.diKernel recognition on the DDG

2. Identification of flow dependences

**3.**Hierarchy of execution scopes reflecting the computational stages & diKernel classification

J.M. Andión et al. A Novel Compiler Support for Automatic Parallelization on Multicore Systems. Parallel Computing, 39(9), 2013.

< UNIVERSIDADE DA CORUÑA

shaded to be omitted in the discovering of parallelism

#### Example of KIR: CONV3D

```
1.int i, j, k, size x, size y, size z;
2.float coefx, coefy, coefz, *input, *output;
3.
4.for (i = 0; i < size x; i++) {
5. for (j = 0; j < size y; j++) {
6. for (k = 0; k < size z; k++) {
7.
    float tempx = input[i][j][k]+coefx*
8.
9.
       input[i-1][j][k]+input[i+1][j][k]+
10.
       input[i-2][j][k]+input[i+2][j][k]+
11.
       input[i-3][j][k]+input[i+3][j][k]+
12.
       input[i-4][j][k]+input[i+4][j][k]
13.
      );
     float tempy = input[i][j][k]+coefy*
14.
15.
16.
       input[i][j-1][k]+input[i][j+1][k]+
17.
       input[i][j-2][k]+input[i][j+2][k]+
18.
       input[i][j-3][k]+input[i][j+3][k]+
19.
       input[i][j-4][k]+input[i][j+4][k]
20.
      );
21.
     float tempz = input[i][j][k]+coefz*
22.
23.
       input[i][j][k-1]+input[i][j][k+1]+
24.
       input[i][j][k-2]+input[i][j][k+2]+
25.
       input[i][j][k-3]+input[i][j][k+3]+
26.
       input[i][j][k-4]+input[i][j][k+4]
27.
      );
28.
     output[i][j][k] =
29.
      output[i][j][k]+tempx+tempy+tempz;
30. }
31. }
32.}
```



C UNIVERSIDADE DA CORUÑA

- Motivation: General Purpose Computation with GPUs
- GPGPU with CUDA & OpenHMPP
- The KIR: an IR for the Detection of Parallelism
- Locality-Aware Generation of Efficient GPGPU Code

CORUÑA UNIVERSIDADE DA CORUÑA

- Case Studies: CONV3D & SGEMM
- Performance Evaluation
- Conclusions & Future Work

#### GPU Programming Features addressed by our Automatic Technique **Threadification** 1 2 Thread grouping 3 Minimization of CPU-GPU data transfers Coalescing 4 5 Maximization of the usage of registers and shared memory 6 Divergency 7 Occupancy 8 Threads per block



## Detection of coalesced accesses to the GPU global memory



< UNIVERSIDADE DA CORUÑA

## Detection of coalesced accesses to the GPU global memory



< UNIVERSIDADE DA CORUÑA

## Locality-Aware Generation of Efficient GPGPU Code (and III)

- 2.Usage of registers to store reused data within a GPU thread
- **3.**Usage of the GPU shared memory for data shared between the threads of a warp
- 4. Increase the computational load of a GPU thread (loop tiling preserving coalescing)



- Motivation: General Purpose Computation with GPUs
- GPGPU with CUDA & OpenHMPP
- The KIR: an IR for the Detection of Parallelism
- Locality-Aware Generation of Efficient GPGPU Code

CORUÑA UNIVERSIDADE DA CORUÑA

- Case Studies: CONV3D & SGEMM
- Performance Evaluation
- Conclusions & Future Work

shaded to be omitted in the discovering of parallelism

### Case Study: CONV3D (I)

```
1.int i, j, k, size x, size y, size z;
2.float coefx, coefy, coefz, *input, *output;
3.
4.for (i = 0; i < size x; i++) {
5. for (j = 0; j < size y; j++) {
6. for (k = 0; k < size z; k++) {
7.
    float tempx = input[i][j][k]+coefx*
8.
9.
       input[i-1][j][k]+input[i+1][j][k]+
10.
       input[i-2][j][k]+input[i+2][j][k]+
11.
       input[i-3][j][k]+input[i+3][j][k]+
12.
       input[i-4][j][k]+input[i+4][j][k]
13.
      );
14.
     float tempy = input[i][j][k]+coefy*
15.
16.
       input[i][j-1][k]+input[i][j+1][k]+
17.
       input[i][j-2][k]+input[i][j+2][k]+
18.
       input[i][j-3][k]+input[i][j+3][k]+
19.
       input[i][j-4][k]+input[i][j+4][k]
20.
      );
21.
     float tempz = input[i][j][k]+coefz*
22.
23.
       input[i][j][k-1]+input[i][j][k+1]+
24.
       input[i][j][k-2]+input[i][j][k+2]+
25.
       input[i][j][k-3]+input[i][j][k+3]+
26.
       input[i][j][k-4]+input[i][j][k+4]
27.
      );
     output[i][j][k] =
28.
29.
      output[i][j][k]+tempx+tempy+tempz;
30. }
31. }
32.}
```



C UNIVERSIDADE DA CORUÑA

### Case Study: CONV3D (II)

- conv3d-cpu
- conv3d-hmpp1: Coalescing

1.int i, j, k, size\_x, size\_y, size\_z; 2.float coefx,coefy,coefz,\*input,\*output; 3. 4.for (i = 0; i < size\_x; i++) { 5. for (j = 0; j < size\_y; j++) { 6. for (k = 0; k < size\_z; k++) { 7. float tempx = input[i][j][k]+coefx\* ...

Default OpenHMPP policy

CHRECS\_input1<sup>T0</sup> =
[{0}][{0}][{0,+,1}]

CHRECS\_input<sub>1</sub><sup>T1</sup> =
[{0}][{1}][{0,+,1}]

< UNIVERSIDADE DA CORUÑA

Loop nest is permuted to forj, fork, fori (permute directive)



#### Case Study: CONV3D (III)

#### conv3d-hmpp2: Registers

```
4.for (i = 0; i < size_x; i++) {
5. for (j = 0; j < size_y; j++) {
6. for (k = 0; k < size_z; k++) {
7. float tempx = input[i][j][k]+coefx*
8. (
9. input[i-1][j][k]+input[i+1][j][k]+
...</pre>
```



J.M. Andión et al. Locality-Aware Automatic Parallelization for GPGPU with OpenHMPP Directives. HLPP 2014.

tives. HLPP 2014. GAC UNIVERSITY OF A CORUNA

< UNIVERSIDADE DA CORUÑA

#### Case Study: CONV3D (and IV)

conv3d-hmpp3: Shared memory

$4 \text{ for } (i = 0 \cdot i \leq size x \cdot i + 1) $			TO		<i>T1</i>	
5. for (j = 0; j < size_y; j++) {		1 <sup>st</sup> dim	2 <sup>nd</sup> dim	$3^{rd}dim \mid 1^{st}dim$	2 <sup>nd</sup> dim	3 <sup>rd</sup> dim
6. for (k = 0; k < size_z; k++) {	CHRECS_input <sub>19</sub>	$  \{0,+,1\}$	$\{0\}$	$\{0\} \mid \{0,+,1\}$	$\{0\}$	{1} {0}
<pre>21. float tempz = input[i][j][k]+coefz* 22. (</pre>	CHRECS_input <sub>20</sub> CHRECS_input <sub>21</sub>	$\{0, +, 1\}$ $\{0, +, 1\}$	$\{0\}$ $\{0\}$	$ \begin{cases} -1 \\ \{1\} \\ \{0, +, 1\} \\ \{0, +, 1\} \end{cases} $	$\{0\}$ $\{0\}$	{0} {2}
<pre>23. input[i][j][k-1]+input[i][j][k+1]+ 24. input[i][j][k-2]+input[i][j][k+2]+</pre>	CHRECS_input <sub>22</sub> CHRECS_input <sub>23</sub>	$ \begin{array}{c c} \{0,+,1\} \\ \{0,+,1\} \end{array} $	$\{0\}\ \{0\}$	$ \begin{array}{c c} \{-2\} \\ \{2\} \\ \end{array} \begin{array}{c} \{0,+,1\} \\ \{0,+,1\} \end{array} $	$\{0\}$ $\{0\}$	$\{-1\}$ $\{3\}$
<pre>25. input[i][j][k-3]+input[i][j][k+3]+ 26. input[i][i][k+4]+input[i][i][k+4]</pre>	CHRECS_input <sub>24</sub> CHRECS_input <sub>25</sub>	$\{0,+,1\}$ $\{0,+,1\}$	$\{0\}$ $\{0\}$	$\begin{cases} -3 \\ \{3\} \end{cases} = \begin{cases} 0, +, 1 \\ \{0, +, 1\} \end{cases}$	$\{0\}$ $\{0\}$	$\{-2\}$ $\{4\}$
27. );	CHRECS_input <sub>26</sub> CHRECS_input <sub>27</sub>	$\begin{cases} \{0,+,1\} \\ \{0,+,1\} \end{cases}$	$\{0\}$ $\{0\}$	$ \begin{cases} -4 \\ \{4\} \end{cases}  \begin{cases} 0, +, 1 \\ \{0, +, 1\} \\ \{0, +, 1\} \end{cases} $	$\{0\}$ $\{0\}$	$\{-3\}$ $\{5\}$

#### shared clause of the gridify directive

< UNIVERSIDADE DA CORUÑA

#### Case Study: SGEMM (I)

```
1.int i, j, l, m, n, k;
2.float A[m][k], B[k][n], C[m][n];
3.float alpha, beta, prod;
4.
5.for (i = 0; i < m; i++) {
6. for (j = 0; j < n; j++) {
7. prod = 0;
8.
    for (1 = 0; 1 < k; 1++) {
       prod += A[i][l] * B[l][j];
9.
10.
   }
11.
     C[i][j] = alpha * prod + beta * C[i][j];
12. }
13.}
```





#### Case Study: SGEMM (II)

- sgemm-cpu
- sgemm-mkl: Intel MKL
- sgemm-hmpp1: Offloading (and check coalescing)

```
1.int i, j, l, m, n, k;
2.float A[m][k], B[k][n], C[m][n];
3.float alpha, beta, prod;
4.
5.for (i = 0; i < m; i++) {
6. for (j = 0; j < n; j++) {
7.
   prod = 0;
    for (l = 0; l < k; l++) {
8.
9.
       prod += A[i][l] * B[l][j];
10.
     }
11.
     C[i][j] = alpha * prod + beta * C[i][j];
12. }
13.}
```

	not instantiated		<i>T0</i>		<i>T1</i>	
	1 <sup>st</sup> dim	$2^{nd} dim$	1 <sup>st</sup> dim	2 <sup>nd</sup> dim	$1^{st}dim$	2 <sup>nd</sup> dim
CHRECS_A CHRECS_B CHRECS_C	$\{ 0, +, 1 \} \\ \{ 0, +, 1 \} \\ \{ 0, +, 1 \}$	$ \begin{array}{c c} \{0,+,1\} \\ \{0,+,1\} \\ \{0,+,1\} \\ \{0,+,1\} \end{array} $	$\{ 0 \} \\ \{ 0, +, 1 \} \\ \{ 0 \}$	$ \begin{array}{c} \{0,+,1\} \\ \{0\} \\ \{0\} \\ \{0\} \end{array} \right $	$\{ \begin{matrix} \{0\} \\ \{0,+,1\} \\ \{0\} \end{matrix}$	$\{ \begin{matrix} 0,+,1 \\ \{1 \\ \{1 \} \\ \{1 \} \end{matrix}$

< UNIVERSIDADE DA CORUÑA

#### Case Study: SGEMM (and III)

sgemm-hmpp2: Tiling preserving coalescing

```
1.int i, j, l, m, n, k;
2.float A[m][k], B[k][n], C[m][n];
3.float alpha, beta, prod;
4.
5.for (i = 0; i < m; i++) {
  for (j = 0; j < n; j++) {
6.
7.
     prod = 0;
8.
     for (1 = 0; 1 < k; 1++) {
       prod += A[i][l] * B[l][j];
9.
10.
     }
11.
     C[i][j] = alpha * prod + beta * C[i][j];
12. }
13.}
```

#### Algorithm 4 Increase the computational load of a GPU thread

1: procedure INCREASELOAD Input: access  $x_k[i_{k,1}][i_{k,2}] \dots [i_{k,n}]$  to an *n*-dimensional array *x* stored in row-major order Input: loop nest  $L = L_1, L_2, \dots, L_l$  where both  $L_1, L_2$  are threadified Input: amount of data  $\Delta$  to be processed by a GPU thread 2: increment the step of the outer loop  $L_1$  to  $\Delta$ 3: for each scalar variable *s* in *L* do

GAC UDC.ES

🗲 UNIVERSIDADE DA CORUÑA

- 4: promote *s* to an array  $s[\Delta]$
- 4. promote s to an array  $s[\Delta]$
- 5: transform reads and writes to s into loops of  $\Delta$  iterations
- 6: end for
- 7: end procedure
- sgemm-hmpp3: Let the compiler use the registers (unroll)
- sgemm-hmpp4: Use the shared memory for B
- sgemm-cublas: NVIDIA CUBLAS library

- Motivation: General Purpose Computation with GPUs
- GPGPU with CUDA & OpenHMPP
- The KIR: an IR for the Detection of Parallelism
- Locality-Aware Generation of Efficient GPGPU Code

CORUÑA UNIVERSIDADE DA CORUÑA

- Case Studies: CONV3D & SGEMM
- Performance Evaluation
- Conclusions & Future Work

#### Performance Evaluation: CONV3D

#### size<sub>X</sub>, size<sub>V</sub> and size<sub>Z</sub> in 128, 256, 384, 512, 640 and 768



GAC COMPUTER ARCHITECTURE GROUP UNIVERSITY OF A CORUNA

🗲 UNIVERSIDADE DA CORUÑA

#### Performance Evaluation: SGEMM

*m, n and k* in 128, 256, 384, 512, 640, 768, 896, 1024, 1152, 1280, 1408, 1536, 1664, 1792, 1920, 2048, 4096, 6144 and 8192



- Motivation: General Purpose Computation with GPUs •
- GPGPU with CUDA & OpenHMPP •
- The KIR: an IR for the Detection of Parallelism
- Locality-Aware Generation of Efficient GPGPU Code •
- Case Studies: CONV3D & SGEMM
- Performance Evaluation
- Conclusions & Future Work



#### Conclusions

KIR-based locality-aware automatic parallelization technique that targets GPU-based heterogeneous systems

- exploit data locality in the complex GPU memory hierarchy
- two representative case studies: CONV3D & SGEMM
- chains of recurrences model accesses to n-dimensional arrays
- OpenHMPP directives enabled a great understandability and portability of the generated GPU code

< UNIVERSIDADE DA CORUÑA

- New automatic partitioning algorithm of the KIR to handle the interactions between computations in full-scale applications
- Auto-tuning approaches to select the best performance on a given hardware architecture
- Test with larger benchmark suite and on other manycore accelerators

< UNIVERSIDADE DA CORUÑA

# Locality-Aware Automatic Parallelization for GPGPU with OpenHMPP Directives

José M. Andión, Manuel Arenaz, François Bodin, Gabriel Rodríguez and Juan Touriño

7th International Symposium on High-Level Parallel Programming and Applications (HLPP 2014) July 3-4, 2014 — Amsterdam, Netherlands

