

# Using a CUDA-Accelerated PGAS Model on a GPU Cluster for Bioinformatics

Jorge González-Domínguez

Parallel and Distributed Architectures Group  
Johannes Gutenberg University of Mainz, Germany  
j.gonzalez@uni-mainz.de

GTC 2015

- 1 Overview of the Problem
- 2 Intra-GPU Parallelization with CUDA
- 3 Inter-GPU Parallelization with UPC++
- 4 Experimental Evaluation
- 5 Conclusions

- 1 Overview of the Problem
- 2 Intra-GPU Parallelization with CUDA
- 3 Inter-GPU Parallelization with UPC++
- 4 Experimental Evaluation
- 5 Conclusions

# Genome-Wide Association Studies (I)

Analyses of genetic influence  
on diseases

# Genome-Wide Association Studies (I)

Analyses of genetic influence  
on diseases

- $M$  individuals



# Genome-Wide Association Studies (I)

Analyses of genetic influence  
on diseases

- $M$  individuals
  - $K$  cases



# Genome-Wide Association Studies (I)

Analyses of genetic influence  
on diseases

- $M$  individuals
  - $K$  cases
  - $C$  controls



# Genome-Wide Association Studies (I)

Analyses of genetic influence  
on diseases

- $M$  individuals
  - $K$  cases
  - $C$  controls
- $N$  genetic markers, Single Nucleotide Polymorphisms (SNPs). 3 genotypes:
  - Homozygous Wild (w, AA, 0)
  - Heterozygous (h, Aa, 1)
  - Homozygous Variant (v, aa, 2)



# Genome-Wide Association Studies (II)

	Cases						Controls									
SNP 1	0	1	2	0	1	2	0	1	2	0	1	2	0	1	2	1
SNP 2	0	1	1	0	2	0	0	0	1	2	2	1	0	1	1	2
SNP 3	0	0	0	0	0	0	0	0	1	2	1	1	1	2	1	1
SNP 4	0	1	0	1	0	1	0	1	2	2	2	2	1	1	1	1
SNP 5	0	2	2	2	0	1	1	1	1	0	0	1	1	0	2	2
SNP 6	1	0	1	0	1	0	1	0	1	2	1	2	1	2	2	1

# Genome-Wide Association Studies (II)

	Cases						Controls									
SNP 1	0	1	2	0	1	2	0	1	2	0	1	2	0	1	2	1
SNP 2	0	1	1	0	2	0	0	0	1	2	2	1	0	1	1	2
SNP 3	0	0	0	0	0	0	0	0	1	2	1	1	1	2	1	1
SNP 4	0	1	0	1	0	1	0	1	2	2	2	2	1	1	1	1
SNP 5	0	2	2	2	0	1	1	1	1	0	0	1	1	0	2	2
SNP 6	1	0	1	0	1	0	1	0	1	2	1	2	1	2	2	1

# Genome-Wide Association Studies (II)

	Cases						Controls									
SNP 1	0	1	2	0	1	2	0	1	2	0	1	2	0	1	2	1
SNP 2	0	1	1	0	2	0	0	0	1	2	2	1	0	1	1	2
SNP 3	0	0	0	0	0	0	0	0	1	2	1	1	1	2	1	1
SNP 4	0	1	0	1	0	1	0	1	2	2	2	2	1	1	1	1
SNP 5	0	2	2	2	0	1	1	1	1	1	0	0	1	1	0	2
SNP 6	1	0	1	0	1	0	1	0	1	2	1	2	1	2	2	1

# Genome-Wide Association Studies (and III)

## Definition

Two SNPs present epistasis or interaction if:

- Their joint genotype frequencies show a statistically significant difference between cases and controls which potentially explains the effect of the genetic variation leading to disease.
- The difference between cases and controls shown by the joint values is significantly higher than using only the individual SNP values.

# BOOST

## BOolean Operation-based Screening and Testing

- Binary traits
- Exhaustive search
- Statistical regression
- Good accuracy (used by biologists)
- Returns a list of SNP pairs with high interaction probability
- Fastest available tool. Intel Core i7 3.20GHz:
  - 40,000 SNPs and 3,200 individuals
    - About 800 million pairs
    - 51 minutes
  - 500,000 SNPs and 5,000 individuals
    - About 125 billion pairs (moderated size)
    - Estimated 7 days

# GBOOST

## CUDA version for GPUs

- Same accuracy as BOOST
- 40,000 SNPs and 6,400 individuals
  - About 800 million pairs
  - 28 seconds on a GTX Titan
- 500,000 SNPs and 5,000 individuals
  - About 125 billion pairs (moderated size)
  - 1 hour on a GTX Titan

# GBOOST

## CUDA version for GPUs

- Same accuracy as BOOST
- 40,000 SNPs and 6,400 individuals
  - About 800 million pairs
  - 28 seconds on a GTX Titan
- 500,000 SNPs and 5,000 individuals
  - About 125 billion pairs (moderated size)
  - 1 hour on a GTX Titan

High-throughput genotyping technologies collect few million SNPs of an individual within a few minutes → Expected datasets with 5M SNPs and 10,000 individuals

- 1 Overview of the Problem
- 2 Intra-GPU Parallelization with CUDA**
- 3 Inter-GPU Parallelization with UPC++
- 4 Experimental Evaluation
- 5 Conclusions

## Calculation of Contingency Tables (I)

For each SNP-pair → Number of occurrences of each combination of genotypes

<b>Cases</b>	<b>SNP2=0</b>	<b>SNP2=1</b>	<b>SNP2=2</b>
<b>SNP1=0</b>	$n_{000}$	$n_{010}$	$n_{020}$
<b>SNP1=1</b>	$n_{100}$	$n_{110}$	$n_{120}$
<b>SNP1=2</b>	$n_{200}$	$n_{210}$	$n_{220}$
<b>Controls</b>	<b>SNP2=0</b>	<b>SNP2=1</b>	<b>SNP2=2</b>
<b>SNP1=0</b>	$n_{001}$	$n_{011}$	$n_{021}$
<b>SNP1=1</b>	$n_{101}$	$n_{111}$	$n_{121}$
<b>SNP1=2</b>	$n_{201}$	$n_{211}$	$n_{221}$

## Calculation of Contingency Tables (II)

SNP 4   0   1   0   1   0   1   0   1   2   2   2   2   1   1   1   1

SNP 6   1   0   1   0   1   0   1   0   1   2   1   2   1   2   2   1

Casos	SNP6=0	SNP6=1	SNP6=2
SNP4=0	0	4	0
SNP4=1	4	0	0
SNP4=2	0	0	0
Controles	SNP6=0	SNP6=1	SNP6=2
SNP4=0	0	0	0
SNP4=1	0	2	2
SNP4=2	0	1	2

## Calculation of Contingency Tables (II)

SNP 4   0   1   0   1   0   1   0   1   2   2   2   2   1   1   1   1

SNP 6   1   0   1   0   1   0   1   0   1   2   1   2   1   2   2   1

Casos	SNP6=0	SNP6=1	SNP6=2
SNP4=0	0	4	0
SNP4=1	4	0	0
SNP4=2	0	0	0
Controles	SNP6=0	SNP6=1	SNP6=2
SNP4=0	0	0	0
SNP4=1	0	2	2
SNP4=2	0	1	2

## Calculation of Contingency Tables (II)

SNP 4   0   1   0   1   0   1   0   1   2   2   2   2   1   1   1   1

SNP 6   1   0   1   0   1   0   1   0   1   2   1   2   1   2   2   1

Casos	SNP6=0	SNP6=1	SNP6=2
SNP4=0	0	4	0
SNP4=1	4	0	0
SNP4=2	0	0	0
Controles	SNP6=0	SNP6=1	SNP6=2
SNP4=0	0	0	0
SNP4=1	0	2	2
SNP4=2	0	1	2

## Filtering Stage

- Epistatic interaction measured via log-linear models
- All SNP-pairs analyzed
- The measure is obtained with numerical calculations from the values of the contingency table
- Pairs with measure higher than a threshold pass the filter
  - They are included in the output file
- multiEpistSearch uses a faster filter than GBOOST (out of the scope)

# CUDA Implementation

## CUDA Kernel

- Genotyping information loaded in device memory through pinned copies
- Each thread performs the whole calculation of independent SNP-pairs
- Only one kernel for the whole computation
- Each call to the kernel analyzes a batch of SNP-pairs

# CUDA Implementation

## CUDA Kernel

- Genotyping information loaded in device memory through pinned copies
- Each thread performs the whole calculation of independent SNP-pairs
- Only one kernel for the whole computation
- Each call to the kernel analyzes a batch of SNP-pairs

## Optimization Techniques

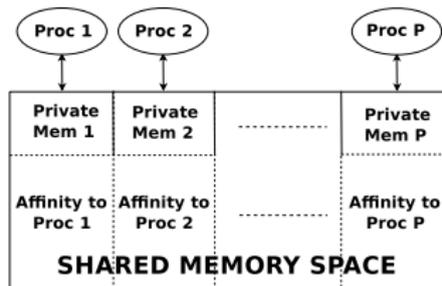
- Boolean representation of genotyping information
- Increase of coalescence
- Exploitation of shared memory

- 1 Overview of the Problem
- 2 Intra-GPU Parallelization with CUDA
- 3 Inter-GPU Parallelization with UPC++**
- 4 Experimental Evaluation
- 5 Conclusions

# UPC++ (I)

- Unified Parallel C++
- Novel extension of ANSI C++
  - Y Zheng, A Kamil, M Driscoll, H Shan, and K Yelick.  
UPC++: a PGAS Extension for C++. *In Proc. 28th IEEE Intl. Parallel and Distributed Processing Symp. (IPDPS'14)*, Phoenix, AR, USA, 2014.
- Follows the Partitioned Global Address Space (PGAS) programming model
- Single Program Multiple Data (SPMD) execution model
- Works on shared and distributed memory systems

# UPC++ (and II)

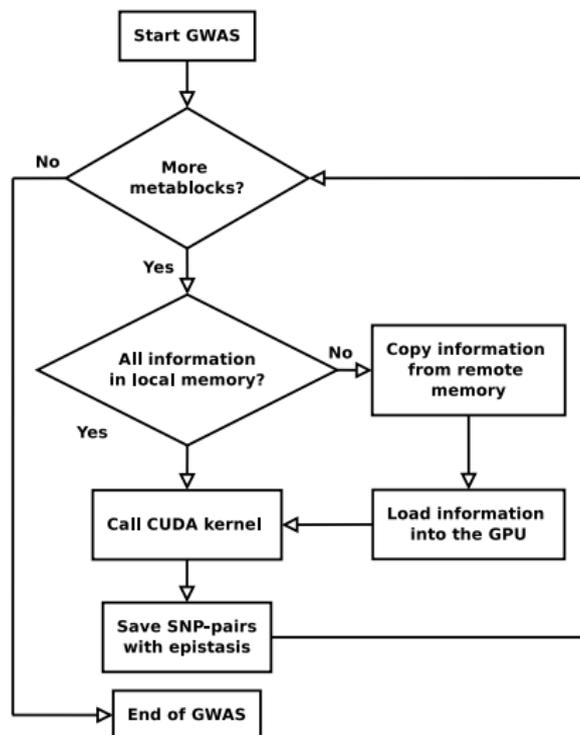


- Global memory logically partitioned among processes
- Processes can directly access (read/write) any part of the global memory
- Memory with affinity usually mapped in the same node (faster accesses)

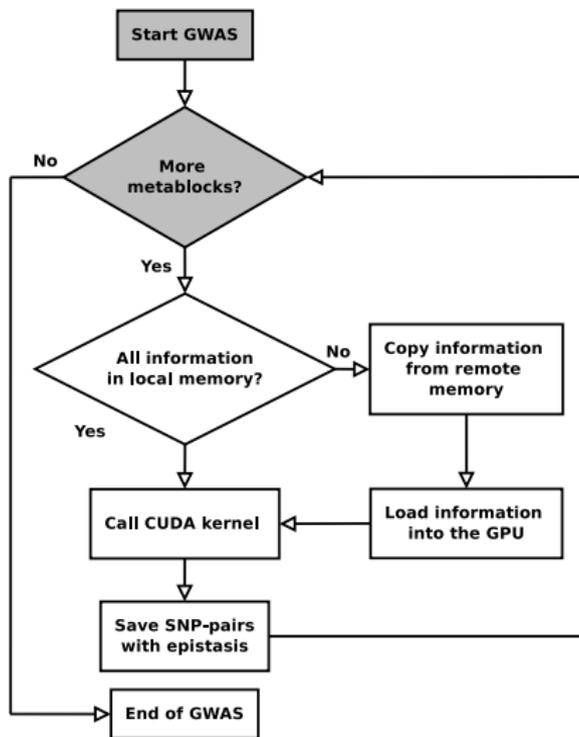
## Multi-GPU Approach (I)

- One UPC++ process per GPU
- SNP data distributed among the parts of the global memory
  - All the information of the same SNP in the same part
- Each GPU (UPC++ process) analyzes different SNP-pairs
  - Creation of contingency table
  - Filtering
- The data of the SNPs to analyze might be in remote memory

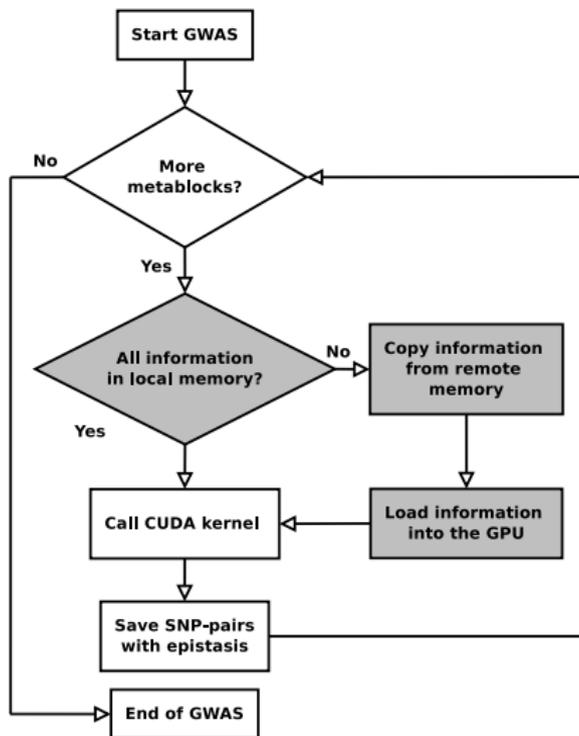
# Multi-GPU Approach (II)



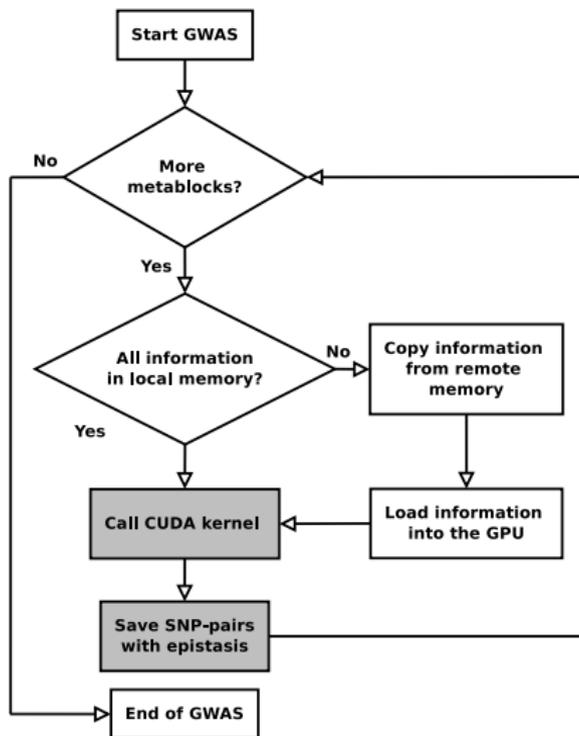
# Multi-GPU Approach (II)



# Multi-GPU Approach (II)



# Multi-GPU Approach (II)



## Multi-GPU Approach (VI)

### Static distribution

- Workload distributed at the beginning
  - Metablocks that will be analyzed by each GPU
- The distribution does not change during the execution
- Balance of the number of metablocks per GPU
  - Similar workload for each GPU
  - Good distribution for systems with similar GPUs
- Minimization of remote copies

## Multi-GPU Approach (and VII)

### On-demand distribution

- The metablocks computed by each GPU initially unknown
- Table with one binary value per metablock that indicates if it has been computed
- When one GPU finishes with one metablock → Looks for the next one that has not been analyzed
- Locks or semaphores necessary for the concurrent accesses to the table
  - Easy with UPC++ support
  - Synchronizations include performance overhead
- GPUs might compute different number of metablocks
  - Faster GPUs analyze more SNP-pairs
  - Good distribution for systems with different GPUs

- 1 Overview of the Problem
- 2 Intra-GPU Parallelization with CUDA
- 3 Inter-GPU Parallelization with UPC++
- 4 Experimental Evaluation**
- 5 Conclusions

# Evaluation with Homogeneous GPUs (I)

## Platform

- Mogon cluster
- Johannes Gutenberg Universität
- 8 nodes with 3 GTX Titan GPUs
  - One of the most powerful GPUs
- Infiniband network

## Evaluation with Homogeneous GPUs (I)

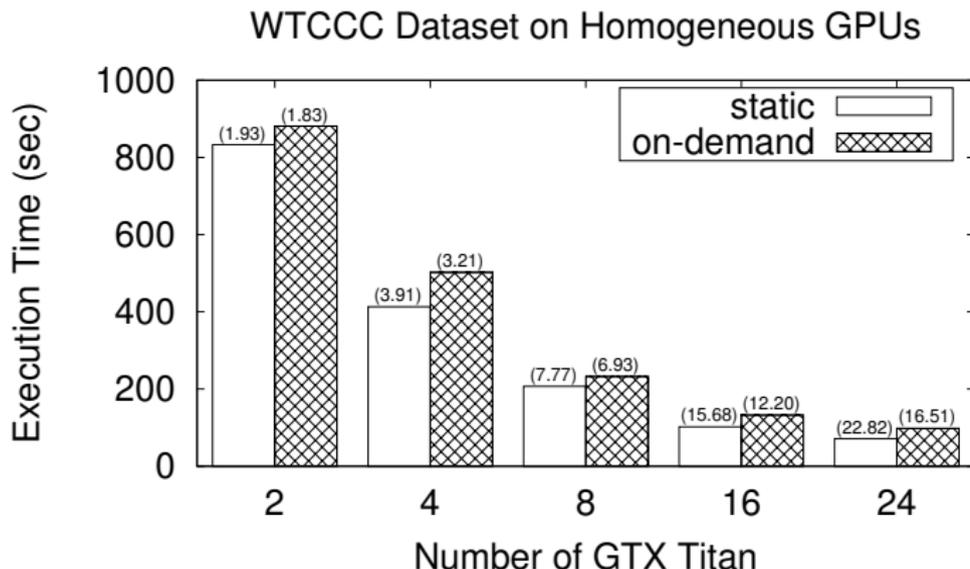
### Platform

- Mogon cluster
- Johannes Gutenberg Universität
- 8 nodes with 3 GTX Titan GPUs
  - One of the most powerful GPUs
- Infiniband network

### Dataset

- Real-world data from the WTCCC database
- Moderately-sized
  - 500,568 SNPs
  - 2,005 cases with bipolar disorder
  - 3,004 controls

## Evaluation with Homogeneous GPUs (II)



- Static 1.38 times faster for 24 GPUs
- Static always > 95 % parallel efficiency

## Evaluation with Homogeneous GPUs (and III)

Design	Architecture	Runtime	Speed ( $10^6$ pairs/s)
multiEpistSearch	24 GTX Titan	1 m 11 s	1764.56
multiEpistSearch	1 GTX Titan	27 m	77.34
GBOOST	1 GTX Titan	1 h 15 m	34.23
EpiGPU*	1 GTX 580	2 h 55 m	11.90
SHEsisEPI*	1 GTX 285	27 h	1.29
BOOST**	Intel Core i7	7 d	0.21

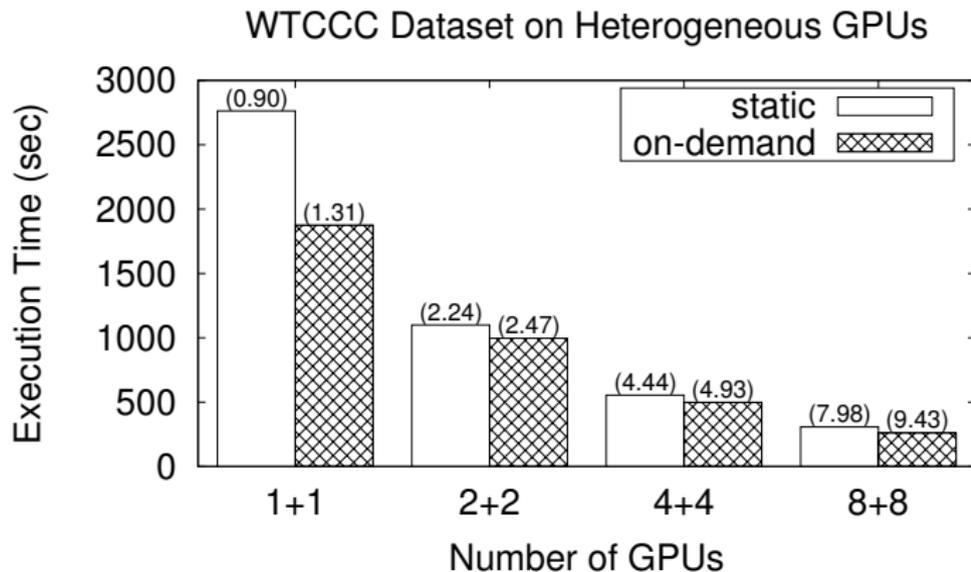
- Speedups for one GPU:
  - 2.77 over GBOOST
  - $> 373$  over estimation for BOOST on a 3GHz Intel Core i7
- With 24 Titan 54.93 and  $> 8,500$  times faster than GBOOST and BOOST, respectively

# Evaluation with Heterogeneous GPUs (I)

## Platform

- Pluton cluster
- Universidade da Coruña (Spain)
- 8 nodes with 1 GTX Tesla K20m
- 4 nodes with 2 Tesla 2050
  - Less cores
- Gigabit Ethernet network

## Evaluation with Heterogeneous GPUs (II)



- On demand 1.18 times faster for 16 GPUs

## Evaluation with Heterogeneous GPUs (and III)

Design	Architecture	Runtime	Speed ( $10^6$ pairs/s)
multiEpistSearch	8 Tesla K20m + 8 2050	4 m 20 s	481.86
multiEpistSearch	8 Tesla K20m	5 m 40 s	348.01
multiEpistSearch	8 Tesla 2050	10 m 12 s	204.71
multiEpistSearch	1 Tesla K20m	41 m	50.93
multiEpistSearch	1 Tesla 2050	1 h 1 m	34.23
GBOOST	1 Tesla K20m	1 h 26 m	24.28
GBOOST	1 Tesla 2050	2 h 17 m	15.22

- With 1 GPU 2.10 and 2.25 times faster than GBOOST
- 1.31 times faster using the whole cluster (on-demand) than only the 8 Tesla K20m

# Evaluation of a Large-Scale Dataset

- Simulated dataset
  - 5M SNPs
  - 5,000 cases
  - 5,000 controls
- **2 hours and 45 minutes on Mogon (24 GTX Titan)**
- Estimation of more than 2 days and 14 hours on 1 GPU
- GBOOST is not able to analyze it
  - Out-of-bound problems in the arrays

- 1 Overview of the Problem
- 2 Intra-GPU Parallelization with CUDA
- 3 Inter-GPU Parallelization with UPC++
- 4 Experimental Evaluation
- 5 Conclusions**

# Conclusions

- multiEpistSearch looks for epistatic interactions on GPU clusters
- Hybrid CUDA&UPC++ implementation
- On only one GPU always speedups higher than 2 over GBOOST
- Two inter-GPU data distributions
  - Static for homogeneous clusters
  - Dynamic for heterogeneous clusters
- High scalability
  - 95% Parallel efficiency with 24 GTX Titans and WTCCC dataset
- 2 hours and 45 minutes for 5M SNPs and 10K samples on 24 GTX Titans

# Bibliography

- **First version of the GPU kernel**

J. González-Domínguez, B. Schmidt, J. C. Kässens, and L. Wienbrandt.

Hybrid CPU/GPU Acceleration of Detection of 2-SNP Epistatic Interactions in GWAS.

In *Proc. 20th Intl. European Conf. on Parallel and Distributed Computing (Euro-Par'14)*, Porto, Portugal.

- **multiEpistSeach (minor revision)**

J. González-Domínguez, J. C. Kässens, L. Wienbrandt, and B. Schmidt.

Large-Scale Genome-Wide Association Studies on a GPU Cluster Using a CUDA-Accelerated PGAS Programming Model.

*Intl. Journal of High Performance Computing Applications (IJHPCA)*.

# Using a CUDA-Accelerated PGAS Model on a GPU Cluster for Bioinformatics

Jorge González-Domínguez

Parallel and Distributed Architectures Group  
Johannes Gutenberg University of Mainz, Germany  
j.gonzalez@uni-mainz.de

GTC 2015