Design and Implementation of Parallel Algorithms for Modern Heterogeneous Platforms Based on Functional Performance Models

#### David Clarke

#### Heterogeneous Computing Laboratory, UCD, Dublin

#### May 11, 2012



**Platform and Application** 

#### Generalised Heterogeneous Platform





**Data Parallel Application** 

while(...) {

}

compute\_parallel(data, size);
syncronise (data);

Background

**Platform and Application** 

#### Generalised Heterogeneous Platform



**Data Parallel Application** 

while(...)  $\{$ 

}

compute\_parallel(data, size);
syncronise (data);

Background

**Platform and Application** 

#### Generalised Heterogeneous Platform



#### Data Parallel Application

while(...)  $\{$ 

compute\_parallel(data, size);
syncronise (data);

Background

# Traditional load Balancing

Background

- Traditionally, processor performance is defined by a constant number.
- Computational units are partitioned as

$$d_i = N imes rac{s_i}{\sum_{j=1}^p s_j}$$

 In reality, speed is a function of problem size.



Problem Size N

#### Load Balanced when:

$$t_1(d_1) \approx t_2(d_2) \approx \ldots \approx t_p(d_p)$$
$$\frac{d_1}{s_1(d_1)} \approx \frac{d_2}{s_2(d_2)} \approx \ldots \approx \frac{d_p}{s_p(d_p)}$$
$$d_1 + d_2 + \ldots + d_p = N$$

Developed two algorithms:

#### • Geometric Partitioning Algorithm

- Convergence guaranteed
- Restriction on shape of speed functions

#### • Numerical Partitioning Algorithm

- No restriction on shape
- No guarantee of convergence (in practice converges for realistic functions)

Input: FPMs and *N*. Output: 
$$d_1, d_2, \ldots, d_p$$
.

# Geometric Partitioning Algorithm

- Points  $(d_i, s_i(d_i))$  lie on a line passing through the origin when  $\frac{d_i}{s_i(d_i)} = constant.$
- Value of *N* determines the slope.
- Algorithm iteratively bisects solution space to find values *d<sub>i</sub>*.



# Geometric Partitioning Algorithm

- Points  $(d_i, s_i(d_i))$  lie on a line passing through the origin when  $\frac{d_i}{s_i(d_i)} = constant.$
- Value of *N* determines the slope.
- Algorithm iteratively bisects solution space to find values *d<sub>i</sub>*.



### Numerical Partitioning Algorithm

Optimal partitioning can be formulated as a system of nonlinear equations, F(x) = 0

$$F(x) = \begin{cases} n - \sum_{i=1}^{p} x_i \\ \frac{x_i}{s_i(x_i)} - \frac{x_1}{s_1(x_1)}, & 2 \le i \le p \end{cases}$$
(1)

Powell's Hybrid method used to solve:

$$x_{k+1} = x_k - J(x_k)F(x_k)$$
 (2)

With Jacobian matrix:

$$J(x) = \begin{pmatrix} -1 & -1 & \dots & -1 \\ -\frac{s_1(x_1) - x_1 s_1'(x_1)}{s_1^2(x_1)} & \frac{s_2(x_2) - x_2 s_2'(x_2)}{s_2^2(x_2)} & 0 & 0 \\ \dots & 0 & \dots & 0 \\ -\frac{s_1(x_1) - x_1 s_1'(x_1)}{s_1^2(x_1)} & 0 & 0 & \frac{s_p(x_p) - x_p s_p'(x_p)}{s_p^2(x_p)} \end{pmatrix}$$
(3)

## Fitting Continuous Functions to Discreet Data



Functional Performance Model Partitioning

Case study: Matrix Partitioning based on FPMs

#### Simple Partitioning



#### 2D Partitioning



### Two-Dimensional Matrix Partitioning with 1D FPMs

- Height and width combined into one parameter, area  $d_i = m_i \times n_i$ .
- Square areas are benchmarked  $m = n = \sqrt{d}$ .
- Partition with 1D FPM algorithm to find area of rectangles (geometric or numerical).
- Use communication volume minimising algorithm\* to compute ordering and shape of these rectangles.



\* Beaumont, O., Boudet, V., Rastello, F., Robert, Y.: Matrix Multiplication on Heterogeneous Platforms. 2001



- Functional performance models are different for each application and each platform.
- Building these models for all conceivable problem sizes is computationally expensive.
- Building full models is not an option for a self adaptive algorithm.
- Present an algorithm that dynamically builds the models at relevant problem sizes





































### Experimental Results: FPM based Partitioning



Function

David Clarke (HCL/UCD)





# Target Hierarchical Heterogeneous Platform



- Heterogeneity between nodes
- Heterogeneity between devices within a node
- Partition matrix between nodes
- Sub-partition between devices within a node
- eg. Grid'5000 Grenoble site.

Hierarchical Partitioning Algorithm

# **Hierarchical Partitioning Algorithm**



• Hierarchy in platform  $\rightarrow$  hierarchy in partitioning

- Nested parallelism
- inter-node partitioning algorithm (INPA)
- inter-device partitioning algorithm (IDPA)
- IDPA is nested inside INPA

















### Experimental Results: Reference full models



Hierarchical Partitioning Algorithm

## **Experimental Results**



90 nodes, 432 CPUs, 12 GPUs from Grid5000 Grenoble site

Two-Dimensional Matrix Partitioning Based on 2D Models

- Numerical solution
- 3-step column based algorithm
- Both algorithms use column based partitioning. Arrangement of processors into a 2D grid is an input.



# Numerical 2D partitioning algorithm

Solve as a constrained non-linear minimisation problem

Function to be minimised:

 $\sum_{j=1}^{q} \sum_{i=1}^{p_j} t_{ij}$ 

with q + 1 constraints:

$$n_1 + n_2 + \ldots + n_q = N$$
  
 $m_{1j} + m_{2j} + \ldots + m_{p_j} = M$ , for  $1 \le j \le q$ 

- Implemented with the NLopt library.
- However, convergence and reliability issues still need to be solved
- Possibly posing the minimisation function and constraints differently will aid convergence.

- Extract 1D model from 2D with P1 P4 n = c for  $1 \le c \le N$ . Partition to find optimum time, P2 P5 add point  $s_{col}(c)$  to model.
- Use 1D partitioning find optimum column widths.
- Find optimum heights for each processor within the columns.



- Extract 1D model from 2D with n = c for  $1 \le c \le N$ . Partition to find optimum time, add point  $s_{col}(c)$  to model.
- Use 1D partitioning find optimum column widths.
- Find optimum heights for each processor within the columns.



- Extract 1D model from 2D with n = c for  $1 \le c \le N$ . Partition to find optimum time, add point  $s_{col}(c)$  to model.
- Use 1D partitioning find optimum column widths.
- Find optimum heights for each processor within the columns.



- Extract 1D model from 2D with n = c for  $1 \le c \le N$ . Partition to find optimum time, add point  $s_{col}(c)$  to model.
- Use 1D partitioning find optimum column widths.
- Find optimum heights for each processor within the columns.



- Extract 1D model from 2D with n = c for  $1 \le c \le N$ . Partition to find optimum time, add point  $s_{col}(c)$  to model.
- Use 1D partitioning find optimum column widths.
- Find optimum heights for each processor within the columns.



- Extract 1D model from 2D with n = c for  $1 \le c \le N$ . Partition to find optimum time, add point  $s_{col}(c)$  to model.
- Use 1D partitioning find optimum column widths.
- Find optimum heights for each processor within the columns.



- Extract 1D model from 2D with n = c for  $1 \le c \le N$ . Partition to find optimum time, add point  $s_{col}(c)$  to model.
- Use 1D partitioning find optimum column widths.
- Find optimum heights for each processor within the columns.



#### Applications

- Signal Processing Systems group, INESC-ID, Lisbon, Portugal
  - Using FPMs to partition database requests on heterogeneous platform.
  - Extended FPMs to overlap communications, ×4 speedup for FFT.
- Division of Scientific Computing, Uppsala University, Sweden
  - Upcoming collaboration to load balance multiphase flow simulations.

#### Experimental Platforms

- Grid'5000, 10 sites, 1260 nodes, France
- HCL Cluster, local 16 node experimental cluster.
- SiPS Cluster, 4 node CPU+GPU, Lisbon, Portugal

## Software

### FuPerMod

- FuPerMod package developed at HCL
- Based on system and mathematical software: C/C++, MPI, Autotools, GNU Scientific Library, Boost C++ libraries, BLAS, CUDA Toolkit
- Contains all presented algorithms.
- Designed for easy integration with existing heterogeneous applications.

#### Conclusion

# Publications

Lastovetsky, A., and R. Reddy, "Distributed Data Partitioning for Heterogeneous Processors Based on Partial Estimation of their Functional Performance Models", HeteroPar'2009, Delft, Netherlands, LNCS, vol. 6043, Springer, pp. 91-101, 2010.

The stand and the stand the stand

- Lastovetsky, A., and R. Reddy, "Two-dimensional Matrix Partitioning for Parallel Computing on Heterogeneous Processors Based on their Functional Performance Models", HeteroPar'2009, Delft, Netherlands, LNCS, vol. 6043, Springer, pp. 112-121, 2010.
- Clarke, D., A. Lastovetsky, and V. Rychkov, "Dynamic Load Balancing of Parallel Computational Iterative Routines on Platforms with Memory Heterogeneity", Europar 2010 / Heteropar'2010, Ischia-Naples, Italy, LNCS 6586, Springer, pp. 41-50, 2011.
- Clarke, D., A. Lastovetsky, and V. Rychkov, "Dynamic Load Balancing of Parallel Computational Iterative Routines on Highly Heterogeneous HPC Platforms", PPL, vol. 21, issue 2: World Scientific, pp. 195-217, 06/2011.
- Rychkov, V., D. Clarke, and A. Lastovetsky, "Using Multidimensional Solvers for Optimal Data Partitioning on Dedicated Heterogeneous HPC Platforms", PaCT-2011, LNCS 6873, Kazan, Russia, Springer, pp. 332-346, 19/09/2011.
- Lastovetsky, A., R. Reddy, V. Rychkov, and D. Clarke, "Design and implementation of self-adaptable parallel algorithms for scientific computing on highly heterogeneous HPC platforms", PARCO (under review).
- Clarke, D., A. Lastovetsky, and V. Rychkov, "Column-Based Matrix Partitioning for Parallel Matrix Multiplication on Heterogeneous Processors Based on Functional Performance Models", HeteroPar'2011, Bordeaux, France, LNCS 7155, Springer, pp. 450-459, 2012.
- Clarke, D., A. Ilic, A. Lastovetsky, L. Sousa, "Hierarchical Partitioning Algorithm for Scientific Computing on Highly Heterogeneous CPU + GPU Clusters" EuroPar'2012 (under review).