Hierarchical Divisible Load Scheduling on Heterogeneous Systems with Limited Memory

Leonel Sousa and Aleksandar Ilić



INESC-ID/ IST, TU Lisbon, Portugal

High Performance Computing, Grids and Clouds Workshop, Cetraro, 2011





COMMODITY COMPUTERS = HETEROGENEOUS SYSTEMS

- Multi-core General-Purpose Processors (CPUs)
- Many-core Graphic Processing Units (GPUs)
- Special accelerators, co-processors, FPGAs,...

CLUSTERS OF COMMODITY COMPUTERS = HIGHLY HETEROGENEOUS SYSTEMS

- => SIGNIFICANT COMPUTING POWER
- Not yet completely explored for COLLABORATIVE COMPUTING
- TO USE THE AVAILABLE RESOURCES AND IMPROVE PERFORMANCE/WATT

HETEROGENEITY MAKES PROBLEMS MUCH MORE COMPLEX!

- Scheduling, performance modeling and load balancing
- Different programming models, languages and implementations

Instituto de Engenharia de Sistemas e Computadores Investigação e Desenvolvimento em Lisboa

technology



DISCRETELY DIVISIBLE LOAD (DDL) PROCESSING

HETEROGENEOUS DISTRIBUTED SYSTEMS WITH HETEROGENEOUS CPU+GPU DESKTOPS AS COMPUTING NODES

Performance Modeling and 3-Step Hierarchical DDL Scheduling Algorithm

CASE STUDY: 2D FFT Batch Execution



Discretely Divisible Load Processing



technology

from seed

- DISCRETELY DIVISIBLE LOAD (DDL) APPLICATIONS
 - Computations divisible into pieces of arbitrary sizes (integers)
 - Fractions independently processed in parallel with no precedence constraints
- APPLICABLE TO A WIDE RANGE OF SCIENTIFIC PROBLEMS
 - Linear algebra, digital signal and image processing, database applications ...
- STATE OF THE ART DDL APPROACHES IN HETEROGENEOUS DISTRIBUTED

COMPUTING

- Very **recent**, but only for traditional distributed systems (no hierarchy, CPU-only)
- Assume symmetric bandwidth and an one-port model for communication links
- Limited memory: only the size of input load is considered, or the exceeding load is simply redistributed among the nodes with available memory



Heterogeneous Desktop Systems



HETEROGENEOUS STAR NETWORK (MASTER-WORKER)

- MULTI-CORE CPU (Master)
 - Global execution controller; access the whole global memory
 - All cores employed for execution
- INTERCONNECTION BUSES
 - Bidirectional full-duplex asymmetric communication
 - Different concurrency levels
 - Potential execution bottleneck
- DEVICES (Distant workers)
 - Different architectures and programming models
 - Computation performed using local memories



technology

from seed

Commodity Cluster Systems



HETEROGENEOUS STAR NETWORK

- MASTER NODE
 - Heterogeneous desktop system
 - Positioned at the center; employed for execution as whole
- INTERCONNECTION NETWORK
 - Limited asymmetric one-port communication
 - Potential execution bottleneck
- WORKER NODES
 - Heterogeneous desktop systems
 - All available devices employed in each



technology

from seed





_

technology



_

Proposed Algorithm Outline (4)



3-STEP HIERARCHICAL DDL SCHEDULING:

- STEP 1 SYSTEM-LEVEL LOAD BALANCING
 - How many load units to send to each node?
- STEP 2 NODE-LEVEL LOAD BALANCING
 - How many chunks to send to each device from a given node's load?



Proposed Algorithm Outline (5)



3-STEP HIERARCHICAL DDL SCHEDULING:

- STEP 1 SYSTEM-LEVEL LOAD BALANCING
 - How many load units to send to each node?
- STEP 2 NODE-LEVEL LOAD BALANCING
 - How many chunks to send to each device from a given node's load?



Proposed Algorithm Outline (6)



3-STEP HIERARCHICAL DDL SCHEDULING:

- STEP 1 SYSTEM-LEVEL LOAD BALANCING
 - How many load units to send to each node?
- STEP 2 NODE-LEVEL LOAD BALANCING
 - How many chunks to send to each device from a given node's load?
- STEP 3 DEVICE-LEVEL SCHEDULING
 - How to sub-partition the device load to:
 - Reduce delays when distributing and retrieving
 - Overlap computation and communication
 - Efficiently use the bidirectional asymmetric bandwidth of buses
 - Respect the amount of supported concurrency
 - Fit into device limited memory



Proposed Algorithm Outline (7)



3-STEP HIERARCHICAL DDL SCHEDULING:

- STEP 1 SYSTEM-LEVEL LOAD BALANCING
 - How many load units to send to each node?
- STEP 2 NODE-LEVEL LOAD BALANCING
 - How many chunks to send to each device from a given node's load?
- STEP 3 DEVICE-LEVEL SCHEDULING
 - How to sub-partition the device load to:
 - Reduce delays when distributing and retrieving
 - Overlap computation and communication
 - Efficiently use the bidirectional asymmetric bandwidth of buses
 - Respect the amount of supported concurrency
 - Fit into device limited memory



Performance Modeling



FUNCTIONAL PERFORMANCE MODELS

- Built from the real application execution
 - No assumptions being made to ease modeling!
- NODE-LEVEL PERFORMANCE MODELS
 - Computation performance models (ψ_w)
 - For each master core and distant worker
 - Full-duplex communication bandwidth ($\sigma_{\iota}, \sigma_{o}$)
 - Bidirectional and asymmetric for each link
 - Total performance ($\psi_ au$) of each device
 - Including computation and communication
- SYSTEM-LEVEL PERFORMANCE MODELS
 - Computation performance models (Ψ_w)
 - For each node (heterogeneous desktop system)
 - Communication bandwidth (Σ_{ι} , Σ_{o})
 - Bidirectional, asymmetric for each network link
 - Total performance ($\Psi_{ au}$) of each node





technology

from seed

• STEP 1 – SYSTEM-LEVEL LOAD BALANCING



- The optimal distribution between nodes lies on a straight line passing through the origin of coordinate system and intersecting total performance curves $(\Psi_{\tau})^*$:

$$\frac{\alpha_1}{\Psi_{t_1}(\alpha_1)} = \dots = \frac{\alpha_k}{\Psi_{t_k}(\alpha_k)} = \dots = \frac{\alpha_n}{\Psi_{t_n}(\alpha_n)} \quad \sum_{i=1}^n \alpha_i = N$$

• STEP 2 – NODE-LEVEL LOAD BALANCING



- Each load fraction α_i is sub-partitioned between node's devices; the optimal distribution that lies on a straight line passing through the origin of coordinate system and intersecting communication-aware total performance curves (ψ_{τ}), such that*:

$$\frac{\beta_1}{\psi_{\tau_1}(\beta_1)} = \dots = \frac{\beta_k}{\psi_{\tau_k}(\beta_k)} = \dots = \frac{\beta_n}{\psi_{\tau_n}(\beta_n)} \qquad \sum_{j=1}^n \beta_j = \alpha_i$$

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

Determination of Load Fractions (2)



- STEP 3 DEVICE-LEVEL SCHEDULING
 - Per-device distributions β_j are allowed to exceed the device memory limits, b_j
 - Device-level multi-installment processing with multi-distributions
 - Γ_k sub-distributions with $\gamma_{k,l}$ sub-load fractions



technology

Determination of Load Fractions (3)



- STEP 3 DEVICE-LEVEL SCHEDULING
 - Per-device distributions β_j are allowed to exceed the device memory limits, b_j
 - Device-level multi-installment processing with multi-distributions



- Γ_k sub-distributions with $\gamma_{k,l}$ sub-load fractions

Instituto de Engenharia de Sistemas e Computadores Investigação e Desenvolvimento em Lisboa

technology

Determination of Load Fractions (4)



- STEP 3 DEVICE-LEVEL SCHEDULING
 - Per-device distributions β_j are allowed to exceed the device memory limits, b_j
 - Device-level multi-installment processing with multi-distributions



- Γ_k sub-distributions with $\gamma_{k,l}$ sub-load fractions

Instituto de Engenharia de Sistemas e Computadores Investigação e Desenvolvimento em Lisboa

technology

Determination of Load Fractions (5) - Limited Memory -



- STEP 3 DEVICE-LEVEL SCHEDULING
 - Per-device distributions β_j are allowed to exceed the device memory limits, b_j
 - Device-level multi-installment processing with multi-distributions
 - Γ_k sub-distributions with $\gamma_{k,l}$ sub-load fractions
 - Application memory requirements are modeled with three functions of load size:
 - Input memory requirements, $\mu_{\iota}(x)$
 - Output memory requirements, $\mu_o(x)$
 - Execution memory requirements, $\mu_w(x, P)$
 - Different implementations of the same problem might have different memory requirements!

 \Rightarrow In each Γ_k sub-distribution the whole amount of memory may be consumed, such that:

$$\sum_{l=1}^{|\Gamma_k|} \left(\mu_\iota(\gamma_{k,l}) + \mu_w(\gamma_{k,l}, p_j) + \mu_o(\gamma_{k,l}) \right) \le b_j$$
$$\sum_{k=1}^{|\Gamma|} \sum_{l=1}^{|\Gamma_k|} \gamma_{k,l} = \beta_j$$

Instituto de Engenharia de Sistemas e Computadores Investigação e Desenvolvimento em Lisboa

technology

Determination of Load Fractions (6)- Computation/Communication Overlapping



- STEP 3 DEVICE-LEVEL SCHEDULING
 - For each Γ_k sub-distribution, $\gamma_{k,l}$ sizes are carefully chosen to allow as best as possible **overlapping of computation and communication** between subsequent sub-fractions



Instituto de Engenharia de Sistemas e Computadores Investigação e Desenvolvimento em Lisboa

technology

Determination of Load Fractions (7)- Computation/Communication Overlapping



- STEP 3 DEVICE-LEVEL SCHEDULING
 - For each Γ_k sub-distribution, $\gamma_{k,l}$ sizes are carefully chosen to allow as best as possible **overlapping of computation and communication** between subsequent sub-fractions



technology

Determination of Load Fractions (8)- Computation/Communication Overlapping



- STEP 3 DEVICE-LEVEL SCHEDULING
 - For each Γ_k sub-distribution, $\gamma_{k,l}$ sizes are carefully chosen to allow as best as possible **overlapping of computation and communication** between subsequent sub-fractions



Instituto de Engenharia de Sistemas e Computadores Investigação e Desenvolvimento em Lisboa

technology

Determination of Load Fractions (9) - Computation/Communication Overlapping



- STEP 3 DEVICE-LEVEL SCHEDULING
 - For each Γ_k sub-distribution, $\gamma_{k,l}$ sizes are carefully chosen to allow as best as possible **overlapping of computation and communication** between subsequent sub-fractions
 - The decisions are made according to the **amount of overlapping concurrency** supported by the device:

$$\begin{array}{c|c} Yk1 & \hline t_{v}(\gamma_{k1}) & \hline t_{w}(\gamma_{k1}) & \hline t_{o}(\gamma_{k1}) \\ \hline Yk2 & \hline t_{v}(\gamma_{k2}) & \hline t_{o}(\gamma_{k2}) & \hline t_{o}(\gamma_{k2}) & \hline t_{v}(\gamma_{k4}) & \hline t_{w}(\gamma_{k4}) \\ \hline Yk3 & \hline t_{v}(\gamma_{k3}) & \hline t_{w}(\gamma_{k3}) & \hline t_{o}(\gamma_{k3}) \\ \hline \end{array} \right)$$

(a) Overlap of a single communication with computation at the time



(b) Complete concurrency between communication and computation

Instituto de Engenharia de Sistemas e Computadores Investigação e Desenvolvimento em Lisboa

technology

Determination of Load Fractions (10) - DEVICE-LEVEL SCHEDULING ALGORITHM -



technology

from seed

STEP 3 – DEVICE-LEVEL SCHEDULING ALGORITHM*/**



- According to the performance models for device computation (ψ_w) and bidirectional asymmetric full-duplex communication links (σ_ι , σ_o)

^{*} Ilić, A., and Sousa, L., "Scheduling Divisible Loads on Heterogeneous Desktop Systems with Limited Memory" (submitted to HeteroPar 2011)

^{**} Ilić, A., and Sousa, L., "Algorithm For Divisible Load Scheduling on Heterogeneous Systems with Realistic Performance Models", Tech. rep., INESC-ID (May 2011)

Determination of Load Fractions (11) - DEVICE-LEVEL SCHEDULING ALGORITHM -



technology

from seed

• STEP 3 – DEVICE-LEVEL SCHEDULING ALGORITHM*/**



- According to the performance models for device computation (ψ_w) and bidirectional asymmetric full-duplex communication links (σ_ι , σ_o)
 - Step 3-I. Determination of the initial optimal distribution with three load fractions.

* Ilić, A., and Sousa, L., "Scheduling Divisible Loads on Heterogeneous Desktop Systems with Limited Memory" (submitted to HeteroPar 2011)

Determination of Load Fractions (12) - DEVICE-LEVEL SCHEDULING ALGORITHM -



technology

from seed

STEP 3 – DEVICE-LEVEL SCHEDULING ALGORITHM*/**



- According to the performance models for device computation (ψ_w) and bidirectional asymmetric full-duplex communication links (σ_ι , σ_o)
 - Step 3-I. Determination of the initial optimal distribution with three load fractions.
 - Step 3-II. Generate additional three-fraction distributions.

* Ilić, A., and Sousa, L., "Scheduling Divisible Loads on Heterogeneous Desktop Systems with Limited Memory" (submitted to HeteroPar 2011)

Determination of Load Fractions (13) - DEVICE-LEVEL SCHEDULING ALGORITHM -



technology

from seed

• STEP 3 – DEVICE-LEVEL SCHEDULING ALGORITHM*/**



- According to the performance models for device computation (ψ_w) and bidirectional asymmetric full-duplex communication links (σ_ι , σ_o)
 - Step 3-I. Determination of the initial optimal distribution with three load fractions.
 - Step 3-II. Generate additional three-fraction distributions.
 - Step 3-III. Insert additional load fractions into existing sub-distributions (iterative).

* Ilić, A., and Sousa, L., "Scheduling Divisible Loads on Heterogeneous Desktop Systems with Limited Memory" (submitted to HeteroPar 2011)

Determination of Load Fractions (14) - DEVICE-LEVEL SCHEDULING ALGORITHM -



technology

from seed

• STEP 3 – DEVICE-LEVEL SCHEDULING ALGORITHM*/**



- According to the performance models for device computation (ψ_w) and bidirectional asymmetric full-duplex communication links (σ_ι , σ_o)
 - Step 3-I. Determination of the initial optimal distribution with three load fractions.
 - Step 3-II. Generate additional three-fraction distributions.
 - Step 3-III. Insert additional load fractions into existing sub-distributions (iterative).
 - Step 3-IV. Generate new sub-distributions by restarting.

* Ilić, A., and Sousa, L., "Scheduling Divisible Loads on Heterogeneous Desktop Systems with Limited Memory" (submitted to HeteroPar 2011)

Determination of Load Fractions (15) - DEVICE-LEVEL SCHEDULING ALGORITHM -



technology

from seed

• STEP 3 – DEVICE-LEVEL SCHEDULING ALGORITHM*/**



- According to the performance models for device computation (ψ_w) and bidirectional asymmetric full-duplex communication links (σ_ι , σ_o)
 - Step 3-I. Determination of the initial optimal distribution with three load fractions.
 - Step 3-II. Generate additional three-fraction distributions.
 - Step 3-III. Insert additional load fractions into existing sub-distributions (iterative).
 - Step 3-IV. Generate new sub-distributions by restarting.
 - Step 3-V. Expand all sub-distributions.

* Ilić, A., and Sousa, L., "Scheduling Divisible Loads on Heterogeneous Desktop Systems with Limited Memory" (submitted to HeteroPar 2011)

Determination of Load Fractions (16) - DEVICE-LEVEL SCHEDULING ALGORITHM -



technology

from seed

• STEP 3 – DEVICE-LEVEL SCHEDULING ALGORITHM*/**



- According to the performance models for device computation (ψ_w) and bidirectional asymmetric full-duplex communication links (σ_ι , σ_o)
 - Step 3-I. Determination of the initial optimal distribution with three load fractions.
 - Step 3-II. Generate additional three-fraction distributions.
 - Step 3-III. Insert additional load fractions into existing sub-distributions (iterative).
 - Step 3-IV. Generate new sub-distributions by restarting.
 - Step 3-V. Expand all sub-distributions.
 - Step 3-VI. Select the distribution with maximum relative performance.

* Ilić, A., and Sousa, L., "Scheduling Divisible Loads on Heterogeneous Desktop Systems with Limited Memory" (submitted to HeteroPar 2011)

Case Study: 2D FFT Batch



DOUBLE FLOATING POINT COMPLEX 2D FFT BATCH EXECUTION

- Size: 1024 times 512 × 512; divisible in the first dimension
- The optimal vendor-provided FFT implementations are used
 - NVIDIA's CUFFT 3.2 for the GPU and Intel MKL 10.3 for the CPU

DISTRIBUTED HETEROGENEOUS SYSTEM WITH 4 CPU+GPU NODES

| Experimental Setup | CPU | GPU | | #CPU Cores | #GPUs |
|--------------------|-------------------|--------------------------|------------|------------|-------|
| | Intel Core 2 Quad | nVIDIA GeForce 285GTX | NODE 1 | 3 | |
| | | | NODE 2 | 2 | |
| Speed/Core (GHz) | 2.83 | 1.476 | NODE 3 | 1 | I |
| Global Memory (MB) | 4096 | 1024 | NODE 4 | - | |

ITERATIVE PROCEDURE FOR ONLINE PERFORMANCE MODELING

- PERFORMANCE ESTIMATION of all heterogeneous devices DURING THE EXECUTION
 - No prior knowledge on the performance of an application is available on any of the devices
- **INITIALLY,** the load is distributed among nodes/devices using FACTORING-BY-TWO STRATEGY*
 - Limited Memory: Factoring-by-two partitioning of the largest loads into new sub-distributions until satisfying the memory limitations
- IN EACH FOLLOWING ITERATION, the load is distributed using the PRESENTED APPROACH

* Ilić, A., and Sousa, L., "Algorithm For Divisible Load Scheduling on Heterogeneous Systems with Realistic Performance Models", Tech. rep., INESC-ID (May 2011)

technology

Case Study: 2D FFT Batch ITERATION 1: FACTORING-BY-TWO STRATEGY



technology

from seed





Instituto de Engenharia de Sistemas e Computadores Investigação e Desenvolvimento em Lisboa

6/28/2011

High Performance Computing, Grids and Clouds - HPC 2011





Case Study: 2D FFT Batch ITERATION 2: PERFORMANCE



technology





Instituto de Engenharia de Sistemas e Computadores Investigação e Desenvolvimento em Lisboa

Node 4:

 $\alpha_{4} = 445$

technology

from seed

Case Study: 2D FFT Batch ITERATION 3





Instituto de Engenharia de Sistemas e Computadores Investigação e Desenvolvimento em Lisboa

6/28/2011

Case Study: 2D FFT Batch ITERATION 4





Conclusions



DYNAMIC LOAD BALANCING

- OBTAINED IN 4 ITERATIONS AND 6.1 SECONDS

TRADITIONAL APPROACHES FOR PERFORMANCE MODELING

- Approximate the performance using number of points equal to the number of iterations
- In this case, **40** POINTS in total

PRESENTED DDL SCHEDULING APPROACH

- Models the performance using 416 POINTS, in this case $\sim 10x$ more than with traditional modeling
- Load balancing solution is 3x faster than the current state of the art approaches
- IN THIS CASE, OBTAINED GPU PERFORMANCE IS AT LEAST 4.1X BETTER THAN THE "OPTIMAL" CUFFT EXECUTION

Instituto de Engenharia de Sistemas e Computadores Investigação e Desenvolvimento em Lisboa

technology

Questions?

Thank you

