A Multi-objective Approach for the Multi-level Scheduling of Large Workloads in Multicore Distributed Systems

Santiago Iturriaga¹, <u>Bernabé Dorronsoro</u>², Andrei Tchernykh³, Sergio Nesmachnow¹, Pascal Bouvry⁴

¹ Univ. de la República - UY
² University of Cadiz - ES
³ CICESE Ensenada - MX
⁴ University of Luxembourg - LU

B. Dorronsoro, S. Nesmachnow, J. Taheri, A. Y. Zomaya, E.-G. Talbi, P. Bouvry, A hierarchical approach for energy-efficient scheduling of large workloads in multicore distributed systems, *Sustainable Computing: Informatics and Systems*, 4(4):252-261, 2014.

Introduction

- Resource management in large-scale geographically distributed computing centers
 - Big supercomputers
 - High performance computing centers
 - Cloud infrastructures
- Minimize energy consumption in conflict with
 - Performance
 - QoS (overdue deadlines)
- Two-level approach
 - Higher level: decides the mapping between jobs and data-centers
 - Lower level: schedule jobs within each data-center with energy consumption considerations
- Use of multi-objective algorithms to analyze the quality of our two-level schedulers

The Problem



Makespan Minimization

$$f_M(\vec{x}) = \max_{0 \le p \le k} CT_p$$

 \vec{x} represents an allocation k is the number of available CNs CT_p is the completion time of CN_p

Energy Minimization



S. Nesmachnow, B. Dorronsoro, J. Pecero, P. Bouvry, Energy-aware scheduling on multicore heterogeneous grid computing systems, The Journal of Grid Computing 11(4):653-680, 2013.

Penalization Cost Minimization

$$f_C(\vec{x}) = \sum_{0 \le q \le n} Cost(\vec{x}, q)$$

$$Cost(\vec{x},q) = \begin{cases} 0 & \text{if } F_q < D_q \\ Penalty_q(F_q) & \text{otherwise} \end{cases}$$

$$Penalty_q(F_q) = \begin{cases} SQRT = \sqrt{F_q - D_q} \\ LIN = F_q - D_q \\ SQR = (F_q - D_q)^2 \end{cases}$$

Multi-objective Optimization

- The goal is to find the Pareto front
- Two key features to measure the quality of solutions
 - Convergence
 - Diversity



Algorithms

NSGA-II

- Reference algorithm
- Panmictic population
- Selection of solutions
 - Ranking
 - Crowding

MOCell

- Cellular population
 - Only next individuals can interact
- External archive
 - Feedback to population

NSGA-II

- Non-dominated Sorting Genetic Algorithm
- Proposed by K. Deb (2002)
- The most popular metaheuristic for multi-objective optimization
- Features
 - Ranking using non-dominated sorting
 - Crowding distance as density estimator



Point B is in a less crowded region than point A 9

NSGA-II



MOCell







- I25 randomly generated workflow batches
 - More than 1000 jobs
 - CNs: between 8 and 64 processors; I to 6 cores each



Fig 3. Jobs' shapes: (a) Serious-Parallel, (b Homogenous-Parallel, (c) Heterogeneous-Parallel, and (d) Single-Task

Experiments

- NSGAII and MOCell
 - Map jobs into CNs and establish scheduling order
 - Lower level: EFTH
 - From previous work
- Average results on 25 instances
- Number of independent runs: 30
- Statistical test: Wilcoxon

processor	frequency	cores	GFLOPS	E_{IDLE}	E_{MAX}	GFLOPS/core
Intel Celeron 430	$1.80~\mathrm{GHz}$	1	7.20	75.0W	94.0W	7.20
Intel Pentium E5300	$2.60~\mathrm{GHz}$	2	20.80	68.0W	109.0W	10.40
Intel Core i7 870	$2.93~\mathrm{GHz}$	4	46.88	76.0W	214.0W	11.72
Intel Core i5 661	$3.33~\mathrm{GHz}$	2	26.64	74.0W	131.0W	13.32
Intel Core i7 980 XE	$3.33~\mathrm{GHz}$	6	107.60	102.0W	210.0W	17.93

Lower Level Scheduler: EFTH

- Based on HEFT
 - Adapted for multicore case: backfilling
- The whole that minimizes task finish time is selected



14

Problem Resolution

- Problem:
 - Allocate jobs to computing nodes
 - Specify the jobs execution order in every CN
 - Evaluation based on EFTH
- Problem representation
 - Permutation of size N_{jobs} + N_{CN} I
 - Job IDs \in [0, N_{jobs}-I]
 - Splitters \in [N_{jobs}, N_{CN}-2]



Operators

• Partially Matched Crossover (PMX)

Parent 1	Offspring 1
4 - 5 - 2 - 10 - <mark>0</mark> - 3 - 1 - 12 - 7 - 8 - 9 - <mark>11</mark> - 6	4 - 5 - 10 - 8 - <mark>0</mark> - 7 - 2 - 9 - 12 - 1 - 3 - <mark>11</mark> - 6
Parent 2	Offspring 2
5 - 0 - 6 - 4 - 10 - 7 - 2 - 9 - 12 - 1 - 3 - 8 - 11	5 - 0 - 6 - 4 - <mark>10</mark> - 3 - 1 - 12 - 7 - 8 - 9 - <mark>8</mark> - 11

• Swap Mutation



- Solutions Repair: can all CNs execute all assigned jobs?
 - Number of cores

Comparison of the Performance of MOEAs

Instance	HV		spread		IGD	
	MOCell	NSGA-II	MOCell	NSGA-II	MOCell	NSGA-II
Heterogeneous	$0,43\pm0,06$	$0,59{\pm}0,07$	$1,00\pm0,02$	$1,32{\pm}0,03$	$0,06\pm0,03$	$0,04{\pm}0,01$
Homogeneous	$0,57\pm0,07$	$0,65{\pm}0,06$	$0,80{\pm}0,05$	$0,90{\pm}0,02$	$0,04{\pm}0,01$	$0,03{\pm}0,01$
Serious Par.	$0,57\pm0,07$	$0,65{\pm}0,06$	$0,95{\pm}0,03$	$1,35{\pm}0,02$	$0,05\pm0,02$	$0,04{\pm}0,01$
Single Task	$0,57\pm0,07$	$0,65{\pm}0,06$	$0,87\pm0,02$	$0,98{\pm}0,04$	$0,04{\pm}0,01$	$0,03{\pm}0,01$
Mix	$0,49\pm0,01$	$0,58{\pm}0,09$	$0,96{\pm}0,04$	$1,08{\pm}0,04$	$0,04{\pm}0,01$	$0,03{\pm}0,01$

Sample Solutions for Heterogeneous Inst.



Comparison of MOEAs vs Heuristic



Conclusions

- Resource management in large geographically distributed datacenters with multi-core architectures
- Multiple objectives: Energy; Makespan; QoS
- Solution techniques
 - Two-levels schedule
 - Two MOEAs
- Results
 - NSGAII outperforms MOCell (accuracy)
 - Both outperform best known heuristic
 - Slightly worse in makespan (2.5%)
 - Better in energy (5%) and penalizations cost (50%)
- Future work
 - Enhance MOEAs with population initialization and LS
 - Design better heuristics

Thank you