



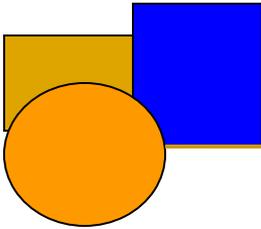
**PUC Minas**  
**Guanhanes**



# **Análise de Algoritmos de Agrupamento e Classificação na Modelagem de Comportamentos de Tarefas Paralelas**

**Lesandro Ponciano do Santos**  
[lesandrop@yahoo.com.br](mailto:lesandrop@yahoo.com.br)

*15 de outubro de 2008*



# Apresentação

---

- PUC Minas em Guanhães
- Curso Bacharelado em Sistemas de Informação
  - Aluno: Lesandro Ponciano dos Santos
  - Orientador: Prof. João Paulo D. Silva
  - Co-orientador: Prof. Luís Fabrício W. Góes
- Projeto voluntário, sem financiamento

- Introdução
  - Contexto
  - Problema
  - Proposta
  - Objetivos
- Trabalhos Relacionados
- Algoritmos de Classificação e Agrupamento
- Planejamento dos Experimentos
- Resultados
- Conclusão

## ■ Contexto

- Aplicações que exigem alto poder de processamento
- *Reconfigurable Gang Scheduling Algorithm (RGSA)* (Góes e Martins, 2005)
  - Tempo de submissão, tempo de execução e número de processos
- Caracterização de Cargas de Trabalho
- Algoritmo baixo (L - *low*) e alto (H - *high*), 4 grupos possíveis (HH, HL, LL, LH)

## ■ Contexto

### Algoritmo – Classificação *Low* e *high\_*

```
se job.n_processos ≤ mediana_n_processos
  se job.tempo_exe ≤ mediana_tempo_exe
    então job.classe = LL;
  senão job.classe = LH;
```

```
senão se job.tempo_exe ≤ mediana_tempo_exe
  então job.classe = HL;
  senão job.classe = HH;
```

## ■ Problema

- O algoritmo de Classificação *Low* e *high* é muito sensível à variações da mediana e se preocupa em minimizar o desvio entre as tarefas de um grupo

## ■ Proposta

- Utilização de algoritmos clássicos de classificação e agrupamento por similaridade

## ■ Objetivos

- Analisar a aplicação dos algoritmos, de agrupamento, *k-means* e, de classificação, *j48*

- **Góes e Martins, 2005**

- Proposta e desenvolvimento do escalonador RGSA em ambiente ideal

- **Santos e Góes, 2007**

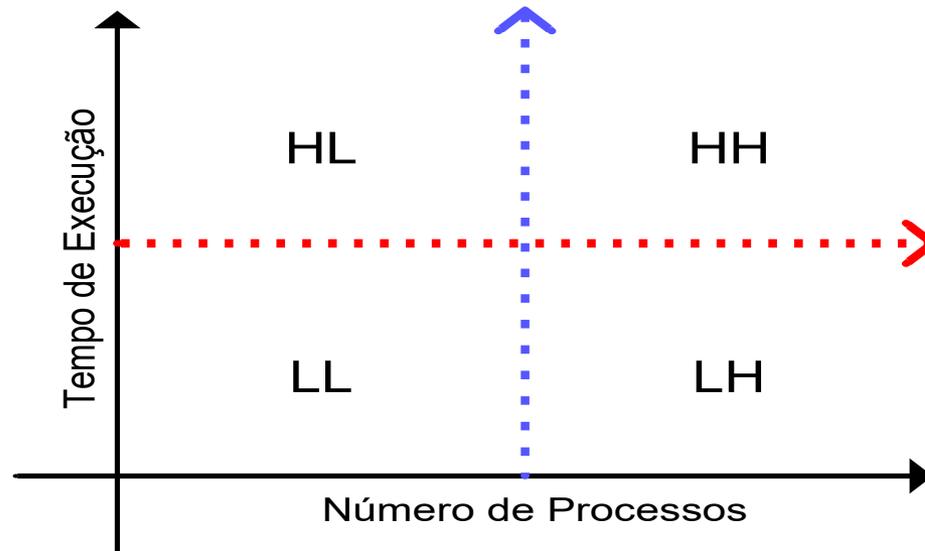
- Técnica de Caracterização de log para extrair informações confiáveis para o escalonador RGSA. (Utilizando classificação *low* e *high*)

- **Feitelson, 2007**

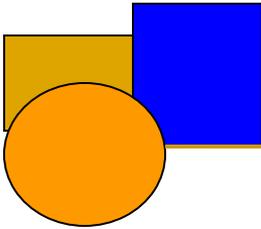
- Constatação que existem padrões na localidade e demanda de serviços em requisições a máquinas paralelas

# Algoritmos de Classificação e Agrupamento

## Agrupamento *low e high*



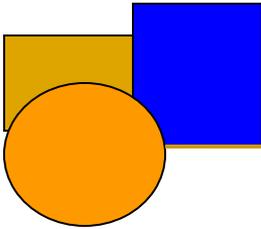
- > Mediana (Tempo de Execução)
- > Mediana (Número de Processos)



# Algoritmos de Classificação e Agrupamento

## Agrupamento *K-means* e Classificação *J48*

- Dado um rastro, com dados de tarefas paralelas, executa-se o algoritmo *k-means* com  $k=4$  e identificam-se quatro grupos (G0, G1, G2 e G3), gerados pela similaridade entre as tarefas
- Executa-se o algoritmo *j48* para identificar as regras que definem os agrupamentos
- As regras obtidas são utilizadas na classificação e predição do comportamento de tarefas futuras



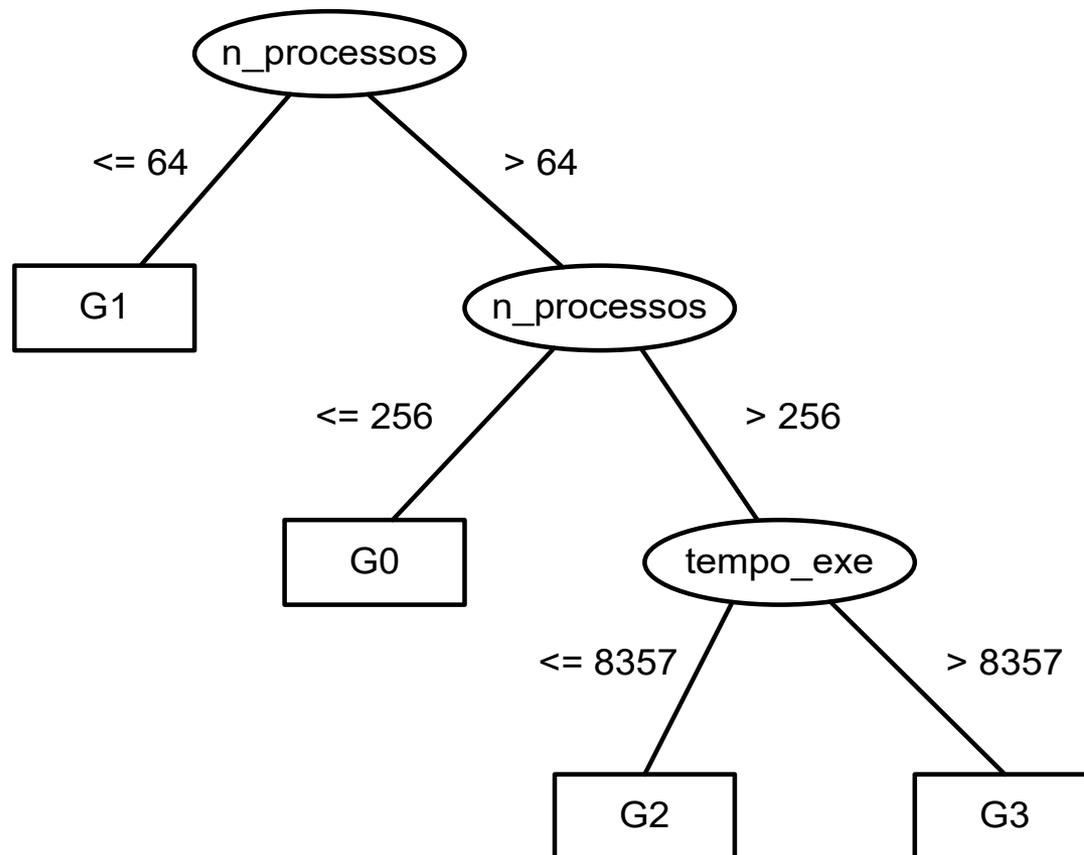
# Algoritmos de Classificação e Agrupamento

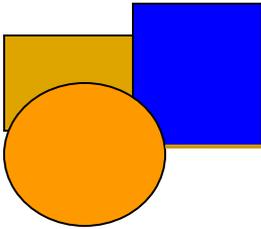
## Exemplo de Regras de classificação (LANCM5 outubro/1995)

```
se job.n_processos <= 64
    então job.classe = G1;
senão
    se job.n_processos <= 256
        então job.classe = G0;
    senão
        se job.tempo_exe <= 8357
            então job.classe = G2;
            senão job.classe = G3;
```

# Algoritmos de Classificação e Agrupamento

## Exemplo de árvore de decisão (LANCM5 outubro/1995)





# Planejamento dos Experimentos

- São utilizadas tarefas submetidas, ao longo de sete meses, em quatro supercomputadores reais (HPC2N, SDSC, SHARCNET, LANCM5)
- Para predição de comportamentos, caracteriza-se um mês e testa-se a capacidade de predição para o mês subsequente

## Análise do número de tarefas caracterizadas

	Número de Tarefas						
	Abr.	Mai.	Jun.	Jul.	Ago.	Set.	Out.
<b>HPC2N</b>	29970	17817	20545	24121	7241	5159	11444
<b>SDSC</b>	5259	8395	8358	5339	7204	8782	9147
<b>SHARC.</b>	87505	76813	52551	95118	90225	63162	158681
<b>LANCM5</b>	5790	6831	6216	6969	6126	6303	5271

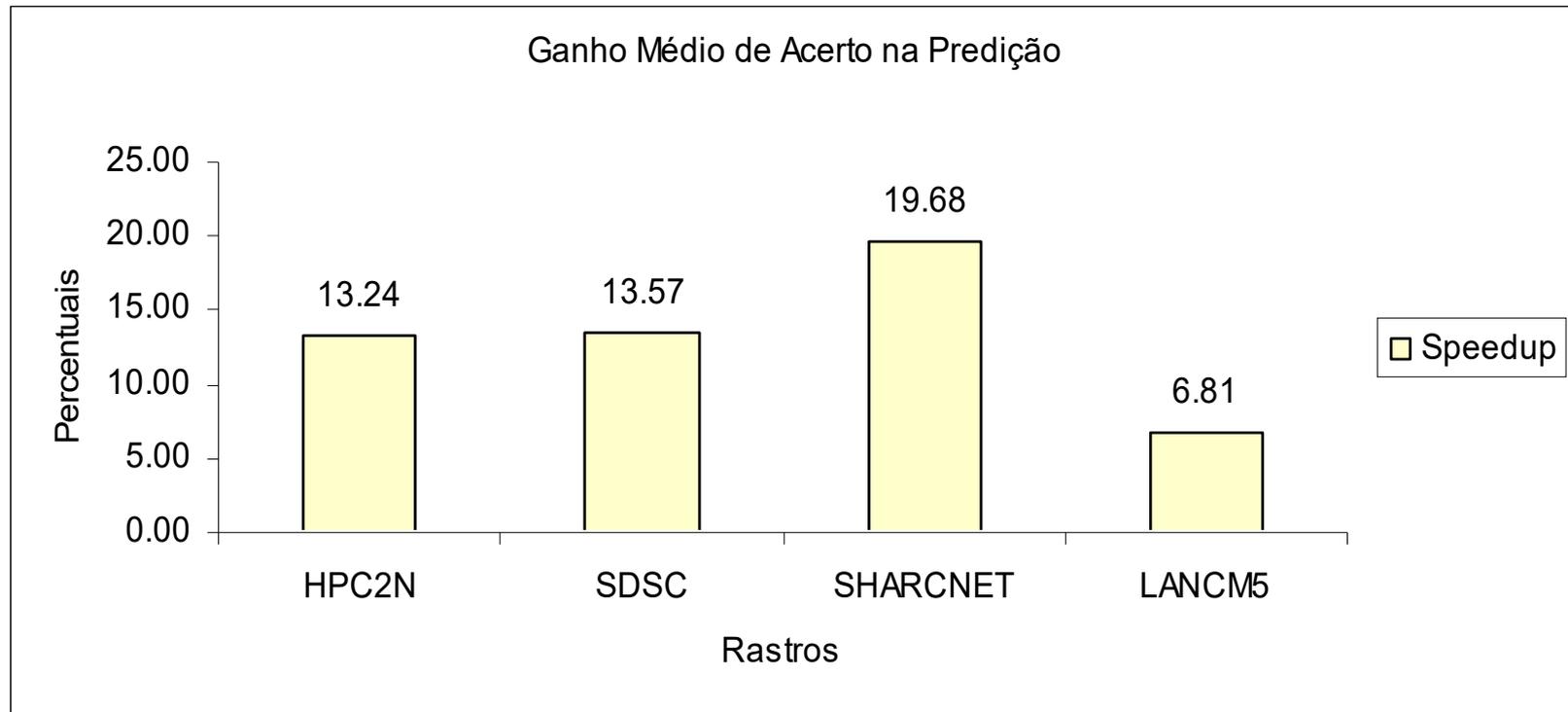
## Teste de Predição SDSC 2005 (algoritmo *low* e *high*)

Meses	Dias da Semana						
	D	S	T	Q	Q	S	S
mai.	73	56	43	65	63	58	65
jun.	77	62	71	71	71	66	71
jul.	46	57	53	56	69	49	55
ago.	55	51	58	57	52	49	58
set.	65	61	63	57	56	66	60
out.	50	63	76	74	67	61	64

## Teste de Predição SDSC 2005 (algoritmo *k-means* e *j48*)

Meses	Dias da Semana						
	D	S	T	Q	Q	S	S
mai.	76	84	70	82	76	82	82
jun.	76	86	90	92	94	94	90
jul.	66	76	76	80	80	82	74
ago.	76	84	82	66	70	70	70
set.	90	86	86	68	70	76	74
out.	78	90	82	92	90	88	86

## Ganho médio do algoritmo *k-means* e *j48* em relação ao algoritmo *low* e *high*, nos testes de precisão

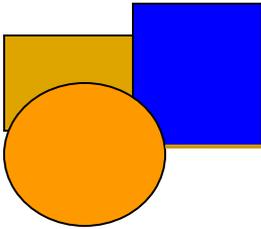


- Os algoritmos *k-means* e *j48* geram, na média, maior percentual de acerto na predição
- A implementação e a estrutura de informações do algoritmo *low* e *high* é mais simples que a do algoritmos *k-means* e *j48*
- Os algoritmos *k-means* e *j48* extraem mais informação (regras), das bases de dados, em relação ao algoritmo *low* e *high* (Mediana)

## ■ Contribuições

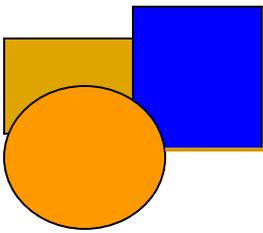
- Proposta, implementação e análise de um método de modelagem de carga de trabalho de máquinas paralelas por atributos de similaridade entre tarefas

- **Publicações de resultados preliminares**
  - Artigo de iniciação científica publicado e premiado como um dos três melhores trabalhos apresentados ao WSCAD-CTIC. Gramado/RS, outubro de 2007
  - Artigo aceito para publicação no WSCAD-CTIC. Campo Grande/MS, outubro, novembro de 2008
  - Dois artigos em processo de avaliação em revistas e congressos



# Principais Referências

- Feitelson, D. **Locality of Sampling and Diversity in Parallel System Workloads**. ACM SIGMETRICS 2007 pp. 53-63..
- GÓES, L. F. W. e MARTINS, C. A. P. S., **Reconfigurable Gang Scheduling Algorithm**, 10th Workshop on Job Scheduling Strategies for Parallel Processing, Lecture Notes in Computer Science, New York. 2005.
- SANTOS, L. P. e GÓES, L. F. W. **Técnica de Caracterização de Cargas de Trabalho para Extração de Informações Utilizadas pelo Escalonador Reconfigurável de Tarefas**. Workshop de Sistemas Computacionais de Alto Desempenho (WSCAD), 2007.
- SANTOS, L. P. e GÓES, L. F. W. **Descoberta e Predição de Comportamento de Tarefas Paralelas através da Caracterização de Padrões em Cargas de Trabalho**. Workshop de Sistemas Computacionais de Alto Desempenho, Concurso de Trabalhos de Iniciação Científica WSCAD-CTIC 2008. Anais Digitais, Campo Grande/MS, 2008. (Artigo aceito, publicação prevista para outubro 2008).
- **Download Workloads** <http://www.cs.huji.ac.il/labs/parallel/workload/logs.html> último acesso em: 14 out. 2008.
- Lee, C.B., et al. **“Are user runtime estimates inherently inaccurate?”**, 10th Workshop on Job Scheduling Strategies for Parallel Processing, 2004.



# Perguntas...