



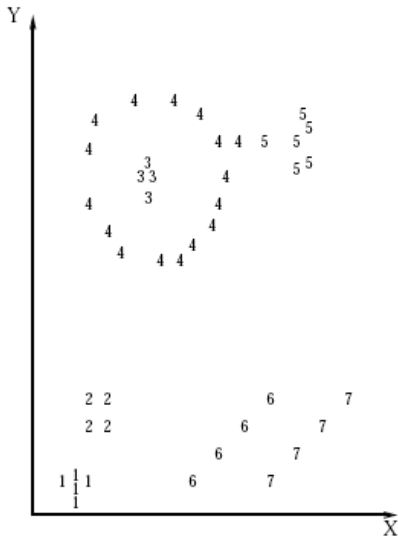
Universidade Estadual de Campinas - UNICAMP
Instituto de Computação - IC

Busca por Similaridade em Espaços Métricos utilizando Técnicas de Agrupamento de Dados

Jurandy Gomes de Almeida Junior

Campinas, 8 de outubro de 2008.

Problema



- Busca seqüencial
- Métodos de acesso
 - Espaciais
 - Métricos

- Busca seqüencial
- Métodos de acesso
 - Espaciais
 - Métricos
 - Estáticos
 - Dinâmicos

- Busca seqüencial
- Métodos de acesso
 - Espaciais
 - Métricos
 - Estáticos
 - Dinâmicos

- Busca seqüencial
- Métodos de acesso
 - Espaciais
 - Métricos
 - Estáticos
 - Dinâmicos

- Busca seqüencial
- Métodos de acesso
 - Espaciais
 - Métricos
 - Estáticos
 - Dinâmicos

- Busca sequencial
- Métodos de acesso
 - Espaciais
 - Métricos
 - Estáticos
 - **Dinâmicos**

- Espaço métrico
 - Positividade
 - Simetria
 - Identidade
 - Desigualdade triangular
- Funções de distância
 - Discreta
 - Contínua
- Consultas por similaridade
 - Abrangência
 - Vizinhos próximos
 - Vizinhos próximos reverso
 - Junção
 - Combinação

- Espaço métrico
 - **Positividade**
 - Simetria
 - Identidade
 - Desigualdade triangular
- Funções de distância
 - Discreta
 - Contínua
- Consultas por similaridade
 - Abrangência
 - Vizinhos próximos
 - Vizinhos próximos reverso
 - Junção
 - Combinação

- Espaço métrico
 - Positividade
 - **Simetria**
 - Identidade
 - Desigualdade triangular
- Funções de distância
 - Discreta
 - Contínua
- Consultas por similaridade
 - Abrangência
 - Vizinhos próximos
 - Vizinhos próximos reverso
 - Junção
 - Combinação

- Espaço métrico
 - Positividade
 - Simetria
 - **Identidade**
 - Desigualdade triangular
- Funções de distância
 - Discreta
 - Contínua
- Consultas por similaridade
 - Abrangência
 - Vizinhos próximos
 - Vizinhos próximos reverso
 - Junção
 - Combinação

- Espaço métrico
 - Positividade
 - Simetria
 - Identidade
 - Desigualdade triangular
- Funções de distância
 - Discreta
 - Contínua
- Consultas por similaridade
 - Abrangência
 - Vizinhos próximos
 - Vizinhos próximos reverso
 - Junção
 - Combinação

- Espaço métrico
 - Positividade
 - Simetria
 - Identidade
 - Desigualdade triangular
- Funções de distância
 - Discreta
 - Contínua
- Consultas por similaridade
 - Abrangência
 - Vizinhos próximos
 - Vizinhos próximos reverso
 - Junção
 - Combinação

- Espaço métrico
 - Positividade
 - Simetria
 - Identidade
 - Desigualdade triangular
- Funções de distância
 - Discreta
 - Contínua
- Consultas por similaridade
 - Abrangência
 - Vizinhos próximos
 - Vizinhos próximos reverso
 - Junção
 - Combinação

- Espaço métrico
 - Positividade
 - Simetria
 - Identidade
 - Desigualdade triangular
- Funções de distância
 - Discreta
 - **Contínua**
- Consultas por similaridade
 - Abrangência
 - Vizinhos próximos
 - Vizinhos próximos reverso
 - Junção
 - Combinação

- Espaço métrico
 - Positividade
 - Simetria
 - Identidade
 - Desigualdade triangular
- Funções de distância
 - Discreta
 - Contínua
- Consultas por similaridade
 - Abrangência
 - Vizinhos próximos
 - Vizinhos próximos reverso
 - Junção
 - Combinação

- Espaço métrico
 - Positividade
 - Simetria
 - Identidade
 - Desigualdade triangular
- Funções de distância
 - Discreta
 - Contínua
- Consultas por similaridade
 - **Abrangência**
 - Vizinhos próximos
 - Vizinhos próximos reverso
 - Junção
 - Combinação

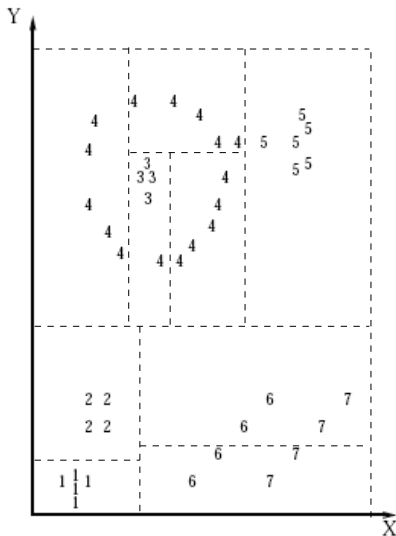
- Espaço métrico
 - Positividade
 - Simetria
 - Identidade
 - Desigualdade triangular
- Funções de distância
 - Discreta
 - Contínua
- Consultas por similaridade
 - Abrangência
 - Vizinhos próximos
 - Vizinhos próximos reverso
 - Junção
 - Combinação

- Espaço métrico
 - Positividade
 - Simetria
 - Identidade
 - Desigualdade triangular
- Funções de distância
 - Discreta
 - Contínua
- Consultas por similaridade
 - Abrangência
 - Vizinhos próximos
 - Vizinhos próximos reverso
 - Junção
 - Combinação

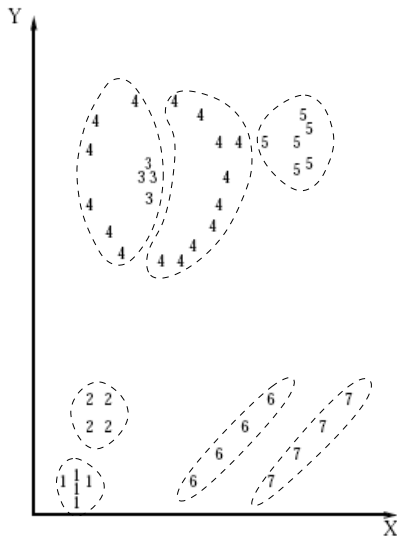
- Espaço métrico
 - Positividade
 - Simetria
 - Identidade
 - Desigualdade triangular
- Funções de distância
 - Discreta
 - Contínua
- Consultas por similaridade
 - Abrangência
 - Vizinhos próximos
 - Vizinhos próximos reverso
 - **Junção**
 - Combinação

- Espaço métrico
 - Positividade
 - Simetria
 - Identidade
 - Desigualdade triangular
- Funções de distância
 - Discreta
 - Contínua
- Consultas por similaridade
 - Abrangência
 - Vizinhos próximos
 - Vizinhos próximos reverso
 - Junção
 - Combinação

Estruturas de indexação



Agrupamento de dados



- Particional
- Hierárquico
 - Aglomerativo
 - Divisivo

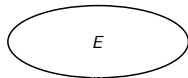
- Particional
- Hierárquico
 - Aglomerativo
 - Divisivo

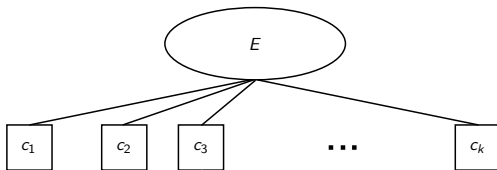
- Particional
- Hierárquico
 - Aglomerativo
 - Divisivo

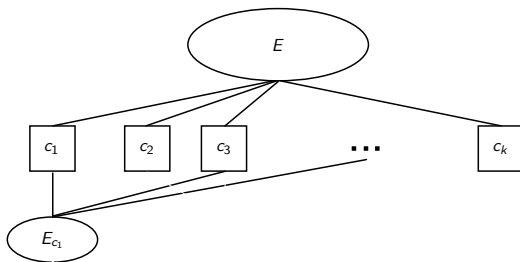
- Particional
- Hierárquico
 - Aglomerativo
 - Divisivo

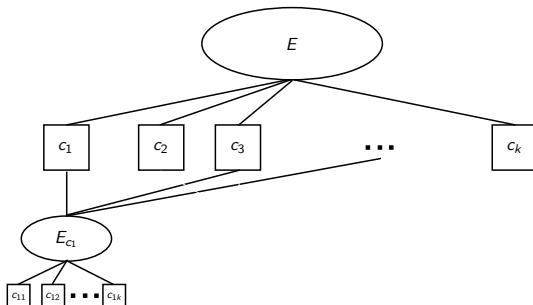
DAHC = Hierárquico Divisivo +
Hierárquico Aglomerativo +
Fator de Reagrupamento

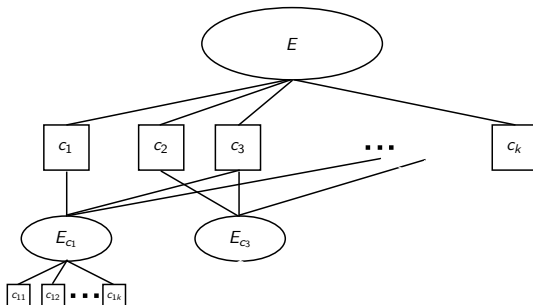
Símbolo	Significado
c, c_{rep}, c_{child}	Um grupo, o elemento representativo do grupo c e um apontador para o nível inferior de c na hierarquia de grupos.
C_i	O conjunto de grupos no i -ésimo nível.
k	O número de grupos em cada tarefa de agrupamento.
$f \in [0, 1)$	O fator de reagrupamento.
E_i	O conjunto de elementos sob análise no i -ésimo nível.
D	Uma medida para avaliar as dissimilaridades entre os elementos em E .

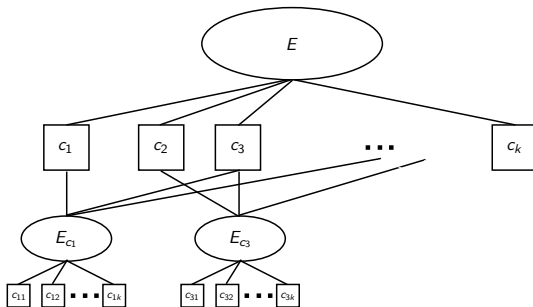


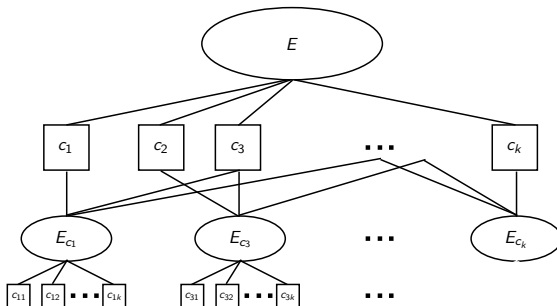


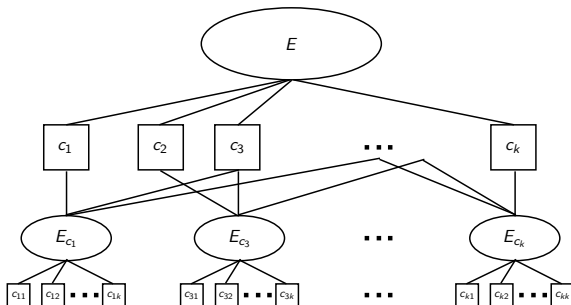


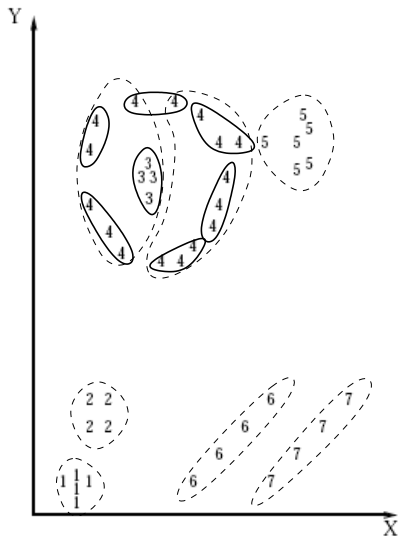












Entrada: O número de grupos k , o fator de reagrupamento $f \in [0, 1)$, o conjunto de elementos E e a função de distância D .

```

1: procedimento DAHC( $k, f, E, D$ )
2:    $C \leftarrow \text{CLUSTER}(k, E, D)$                                 ▷ Etapa divisiva
3:   para  $c \in C$  faça
4:      $C^* \leftarrow \lfloor f \times k \rfloor$  grupos mais próximos de  $c \in C$ 
5:      $E^* \leftarrow \{\}$ 
6:     para  $c^* \in C^*$  faça                                       ▷ Etapa aglomerativa
7:        $E^* \leftarrow E^* \cup c^*$ 
8:     fim para
9:     se  $\text{cardinalidade}(E^*) > k$  então                             ▷ Aprofundamento
10:       $c_{child} \leftarrow \text{DAHC}(k, f, E^*, D)$ 
11:    fim se
12:  fim para
13: fim procedimento

```

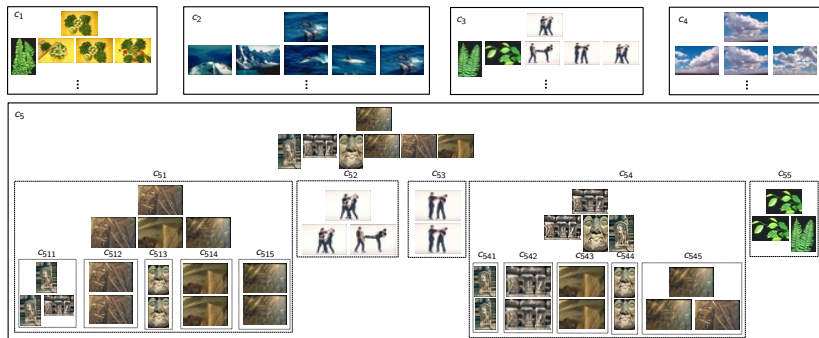



Figura: Exemplo da estrutura hierárquica criada pelo DAHC.

- Consultas por exemplo
- Bancos de dados
 - Corel RRsets
 - FreeFoto
- Validação cruzada
- Medidas
 - Precisão após 30 imagens
 - Número de comparações

- Consultas por exemplo
- Bancos de dados
 - Corel RRsets
 - FreeFoto
- Validação cruzada
- Medidas
 - Precisão após 30 imagens
 - Número de comparações

Experimentos

- Consultas por exemplo
- Bancos de dados
 - Corel RRSets
 - FreeFoto
- Validação cruzada
- Medidas
 - Precisão após 30 imagens
 - Número de comparações



Parada



Locomotivas



Montanhas



Surf



Peixes



Casas



Cervos



Árvores

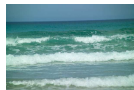


Cogumelos

- Consultas por exemplo
- Bancos de dados
 - Corel RRsets
 - FreeFoto
- Validação cruzada
- Medidas
 - Precisão após 30 imagens
 - Número de comparações



Árvores



Praias



Nuvens



Flores



Folhas



Poente



Luar



Montanhas



Rios

- Consultas por exemplo
- Bancos de dados
 - Corel RRsets
 - FreeFoto
- Validação cruzada
- Medidas
 - Precisão após 30 imagens
 - Número de comparações

- Consultas por exemplo
- Bancos de dados
 - Corel RRsets
 - FreeFoto
- Validação cruzada
- **Medidas**
 - Precisão após 30 imagens
 - Número de comparações

- Consultas por exemplo
- Bancos de dados
 - Corel RRsets
 - FreeFoto
- Validação cruzada
- Medidas
 - Precisão após 30 imagens
 - Número de comparações

- Consultas por exemplo
- Bancos de dados
 - Corel RRsets
 - FreeFoto
- Validação cruzada
- Medidas
 - Precisão após 30 imagens
 - Número de comparações

(a) Precisão após 30 imagens para cada descritor.

	BIC	GCH	CCV
<i>Corel RRsets</i>	53.0%	41.8%	40.7%
<i>FreeFoto</i>	68.1%	55.3%	59.9%

(b) Número de comparações.

<i>Corel RRsets</i>	<i>FreeFoto</i>
$\approx 21,000$	$\approx 48,000$

Tabela: Valores de referência utilizando uma busca sequencial.

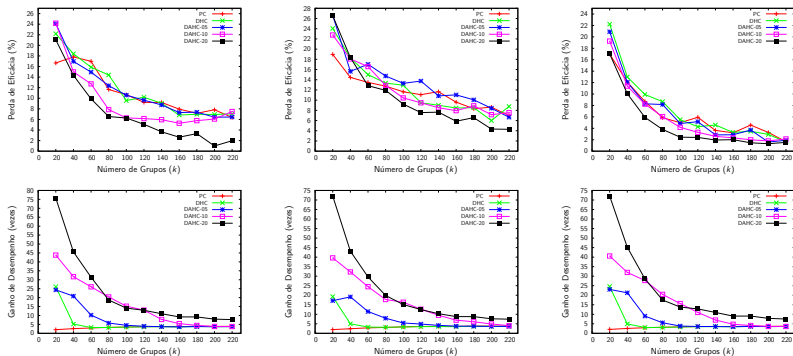


Figura: Perda de eficácia (acima) e ganho de desempenho (abaixo) para o banco *Corel RRsets*. Resultados dos descritores GCH, CCV e BIC, são mostrados, respectivamente, da esquerda para direita.

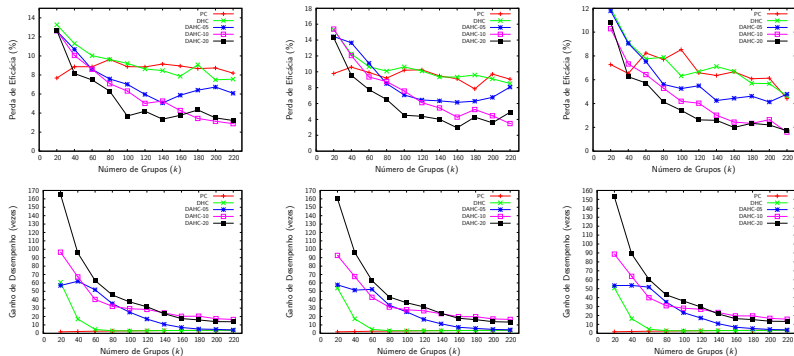


Figura: Perda de eficácia (acima) e ganho de desempenho (abaixo) para o banco *FreeFoto*. Resultados dos descritores GCH, CCV e BIC, são mostrados, respectivamente, da esquerda para direita.

- 1 Suporte a inserção e a remoção de dados
- 2 Suporte aos diversos tipos de consultas por similaridade
- 3 Suporte a otimização da estrutura depois de criada

Problemas sendo tratados

- 1 Suporte a inserção e a remoção de dados
- 2 Suporte aos diversos tipos de consultas por similaridade
- 3 Suporte a otimização da estrutura depois de criada

- 1 Suporte a inserção e a remoção de dados
- 2 Suporte aos diversos tipos de consultas por similaridade
- 3 Suporte a otimização da estrutura depois de criada

Dúvidas?

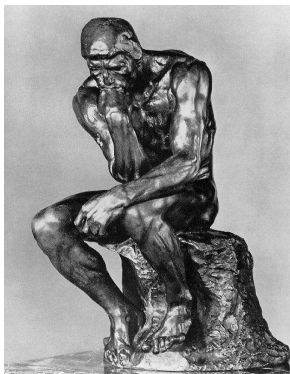


Figura: *O pensador* - Rodin

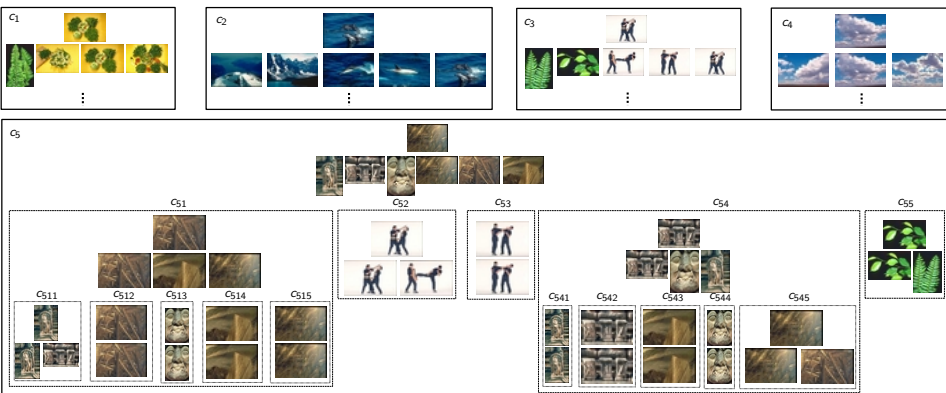


Figura: Exemplo da estrutura hierárquica criada pelo DAHC.



Parada



Locomotivas



Montanhas



Surf



Peixes



Casas



Cervos



Árvores



Cogumelos



Árvores



Praias



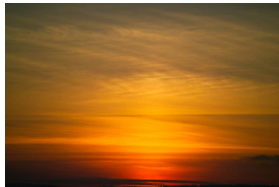
Nuvens



Flores



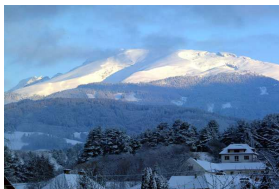
Folhas



Poente



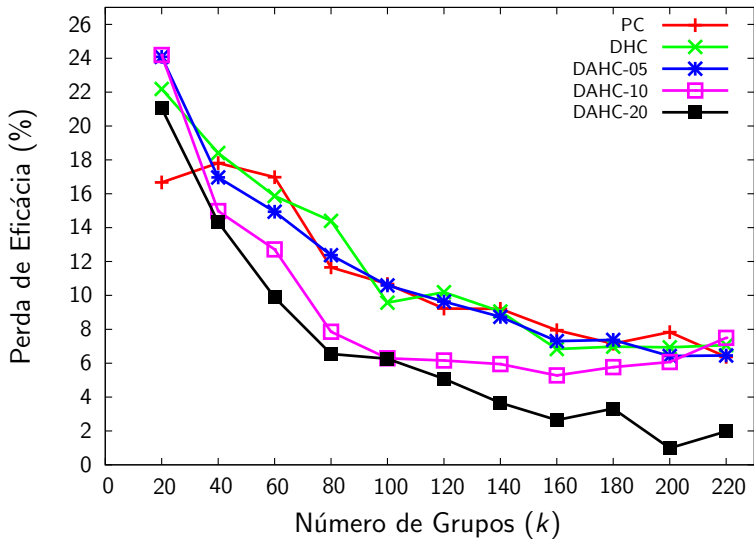
Luar



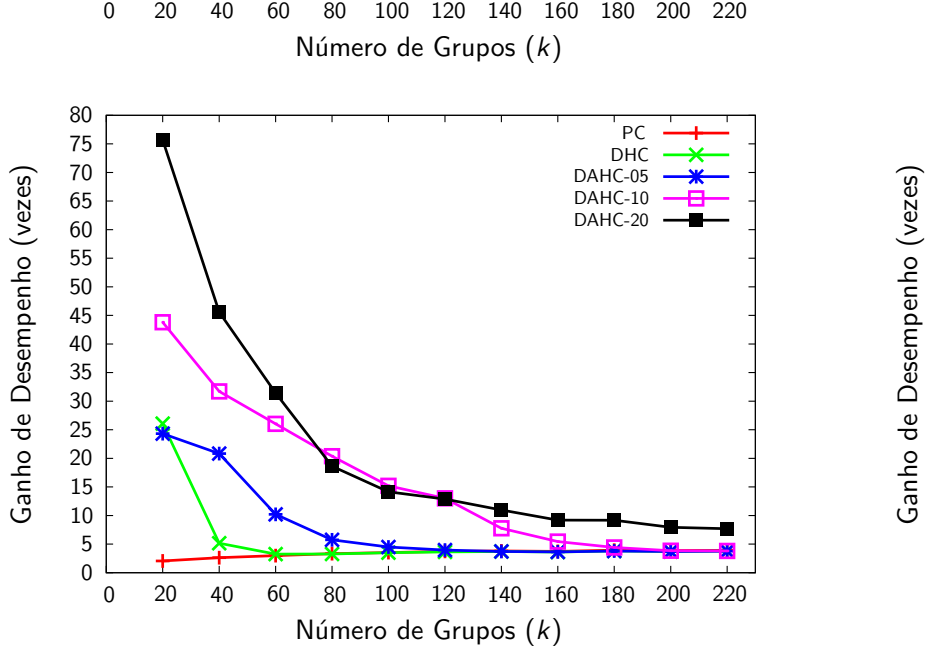
Montanhas

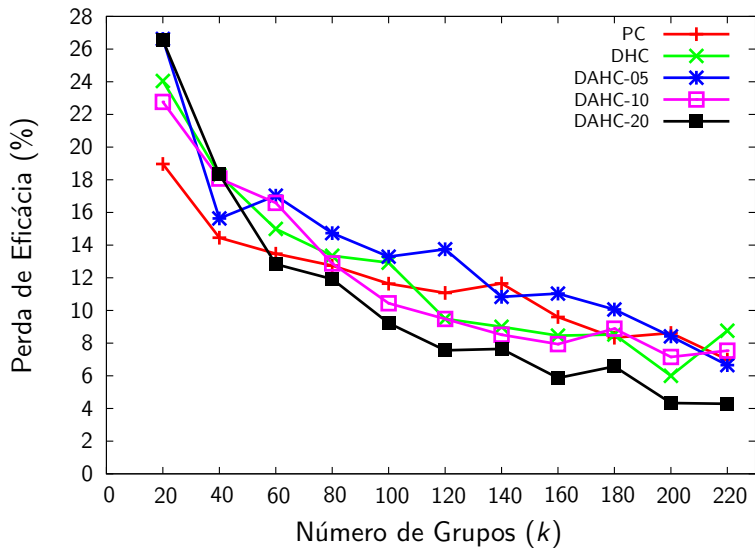


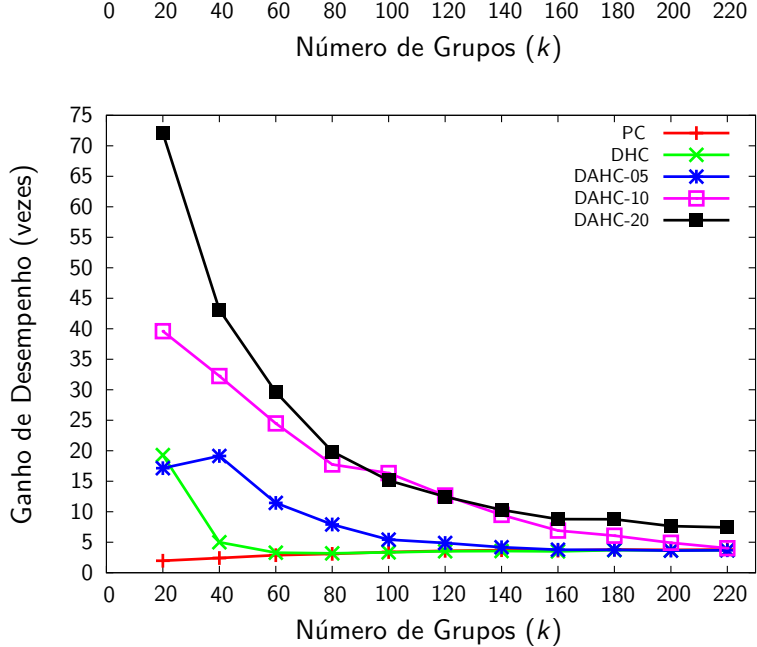
Rios

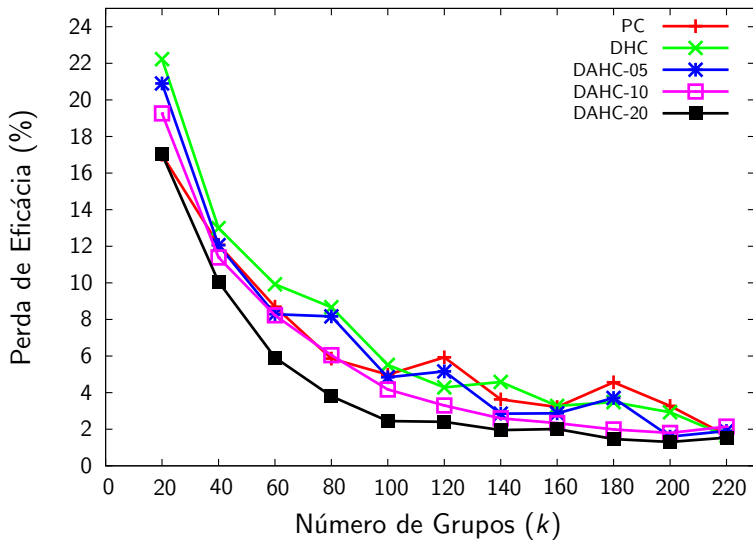


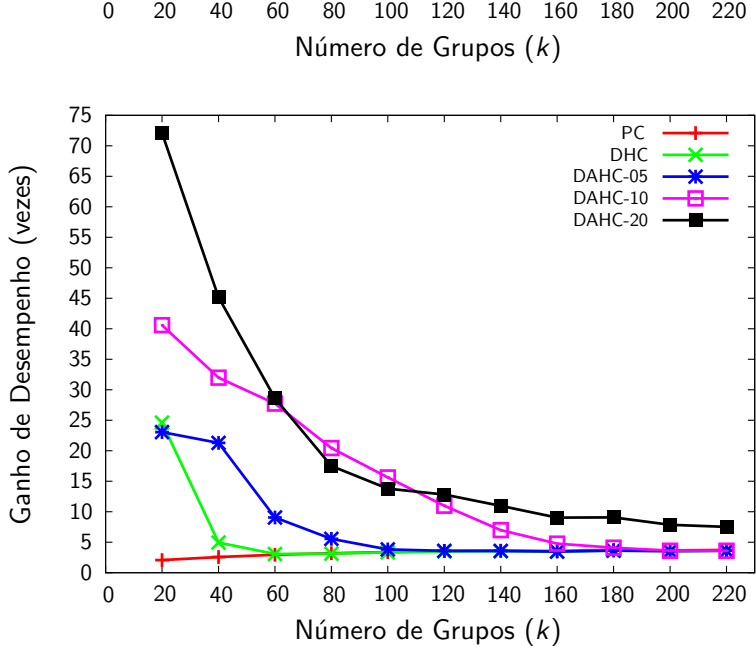
Perda de Eficácia (%)

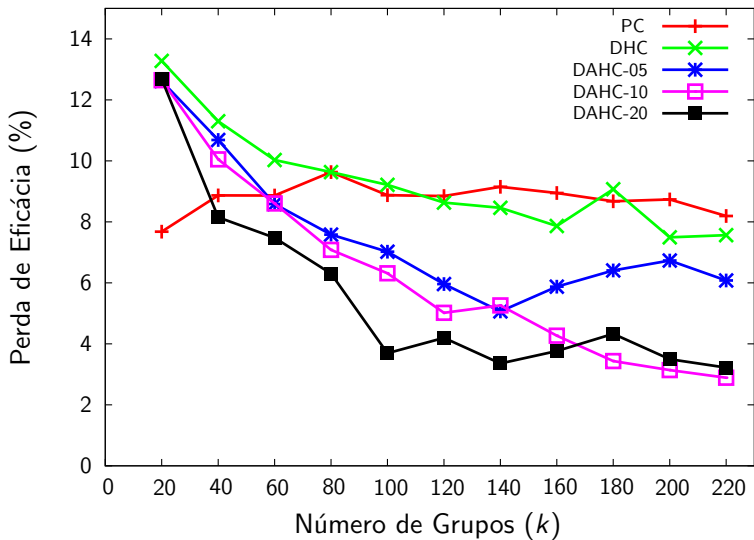












Perda de Eficácia (%)

