

TRABALHO PRÁTICO 4:

Problemas NP-Completo e Programação paralela

Sandro Miccoli - 2009052409 - smiccoli@dcc.ufmg.br

¹Departamento de Ciência da Computação – Universidade Federal de Minas Gerais (UFMG)

29 de dezembro de 2012

Resumo. *Este relatório descreve como foi modelado e implementado um problema de alto custo computacional. Será descrito como foi modelado o problema, as estruturas utilizadas e porquê não é possível implementar uma solução que consiga resolver entradas não triviais de maneira ótima. Finalmente será detalhado a análise de complexidade dos algoritmos e, por último, uma breve conclusão do trabalho implementado.*

1. INTRODUÇÃO

O trabalho propõe resolver um problema de uma rede social exclusiva para cozinheiros, o Facecook; para isso foi utilizado a teoria dos grafos para modelar o problema. Neste grafo, cada vértice representa um membro da rede, e as arestas representam uma conexão de amizade entre dois membros.

O que devemos encontrar nessa lista de membros, é o maior grupo de amigos registrados na rede social. Podemos perceber que isso se trata de um problema conhecido na teoria dos grafos, o problema da clique. Esse problema refere-se a qualquer problema que possui como objetivo encontrar subgrafos completos ("cliques") em um grafo. Um subgrafo é completo se existir uma aresta para todos os pares de vértices. Como exemplo, o problema de encontrar conjuntos de nós em que todos os elementos estão conectados entre si. [Wikipedia]

Para resolver a questão da clique utilizamos o algoritmo Bron-Kerbosch, que é utilizado para encontrar cliques máximos em grafos não direcionados. Iremos detalhar mais sobre seu funcionamento na Seção 3.

Esse problema é **NP-Completo**, pois não pode ser resolvido de forma ótima para entradas grandes. Ou seja, não existe um algoritmo de tempo polinomial para resolver esse problema, e isso que iremos mostrar neste trabalho.

Além disso, para entradas muito grandes o algoritmo sequencial não resolveria em tempo hábil, então foi proposto que construíssemos uma versão paralelizada do algoritmo.

O restante deste relatório é organizado da seguinte forma. A Seção 2 descreve como foi feita a modelagem do problema e o armazenamento das conexões do grafo. A Seção 3 descreve as estruturas e como resolvido o problema do clique sequencialmente. A Seção 4 descreve como seria feito a paralelização do problema. A Seção 5 trata de detalhes específicos da implementação do trabalho: quais os arquivos utilizados; como é feita a compilação e execução; além de detalhar o formato dos arquivos de entrada e saída. A Seção 6 contém a análise de desempenho do problema para entradas triviais e não-triviais. A Seção 7 conclui o trabalho.

2. MODELAGEM

Modelamos o trabalho com grafos implementados utilizando matriz de adjacência, além disso, para trabalhar com o problema da clique, foi necessário implementar um tipo de conjunto, para trabalhar melhor com os nós do grafo.

```
typedef struct Matriz{
    int col, lin;
    int ** matriz;
} Matriz;

typedef struct grafo {
    Matriz matrizAdj;
    int N; // Quantidade de vertices no grafo
} Grafo;

typedef struct conjunto {
    int * elementos;
    int tam;
    int qtde;
} Conjunto;
```

Essas três estruturas foram necessárias para implementar todo o trabalho. Podemos dizer que a complexidade espacial da estrutura é de $O(n^2)$, sendo n o número de vértices que o grafo tem. A estrutura de conjunto não passa de $O(n)$, pois apenas armazena quais vértices estão contidos ou não nos conjuntos, não necessariamente precisando armazenar as arestas que conectam cada um.

3. SOLUÇÃO SEQUENCIAL PROPOSTA

Para nossa solução sequencial, utilizamos uma forma básica do **Bron-Kerbosch**, que se trata de um algoritmo recursivo que utiliza *backtracking* para encontrar todos os cliques máximos de um grafo. A base do algoritmo é a busca exaustiva, mas utiliza-se um conhecimento maior sobre a natureza do problema. Ele gera apenas cliques maximais, evitando assim que cada conjunto gerado tenha que ser comparada com os previamente testados. [Bron and Kerbosch]

A estratégia de backtracking constrói incrementalmente candidatos para solução do problema, e na medida que determina que não tem possibilidade do candidato fazer parte da solução válida ele o abandona. [Backtracking]

A estrutura básica do algoritmo utiliza três conjuntos: um para armazenar os vértices que já foram definidos como parte da clique; outro que armazena os vértices que são candidatos a fazer parte da clique; e o último armazena os vértices que já foram analisados e que não irão fazer parte da clique.

Após o algoritmo fazer os cálculos para cada vértice, no final é possível saber qual o clique maximal do grafo. Esse clique máximo representa o maior grupo de usuários da rede social que são amigos.

No pior caso, o algoritmo apresenta complexidade exponencial de $O(2^n)$, sendo n o número de vértices no grafo.

4. SOLUÇÃO PARALELIZADA PROPOSTA

Aqui será detalhado nossa proposta de solução paralelizada, apesar dela não ter sido implementada.

4.1. Algoritmos implementados

- *void intersecaoVizinhos(Conjunto * P, Grafo * G, int vertice)*

Descrição: Função que calcula a interseção entre o conjunto de vértices P para descobrir os vizinhos do vértice em questão.

Parâmetros: Estrutura de conjunto P, estrutura de grafo G e um inteiro representando o vértice.

Complexidade: $O(n)$, onde n é a quantidade de vértices do grafo G.

- *void BK(Conjunto * C, Conjunto * P, Conjunto * S, Grafo * G)*

Descrição: Função que faz o cálculo do Bron-Kerbosch.

Parâmetros: Três conjuntos C, P e S e um grafo G. Dos três conjuntos, um é utilizado para armazenar os vértices que já foram definidos como parte da clique; outro que armazena os vértices que são candidatos a fazer parte da clique; e o último armazena os vértices que já foram analisados e que não irão fazer parte da clique.

Complexidade: $O(2^n)$, onde n é o número de vértices no grafo.

4.2. COMPLEXIDADE GERAL

Cada complexidade separada dos algoritmos será $O(n)$, onde n é o número de páginas na memória, igual para cada um deles. Porém, a complexidade temporal do programa deve levar em consideração a quantidade instâncias, de páginas e de acessos.

Assim, nosso cálculo de complexidade poderia ser definido como *número de páginas * número de acessos * 3*. Nesses três parâmetros, o que tem um maior peso seria os acessos, pois ele que irá dominar todos os outros. Por isso, podemos dizer então, que a complexidade temporal final do programa será de $O(n)$, porém n aqui será a quantidade de acessos feitos.

5. IMPLEMENTAÇÃO

5.1. Código

5.1.1. Arquivos .c

- **tp4-seq.c** Arquivo principal do programa. Lê os arquivos de entrada, calcula o maior clique de cada instância e escreve resultado em um arquivo de saída.
- **grafos_matriz.c** Contém todas as funções de manipulação, leitura e escrita de grafos. Utiliza o TAD de matriz como implementação da matriz de adjacências.
- **clique.c** Contém as funções necessárias para o cálculo do clique.
- **arquivos.c** Um tipo abstrato de dados de manipulação de arquivos, contendo funções de abertura, leitura, escrita e fechamento.

5.1.2. Arquivos .h

- **grafos_matriz.h** Além de definir a estrutura de grafos, contém todas as funções de manipulação, leitura e escrita de grafos. Utiliza o TAD de matriz como implementação da matriz de adjacências.
- **clique.h** Contém as definições das funções necessárias para o cálculo do clique. Além disso, contém a definição da estrutura de conjuntos.
- **arquivos.h** Definição das funções utilizadas para ler, escrever e fechar corretamente um arquivo.

5.2. Compilação

O programa deve ser compilado através do compilador GCC através dos seguintes comandos

Para programação dinâmica:

```
gcc -Wall -Lsrc src/tp4-seq.c src/grafos\_matriz.c src/arquivos.c src/clique.c -o tp4-seq
```

Ou através do comando *make*.

5.3. Execução

A execução do programa tem como parâmetros:

- Um arquivo de entrada contendo várias instâncias a serem simuladas.
- Um arquivo de saída que irá receber a quantidade de nós do maior clique de cada instância.

O comando para a execução do programa é da forma:

```
./tp4-seq <arquivo_de_entrada> <arquivo_de_saida>
```

5.3.1. Formato da entrada

A primeira linha do arquivo de entrada contém o valor k de instâncias que o arquivo contém. As k instâncias são definidas da seguinte forma. A primeira linha contém um inteiro n indicando quantos nós (que representa o número de usuários da rede) o grafo possui. As n linhas seguintes contém a matriz de adjacência.

A seguir um arquivo de entrada de exemplo:

```
1
7
0 1 1 0 0 1 0
1 0 1 1 1 0 0
1 1 0 1 1 0 0
0 1 1 0 1 0 0
0 1 1 1 0 0 1
1 0 0 0 0 0 1
0 0 0 0 1 1 0
```

5.3.2. Formato da saída

O arquivo de saída consiste em k linhas, cada uma representando o resultado de uma instância. Cada linha contém um inteiro que representa o número de convidados para a festa, ou, em outras palavras, o maior clique encontrado no grafo.

4

6. AVALIAÇÃO EXPERIMENTAL

As análises feitas para para este trabalho foram todas detalhadas na especificação. O que foi pedido foi o seguinte:

- Calcular a localidade de referência temporal.
- Calcular a localidade de referência espacial.
- Gerar o histograma das distâncias de acessos.
- Gerar o histograma das distâncias de pilha.
- Gerar um gráfico "Tamanho da página"x "Bytes movimentados".
- Gerar um gráfico "Tamanho da memória"x "Falhas".

Para realizar essas análise foi disponibilizado um arquivo com uma configuração diferente no fórum da disciplina. Esse arquivo possuía apenas acessos à memória, sem informações de tamanho da memória ou da página. Foi criado dois scripts que auxiliaram na execução destas análises. O primeiro de localidade espacial e outro de localidade temporal.

Nas próximas seções iremos descrever a máquina utilizada para os testes e o resultado de cada item descrito acima.

6.1. Máquina utilizada

1

Segue especificação da máquina utilizada para os testes:

```
model name:      Intel(R) Core(TM) i3 CPU           M 330  @ 2.13GHz
cpu MHz:         933.000
cache size:      3072 KB
MemTotal:        3980124 kB
```

6.2. Localidade de referência temporal

A Tabela 1 mostra os resultados que obtivemos para o cálculo da localidade de referência temporal.

Instância	Temporal
1	19.67
2	14.10
3	21.08
4	10.67

Tabela 1. Localidade de referência temporal

6.3. Localidade de referência espacial

A Tabela 2 mostra os resultados que obtivemos para o cálculo da localidade de referência espacial. É gritante a diferença entre o resultado da instância 2 pras outras, mas isso ocorreu pois ela tem um padrão bem peculiar. A grande maioria dos acessos são sequenciais na memória (posição 8, depois 7, depois 6, depois 5), isso caracteriza uma boa localidade de referência espacial.

Instância	Espacial
1	16.69
2	2.55
3	10.46
4	19.71

Tabela 2. Localidade de referência temporal

6.4. Histogramas de Distância de Acessos

Na Figura 1 agrupamos os quatro histogramas que foram pedidos. Logo abaixo é possível ver o resultado para cada uma das instâncias.

Na instância 1 e 4 ocorre uma distribuição maior dos acessos. Na primeira os ápices se encontram nas menores distâncias, já na outra o ápice ocorre na distância 32. As distancias de acesso da instância 1 se encontram muito melhor distribuídas que a instância 4.

Já as instâncias 2 e 3 possuem uma similaridade: quase a totalidade dos acessos se encontram à distância 1. Isso indica que a localidade espacial deles é melhor que a das outras instâncias, por não terem que caminhar tanto na memória para acessar o próximo valor.

6.5. Histogramas de Distância de Pilha

Na Figura 2 foi agrupado os histogramas referentes ao calculo de distância de pilha de cada instância.

Nessa situação a instância que apresentou o melhor resultado foi a instância 2, pois a distância na pilha que mais correu foi o valor 3. Todas as outras distâncias ficaram bem distribuídas no histograma.

A segunda melhor instância nesse quesito seria a instância 4, pois muitos acessos se concentram em distâncias menores que 5. O resto dos acessos também se distribuíram bem pelo resto das distâncias.

As piores instâncias nessa análise foram as instâncias 1 e 3. A instância 1 possui muitos acessos de distâncias variadas, desde as menores até as maiores. Já a instância 3, apesar de possuir o ápice de acessos com distância de pilha igual a 1, possui dezenas de acessos com distâncias maiores que 40, o que prejudica muito seu resultado final.

6.6. Tamanho da página x Bytes Movimentados

Na Figura 3 temos a relação de "Tamanho da página x Bytes Movimentados". Com essas informações podemos analisar o *trade-off* entre o tamanho da página e o volume transferido entre disco e memória. Como foi dito na especificação do trabalho "Páginas maiores

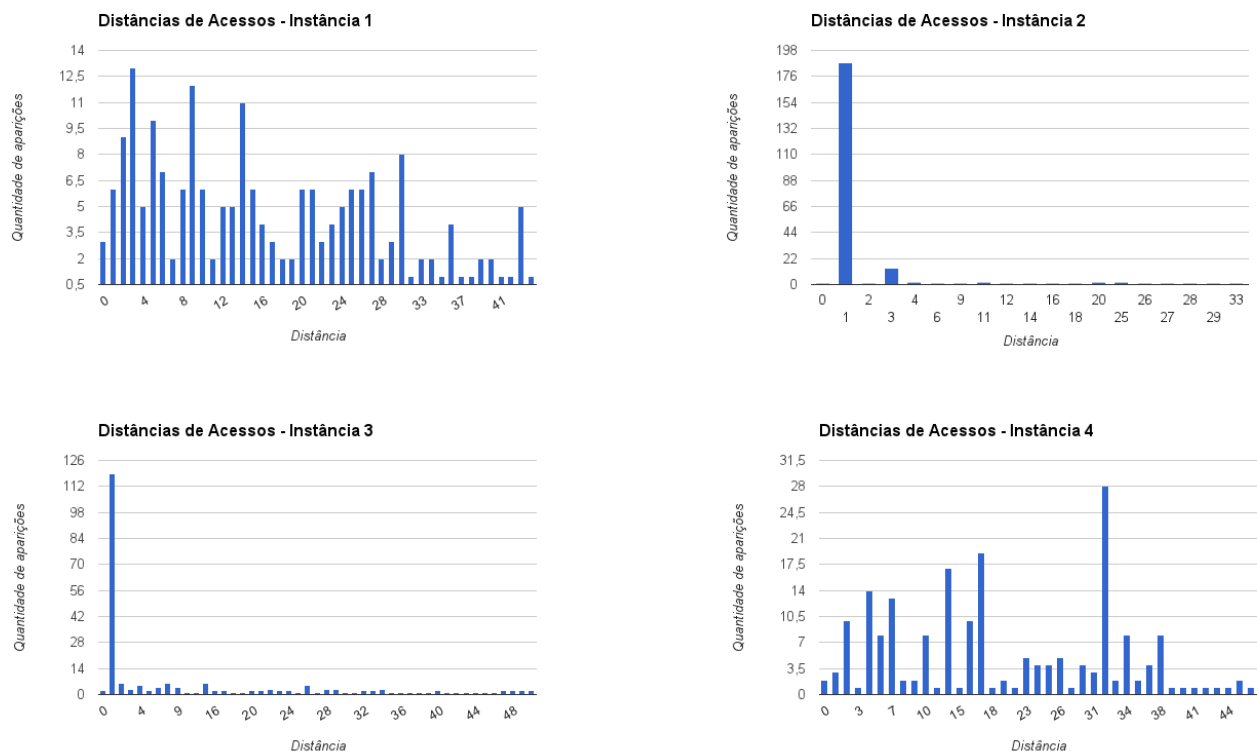


Figura 1. Distâncias de Acessos de todas as instâncias

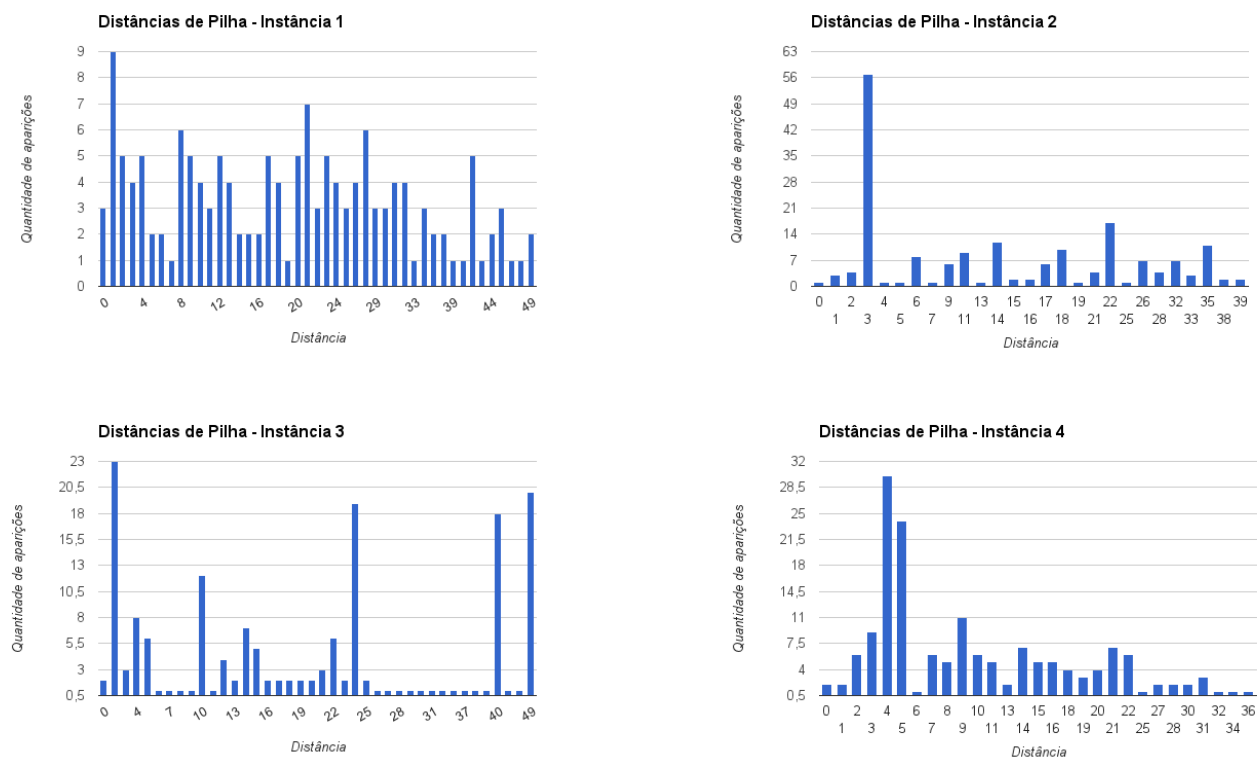


Figura 2. Distâncias de Pilha de todas as instâncias

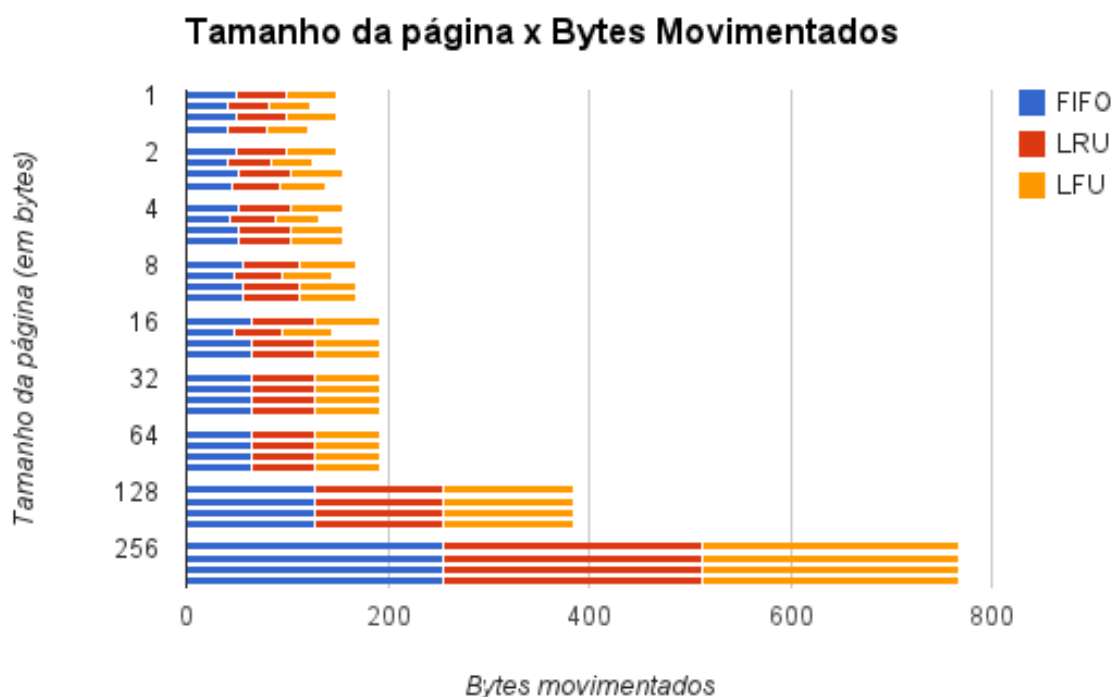


Figura 3. Tamanho da página x Bytes Movimentados (Memória fixa em 256 bytes)

irão certamente diminuir o número de falhas, porém vão causar o carregamento desnecessário de muitos dados”, podemos verificar isso ocorrendo de fato em nossos testes.

A política adotada foi a seguinte: foi fixado o tamanho da memória em 256 bytes, e então variamos o tamanho de cada página (sempre na potência de 2), de 1 até 256. Para páginas pequenas (1, 2, 4, 8 bytes), podemos ver que poucas falhas ocorreram, então consequentemente, poucos dados foram movimentados. Já no caso de páginas maiores (128, 256 bytes), a quantidade de falhas reduziu drasticamente, porém a quantidade de bytes movimentados também cresceu vertiginosamente.

A performance de cada uma das políticas de reposição nessa situação foi igual para todas. Todas obtiveram o mesmo número de falhas e, consequentemente, a mesma quantidade de bytes movimentados.

O tamanho da página, para memória fixa em 256 bytes, só começou a ser significativo quando chegou em 128 bytes, metade do tamanho da memória. Até então, a quantidade de bytes movimentados tinha crescido, porém de forma lenta. Então podemos dizer que o tamanho ideal da página pode ser definido como no máximo $1/4$ do tamanho da memória física disponível.

6.7. Tamanho da memória x Page Faults

Na Figura 4 temos a relação de “Tamanho da memória x Page Faults”. O objetivo desta análise seria estimar o tamanho ideal de memória física disponível. Obviamente, quanto mais memória menos falhas ocorrerão, porém este teste deixou visível a partir de qual

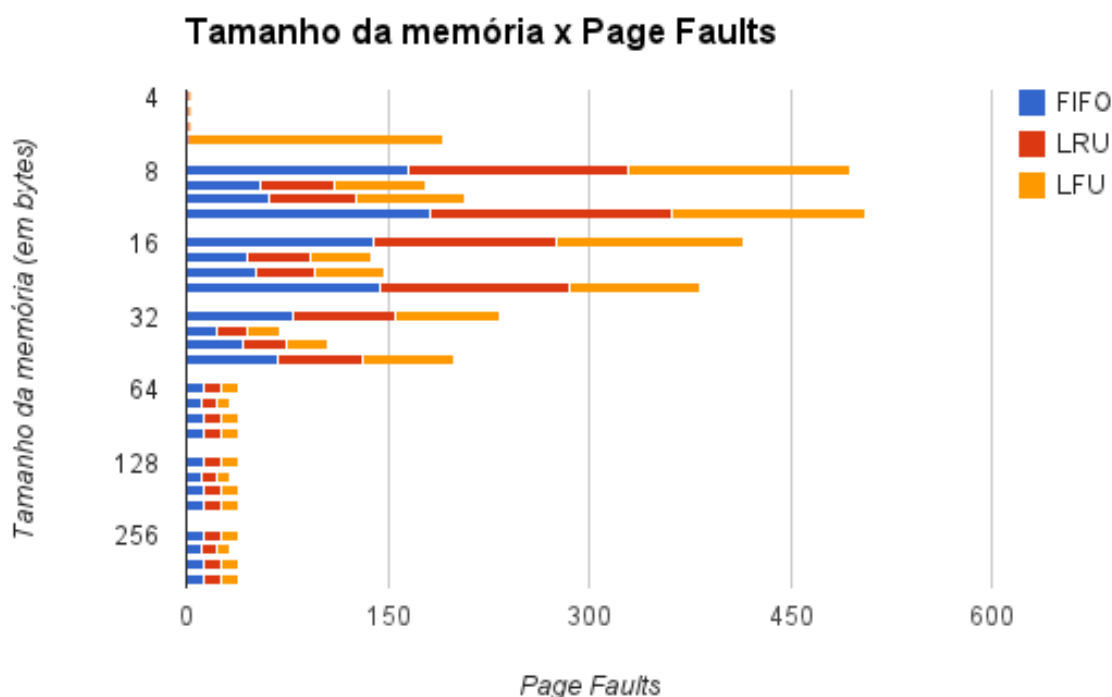


Figura 4. Tamanho da memória x Page Faults (Página fixa em 4 bytes)

momento já é interessante determinar qual a melhor quantidade de memória necessária para obter menos falhas.

A política adotada neste teste foi o seguinte: foi fixado o tamanho da página em 4 bytes, após isso, variamos o tamanho da memória (sempre na potência de 2), de 4 até 256. E podemos perceber, como já era esperado, que quanto maior a memória disponível, menor a quantidade de *page faults*.

Nessa situação, com páginas fixas em 4 bytes, um tamanho de memória ideal poderia ser 32 bytes, porém para sequências de acesso com uma localidade de referência espacial boa. Se não for possível saber o comportamento da sequência de acessos, o tamanho ideal da memória seria 64 bytes, uma vez que todas as políticas e todas as instâncias obtiveram um número pequeno de *page faults*.

6.8. Resultado

De acordo com as nossas análises feitas, podemos estimar os valores ideais para o tamanho da memória física e o tamanho das páginas, para obter o menor número de falhas possível para cada uma das quatro sequências de acessos disponibilizadas para os testes.

Na Figura 3 o desempenho de cada uma das políticas é igual para as quatro instâncias, obtendo a mesma quantidade de *page faults* e, consequentemente, de bytes movimentados. Isso ocorreu pois com a quantidade de memória estabelecida, independente da política escolhida, a mesma quantidade de *page faults* ocorreria.

Na Figura 4 é possível ver que a Instância 1 e 4 obtiveram os piores resultados.

Isso ocorreu justamente pela característica de cada uma. A sequência de acessos nessas instâncias ocorreram de uma forma mais caótica, quando posições muito distantes umas das outras foram acessadas em sequência. Já no caso das Instâncias 2 e 3 ocorre o contrário, as sequências de acessos ocorrem de uma maneira mais sequencial, com uma localidade de referência espacial melhor, pois dados próximos são acessados em um curto espaço de tempo.

7. CONCLUSÃO

No geral, a relação entre cada uma das políticas se manteve praticamente constante, nos gráficos pode não estar totalmente perceptível, mas é possível perceber que a política de reposição de remover as páginas com a menor quantidade de acessos (LFU) obteve uma quantidade de page faults melhor do que as outras políticas.

Foi possível perceber também que sequências de acessos com perfis mais sequenciais, ou seja, com localidade de referência espacial melhores, possuem um resultado melhor, independente do tamanho da memória disponibilizado.

Para perfis mais desordenados, com localidade de referência espacial pior, a quantidade de *page faults* para qualquer uma das políticas aumenta consideravelmente.

Acredito que os objetivos do trabalho foram concluídos com sucesso, uma vez que o sistema de memória virtual (SMV) foi implementado com sucesso e foi possível exercitar, mais uma vez, os conceitos de gerenciamento de memória. Além disso, na parte experimental, foi possível ver, na prática, a localidade de referência de diferentes tipos de sequências de acessos e como cada política trabalha com cada um desses perfis.

Referências

Backtracking, W. Backtracking. <http://en.wikipedia.org/wiki/Backtracking>.

Bron, C. and Kerbosch, J. *Algorithm 457: finding all cliques of an undirected graph*.

Wikipedia. Problema do clique. http://pt.wikipedia.org/wiki/Problema_do_clique.