

Introdução à Programação Paralela
2017/02

Trabalho Prático Parte 01:
Melhora de Performance em uma Multiplicação de
Matrizes simples

Gustavo Borba
Belo Horizonte, September 3, 2017

1 Introdução

O presente trabalho tem o objetivo de, em sua primeira etapa, experimentar otimizações de código em um algoritmo de multiplicação de matrizes, de forma a melhorar a performance do mesmo em observância à plataforma que será executado. Ao final desta etapa, espera-se que o código tenha um desempenho melhor, comparado ao seu original, e que esteja mais fácil de implementar sua paralelização.

2 Metodologia

De acordo com orientações descritas por BILMES[1], o código de duas das três implementações disponibilizadas foi alterado para explicitar características que podem ser utilizadas pelo otimizador do compilador.

Os arquivos com os códigos originais foram mantidos, para fins de comparação com os algoritmos alterados. Para estes últimos, uma nova linha no makefile foi criada para cada um, e na pasta do projeto pode-se verificar a presença do sufixo *borba-*, indicando uma versão alterada do arquivo.

2.1 Remover explicitamente falsas dependências

Para garantir a corretude das operações, os compiladores não assumem que, ao acessar um ponteiro que já fora acessado anteriormente, os mesmos dados estarão lá (afinal de contas, outra função pode ter acesso a esse ponteiro e alterar seu valor), de modo que as operações sobre ponteiros exigem, em sua maioria, acessos constantes à memória.

Como o programa em questão trata de três matrizes distintas, nenhuma das matrizes será alterada por uma função externa durante a execução das multiplicações. Dessa forma, é possível adicionar a cláusula **restrict** aos ponteiros recebidos pela função, que indica que apenas aqueles ponteiros ou valores diretamente derivados dele serão utilizados para acessar aquela área de memória.

```
void square_dgemm (int n,  
                  double* restrict A, double* restrict B, double* restrict C)
```

O compilador saberá, dessa forma, que não é necessário realizar um fetch dos dados, à menos que haja um cache miss, pois apenas aquele ponteiro referencia a matriz em questão.

2.2 Explorando melhor os Registradores

Realizando o profiling do código original com a ferramenta **perf**, verifica-se que o tempo de acesso à memória é o principal gargalo do programa em questão. Ao verificar os *misses*, obtemos um resultado como o apresentado na Figura 1, ficam evidentes dois eventos onde se têm um custo de performance elevado

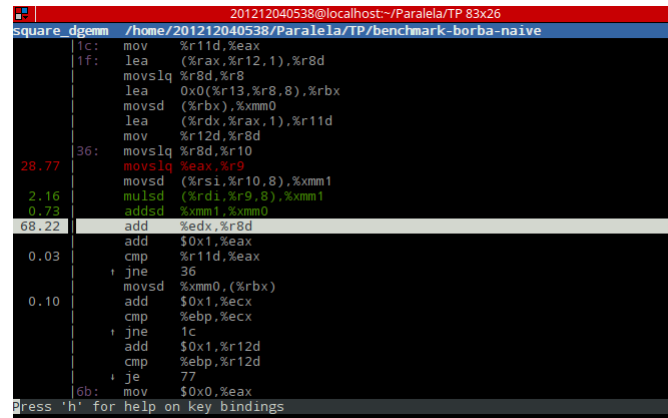


Figure 1: Profiling realizado pelo programa `perf` no algoritmo de multiplicação de matrizes ingênuo, em termos de *misses* totais. Percebe-se duas instruções com maior custo em performance

Analisando o código exibido, percebe-se que partes mais onerosas em termos de memória são um `load` e uma soma, cujo armazenamento se dá em endereço de memória.

```
void square_dgemm (int n, double* A, double* B, double* C){
    for (int i = 0; i < n; ++i)
        for (int j = 0; j < n; ++j){
            C[i+j*n] = C[i+j*n]; // aqui
            for( int k = 0; k < n; k++)
                C[i+j*n] += A[i+k*n] * B[k+j*n]; // aqui
        }
}
```

Para reduzir essa demanda por acessos à memória, transforma-se o ponteiro em uma variável local. Como trata-se de um valor muito utilizado, inclusive, a utilização da cláusula **register** é ainda mais eficiente. Assim, por completo, o código ingênuo se torna o seguinte:

```
void square_dgemm (int n, double* restrict A, double* restrict B, double* restrict
    C) {
    for (int i = 0; i < n; ++i)
        for (int j = 0; j < n; ++j) {
            int jn = j*n;
            register double cij = C[i+jn];
            for( int k = 0; k < n; k++) {
                double Aikn=A[i+k*n];
                double Bkjn=B[k+jn];
                cij += Aikn * Bkjn;
            }
            C[i+jn] = cij;
        }
}
```

repetindo o processo com a ferramenta **perf**, mas agora com tempo de CPU, verifica-se na Figura 2 que o gargalo de processamento ainda ocorre nestas duas operações, sendo dispensável a criação de mais variáveis temporárias para a otimização de índices, por exemplo, pois estes são pouco utilizados com relação aos anteriormente citados.

```

201212040538@localhost:~/Paralela/TP 168x26
square_dgemm /home/201212040538/Paralela/TP/benchmark-borba-naive
    jle 77
    mov %rdx,%rdi
    mov %rcx,%r13
    mov $0x0,%r12d
    mov %ebp,%edx
    jmp 6b
1c: mov %r11d,%eax
1f: lea (%rax,%r12,1),%r8d
0.12 movslq %r8d,%r8
0.00 lea 0x0(%r13,%r8,8),%rbx
0.00 movsd (%rbx),%xmm0
0.77 lea (%rdx,%rax,1),%r11d
0.09 mov %r12d,%r8d
36: movslq %r8d,%r10
25.52 movslq %eax,%r9
0.05 movsd (%rsi,%r10,8),%xmm1
13.86 mulsd (%rdi,%r9,8),%xmm1
6.50 addsd %xmm1,%xmm0
52.43 add %edx,%r8d
    add $0x1,%eax
0.51 cmp %r11d,%eax
    jne 36
    movsd %xmm0,(%rbx)
    add $0x1,%ecx
    Press 'h' for help on key bindings

```

Figure 2: Profiling realizado pelo programa perf no algoritmo modificado, em termos de processamento.

3 Resultados Computacionais

Todos os testes aqui descritos foram executados na máquina Katrina, que tem como processador um Intel Haswell-EP de 6 cores (12 threads) executando a 2.4 GHz e Cache de 15 M, resultando em:

$$8 \text{ double-precision (64-bit) flops por pipeline} * 6 \text{ pipelines} * 2.4 \text{ GHz} = 115.2 \text{ Gflops}$$

Para a execução dos testes computacionais foi utilizada a função de benchmark disponibilizada para este mesmo fim. Essa função utiliza um conjunto representativo de tamanhos de matrizes (múltiplos de 32 ± 1) para executar a função de multiplicação de matrizes capturando o tempo de execução. A partir das características da máquina, a função calcula o desempenho do programa em Mflops e em porcentagem da capacidade total de processamento da máquina.

201212040538@localhost:~/Paralela/TP 83x37				201212040538@localhost:~/Paralela/TP 83x37			
Description: Naive, three-loop dgemv.				Description: Naive, three-loop dgemv.			
Size: 31	Mflop/s: 1530.23	Percentage: 8.07		Size: 31	Mflop/s: 1530.27	Percentage: 8.07	
Size: 32	Mflop/s: 1538.2	Percentage: 8.12		Size: 32	Mflop/s: 1552.67	Percentage: 8.09	
Size: 96	Mflop/s: 2217.75	Percentage: 11.55		Size: 96	Mflop/s: 2218.34	Percentage: 11.55	
Size: 97	Mflop/s: 2220.25	Percentage: 11.56		Size: 97	Mflop/s: 2220.16	Percentage: 11.56	
Size: 127	Mflop/s: 2195.01	Percentage: 11.43		Size: 127	Mflop/s: 2194.94	Percentage: 11.43	
Size: 128	Mflop/s: 1974.6	Percentage: 10.28		Size: 128	Mflop/s: 2056.18	Percentage: 10.71	
Size: 129	Mflop/s: 2194.08	Percentage: 11.43		Size: 129	Mflop/s: 2194.11	Percentage: 11.43	
Size: 191	Mflop/s: 2108.39	Percentage: 10.98		Size: 191	Mflop/s: 2110.38	Percentage: 10.99	
Size: 192	Mflop/s: 1983.99	Percentage: 10.23		Size: 192	Mflop/s: 1997.06	Percentage: 10.45	
Size: 229	Mflop/s: 2089.14	Percentage: 10.93		Size: 229	Mflop/s: 2097.89	Percentage: 10.93	
Size: 255	Mflop/s: 2079.26	Percentage: 10.83		Size: 255	Mflop/s: 2079.47	Percentage: 10.83	
Size: 256	Mflop/s: 956.254	Percentage: 4.98		Size: 256	Mflop/s: 957.788	Percentage: 4.99	
Size: 257	Mflop/s: 2073.71	Percentage: 10.80		Size: 257	Mflop/s: 2073.83	Percentage: 10.80	
Size: 319	Mflop/s: 2082.73	Percentage: 10.85		Size: 319	Mflop/s: 2078.65	Percentage: 10.83	
Size: 320	Mflop/s: 1621.97	Percentage: 8.45		Size: 320	Mflop/s: 1980.53	Percentage: 10.32	
Size: 321	Mflop/s: 2081.65	Percentage: 10.84		Size: 321	Mflop/s: 2079.71	Percentage: 10.83	
Size: 417	Mflop/s: 2078.63	Percentage: 10.81		Size: 417	Mflop/s: 2065.21	Percentage: 10.76	
Size: 479	Mflop/s: 2078.86	Percentage: 10.83		Size: 479	Mflop/s: 2083.45	Percentage: 10.85	
Size: 480	Mflop/s: 1988.93	Percentage: 10.36		Size: 480	Mflop/s: 1818.23	Percentage: 9.47	
Size: 511	Mflop/s: 2083.88	Percentage: 10.91		Size: 511	Mflop/s: 2093.1	Percentage: 10.90	
Size: 512	Mflop/s: 948.632	Percentage: 4.94		Size: 512	Mflop/s: 950.895	Percentage: 4.95	
Size: 639	Mflop/s: 2065.74	Percentage: 10.76		Size: 639	Mflop/s: 2051.41	Percentage: 10.68	
Size: 640	Mflop/s: 954.185	Percentage: 4.97		Size: 640	Mflop/s: 958.274	Percentage: 4.99	
Size: 767	Mflop/s: 2069.25	Percentage: 10.78		Size: 767	Mflop/s: 2047.37	Percentage: 10.66	
Size: 768	Mflop/s: 962.121	Percentage: 5.01		Size: 768	Mflop/s: 961.029	Percentage: 5.01	
Size: 769	Mflop/s: 2064.16	Percentage: 10.75		Size: 769	Mflop/s: 2050.06	Percentage: 10.68	
Average percentage of Peak = 9.47504				Average percentage of Peak = 9.7198			
[201212040538@localhost:TP] █				[201212040538@localhost:TP] █			

Figure 3: Execução das multiplicações ingênuas de matrizes, original (à esquerda) e modificado (à direita)

Executando os processos em sequência (para evitar a concorrência), tal como exibido na figura 5 e comparando unicamente os valores obtidos em MFlops/s foi possível montar dois gráficos, que podem ser vistos nas duas figuras abaixo. O primeiro gráfico compara a porcentagem de uso da capacidade total da CPU do servidor Katrina, enquanto o segundo gráfico explora as diferenças em percentual dos dois algoritmos.

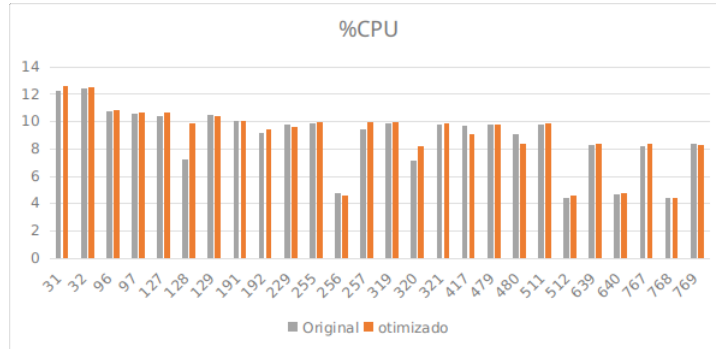


Figure 4: Rendimento em porcentagem de Capacidade total da CPU para os algoritmos

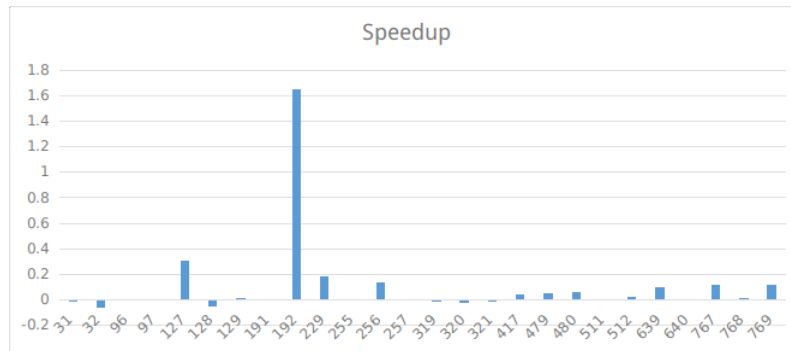


Figure 5: Diferença percentual entre o algoritmo utilizado e o algoritmo

Percebe-se uma grande otimização para o tamanho de matrizes de 192x192. Esse valor é um divisor do tamanho da cache, o que indica que armazenar o valor mais utilizado em um registrador pode realmente valer a pena:

$$15M = 1024 \times 15 = 15360.$$

$$\frac{15360}{192} = 80 \text{ páginas}$$

Utilizando a mesma estrutura de testes, agora compara-se o código escrito com otimizações e seu código original compilado com as diretivas *-O3 fstrict-aliasing -std=gnu99*, obtivemos o gráfico da figura 6, que representa o ganho (em porcentagem) da otimização realizada manualmente ante à otimização do compilador.

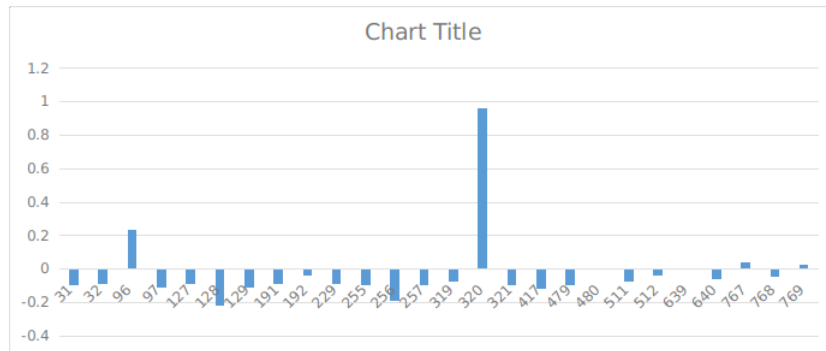


Figure 6: Resultados do ganho da otimização escrita, em comparação com a realizada pelo compilador

Nota-se que os resultados do otimizador do compilador são ligeiramente melhores (valores negativos) enquanto as otimizações manuais se mostraram mais efetivas nos tamanhos de matrizes 96 e 320 (chegando a 1% a mais de utilização da capacidade total da CPU). Tal resultado pode ser explicado pelas características do processador (que tem um processador de 3.20 GHz) em conjunto com a otimização da gerência de memória, utilizando um registrador para armazenar o valor mais comumente utilizado, evitando cache misses e aproveitando melhor o tempo para realizar operações de ponto flutuante.

4 Conclusão

As otimizações realizadas foram apenas ligeiramente eficientes, chegando próximas ao otimizador do compilador, mas perdendo um pouco para ele. Estima-se, no entanto, que estas otimizações ajudem a melhorar significativamente o paralelismo que será futuramente feito.

Os próximos passos envolvem Loop Unrolling, e a pesquisa e implementação de um método mais eficiente de multiplicação de matrizes.

References

- [1] BILMES, J. Optimizing Matrix Multiply using PHiPAC: a Portable, High-Performance, ANSI C Coding Methodology
Disponível em https://people.eecs.berkeley.edu/~krstepapers/shipac_ics97.pdf
Acesso em 2 de Setembro de 2017