

UNIVERSIDADE FEDERAL DE VIÇOSA  
INF310 – PROGRAMAÇÃO CONCORRENTE E DISTRIBUÍDA  
Lista de Exercícios 3

1. Um sistema linear pode ser resolvido através de operações executadas entre as linhas com o objetivo de organizá-lo como um sistema linear triangular, permitindo sua resolução de forma fácil. Por exemplo, no sistema linear:

$$\begin{array}{rcl} 2x_0 - 3x_1 & & = 3 \\ 4x_0 - 5x_1 + x_2 & = & 7 \\ 2x_0 - x_1 - 3x_2 & = & 5 \end{array}$$

ao se multiplicar a primeira linha por -2 e somar com a segunda, multiplicar a primeira por -1 e somar com a terceira e, finalmente, subtrair a nova linha 3 do dobro da nova linha 2, obtém-se:

$$\begin{array}{rcl} 2x_0 - 3x_1 & = & 3 \\ x_1 + x_2 & = & 1 \\ -5x_2 & = & 0 \end{array}$$

A resolução de um sistema linear triangular superior formado por n equações pode ser realizada através do trecho de código a seguir, onde A armazena os coeficientes do sistema, b armazena os termos independentes, e x armazena a solução do sistema:

```
for (row = n-1; row >= 0; row--) {
    x[row] = b[row];
    for (col = row+1; col < n; col++)
        x[row] -= A[row][col] * x[col];
    x[row] /= A[row][row];
}
```

- a) O loop externo pode ser paralelizado? Como ou por que não?  
b) O loop interno pode ser paralelizado? Como ou por que não?
2. Como o uso do OpenMP aumenta a portabilidade do código em relação ao uso de Pthreads?
3. Considerando que o escalonamento static do OpenMP apresenta um overhead menor que os demais escalonamentos, em que situação é melhor abrir mão deste modo de escalonamento?
4. Qual é o resultado do trecho de código mostrado a seguir?
- ```
...
int x=10;
#pragma omp parallel reduction(*:x) num_threads(3)
x+=2;
printf("x = %d\n",x);
...
```
5. Como funciona uma arquitetura SIMD? Pode-se dizer que CUDA segue essa arquitetura? Justifique.
6. Explique por que é possível obter maior eficiência com o uso da memória compartilhada no CUDA. Explique também em que situações o uso dessa memória oferece maior ganho de eficiência e quais são as particularidades em seu uso (considerando um volume grande de dados).

7. Quais são as principais diferenças entre tecnologias como OpenMP e o MPI? Descreva como a arquitetura do sistema influencia no funcionamento das mesmas.
8. Considerando as duas versões de um kernel CUDA apresentadas a seguir responda porque o acesso aos dados é mais eficiente no kernel da direita. Dica: Considere o conceito de coalescing.

```

struct myStruct {
    int x;
    int y;
};

__global__
void kernel(myStruct *a, int *r) {
    int tid=blockIdx.x*blockDim.x +
        threadIdx.x;

    int valor_x = a[tid].x;
    int valor_y = a[tid].y;

    r[tid]=valor_x - valor_y;
}

```

```

__global__
void kernel(int *x, int *y, int *r) {
    int tid=blockIdx.x*blockDim.x +
        threadIdx.x;

    int valor_x = x[tid];
    int valor_y = y[tid];

    r[tid]=valor_x - valor_y;
}

```

9. Uma imagem preto e branco pode ser representada por uma matriz onde cada pixel da imagem é dado por um número inteiro no intervalo [0,255] armazenado na posição respectiva da matriz. Imagine que desejamos estender uma imagem achatada para que a mesma passe a ter o dobro de sua altura. Para isso, podemos processar a imagem de modo a inserir uma nova linha entre cada par de linhas da imagem original, e preencher cada pixel dessa nova linha com a média dos valores encontrados nos pixels imediatamente acima e imediatamente abaixo. Assim, se a imagem original tem dimensões  $M \times N$ , a imagem processada terá tamanho  $(2*M-1) \times N$ , como ilustrado a seguir.

|     |     |     |     |
|-----|-----|-----|-----|
| 0   | 1   | 1   | 2   |
| 3   | 3   | 4   | 5   |
| 200 | 250 | 230 | 230 |

Representação de uma parte da matriz original

|     |     |     |     |
|-----|-----|-----|-----|
| 0   | 1   | 1   | 2   |
| 1   | 2   | 2   | 3   |
| 3   | 3   | 4   | 5   |
| 101 | 126 | 117 | 117 |
| 200 | 250 | 230 | 230 |

Matriz estendida destacando as duas linhas inseridas contendo a média (truncada) dos valores imediatamente acima e abaixo

Escreva um kernel CUDA que faz esse processamento de forma paralela. Os parâmetros do kernel devem ser:

- uma matriz de entrada representando a imagem original;
- uma matriz de saída que deverá armazenar o resultado do processamento;
- o número de linhas da matriz original;
- o número de colunas da matriz original.