

Jacobi iterative equation solver

$$A x = b$$

Each x_i can be calculated in parallel

$$a_{i0}x_0 + a_{i1}x_1 + \dots + a_{ii}x_i + \dots + a_{in-1}x_{n-1} = b_i$$

$$a_{ii}x_i = b_i - \sum_{j \neq i} a_{ij}x_j$$

update rule

$$x_i = \frac{1}{a_{ii}} \left(b_i - \sum_{j \neq i} a_{ij}x_j \right)$$

Version 1

$$A \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_{n-1} \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_{n-1} \end{bmatrix}$$

thread block: $(1, \text{SIZE})$
 grid: $(1, \frac{\# \text{ rows}}{\text{SIZE}})$ ← covers all output elements of x

Host code:

1. Allocate and copy A, x, b , to device
2. Allocate for x_{new} and ssd
3. while (!done) {
 - $\text{cudaMemset}(\text{ssd}, 0)$ ← src dest
 - launch kernel ($A, b, x, x_{\text{new}}, \text{ssd}$)
 - $\text{cudaMemcpy}(\text{ssd})$
 - if ($\sqrt{\text{ssd}} < \text{eps}$) done = 1
 - flip pointers to x, x_{new}
 - $\text{cudaMemcpy}(x)$

Use double precision for ssd values

Effect of thread block size on performance.

- maximize utilization on GPU

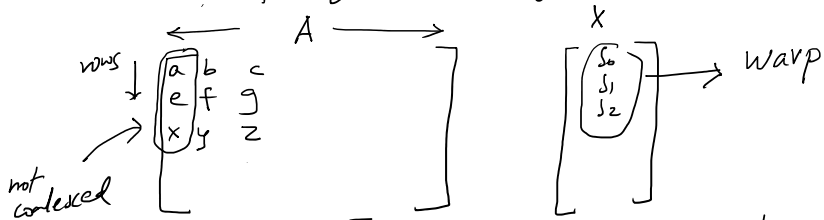
kernel code

1. Allocate shared memory for local_ssd values
2. old_value ← $x[i]$
3. new_value ← update
4. $x_{\text{new}}[i] \leftarrow \text{new_value}$
5. $\text{local_ssd} = (\text{old_value} - \text{new_value})^2$
6. Store local_ssd to shared memory
7. Reduce local_ssd to single value at thread block level
8. if (threadIdx.x == 0)
 - accumulate reduced local ssd value into shared ssd value in GPU global mem

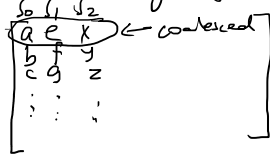
use atomicAdd()

Optimized version:

- Make accesses from global mem of matrix A to be coalesced *



Convert A to $B = A^T$ on CPU s.t. elements of B are laid out in column major form



Rewrite update logic in kernel to perform update assuming column-major layout.

Look at the example of matrix-vector multiplication available on BBLearn