

## Problem Scenario

The task at hand is the development of a sparse matrix-vector product kernel. You are required to implement a kernel capable of computing

$$y \leftarrow Ax$$

where  $A$  is a sparse matrix stored in

1. CSR
2. ELLPACK

storage formats (see below). The kernel shall be parallelized to exploit available computing capabilities. The code shall be developed in both OpenMP and CUDA versions, and shall be tested for correctness against a serial implementation. You shall need to develop auxiliary functions to preprocess the matrix data and represent it in the desired format. The code shall be tested for performance on a specified set of matrices; performance tests shall be carried out on the GRID computing facility.

## Test Matrices

The test matrices shall be obtained from the University of Florida Sparse Matrix Collection at the website <https://www.cise.ufl.edu/research/sparse/matrices/>. It is recommended to download the matrices in the MatrixMarket format; software for reading matrices from file into memory is available at <http://math.nist.gov/MatrixMarket/>. The matrices are characterized by a number of rows  $M$ , number of columns  $N$  and number of nonzero entries  $NZ$ . Typically, for symmetric matrices only the upper (or lower) triangle is stored on file; these should be appropriately expanded to reconstruct the complete matrix.

The test set shall include at least the following matrices:

cage4	Cube_Coup_dt0	FEM_3D_thermal1
mhda416	ML_Laplace	thermal1
mcfe	bcsstk17	thermal2
olm1000	mac_econ_fwd500	thermomech-TK
adder_dcop_32	mhd4800a	nlpkkt80
west2021	cop20k_A	webbase-1M
cavity10	raefsky2	dc1
rdist2	af23560	amazon0302
cant	lung2	af_1_k101
olafu	PR02R	roadNet-PA

although testing with other matrices is also encouraged.

## Performance measurements

All performance data shall be obtained by repeated invocation of the kernel for a certain number of times, resulting in a measurement of the average time per kernel invocation.

For all cases, the measure of performance in MFLOPS or GFLOPS will be obtained as

$$FLOPS = \frac{2 \cdot NZ}{T}$$

where  $NZ$  is the number of nonzero entries in the matrix, adjusted as necessary in the case of symmetric storage on file, and  $T$  is the (average) time per kernel invocation.

The measurements shall only include the time needed to perform the matrix-vector product operation; input/output from file and data structure preprocessing shall *not* be included in the timings.

For the OpenMP version the code shall be tested with varying number of threads, from 1 up to the maximum number of available cores.

## Storage formats

The storage formats will be described in Matlab.

When translating the codes into C/CUDA, you should pay attention to the actual data access pattern from memory and the resulting data structure adjustments that may be necessary, keeping in mind that Matlab employs storage by columns, whereas C and C++ employ storage by rows.

### CSR

Compressed Storage by Rows format. An  $M \times N$  matrix with  $NZ$  non-zero entries is described by the following data:

**M** Rows;

**N** Columns;

**IRP(1:M+1)** Array of pointers to row start;

**JA(1:NZ)** Array of column indices;

**AS(1:NZ)** Array of coefficients;

so that the matrix-vector product in Matlab would be implemented by the following code:

```
for i=1:m
    t=0;
    for j=irp(i):irp(i+1)-1
        t = t + as(j)*x(ja(j));
    end
```

```

    y(i) = t;
end

```

As an example, the matrix

$$\begin{pmatrix} 11 & 12 & 0 & 0 \\ 0 & 22 & 23 & 0 \\ 0 & 0 & 33 & 0 \\ 0 & 0 & 43 & 44 \end{pmatrix}$$

would be stored in CSR (assuming 1-base indexing as in Matlab) as:

$$\begin{aligned} M &= 4 \\ N &= 4 \\ IRP &= [1, 3, 5, 6, 8] \\ JA &= [1, 2, 2, 3, 3, 3, 4] \\ AS &= [11, 12, 22, 23, 33, 43, 44] \end{aligned}$$

## ELLPACK

ELLPACK storage format. An  $M \times N$  matrix with  $NZ$  non-zero entries and such that  $MAXNZ$  is the maximum number of nonzero entries per row across all rows is described by the following data:

**M** Rows;

**N** Columns;

**MAXNZ** Max nonzeros per row

**JA(1:M,1:MAXNZ)** 2D array of column indices

**AS(1:M,1:MAXNZ)** 2D array of coefficients;

so that the matrix-vector product in Matlab would be implemented by the following code:

```

for i=1:m
    t=0;
    for j=1:maxnzs
        t = t + as(i,j)*x(ja(i,j));
    end
    y(i) = t;
end

```

Rows with a number of nonzeros less than  $MAXNZ$  shall be padded with appropriate values, e.g. zeros in  $AS$ . As an example, the matrix

$$\begin{pmatrix} 11 & 12 & 0 & 0 \\ 0 & 22 & 23 & 0 \\ 0 & 0 & 33 & 0 \\ 0 & 0 & 43 & 44 \end{pmatrix}$$

would be stored in ELLPACK (assuming 1-base indexing as in Matlab) as:

$$\begin{aligned} M &= 4 \\ N &= 4 \\ MAXNZ &= 2 \\ JA &= \begin{pmatrix} 1 & 2 \\ 2 & 3 \\ 3 & 3 \\ 3 & 4 \end{pmatrix} \\ AS &= \begin{pmatrix} 11 & 12 \\ 22 & 23 \\ 33 & 0 \\ 43 & 44 \end{pmatrix} \end{aligned}$$

Note that describing the ELLPACK format in Matlab, which has storage by columns, implies that coefficients in the same column are adjacent in memory; reproducing the same data layout in a C-like language means that the matrix would appear to be transposed. This has an effect in the CUDA version: most matrices have a relatively small number of nonzeros per row, so it makes sense to assign only one thread per row, but this means that if we store the matrix by rows and then access it by columns, accesses will be uncoalesced; so, the layout in memory should be adjusted to bring consecutive threads to access consecutive entries, i.e. the rows in the device memory should correspond to the columns of the matrix as described in the Matlab code above, and vice-versa.

Note that for some matrices the amount of padding needed may be so large as to make it impossible to use the ELLPACK format on some computing devices.

## Deliverables

The students shall have to turn in:

- A report detailing the approach they have taken,
- The code implemented,
- A report with analysis of the performance data;
- A short summary design specification and test plan (including performance testing).

## **Marking Scheme**

- 25** Software implementation & performance obtained;
- 40** Project Report;
- 25** Design and test documents;
- 10** Performance analysis;

**Submission deadline: 2:30 pm, April 3rd, 2018**