

ECE408 Spring 2020

Applied Parallel Programming

Lecture 19: Parallel Sparse Methods

©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

1

1

Objective

- To learn to regularize irregular data with
 - Limiting variations with clamping
 - Sorting
 - Transposition
- To learn to write a high-performance SpMV kernel based on JDS transposed format

©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

2

2

Coordinate (COO) format

- Explicitly list the column and row indices for every non-zero element

		Row 0	Row 2	Row 3
Nonzero values	data[7]	{ 3, 1,	2, 4, 1,	1, 1 }
Column indices	col_index[7]	{ 0, 2,	1, 2, 3,	0, 3 }
Row indices	row_index[7]	{ 0, 0,	2, 2, 2,	3, 3 }

©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

3

3

COO Allows Reordering of Elements

		Row 0	Row 2	Row 3
Nonzero values	data[7]	{ 3, 1,	2, 4, 1,	1, 1 }
Column indices	col_index[7]	{ 0, 2,	1, 2, 3,	0, 3 }
Row indices	row_index[7]	{ 0, 0,	2, 2, 2,	3, 3 }

Nonzero values	data[7]	{ 1 1, 2, 4, 3, 1 1 }
Column indices	col_index[7]	{ 0 2, 1, 2, 0, 3, 3 }
Row indices	row_index[7]	{ 3 0, 2, 2, 0, 2, 3 }

©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

4

4

```

1. for (int i = 0; i < num_elem; row++)
2.   y[row_index[i]] += data[i] * x[col_index[i]];

```

a sequential loop that implements SpMV/COO

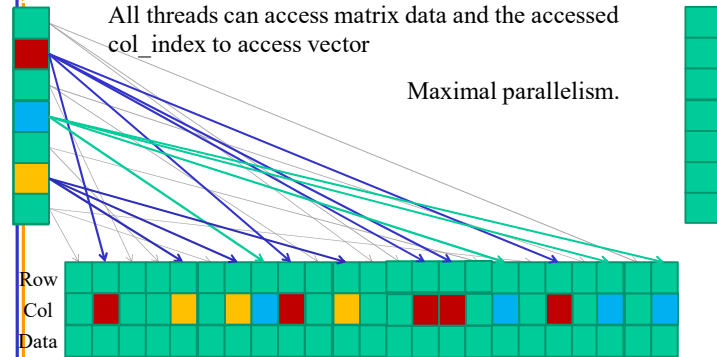
©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

5

COO Kernel Design Accessing Input Matrix and Vector

All threads can access matrix data and the accessed
col_index to access vector

Maximal parallelism.



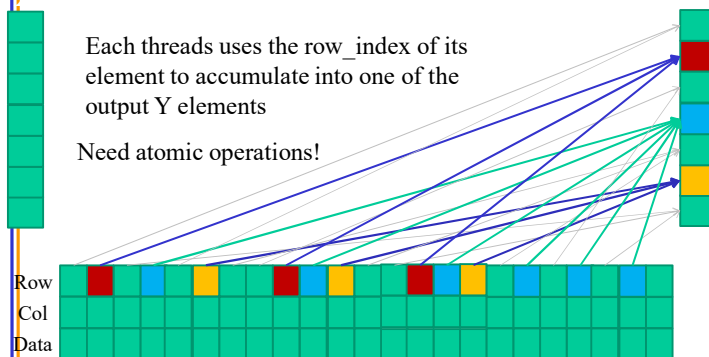
©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

6

COO kernel Design Accumulating into Output Vector

Each threads uses the row_index of its
element to accumulate into one of the
output Y elements

Need atomic operations!



©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

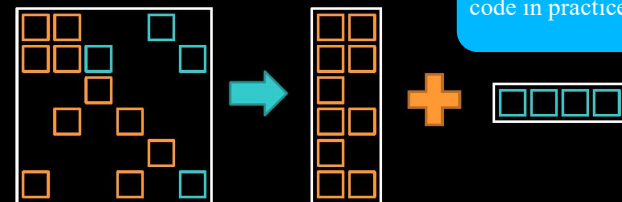
7

Hybrid Format



- ELL handles *typical* entries
- COO handles *exceptional* entries
 - Implemented with segmented reduction

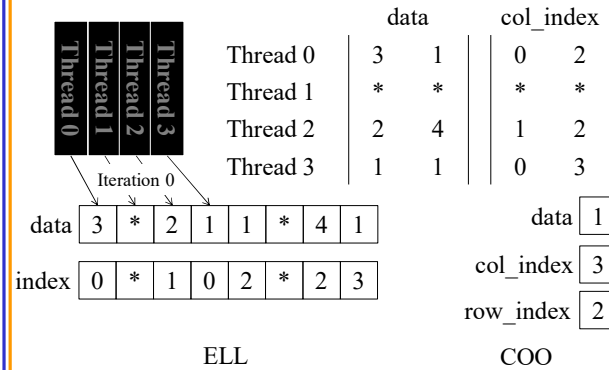
Often implemented
in sequential host
code in practice



© 2010 NVIDIA Corporation

8

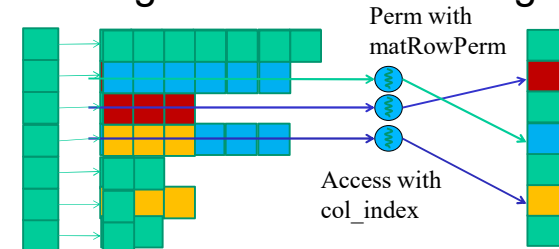
Reduced Padding with Hybrid Format



©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

9

JDS (Jagged Diagonal Sparse) Kernel Design for Load Balancing

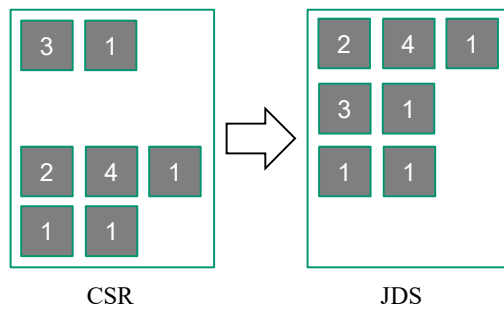


Sort rows into descending order according to number of non-zero. Keep track of the original row numbers so that the output vector can be generated correctly.

©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

10

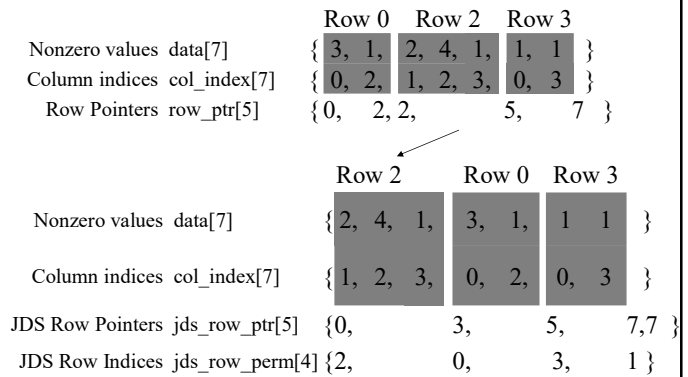
Sorting Rows According to Length (Regularization)



©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

11

CSR to JDS Conversion

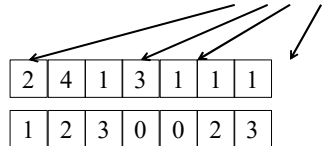


©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

12

JDS Summary

Nonzero values data[7] { 2, 4, 1, 3, 1, 1, 1 }
 Column indices Jds_col_index[7] { 1, 2, 3, 0, 2, 0, 3 }
 JDS row indices Jds_row_perm[4] { 2, 0, 3, 1 }
 JDS Row Ptrs Jds_row_ptr[5] { 0, 3, 5, 7, 7 }



©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

13

13

A Parallel SpMV/JDS Kernel

```
1. __global__ void SpMV_JDS(int num_rows, float *data,
    int *col_index, int *jds_row_ptr, int *jds_row_perm,
    float *x, float *y) {
2.   int row = blockIdx.x * blockDim.x + threadIdx.x;
3.   if (row < num_rows) {
4.     float dot = 0;
5.     int row_start = jds_row_ptr[row];
6.     int row_end = jds_row_ptr[row+1];
7.     for (int elem = row_start; elem < row_end; elem++) {
8.       dot += data[elem] * x[col_index[elem]];
9.     }
10.    y[jds_row_perm[row]] = dot;
11.  }
```

Nonzero values data[7] { 2, 4, 1, 3, 1, 1, 1 }
 Column indices col_index[7] { 1, 2, 3, 0, 2, 0, 3 }
 JDS Row Pointers jds_row_ptr[5] { 0, 3, 5, 7, 7 }
 JDS Row Indices jds_row_perm[4] { 2, 0, 3, 1 }

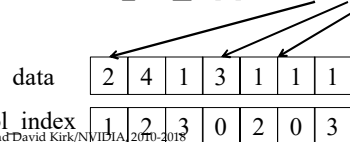
©

14

JDS vs. CSR - Control Divergence

- Threads still execute different number of iterations in the JDS kernel for-loop
 - However, neighboring threads tend to execute similar number of iterations because of sorting.
 - Better thread utilization, less control divergence

Nonzero values data[7] { 2, 4, 1, 3, 1, 1, 1 }
 Column indices col_index[7] { 1, 2, 3, 0, 2, 0, 3 }
 JDS row indices Jds_row_perm[4] { 2, 0, 3, 1 }
 JDS Row Ptrs Jds_row_ptr[5] { 0, 3, 5, 7, 7 }



©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

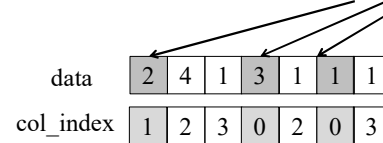
15

15

JDS vs. CSR Memory Divergence

- Adjacent threads still access non-adjacent memory locations

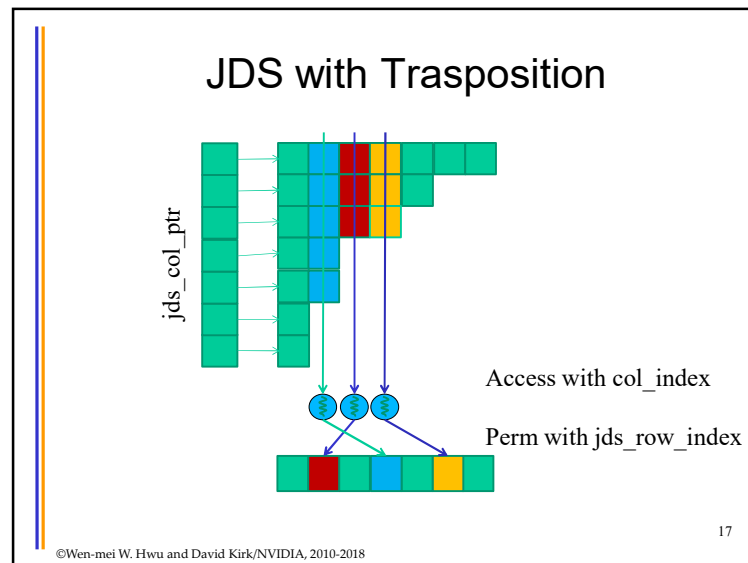
Nonzero values data[7] { 2, 4, 1, 3, 1, 1, 1 }
 Column indices col_index[7] { 1, 2, 3, 0, 2, 0, 3 }
 JDS row indices jds_row_perm[4] { 2, 0, 3, 1 }
 JDS Row Ptrs jds_row_ptr[5] { 0, 3, 5, 7, 7 }



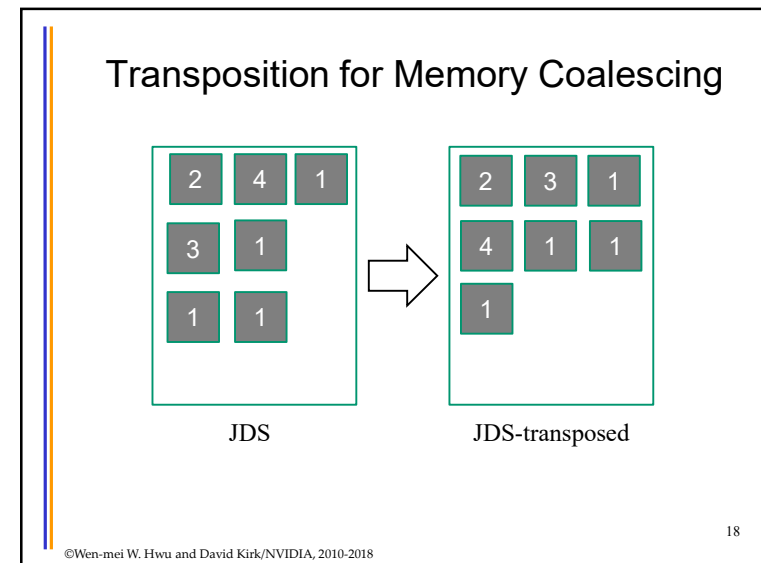
©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

16

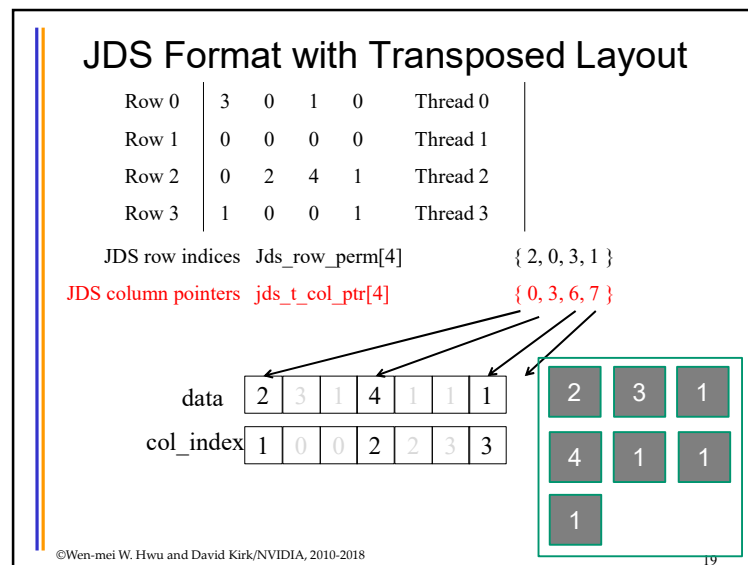
16



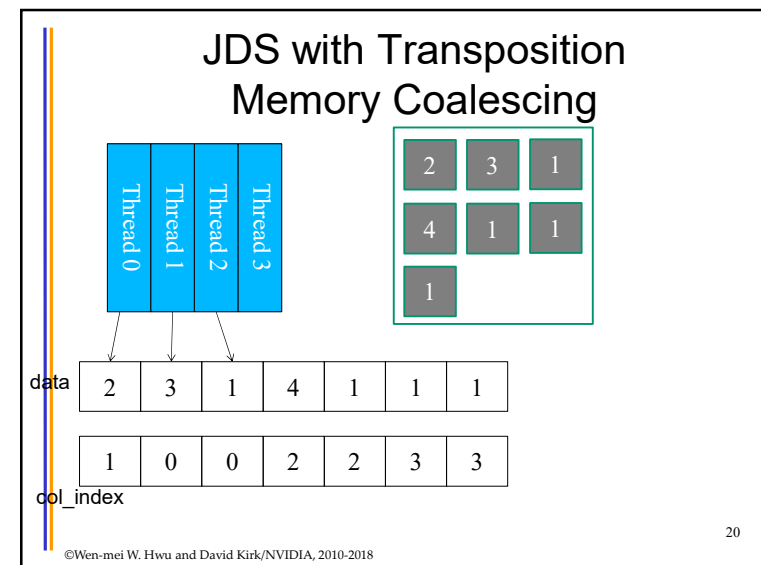
17



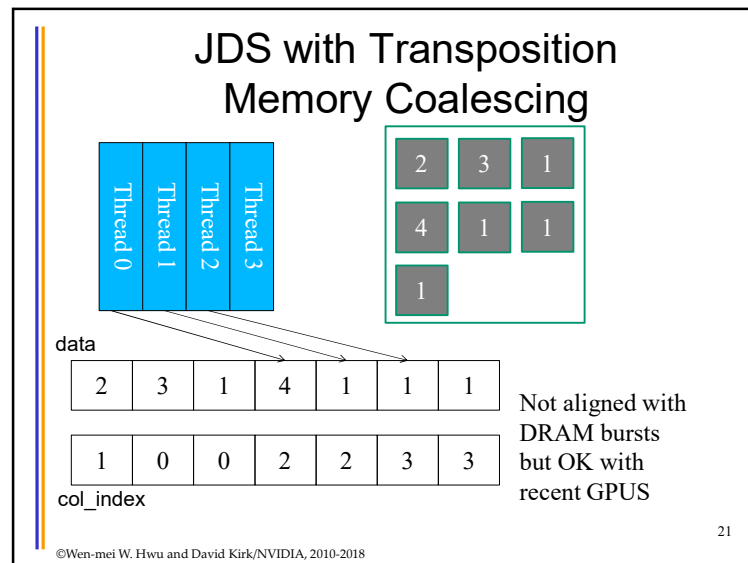
18



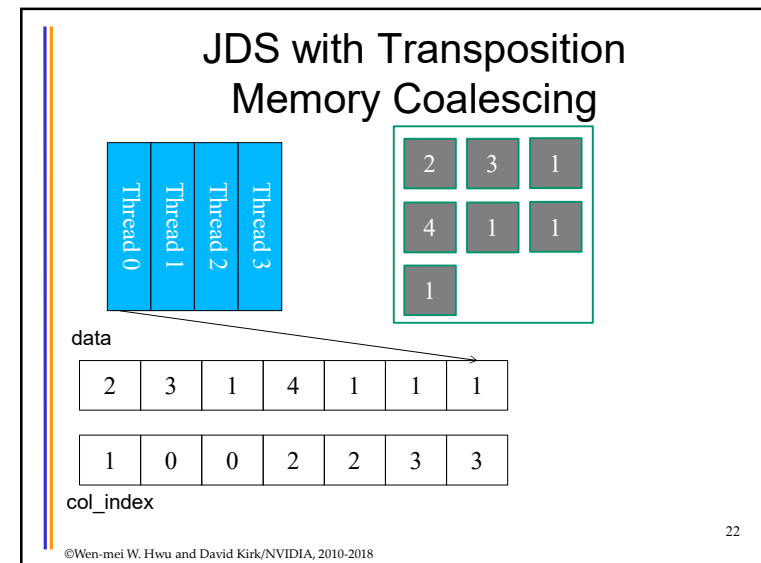
19



20



21



22

A Parallel SpMV/JDS_T Kernel

```

1. __global__ void SpMV_JDS_T(int num_rows, float *data,
   int *col_index, int *jds_t_col_ptr, int *jds_row_perm,
   float *x, float *y) {
2.   int row = blockIdx.x * blockDim.x + threadIdx.x;
3.   if (row < num_rows) {
4.     float dot = 0;
5.     unsigned int sec = 0;
6.     while (jds_t_col_ptr[sec+1]-jds_t_col_ptr[sec] > row) {
7.       dot += data[jds_t_col_ptr[sec]+row] *
8.         x[col_index[jds_t_col_ptr[sec]+row]];
9.       sec++;
10.    }
11.    y[jds_row_perm[row]] = dot;
12.  }
13.}

```

Column indices col_index[7] { 1, 0, 3, 2, 2, 3, 3 }

JDS_T Column Pointers jds_t_col_ptr[5] { 0, 3, 6, 7, 7 }

JDS Row Indices jds_row_perm[4] { 2, 0, 3, 1 }

©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

23

MP7 Variable Names

JDS_T Length of Cols matRows[4] { 3, 2, 2, 0 }

	Sec 0	Sec 1	Sec 2
Nonzero values matData[7]	{ 2, 3, 1 }	{ 4, 1, 1 }	{ 1 }
Column indices matCols[7]	{ 1, 0, 3 }	{ 2, 2, 3 }	{ 3 }
JDS_T Column Pointers matColStart[4]	{ 0, 3, 6, 7 }		
JDS Row Indices matRowPerm[4]	{ 2, 0, 3, 1 }		

©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

24

Sparse Matrices as Foundation for Advanced Algorithm Techniques

- Graphs are often represented as sparse adjacency matrices
 - Used extensively in social network analytics, natural language processing, etc.
- Binning techniques often use sparse matrices for data compaction
 - Used extensively in ray tracing, particle-based fluid dynamics methods, and games
- These will be covered in ECE508/CS508

25

©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

25

ANY QUESTIONS?

26

©Wen-mei W. Hwu and David Kirk/NVIDIA, 2010-2018

26