

15-418/618, Spring 2017

Recitation Week 7

Randal E. Bryant

3 March, 2017

1 Logistics

All files can be found in:

`/afs/cs.cmu.edu/academic/class/15418-s17/recitations/recw7`

2 Sparse matrix times dense vector

All code in subdirectory `mvmul`.

Data structures and declarations

In file `matrix.h`.

Base data types:

```
1 // Underlying data type for representing matrix indices.
2 typedef unsigned index_t;
3 // Type of matrix data
4 typedef float data_t;
```

Vectors:

```
1 // Dense representation of vector
2 typedef struct {
3     index_t length;
4     // Array of values
```

```

5     data_t *value;
6 } vec_t;

```

Compressed Sparse Row (CSR) Representation of matrix. Note introduction of `rindex`. Indicates row number for each nonzero matrix element.

```

1 // CSR representation of square matrix
2 typedef struct {
3     index_t nrow;           // Number of rows (= number of columns)
4     index_t nnz;            // Number of nonzero elements
5     data_t *value;          // Nonzero matrix values, row-major order [nnz]
6     index_t *cindex;        // Column index for each nonzero entry [nnz]
7     index_t *rindex;        // Row index for each nonzero entry [nnz]
8     index_t *rowstart;      // Offset of each row [nrow+1]
9 } csr_t;

```

Row-oriented multiplication

In file `mvmul.cpp`.

Using static parallelism over rows

```

1 // Using OMP static parallel over rows
2 void mul_csr_mrow(csr_t *m, vec_t *x, vec_t *y) {
3     index_t nrow = m->nrow;
4     #pragma omp parallel for schedule(static)
5     for (index_t r = 0; r < nrow; r++) {
6         index_t idxmin = m->rowstart[r];
7         index_t idxmax = m->rowstart[r+1];
8         float val = 0.0;
9         for (index_t idx = idxmin; idx < idxmax; idx++) {
10             index_t c = m->cindex[idx];
11             data_t mval = m->value[idx];
12             data_t xval = x->value[c];
13             val += mval * xval;
14         }
15         y->value[r] = val;
16     }
17 }

```

Measured on `latedays.andrew.cmu.edu`:

- Xeon E5-2620
- 2.40 GHz

- 15MB L3 cache
- 12 cores, each with 2-way hyperthreading

Test for matrix with N rows, for $N = 10^3$ and $N = 10^4$. All matrices have 10% density. The number of nonzero entries is $M = 0.1 \cdot N^2$, i.e., either $M = 10^5$ or $M = 10^7$.

Random data distributed among matrix rows by two different conventions:

Uniform (U): Each row has (almost) the same number of elements

Skewed (S): All data packed into as few rows as possible

Performance expressed in Gigafllops. Matrix with N rows and M non-zero elements requires $2M$ floating-point operations (either 2×10^5 or 2×10^7 .)

Method	$N = 10^3$		$N = 10^4$	
	U	S	U	S
Sequential, row-oriented	2.2	1.9	1.6	1.7
Parallel, row-oriented	4.3	0.7	4.3	1.3

Data-oriented multiplication, updating memory

Idea: compute in parallel over the nonzero data entries, regardless of which rows they are in. Simplest version performs atomic additions to the result vector.

```

1  // Parallel over elements, updates to memory
2  void mul_csr_data_mps(csr_t *m, vec_t *x, vec_t *y) {
3      clear_vector(y);
4      #pragma omp parallel for schedule(static)
5      for (index_t idx = 0; idx < m->nnz; idx++) {
6          index_t r = m->rindex[idx];
7          index_t c = m->cindex[idx];
8          data_t xval = x->value[c];
9          data_t mval = m->value[idx];
10         #pragma omp atomic
11         y->value[r] += mval * xval;
12     }
13 }
```

Method	$N = 10^3$		$N = 10^4$	
	U	S	U	S
Sequential, row-oriented	2.2	1.9	1.6	1.7
Parallel, row-oriented	4.3	0.7	4.3	1.3
Sequential, data-oriented. Update memory	0.7	0.7	0.7	0.7
Parallel, data-oriented. Update memory	1.9	0.1	2.3	2.1

Questions:

1. Why is the sequential data-oriented code slower than row-oriented code?
2. Why is the parallel performance so poor?
3. What does this tell you about the performance of atomic operations?

Data-oriented multiplication, updating registers

Can take advantage of fact that values of `rindex` will be monotonic. Use variable `last_r` to track when transition from one row to the next.

```

1 // Sequential over elements, updates to registers
2 void mul_csr_rdata_seq(csr_t *m, vec_t *x, vec_t *y) {
3     clear_vector(y);
4     index_t last_r = 0;
5     data_t val = 0.0;
6     for (index_t idx = 0; idx < m->nnz; idx++) {
7         index_t r = m->rindex[idx];
8         if (r != last_r) {
9             if (val != 0.0)
10                 y->value[last_r] += val;
11             last_r = r;
12             val = 0.0;
13         }
14         index_t c = m->cindex[idx];
15         data_t xval = x->value[c];
16         data_t mval = m->value[idx];
17         val += mval * xval;
18     }
19     if (val != 0.0)
20         y->value[last_r] += val;
21 }

```

Questions:

1. Function `clear_vector` sets elements of vector to zero. It uses `memset`. Why is this needed?
2. Is the test `val != 0.0` necessary? What would be the effect of removing it?
3. Why is the final update required?
4. What would happen if attempt to parallelize the loop with an OpenMP pragma?
 - (a) Would variable `val` be thread-local or global?
 - (b) Would variable `last_r` be thread-local or global?
 - (c) Would variable `r` be thread-local or global?
 - (d) If the matrix has 1600 nonzero entries, and there are 16 threads:
 - i. What value would thread i initially assign to `idx`?
 - ii. What value would thread i initially assign to `r`?
 - iii. What would happen with the test `r != last_r`?
 - (e) What would be the overall effect of the function?

Understanding OpenMP parallelism:

- Declaring code block to be `parallel` indicates that program should split into multiple threads, each running the code block independently, except when other forms of synchronization are specified.
- Any variable declared within block is thread local.
- Declaring a for loop to be parallelized with static scheduling means that the different threads will divide the loop indices (almost) evenly.
- Declaring a particular read+write operation to be atomic means that only one thread may perform the operation at a time.

```

1 // Parallel over elements, updates to registers
2 // Shared destination vector
3 void mul_csr_rdata_mps(csr_t *m, vec_t *x, vec_t *y) {
4     clear_vector(y);
5     #pragma omp parallel
6     {
7         index_t last_r = 0; // Private to thread
8         data_t val = 0.0;   // Private to thread
9         #pragma omp for schedule(static)
10        for (index_t idx = 0; idx < m->nnz; idx++) {

```

```

11         index_t r = m->rindex[idx];
12         if (r != last_r) {
13             if (val != 0.0)
14                 #pragma omp atomic
15                 y->value[last_r] += val;
16             last_r = r;
17             val = 0.0;
18         }
19         index_t c = m->cindex[idx];
20         data_t xval = x->value[c];
21         data_t mval = m->value[idx];
22         val += mval * xval;
23     }
24     if (val != 0.0)
25         #pragma omp atomic
26         y->value[last_r] += val;
27 }
28 }

```

Questions:

1. What would be the effect of removing the test `val != 0.0`?
2. How many threads will update element $r > 0$ of the result vector?
3. Would variable `r` be thread-local or global?
4. If the matrix has 1600 nonzero entries, and there are 16 threads:
 - (a) What value would thread i initially assign to `idx`?
 - (b) What value would thread i initially assign to `r`?
 - (c) What would happen with the test `r != last_r`?
5. What would be the overall effect of the function?

Method	$N = 10^3$		$N = 10^4$	
	U	S	U	S
Sequential, row-oriented	2.2	1.9	1.6	1.7
Parallel, row-oriented	4.3	0.7	4.3	1.3
Sequential, data-oriented. Update memory	0.7	0.7	0.7	0.7
Parallel, data-oriented. Update memory	1.9	0.1	2.3	2.1
Sequential, data-oriented. Update register	1.7	1.8	1.5	1.7
Parallel, data-oriented. Update register	5.8	4.4	3.5	2.0

Questions:

1. How does the sequential performance compare to the other two sequential versions? Explain.
2. How sensitive is this version to skewed data?
3. Why does the performance drop off for large values of N ?

Avoiding memory contention

Use thread-local vectors to accumulate result values. Track range of rows values computed by each thread to reduce number of updates.

```
1 // Parallel over elements, updates to registers
2 // Separate destination vector
3 void mul_csr_rldata_mps(csr_t *m, vec_t *x, vec_t *y) {
4     clear_vector(y);
5     #pragma omp parallel
6     {
7         index_t nrow = m->nrow;
8         index_t last_r = 0;           // Private to thread
9         data_t val = 0.0;             // Private to thread
10        data_t local_y[nrow];         // Private to thread
11        index_t min_r = 0;            // Private to thread
12        bool first = true;            // Private to thread
13        memset((void *) local_y, 0, nrow * sizeof(data_t));
14        #pragma omp for schedule(static) nowait
15        for (index_t idx = 0; idx < m->nnz; idx++) {
16            index_t r = m->rindex[idx];
17            if (first)
18                min_r = last_r = r;
19            first = false;
20            if (r != last_r) {
21                local_y[last_r] += val;
22                last_r = r;
23                val = 0.0;
24            }
25            index_t c = m->cindex[idx];
26            data_t xval = x->value[c];
27            data_t mval = m->value[idx];
28            val += mval * xval;
29        }
30        local_y[last_r] += val;
31        // Combine local y values
```

```

32     for (index_t r = min_r; r <= last_r; r++) {
33         #pragma omp atomic
34         y->value[r] += local_y[r];
35     }
36 }
37 }

```

Questions

1. Where are the arrays `local_y` allocated?
2. How many calls will the function make to `memset`?
3. How will each thread assign a distinct value to `min_r`?
4. Why is it possible to specify that the for loop can be parallelized with option `nowait`?
5. What would happen if the final loop ranged from 0 to `nrow`?

Method	$N = 10^3$		$N = 10^4$	
	U	S	U	S
Sequential, row-oriented	2.2	1.9	1.6	1.7
Parallel, row-oriented	4.3	0.7	4.3	1.3
Sequential, data-oriented. Update memory	0.7	0.7	0.7	0.7
Parallel, data-oriented. Update memory	1.9	0.1	2.3	2.1
Sequential, data-oriented. Update register	1.7	1.8	1.5	1.7
Parallel, data-oriented. Update register	5.8	4.4	3.5	2.0
Parallel, data-oriented. Local copy of y	6.8	4.9	3.6	2.5

Further improvements

Imagine a scenario where the program must perform many multiplications for a given matrix. This would occur, for example, when computing an iterative solution to a linear system.

1. What data structure could be used for the local values of `y` that would eliminate the need for allocating them on each call?
2. How could the program avoid having to clear out all of the vectors with each call?
3. How could the final summations be performed without requiring any atomic operations?

3 Radix Sorting

All code in subdirectory `radixsort`.

- Sort set of 64-bit values
- Algorithm
 - Make 8 passes, from least significant byte (byte 0) to most (byte 7).
 - On each pass, perform stable sort using byte i as key.
 - Implement each pass with bucket sort.
- Implementation of one pass
 - Bucket has counter for each of the 256 possible key values
 - Read through all data and count number of keys for each bucket
 - Compute starting offsets for each key value
 - Read through all data and place in appropriate location in destination.
- Use two different arrays, alternating between source and destination.

Exercise: Use radix sort to sort the following words in 3 passes:

Input	Pass 1	Pass 2	Pass 3
pig			
pat			
dog			
cat			
boy			

Declarations

From `rsort.h`

```

1 typedef unsigned long data_t;
2 typedef unsigned index_t;
3
4 // How many bits does each pass use
5 #define BASE_BITS 8
6 // What is radix of sort
7 #define BASE (1 << BASE_BITS)
8 // Mask of all 1's over BASE_BITS
9 #define MASK (BASE-1)
10 // Extract sorting key from data
11 #define DIGITS(v, shift) (((v) >> (shift)) & MASK)

```

Sequential code

```

1 // Radix sorting, sequential.
2 void seq_radix_sort(index_t N, data_t *indata,
3                     data_t *outdata, data_t *scratchdata) {
4     data_t *src = indata;
5     // Assume even number of steps, so by toggling between
6     // outdata and scratchdata, result will end up in outdata
7     data_t *dest = scratchdata;
8     index_t count[BASE];
9     index_t offset[BASE];
10    index_t total_digits = sizeof(data_t) * 8;
11
12    for (index_t shift = 0; shift < total_digits; shift+= BASE_BITS) {
13        memset(count, 0, BASE*sizeof(index_t));
14        // Accumulate count for each key value
15        for (index_t i = 0; i < N; i++) {
16            data_t key = DIGITS(src[i], shift);
17            count[key]++;
18        }
19        // Compute offsets
20        offset[0] = 0;
21        for (index_t b = 1; b < BASE; b++)
22            offset[b] = offset[b-1] + count[b-1];
23        // Distribute data
24        for (index_t i = 0; i < N; i++) {
25            data_t key = DIGITS(src[i], shift);
26            index_t pos = offset[key]++;
27            dest[pos] = src[i];
28        }
29        // Find new src & dest

```

```

30         src = dest;
31         dest = (dest == outdata) ? scratchdata : outdata;
32     }
33 }

```

Questions:

1. What is the role of array `scratchdata`?
2. What would be the values of `count` and `offset` when performing the first pass for the exercise example?
3. What kind of cache performance would be seen for the data reads (array `src`) and writes (array `dest`)?

Parallel version

This version is derived from code made available online by Haichuan Wang, based on code he wrote for a class at UIUC. It uses lots of clever tricks.

```

1  // Variation of code due to Haichuan Wang, UIUC
2  void full_par_radix_sort(index_t N, data_t *indata,
3                          data_t *outdata, data_t *scratchdata) {
4      data_t *src = indata;
5      // Assume even number of steps, so by toggling between
6      // outdata and scratchdata, result will end up in outdata
7      data_t *dest = scratchdata;
8      index_t total_digits = sizeof(data_t) * 8;
9      index_t count[BASE];
10     index_t offset[BASE];
11
12     for(index_t shift = 0; shift < total_digits; shift+=BASE_BITS) {
13         memset(count, 0, BASE*sizeof(index_t));
14         #pragma omp parallel
15         {
16             // Per-thread counts and offsets
17             index_t local_count[BASE] = {0};
18             index_t local_offset[BASE];
19             #pragma omp for schedule(static) nowait
20             // Gather counts on per-thread basis
21             for(index_t i = 0; i < N; i++){
22                 data_t key = DIGITS(src[i], shift);
23                 local_count[key]++;
24             }

```

```

25      // Compute global counts based on local ones
26      // Critical faster than each addition being atomic
27      #pragma omp critical
28      for(index_t b = 0; b < BASE; b++) {
29          count[b] += local_count[b];
30      }
31      #pragma omp barrier
32      // Compute global offsets
33      #pragma omp single
34      {
35          offset[0] = 0;
36          for (index_t b = 1; b < BASE; b++)
37              offset[b] = count[b-1] + offset[b-1];
38      }
39      int nthreads = omp_get_num_threads();
40      int tid = omp_get_thread_num();
41      // Compute local offsets
42      // Enforce serialization by triggering each thread in sequence
43      for (int t = 0; t < nthreads; t++) {
44          if(t == tid) {
45              for(index_t b = 0; b < BASE; b++) {
46                  local_offset[b] = offset[b];
47                  offset[b] += local_count[b];
48              }
49          }
50          #pragma omp barrier
51      }
52      #pragma omp for schedule(static)
53      for(index_t i = 0; i < N; i++) {
54          data_t key = DIGITS(src[i], shift);
55          index_t pos = local_offset[key]++;
56          dest[pos] = src[i];
57      }
58  }
59      // Find new src & dest
60      src = dest;
61      dest = (dest == outdata) ? scratchdata : outdata;
62  }
63  }

```

Questions:

1. Why are the updates to `count` placed in critical section (rather than either unsynchronized, or atomic)?

2. What is the purpose of the `barrier` pragma?
3. What is the effect of the `single` pragma?
4. What happens in the loop over thread ids?
 - (a) In what order will the copies of `local_offset` be computed?
 - (b) Why couldn't the computation of `local_offset` be done using a simple critical section?
 - (c) Why are the values of `offset` updated?

Performance

Express in megawords of data sorted per second.

Method	$N = 10^6$	$N = 10^7$
Library quicksort	28.1	22.1
Sequential radix	31.7	26.1
Parallel radix	330.0	67.9

Questions:

1. Why does the performance of quicksort fall off faster than radix sort for large values of N ?
2. What limits the parallel performance for large values of N ?