

ACM/CS 114

Parallel algorithms for scientific applications

Michael A. G. Aïvázis

California Institute of Technology

Winter 2012

Required changes to the sequential solution

- ▶ what is needed
 - ▶ an object to hold the problem information shared among the threads
 - ▶ the per-thread administrative data structure that holds the thread id and the pointer to the shared information
 - ▶ this is the argument to `pthread_create`
 - ▶ a mutex to protect the update of the global convergence criterion
 - ▶ a `pthread_create` compatible worker routine
 - ▶ a change at the top-level driver to enable the user to choose the number of threads
- ▶ and a strategy for managing the thread life cycle
 - ▶ synchronization is trivial if
 - ▶ we spawn our threads to perform the updates of a single iteration
 - ▶ harvest them
 - ▶ check the convergence criterion
 - ▶ stop, or respawn them if another iteration is necessary
 - ▶ can the convergence test be done in parallel?
 - ▶ so we don't have to pay the create/harvest overhead?
 - ▶ if so, how do we guarantee correctness and consistency?

Threaded Jacobi: thread data

```
1 struct Task {
2     // shared information
3     size_t workers;
4     Grid & current;
5     Grid & next;
6     double maxDeviation;
7     // mutex to control access to the convergence criterion
8     pthread_mutex_t lock;
9
10    // constructor
11    Task(size_t workers, Grid & current, Grid & next) :
12        workers(workers), current(current), next(next), maxDeviation(0.0) {
13        pthread_mutex_init(&lock, 0);
14    }
15    // destructor
16    ~Task() {
17        pthread_mutex_destroy(&lock);
18    }
19 };
20
21 struct Context {
22     // thread info
23     size_t id;
24     pthread_t descriptor;
25     Task * task;
26 };
```

Threaded Jacobi: driving the update

```
28 void Jacobi::solve(Problem & problem) {
29     // initialize the problem
30     problem.initialize();
31     // do the actual solve
32     _solve(problem);
33     // compute and store the error
34     std::cout << " computing absolute error" << std::endl;
35     // compute the relative error
36     Grid & error = problem.error();
37     const Grid & exact = problem.exact();
38     const Grid & solution = problem.solution();
39
40     for (size_t j=0; j < exact.size(); j++) {
41         for (size_t i=0; i < exact.size(); i++) {
42             if (exact(i,j) == 0.0) {
43                 error(i,j) = std::abs(solution(i,j));
44             } else {
45                 error(i,j) = std::abs(solution(i,j) - exact(i,j))/exact(i,j);
46             }
47         }
48     }
49     std::cout << " --- done." << std::endl;
50     return;
51 }
```

Threaded Jacobi: the master thread

```
52 void Jacobi::_solve(Problem & problem) {  
53     Grid & current = problem.solution();  
54  
55     // create and initialize temporary storage  
56     Grid next(current.size());  
57     problem.initialize(next);  
58  
59     // shared thread info  
60     Task task(_workers, current, next);  
61     // per-thread information  
62     Context context[_workers];  
63  
64     // let's get going  
65     std::cout << "jacobi: tolerance=" << _tolerance << std::endl;  
66  
67     // put an upper bound on the number of iterations  
68     const size_t max_iterations = (size_t) 1.0e4;
```

Threaded Jacobi: the master thread, part 2

```
69 for (size_t iterations = 0; iterations < max_iterations; iterations++) {
70     if (iterations % 100 == 0) {
71         std::cout << " " << iterations << std::endl;
72     }
73     // reset the maximum deviation
74     task.maxDeviation = 0.0;
75     // spawn the threads
76     for (size_t tid=0; tid < _workers; tid++) {
77         context[tid].id = tid;
78         context[tid].task = &task;
79
80         int status = pthread_create(&context[tid].descriptor, 0, _update, &context[tid]);
81         if (status) {
82             throw ("error in pthread_create");
83         }
84     }
85     // harvest the threads
86     for (size_t tid = 0; tid < _workers; tid++) {
87         pthread_join(context[tid].descriptor, 0);
88     }
89
90     // swap the blocks between the two grids
91     Grid::swapBlocks(current, next);
92     // check convergence
93     if (task.maxDeviation < _tolerance) {
94         std::cout << " ### convergence in " << iterations << " iterations!" << std::endl;
95         break;
96     }
97 }
98 std::cout << " --- done." << std::endl;
99
100 return;
101 }
```

Threaded Jacobi: update in the worker threads

```
102 void * Jacobi::_update(void * arg) {
103     Context * context = static_cast<Context *>(arg);
104
105     size_t id = context->id;
106     Task * task = context->task;
107
108     size_t workers = task->workers;
109     Grid & current = task->current;
110     Grid & next = task->next;
111     pthread_mutex_t lock = task->lock;
112
113     double max_dev = 0.0;
114     // do an iteration step
115     // leave the boundary alone
116     // iterate over the interior of the grid
117     for (size_t j=id+1; j < current.size()-1; j+=workers) {
118         for (size_t i=1; i < current.size()-1; i++) {
119             next(i,j) = 0.25*(current(i+1,j)+current(i-1,j)+current(i,j+1)+current(i,j-1));
120             // compute the deviation from the last generation
121             double dev = std::abs(next(i,j) - current(i,j));
122             // and update the maximum deviation
123             if (dev > max_dev) {
124                 max_dev = dev;
125             }
126         }
127     }
128
129     // grab the lock and update the global maximum deviation
130     pthread_mutex_lock(&lock);
131     if (task->maxDeviation < max_dev) {
132         task->maxDeviation = max_dev;
133     }
134     pthread_mutex_unlock(&lock);
135
136     return 0;
137 }
```

Assessing the threaded implementation

- ▶ the implemented synchronization scheme is very simple
 - ▶ each grid update step spawns some number of workers to update a subset of the cells
 - ▶ the workers are harvested after the grid is updated
 - ▶ the main thread checks for convergence
 - ▶ if another iteration is required, a new set of workers is spawned
- ▶ the simplicity of this strategy comes at a cost
 - ▶ *scalability* suffers when the overhead of creating and harvesting threads is comparable to amount of work done by each thread
 - ▶ for low thread counts, it is still an overall win, since the time to solution decreases and the machine utilization is better
 - ▶ but as the number of threads increases, the program becomes *slower*
 - ▶ timing a 100×100 grid to convergence on a recent MacPro

threads	1	2	4	8	16
time(s)	4.367	2.517	1.918	1.937	3.537

- ▶ and 10,000 iterations of a 1000×1000 grid

threads	1	2	4	8	16
time(s)	413.306	211.050	109.509	98.279	74.087

Improving the update loop

- ▶ the plan is to keep the workers alive and updating the grid while either we converge or `max_iterations` is reached
- ▶ the main thread
 - ▶ loops to spawn all the threads
 - ▶ and immediately enters a loop to harvest them
- ▶ the workers use a condition variable to synchronize among themselves
 - ▶ they iterate, updating the grid
 - ▶ grab a mutex, deposit their local maximum deviation from the last iterations, update a counter that records how many workers have completed their update, and release the lock
 - ▶ enter another critical section with the termination logic
 - ▶ everybody uses a condition variable to wait for the slowest worker
 - ▶ the slowest worker checks the convergence criterion and updates the termination flag, swaps the grid blocks and signals everybody else
 - ▶ if the termination flag is set, or if the maximum number of iterations has been reached, all threads exit

Threaded Jacobi: the main thread

```
1 void Jacobi::_solve(Problem & problem) {
2     Grid & current = problem.solution();
3
4     // create and initialize temporary storage
5     Grid next(current.size());
6     problem.initialize(next);
7
8     // shared thread info
9     Task task(_workers, _tolerance, current, next);
10    // per-thread information
11    Context context[_workers];
12    // spawn the threads
13    std::cout << "jacobi: spawning " << _workers << " workers" << std::endl;
14    for (size_t tid=0; tid < _workers; tid++) {
15        context[tid].id = tid;
16        context[tid].task = &task;
17
18        int status = pthread_create(&context[tid].descriptor, 0, _update, &context[tid]);
19        if (status) {
20            throw ("error in pthread_create");
21        }
22    }
23    // harvest the threads
24    for (size_t tid = 0; tid < _workers; tid++) {
25        pthread_join(context[tid].descriptor, 0);
26    }
27    // done
28    std::cout << "jacobi: done." << std::endl;
29    return;
30 }
```

Threaded Jacobi: updated thread data

```
1 struct Task {
2     // shared information
3     size_t workers; // the number of threads
4     double tolerance; // the convergence tolerance
5     Grid & current;
6     Grid & next;
7
8     bool done; // is there more work?
9     double maxDeviation; // the value
10    size_t contributions; // the number of threads that have deposited contributions
11    pthread_mutex_t gridUpdate_lock; //the mutex
12    pthread_cond_t gridUpdate_check;
13
14    Task(size_t workers, double tolerance, Grid & current, Grid & next) :
15        workers(workers), tolerance(tolerance), current(current), next(next),
16        done(false), maxDeviation(0.0), contributions(0),
17        gridUpdate_lock(), gridUpdate_check() {
18        // initialize the grid update lock
19        pthread_mutex_init(&gridUpdate_lock, 0);
20        pthread_cond_init(&gridUpdate_check, 0);
21    }
22
23    ~Task() {
24        pthread_mutex_destroy(&gridUpdate_lock);
25        pthread_cond_destroy(&gridUpdate_check);
26    }
27 };
```

Threaded Jacobi: workers, part 1

```
31 // the threaded update
32 void * Jacobi::_update(void * arg) {
33     Context * context = static_cast<Context *>(arg);
34
35     size_t id = context->id;
36     Task * task = context->task;
37
38     const size_t workers = task->workers;
39     Grid & current = task->current;
40     Grid & next = task->next;
41
42     size_t maxIterations = (size_t) 1e4;
43     // iterate, updating the grid until done
44     for (size_t iteration = 0; iteration < maxIterations; iteration++) {
45         // thread 0: print an update
46         if (id == 0 && iteration % 100 == 0) {
47             std::cout << " " << iteration << std::endl;
48         }
49
50         double max_dev = 0.0;
51         // do an iteration step
52         // leave the boundary alone
53         // iterate over the interior of the grid
54         for (size_t j=id+1; j < current.size()-1; j+=workers) {
55             for (size_t i=1; i < current.size()-1; i++) {
56                 next(i,j) = 0.25*(current(i+1,j)+current(i-1,j)+current(i,j+1)+current(i,j-1));
57                 // compute the deviation from the last generation
58                 double dev = std::abs(next(i,j) - current(i,j));
59                 // and update the maximum deviation
60                 if (dev > max_dev) {
61                     max_dev = dev;
62                 }
63             }
64         }
65         // done with the grid update
```

Threaded Jacobi: workers, part 2

```
66 // grab the grid update lock
67 pthread_mutex_lock(&task->gridUpdate_lock);
68 // update the global maximum deviation
69 if (task->maxDeviation < max_dev) {
70     task->maxDeviation = max_dev;
71 }
72 // leave a mark
73 task->contributions++;
74 // bookkeeping at the end of the update
75 if (task->contributions == workers) {
76     // if i am the slowest worker
77     // swap the blocks between the two grids
78     Grid::swapBlocks(current, next);
79     // check convergence
80     if (task->maxDeviation < task->tolerance) {
81         std::cout
82             << " +++ thread " << id << ": convergence in " << iteration << " iterations"
83             <<std::endl;
84         task->done = true;
85     }
86     // reset our accounting and signal everybody
87     task->contributions = 0;
88     task->maxDeviation = 0;
89     pthread_cond_broadcast(&task->gridUpdate_check);
90 } else {
91     // all but the slowest wait here
92     pthread_cond_wait(&task->gridUpdate_check, &task->gridUpdate_lock);
93 }
94 // release
95 pthread_mutex_unlock(&task->gridUpdate_lock);
96 // check whether we are done
97 if (task->done) {
98     break;
99 }
100 }
101 return 0;
102 }
```

Assessing the improved implementation

- ▶ the improved threading scheme is not much more complex
 - ▶ we keep track of how many threads have computed their grid update
 - ▶ the slowest worker check the convergence criterion and performs all the necessary bookkeeping
 - ▶ while everybody else waits
 - ▶ use `pthread_cond_broadcast` to wake the other workers
- ▶ here is the performance comparison for 10,000 iterations on a 1000×1000 grid on the same 8-core MacPro

threads	1	2	4	8	16
previous(s)	413.306	211.050	109.509	98.279	74.087
updated(s)	408.636	208.832	107.015	59.043	61.481