

Topic Notes: Scientific Computing

Distributed Data Structures

Recall that in an SPMD parallel environment, cooperating processes do not share memory, and often coordinate with each other through message passing. One of the motivations for programming in this model, such as with MPI, is that not only can these programs take advantage of more processing cores as compute nodes are added, but also the additional memory resources. We will assume the MPI model for this, where even if running on cores of the same physical node that share memory, each has access only to its own virtual address space and any sharing of information among processes required message passing.

A *distributed data structure* can be used in these situations. With a distributed data structure, the global data structure does not exist in its entirety in any one process' memory space. Each allocates and manages some subset of the structure, possibly including some overlap.

So far, we have looked at distributed data structures that use only arrays. Arrays usually are the easiest to distribute; we simply assign a range of subscripts to each process. As a programmer of a distributed array data structure, we would need to take care to understand the difference between a *local* and a *global indexing*. In the case when there is some duplication, we likely would also need to consider which process is considered the *owner* of any duplicated data.

In your Jacobi solver, solution points were distributed evenly through the domain:



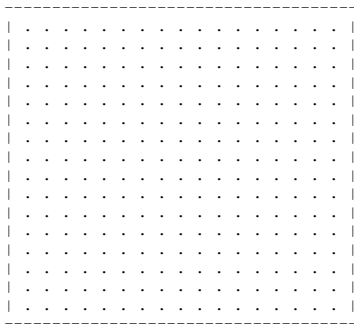
When we distributed this across a team of MPI processes, we assigned a subset of the rows to be assigned to each process, along with a replicated boundary row on each side.

In our programs, we stored only local indices of the distributed rows. For the purposes of the computation, it did not matter what the global index was. The only place it mattered was in the functionality to output the solution data with coordinates.

Adaptivity and Linked Structures

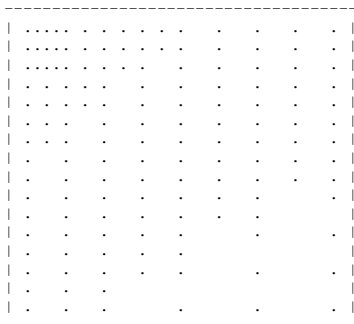
Solution points here are distributed evenly through the domain. That is, each represents the same size subset of the computational domain. The density of these points is what determines how accurately the structure can represent temperatures throughout the domain (since in our simple simulation, we assume temperature is a constant throughout the portion of the domain it models).

If you wanted a more accurate solution, you could add more solution points. But since your program used a uniform distribution of the points, better accuracy requires adding points everywhere.



While what we described above will work to improve our solution accuracy, it increases the total amount of work very quickly.

In some cases, we really only need more points in just parts of the domain. In a heat distribution problem, if there is a heat source in the northwest corner, we may only need more accuracy near that heat source. Fewer points may provide a sufficiently accurate solution further from the source.



This is significantly more efficient in terms of the amount of computation we need to do to obtain a solution of acceptable accuracy.

If we know ahead of time where the extra work is needed, we could assign extra points there at the start of the computation. However, we often do not know this information. After all, we probably wouldn't be solving problems for which we already have a solution handy, so an *adaptive* approach is taken. Periodically, the accuracy of the solution can be checked, and extra points added

as needed. In the context of the Jacobi solver, we may have a threshold for the greatest allowable difference in temperature between adjacent points. If the difference exceeds the threshold, points are added in that vicinity and the solution is recomputed.

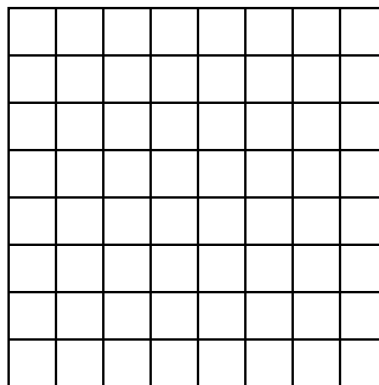
Adaptivity with arrays is difficult, as the completely regular structure is lost. In many cases, a *mesh* structure, often implemented as a linked data structure, is used instead of arrays.

Meshes come in a variety of types. Meshes consisting of quadrilaterals (hexahedra in three dimensions) are sometimes called *structured* meshes. Meshes constructed from triangles (tetrahedra in three dimensions) are called *unstructured* meshes. There are big differences in terms of what happens mathematically when you use them to solve a problem and how hard they are to generate, but for the purpose here, studying how we handle the parallelization of these irregular structures, the issues are similar.

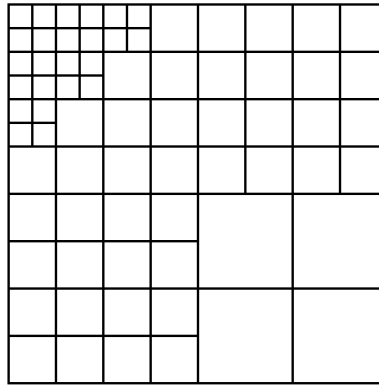
Some terminology: a typical mesh consists of three-dimensional *regions* or *volumes*, and their bounding *faces*, *edges*, and *vertices*. A data structure implementing the mesh will often allow queries such as “what faces bound this region” and “what edges are incident on this vertex” to be made efficiently.

The term “element” often is used to refer to the highest-dimension entity, and in many cases is the entity with which the solution is stored.

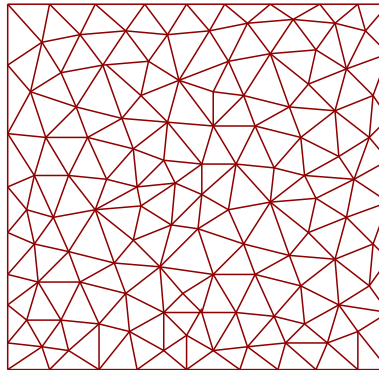
In fact, our Jacobi example is really just a computation on a uniform quadrilateral structured mesh. Here, squares serve as the mesh elements and these contain the solution values.



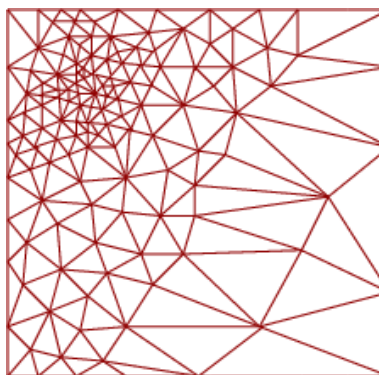
A (hypothetical) *adaptive refinement* of such a mesh might look like:



Here is a simple unstructured mesh:



The use of adaptivity (and some other considerations) necessitate linked data structures for the mesh. Here is the triangular mesh after a hypothetical refinement operation:



Such mesh structures may be stored in memory as arrays of entities or with a full topological hierarchy. We will soon consider an implementation of an adaptive structured mesh using a *quadtree* data structure.

We can see some examples of mesh structures implemented with a full, linked entity hierarchy in the Parallel Unstructure Mesh Infrastructure, from RPI's Scientific Computation Research Center. (See links on lecture page).

A variety of algorithms can be used to determine just how to distribute mesh elements among a set of cooperating processes and several of these *mesh partitioning* algorithms will be topics of study for us very soon.

Computation on Quadrees

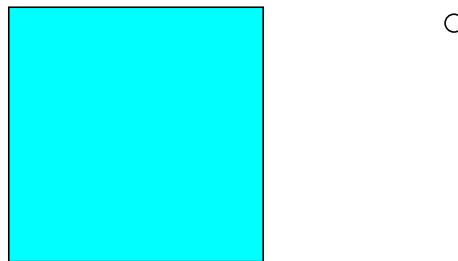
We will consider a relatively simple example of an adaptive computation. We return to our Jacobi Iteration heat solver.

We could approach this problem using an unstructured mesh as described above, but that can get complicated very quickly, even more so than the approach we will take. Instead, we will use a C program to solve Laplace's equation on a square domain using Jacobi iteration that operates on an adaptive *quadtree* structure. We will think about how to parallelize this adaptive Jacobi solver, using the SPMD model with MPI for message passing and finally, how to implement a redistribution procedure to rebalance the load after imbalance is introduced by adaptivity.

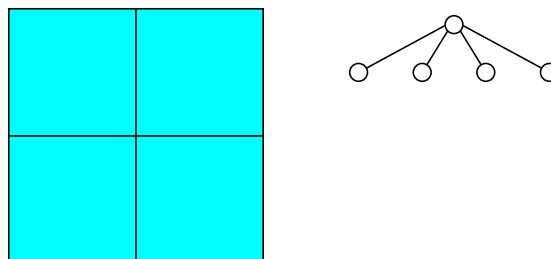
Quadrees

A quadtree is a spatial tree data structure. Each node in the tree is called a *quadrant*. The leaves of the quadtree will be referred to as *leaf quadrants* or *terminal quadrants*. These terminal quadrants will serve as the elements on which we will perform computation.

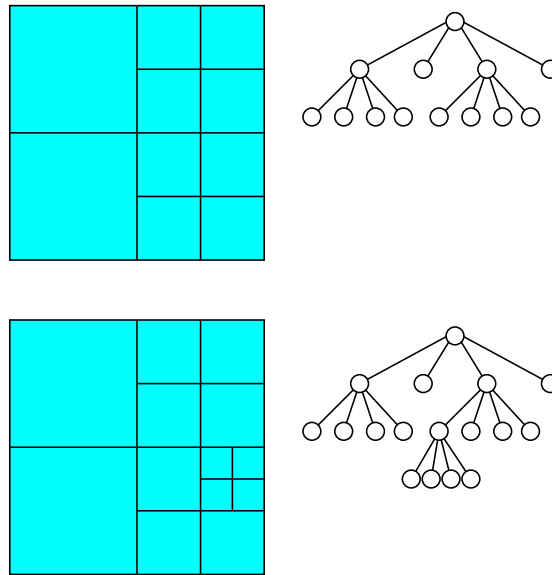
The tree's root represents the entire domain, which, in our case, is a square.



The four children of the root each represent one quarter of the space taken by the root.



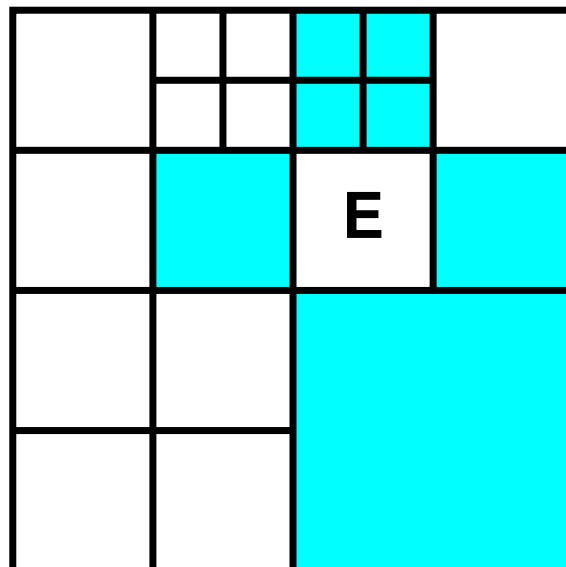
These children can then be divided in four, continuing down as many levels as desired. Different parts of the domain may be refined to different levels.



METAL’s HDX has an algorithm visualization of the construction of a quadtree.

Computing Sequentially on a Quadtree

The Jacobi iteration on a given quadtree is similar to the iteration you’ve done in the array-based versions. However, since there may be a deeper quadtree in some parts of the domain, some leaf quadrants have neighbors more or less refined than themselves.



In the above figure, the element “E” should compute its value based on the values of the shaded neighbor quadrants. Note that E’s north neighbor is actually refined one level further, so it has two immediate north neighbors. Ideally, the north neighbor value should be the average of the two

immediate quadrant neighbors, but it is sufficient to use the value at the shaded quadrant (which is the average of the four leaf quadrants).

One obstacle to overcome is to locate neighboring quadrants efficiently. It is possible to maintain direct neighbor links within your quadtree structure, but you may find it more useful to search based on coordinate information. Each quadrant knows its own bounding box coordinates and from this can easily compute the coordinates of the adjacent quadrants at the same level. A nice feature of the quadtree structure is that the quadrants that contain any point in space can easily be found with a simple traversal from the root, determining the correct child at each step by comparing coordinates. If the neighbor point is outside the root quadrant, you know to apply a boundary condition.

Sequential Program

<https://github.com/SienaCSISParallelProcessing/quadtree-jacobi>

To make the program more interesting in the context of adaptivity, our program will allow a wider range of initial and boundary conditions than your previous implementations.

- Initial conditions are specified by a C function that provides values given (x,y) coordinates of a point in the domain. Each leaf quadrant's solution value is initialized to the value obtained by passing the quadrant's centroid to this function.
- Boundary conditions are also specified through C functions. The left and right boundaries take the y coordinate as a parameter and the top and bottom boundaries take the x coordinate as a parameter, allowing boundary conditions to be functions in addition to simple constants. When a quadrant needs to query a boundary condition when computing its value, it should call the appropriate function.
- Special "internal boundary conditions" are specified by one more C function. The function takes a leaf quadrant and sets its value if the quadrant contains any points that have internal boundary values. For example, if there is a point heat source at (0.25,0.25) keeping the immediate area at a constant temperature of 2, this function will set the solution value of any quadrant containing (0.25,0.25) to 2. This function should be called on each leaf quadrant during each Jacobi iteration step. If the function sets the quadrant's solution value, that value should be used instead of the solution computed based on its neighbors.

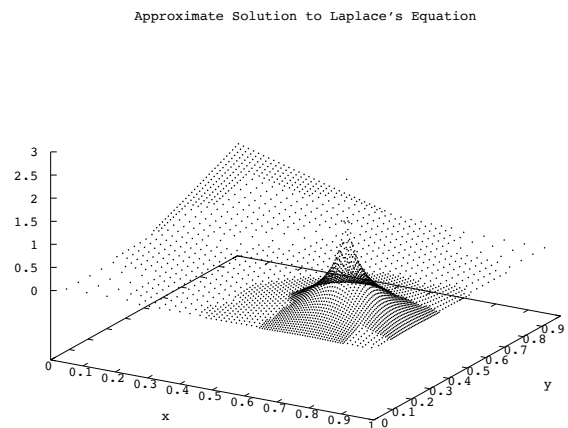
These functions can be hard-coded into our program, but we could implement a system where initial and boundary conditions can be specified through a configuration file instead.

The program will take several parameters:

1. an initial global refinement level, which in turn determines the initial number of leaf quadrants
2. a Jacobi iteration tolerance, similar to the tolerance from the previous implementations
3. a limit on the number of Jacobi iteration steps

4. a refinement threshold, specifying the maximum difference in temperature allowed between adjacent leaf quadrants
5. a limit on the total number of refinement steps

Here is the result of the program when run on the unit square with the $x = 0$ boundary condition set to $2(y - 0.5)$ when $y > 0.5$, the $y = 1$ boundary condition set to $1 - 2x$ when $x < 0.5$, all other side boundaries are fixed at 0, and a special boundary condition setting the point at $(0.8, 0.2)$ to 3. The entire domain is initialized to 0. The initial quadtree is refined 3 levels (64 leaf quadrants), a Jacobi tolerance of 0.0001 with a maximum number of iterations of 2000, and a refinement threshold of 0.05 and a maximum number of refinement levels of 3.



Sequential Program Design and Implementation

The structure of the program looks something like this:

```

create quadtree to initial refinement level
set initial solution values of leaf quadrants to initial conditions
do {
  do {
    foreach leaf quadrant {
      if (!special boundary condition applied)
        do Jacobi iteration step
    }
  } while (error > Jacobi tolerance && jac_steps < jac_max)
  refine quadrants based on refinement threshold
} while (any quadrants were refined && ref_steps < ref_max)

```


Even though our program is written in C instead of an object-oriented language like C++ or Java, we will follow good object-oriented design. One way to do this is to separate the data structures and functionality of the quadtree from the solution process.

We'll start by looking at the main program, then look at quadtree-related structures and functions as we encounter them.

Some things to notice at this point:

- The list of local variables in `main` is very short - just the pointer to the root of our quadtree and several solution parameters and miscellaneous variables like loop indices and file pointers.
- The program requires 6 command-line arguments. In addition to the solution parameters mentioned above, we have an “output level” that we can use to specify how frequently our program will print solution data.
- We create the initial quadtree structure: the parameters are the bounding box (top, left, bottom, right), the initial solution value, and the parent point (NULL for the root).
- Next, we have to do our initial global refinement.
 - We do this with a *visitor function* that will in turn call a *callback function* with each leaf quadrant as a parameter.
 - The function `visit_all_leaf_quadrants` is, as you might suspect, a recursive function. It takes as parameters the quadrant whose leaves are to be visited, a pointer to the function to call on each leaf, and a pointer to some caller-specified data that will also be passed along to the callback function. The function's two options are:
 - * if we are already at a leaf, call the callback function
 - * if we are an interior quadrant, make a recursive call to the visitor function on each child quadrant.
 - In this case, we use a callback function `do_refine`, which performs one level of refinement on each leaf quadrant.
 - Our `do_refine` function just calls a quadtree function `refine_leaf_quadrant`.
 - The `refine_leaf_quadrant` function:
 - * Checks to make sure the quadrant is in fact a leaf. Note the use of the function `is_leaf_quadrant` to check and the use of the macro `ASSERT` to terminate our program with an error condition if this is not a leaf.
 - * Creates 4 new leaf quadrants, each $\frac{1}{4}$ the size of the original leaf, and with the former leaf as their parent.
 - This happens as many times as we specified for the initial refinement level (`init_ref`), and results in a uniform quadtree with 4^{init_ref} leaves.
- Next, we need to set our initial conditions based on the function provided.

- We again use our leaf quadrant visitor (leaves as the only quadrants that have a solution value in our implementation) this time with `set_init_cond` as the callback function.
 - The `set_init_cond` function calls our initial condition function `initial_cond`, with the coordinates of the centroid of the leaf to get the appropriate initial condition value for this leaf, then calls `quadrant_set_value` to assign it to the leaf.
- For the moment, we’ll ignore the solution printing and look at the solution process.
 - Our main loop is a nested `do/while`. The outer loop guides solutions on different refinements of the tree. The inner loop is the Jacobi iteration corresponding to what we had done previously. We’ll consider that loop to start.
 - The solution step is done with another visitor/callback.
 - * Like you did in your implementations earlier, we always do two iterations in succession: one to compute a second solution from the current, then another to compute current again based on the second.
The callback `do_jacobi_iter_phase1` computes what we call `previous` from `value`, and `do_jacobi_iter_phase1` computes `value` from `previous`.
 - * Each of these follows the same general procedure:
 - Check to see if there is a special boundary condition that applies to this leaf. This is done with the `apply_other_bc` function. If a condition was applied, the function returns `true` and we’re done with this leaf.
 - Otherwise, we need to find our 4 neighbor values that we’ll be averaging to get our new value. This is not as simple as changing array indicies by one like we did when the computation was being done on a simple 2D array.
 - Neighbor-finding in a quadtree involves a simple search. In our quadtree code, it is done with the function `neighbor_quadrant`, which takes any quadrant and a direction and returns either a neighboring quadrant in the desired direction, or `NULL` if there is no neighbor (i.e., we are on a boundary).
It is not even always clear what we mean by the “north neighbor” There is potentially a whole hierarchy of quadrants neighboring us in a given direction. What we need for our solution is the neighbor *at our own level* in the tree hierarchy. If our neighbor is not refined as far as us, we’ll go ahead and use the leaf at a higher level. If it is refined more, we want to use the leaf quadrants adjacent to us, but only count them once, even if more than one is adjacent. See the recursive functions `quadrant_side_value` and `quadrant_side_previous` for the details.
We find the appropriate neighbor by finding a point in space that we know will be inside our neighbor in a given direction, then searching for that point in the tree.
Something to think about: we could start this search at the root and work down, but can we do better by searching up the tree to find our nearest ancestor that contains the desired point, then back down to the appropriate leaf? This is faster when our neighbors are most likely our siblings or cousins. Only

in those cases where our neighbor is in a different level 1 quadrant will we need to search all the way back to the root. This is the search we use in the `neighbor_quadrant` function.

- Any neighbor search that returns `NULL` indicates that we've gone off the edge of the universe and should use a boundary condition instead. We do this with the `bc_*` functions.
 - With the 4 neighbor or boundary values to average ready, we compute the new value and store it in the quadrant.
 - In the phase 2 computation, we also check the error value, where the maximum encountered so far being passed in as `maxerr`. Note that this uses the extra callback parameter to pass the pointer to `main`'s local variable `max_jac_diff` to the callbacks. At the end, we'll have the maximum error available in `max_jac_diff`.
- After the 2 iterations, we check to see if we've reached our error tolerance or the maximum number of iterations.

We've ignored an important part here so far. The adaptivity. That's the outer loop.

- Recall that we want to find places where adjacent quadrants have solution values whose difference exceeds a given tolerance. When we find such situations, we want to refine the tree in that area to have a more accurate solution. That's what happening in the outer loop.
- The function `calc_error_and_refine` is the heart of the refinement functionality. It determines where refinement is needed, performs that refinement, then returns the number of refinements that were performed. (If there are 0, we can stop processing, since we'd only recompute the same solution we just computed on the previous grid.)
- The implementation of `calc_error_and_refine` again makes good use of our visitor function in each phase of our refinement procedure:
 - Marking quadrants for refinement: the `check_if_refinement_needed` function. We locate each of our neighbors, then see if any of them have solution values too far from our own. If one is found, we mark ourself for refinement.

Marking for refinement involves a little trick in the quadtree data structure. Since we don't have children when we're a leaf, we use a non-zero value in the third child pointer to indicate the refinement mark.
 - Refining marked quadrants: the `refine_if_marked` function. If the leaf being visited is marked for refinement, we refine it. Notice that the new quadrants inherit the solution value from their parent. This means we use the last solution on one grid as the initial solution on the next.

Finally, a bit about printing the solution. As you learned even in the non-adaptive versions you wrote, it is nearly impossible to understand the solution based only on printing out values from

the grid. This problem becomes even more difficult once we include adaptivity. Our approach is to print solutions to a file in a format that includes the coordinates of each leaf quadrant and its solution value. These values may then be plotted with gnuplot. A script that I used to generate some of the solution plots in these notes is available as `solutionplot.gp`. This is a very simple gnuplot script and I am sure you can do better, but it gets the job done.

This version of the solver also has the ability to write solution files more frequently, and the scripts in the `video` directory can be used to visualize these solutions and paste them together into a solution animation.

Parallelization: Shared Memory

Our ultimate goal is to parallelize this program using either pure message passing (MPI only) or a hybrid of message passing and shared memory (MPI processes, one per node, then a number of threads per node that can take advantage of multiple cores).

There are many choices to make when parallelizing a program such as this. Which parts of the computation are to be performed by each processor? What is the granularity of the units of work to be divided among the processors? What information must be maintained and what communication is needed to support the computation once the work has been divided? Will the workloads need to be adjusted following adaptive steps? How can such a rebalancing be performed?

We can delay some of these concerns at the start by using shared memory first. We do need to think about load balancing but it doesn't matter which threads do which work, since the whole tree and all solution values in it are shared.

The program in the `pthread` directory of the repository is a version parallelized with pthreads.

First, the program needs to know the number of threads. Instead of adding a command-line parameter, this version uses an environment variable `NUM_THREADS`, whose value is retrieved from the environment by a call to `getenv`.

It makes sense here to keep it simple: use a traversal of the quadtree, where each thread gets about the same number of leaf quadrants to compute, and each leaf's new value is computed by exactly one thread.

Much of the work can be done through a pthreads-aware version of the `visit_all_leaf_quadrants` function: `visit_leaf_quadrant_range`.

Each time the quadtree is modified, a new range of leaf quadrants needs to be computed for each thread. This range is then passed to the `visit_leaf_quadrant_range` function.

Study this function to see how it decides which parts of the tree to traverse on each thread.

Note, however, that the refinement step is complicated by the fact that it *modifies* the tree during the traversal. In the serial case, this was just not a problem. We only refine leaves, and when we are done visiting a leaf, whether it got refined or not, the traversal can continue. See the `calc_error_and_refine` function for details on the approach taken.

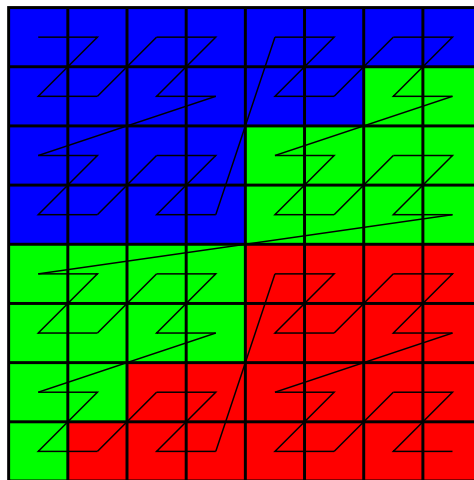
Parallelization: Message Passing

Parallelization of the adaptive quadtree computation using message passing is much more complex.

The following is a possible phased development process to a message passing parallelization. Many details would need to be worked out, but this approach should be feasible.

Phase 1 Parallel computation on a uniform quadtree.

- Build the initial quadtree to the requested refinement level, replicated it completely on each process. This will be referred to as the *global quadtree*.
- Assign a unique owner process to each leaf quadrant, essentially partitioning the quadtree. Each quadrant's solution is computed only by its owner, but quadrants along partition boundaries will need to exchange solution values between iterations.
- There are many possible approaches to this partitioning problem, but it should be fairly straightforward to divide the quadrants into partitions based on a tree traversal. If you have n leaf quadrants and p processes, assign the first $\frac{n}{p}$ to the first partition (process 0), the next $\frac{n}{p}$ to the next partition (process 1), and so on, handling remainder quadrants in some reasonable way. If child quadrants are ordered NW-NE-SW-SE, a partitioning of a base quadtree refined to three levels might be done as follows:



- Message passing is needed to send solution information from owned quadrants on partition boundaries to those other processes (and only to those other processes) that will need the information during the solution process.
- There is no adaptivity for this phase, so the working program should compute only on the initial quadtree.
- Solution output procedures would be needed to compare the solution from the parallel version with those from the sequential version.

Phase 2 Introducing adaptivity.

- Adaptivity would work much as it did in the sequential program, except that when a quadrant is refined, it would only be refined on the process that owns that quadrant (or its ancestor in the global quadtree). Once any quadrant from the global quadtree is refined on some process, the quadtree is no longer fully represented on all other processes.
- Working with the locally-refined quadtree is straightforward on the interiors of partitions, but near partition boundaries, complications will arise.
 - when computing solution values, neighbors can be owned by off-process copies of the global quadtree that may have been refined
 - when checking the difference in solution values for adaptivity control, off-process neighbors might already be refined to different levels.
- possible approaches here include
 - message-passing for the solution-value exchange could include the entire refinement structure for neighbor quadrants from the global quadtree
 - at least some partial refinements could be conducted of the off-process copies of global quadrants that need to be involved in interprocess communication.

Phase 3 Dynamic load balancing.

- Adaptivity will likely introduce a load imbalance. If all or most of the refinement takes place in just a few processes, those processes will have a larger workload during each Jacobi iteration, causing other processes to wait before the boundary exchange.
- After each refinement phase, a rebalancing phase should be conducted, where the partitions of the global quadtree structure are adjusted (and refined parts of the tree *migrated* appropriately) to ensure that each process has approximately the same number of owned leaf quadrants.
- To keep this relatively straightforward, the units of work that are allowed to be migrated can be restricted to be the leaves of the global quadtree.
- A disadvantage of this is that the granularity of the “work objects” that you are partitioning can get large after several refinement steps have occurred, meaning a perfect load balance may not be possible.
- Once adaptivity has been performed, the computational costs of each global quadtree leaf will be different, and it is those costs that should be balanced across the processes, not the number of global quadtree leaf quadrants on each process.
- Message passing to share those costs is straightforward (which of our communication patterns would that be?), and allows each process to compute this new decomposition independently and to determine which parts of its tree need to be sent elsewhere. The actual migration of the refined quadtree structures will introduce some complexity.

An actual implementation is beyond what can reasonably be completed in a short-term project for this course.

Partitioning and Dynamic Load Balancing

We have considered partitioning and dynamic load balancing in some specific situations. Let's now think about it in more general circumstances.

Our assumption here is that we have a computation whose memory and computational requirements are dominated by some set of objects that we distribute among a set of cooperating processors. We will most often think of this as a mesh being used to solve a PDE, but other structures are possible.

Typically, one process is assigned to each processor. Data are distributed among the processes, and each process computes the solution on its local data (its *subdomain*). Inter-process communication provides data that are needed by a process but "owned" by a different process. This model introduces complications including

1. assigning data to subdomains (i.e., *partitioning*, or when the data is already distributed, *dynamic load balancing*)
2. constructing and maintaining distributed data structures that allow for efficient data migration and access to data assigned to other processes, and
3. communicating the data as needed during the solution process.

More details on partitioning and dynamic load balancing have been moved to a separate slide deck.