

Parallel Sudoku Solver

Members:

Kevin Zhu (kzhu2)

Hao Zhang (haozhan2)

URL

<https://github.com/Hao2747/sudoku-parallel-solver.github.io>

SUMMARY

We plan on parallelizing a sudoku solver using multiple algorithms and then comparing overall runtime performance for each approach. Approaches include breadth-first search, distributed work queue, work-stealing distributed queue (hopefully), and a constraints-solver approach (hopefully), and if time permits, a DFS approach. We plan to be coding in OpenMP and if time permits, we also will try an approach in MPI for large sudoku boards.

BACKGROUND

Sudoku varies on the compute intensity based on the approach taken to solve the board. At a classical level, it is a graph based problem, where each node is a possible sudoku board and each connecting vertex indicates which square was filled to move from one possible board to the next. Proper sudoku boards have one solution.

Initially, we need to find the possible numbers for each empty square. This can be done in parallel because we need to traverse the entire board to find the missing numbers in each row, column, and box for each square, which can be done in parallel. In the next step, the cpus can choose which empty square they want to work on.

They can all choose to work on one square or try a divide and conquer approach to work on different squares, which is more 'trial-and-error' focused. Another approach is Constraint Propagation (eliminate the choices of possibility of each empty square using humanic heuristic). The sudoku wiki has a series of these heuristic strategies ranging from simple to extreme strategies to solve the board. Each of these strategies has varying amounts of parallelism. While it may take extra work to check for these strategies, they may also greatly reduce and prune the search space.

THE CHALLENGE

Parallelizing in DFS fashion is inherently challenging. The naive method for solving Sudoku involves using backtracking to solve one square at a time and recursively, but it's not so intuitive how to parallelize this due to the directed nature of the graph and the relational dependencies between them. Trying to parallelize over the nodes is hard because each node is dependent on the previous node. This doesn't mean that it is impossible to do (check references for work done on GPUs, and we will try to make an attempt on this).

As a result, we plan on starting off with a BFS approach. Parallelism is not exactly as intuitive as in Lab 3 with Nbody. Blindly spawning tasks for every empty square leads to work redundancy. For example, in a trial-error approach, where each CPU gets its own square, as it goes down the tree, it will keep filling other empty squares that might overlap with what other cores are working on. This would result in a source of redundancy (e.g. if a square cannot be value 9, two processes might both reach the same square in their tree and try 9 if they do not communicate). Additionally, if one square is solved, should it notify all other working threads to update their board.

To implement these algorithms, we plan to use a global work queue and then progress to a distributed work queue, measuring each approach against one another. How to structure the task queue can also be a challenging task since resolving one task (completing one empty square) can sometimes resolve another task (complete another square). Therefore, careful design needs to be taken in order to avoid task redundancy.

From our project, we hope to have a more concrete understanding of the parallel algorithms, possible communication overheads, and see the realization of the benefits or drawbacks of a distributed work queue structure.

As for workload, because of a tree-like structure, the child is dependent on the parent. For parallelization, as described, the children need to communicate or understand what is not being done. Generally, the memory access can be quite random within the grid and there is a lot of divergent execution. Though a DFS approach will generally have better locality than BFS approaches. Because the Sudoku problem could technically be represented as a graph, we also have to make sure that there is no workload imbalance as some parts of the tree might be deeper than other sections.

What makes this challenging to implement as a work queue is that in OpenMp there is a performance penalty in using critical and having high contention regions. We need to structure and assign tasks in a way that can circumvent this issue.

RESOURCES

We plan on using both the GHC and PSC machines, and may get benchmarks from either machine. The PSC machine will be important for the MPI approach(if we have the time to do it) as well as testing the scaling of our algorithm. As of now, we plan to code the entire thing ourselves starting from scratch.

We are planning to use existing sudoku datasets as our test cases listed here:

<https://github.com/t-dillon/tdoku/blob/master/benchmarks/README.md>

<https://www.sudokuwiki.org/sudoku.htm> (and manually implement some commonly tested puzzles, also possible from different websites)

Paper Reference of approaches to parallel sudoku solver:

<http://individual.utoronto.ca/rafatrashid/Projects/2012/SudokuReport.pdf>

Paper Reference of Parallel DFS:

<https://research.nvidia.com/sites/default/files/publications/nvr-2017-001.pdf>

GOALS AND DELIVERABLES

PLAN TO ACHIEVE

- Sequential/Parallel Breadth-first search solver
- Implemented distributed task queue strategy
- Sequential/Parallel CP (Constraint Problem) and integrate some modern Sudoku techniques

HOPE TO ACHIEVE

- Sequential/ Parallel Depth-first search solver
- Sequential and Parallel extension for MPI for large workloads
- High performance for MPI workload

In our project, we plan to have separate sequential and parallel implementations of various algorithms. We will also be computing the overall speedup for the workload in a separate measurement as well. This way, we think the speedup comparison is more fair for different workloads and show the true achieved speedup for each task.

For example, a depth and breadth first search may not remotely have the same runtime because it may be the farthest depth (favoring depth first) or the 2nd node (favoring breadth).

We also plan to time and measure different sections of the code. For example, the initial phase of finding the possible digits may be parallelized. We will try to see our MPI performance for that section if time permits. For scaling, we plan to use the PSC machines to test the parallelization on large numbers of cores.

We plan to analyze the performance differences and characteristics between each of the parallel models we create and find out where the bottlenecks are in each approach (ex was it communication or work redundancy) and try to further tune our algorithm. For a given workload, we should be able to see what strategies made a positive difference in solving problems.

PLATFORM CHOICE

We have been coding in C++ for MPI and OpenMP tasks and we plan to continue to use them for our project. The Gates and PSC machines, along with the OpenMP and MPI versions, are resources that we have already been using and are familiar with.

Sudoku boards are relatively small so using a shared memory approach makes sense. The naive intuition for creating a sequential algorithm involves backtracking which is recursive by nature. As can be seen in Lab 3, recursive algorithms with little computational complexity are better done on OpenMP as the communication overhead is not worth it.

If the grid is large or if there are many critical sections in the OpenMP code, then using MPI is probably better suited.

SCHEDULE

	Personal Deadlines (End of Week Goals)
Week of March 27	Finish Writing Proposal Initial Research and Outline Done (What is possible, what we want to accomplish, and goals)
Week of April 3	Any additional Research and Consultation necessary A finished sequential version with a bootstrapped code base and data structures <ul style="list-style-type: none">- Agreed Code Style- Ways to import / Validate Designs

	<ul style="list-style-type: none"> - Few Sample Test Cases - Backtracing/Recursive Implementation
Week of April 10	<ul style="list-style-type: none"> - Adding More Test Cases - Completing and finishing an initial openMP implementation and gauging/documenting performance <p>We may also want to start the work queue version and DFS version at this time.</p>
Week of April 17	<p>Finish work queue version and continue to work on DFS version</p> <p>Start adding strategic version ideas in parallel with fall back to regular work Queue Version</p> <p>Start Distributed Queue Version</p>
Milestone April 19th	<p>Finish Initial Code Base + validation</p> <p>Have performance checker</p> <p>Have quite a few testcases ,but a system that can easily be expanded on</p> <p>Finish all sequential Versions</p> <p>Finished Naive OpenMP version (BFS)</p> <p>Working on/nearing completion of WorkQueue Version</p>
Week of April 24	<p>Hopefully complete DFS version</p> <p>Finish Distributed Queue Version</p> <p>If time permits, MPI version</p>
Week of May 1	<p>Final report, summary and presentation.</p> <p>Wrap up any loose ends</p>