

0.1 Project Specifications

The group project for the course Programming Massively Parallel Hardware is about parallelizing a serial implementation of `tridag`, contained in the files `ProjCoreOrig.cpp`, `ProjectMain.cpp`, `ProjHelperFun.cpp`, and `ProjHelperFun.h`. The project consists of three parts: converting to OpenMP parallel, converting to CUDA parallel, and results and comparisons. We will be comparing the CUDA and OpenMP implementation with the original code.

0.2 OpenMP

0.2.1 The Idea

Converting serial loops into parallel loops is as simple as ensuring that there are no loop dependencies. We must check that for every loop we wish to parallelize, that no iteration depends on the results of any other iteration, and that variables which belong to multiple iterations have either the correct scope, so that information does not pass from one iteration to the next, or a separate index for each iteration. After that, we denote the loop ready by adding the appropriate `#pragma ...` command, and compiling with the correct options.

0.2.2 The Work

Converting the original code to OpenMP was a straightforward matter, involving only modifying code in `ProjCoreOrig.cpp`. After ensuring that none of the inner loops would create race conditions, we moved the `REAL strike` and `PrivGlobs globs` declarations inside the main loop, and then simply added the line

```
#pragma omp parallel for default(shared) schedule(static) if(outer>8)
```

right before the same loop declaration. We decided to add the static scheduling line for more than 8 threads, since that was what we were working with. After that, all that was left was to compile with `g++` using the command line option `-fopenmp`. The results are displayed in the table on page 6.

0.2.3 Correctness

The changes we made were very minor, involving moving 2 variable declarations inwards in scope, and adding the `#pragma ...` line. The variable declarations made no difference, since every iteration of the loop resulted in both `strike` and `globs` being immediately assigned new values, but moving them allowed for us to parallelize without worrying if the results of one iteration spilled over into another. The `#pragma ...` declaration signaled that this loop was able to be run in parallel. We can assure that it was by looking through the code. From the outermost loop, there were no loop dependencies from one iteration to the next: the results were saved in separate `res[]` indexes, and no result from a previous iteration was used in a later iteration.

0.3 CUDA

0.3.1 Preparations

The code as it stands can be naïvely converted to CUDA, by replacing the inner loops with the exact same function inside a CUDA kernel. However, we do not intend to merely convert the OpenMP version to CUDA. Instead, we are going to implement an optimized version of TRIDAG

0.3.2 Loop Distribution and Array Expansion

The Idea

The first step was to conglomerate every secondary function, which allowed us to hoist the initialization of the globs (global variables) into glob arrays. That is, instead of initializing a temp variable or array for every iteration of a loop, we instead initialize an array one dimension larger, with size equal to the number of iterations of the loop, before the looping code.

```
for (int i = 0; i<max; i++){
float tmpA = 0.0;
  for (int j = 0; j<max2; j++){
    tmpA += 2*B[j];
    ...
  }
  ...
}
```

Figure 1: A code snippet with tmpA initialized for every iteration.

```
float tmpA[max] = {0.0, ...};
for (int i = 0; i<max; i++){
  for (int j = 0; j<max2; j++){
    tmpA[i] += 2*B[j];
    ...
  }
  ...
}
```

Figure 2: The same code with tmpA hoisted.

The purpose of this is twofold. First, this allows us to serialize a computation which would otherwise be repeated by every thread. By computing it beforehand, the threads can be spared this extra work. Second, this allows us to easily parallelize the inner loops. Since each iteration of the inner j loop

requires access to a single `tmpA` per iteration of the `i` loop, we would have to compute it, then pass it on as a variable to each of the threads. In the hoisted version, `tmpA[]` is copied to device memory, so that the inner loop can just access the appropriate version without much trouble.

The Work

The first step for us was to begin preparing the CPU for kernel parallelization. Before distributing the outer loops in `ProjCoreOrig.cpp`, we began by moving all of the secondary functions (`void updateParams`, `void setPayoff`, `REAL value`, and `void rollback`) into `void run_OrigCPU`. This allowed us to easily see the globs and their relations to one another.

The work at this step is temporary. The code is objectively de-optimized, since more time is spent on initializing arrays, memory usage is larger, and the code runs slower. The purpose of this is to prepare the code for optimal loop distribution.

The `REAL strike` variable has been removed, being placed instead inside of one of the kernels. The variables which have been hoisted at this step are listed in figure 3.

```
void run_OrigCPU(...)
{
    ...
    // Generate vector of globs. Initialize grid and operators onces
    // and make default element of vector
    // Hoisted from "value"
    PrivGlobs globs(numX, numY, numT);
    initGrid(s0, alpha, nu, t, numX, numY, numT, globs);
    initOperator(globs.myX, globs.myDxx);
    initOperator(globs.myY, globs.myDyy);
    vector<PrivGlobs> globArr (outer, globs);
    ...
    //Rollback globs
    vector<vector<vector<REAL> > > u(outer, vector<vector<REAL> >(numY,
        vector<REAL>(numX))); // [outer] [numY] [numX]
    vector<vector<vector<REAL> > > v(outer, vector<vector<REAL> >(numX,
        vector<REAL>(numY))); // [outer] [numX] [numY]
    vector<vector<REAL> > a(outer, vector<REAL>(numZ)),
        b(outer, vector<REAL>(numZ)), c(outer, vector<REAL>(numZ)),
        y(outer, vector<REAL>(numZ)); // [outer] [max(numX, numY)]
    vector<vector<REAL> > yy(outer, vector<REAL>(numZ)); // temporary
        used in tridag // [outer] [max(numX, numY)]
    ...
}
```

Figure 3: Hoisted variables in `ProjCoreOrig.cpp`.

The other necessary modifications were to simply update the relevant vari-

able references, for example, changing `globs.myResult[j][k]` to `globArr[i].myResult[j][k]`.

Correctness

The reason we are allowed to do this is because we are not fundamentally changing anything about the flow of the program. Moving all of the functions together does nothing to program flow, only impeding readability slightly. For the hoisted variables, the extra dimension can be easily compared to a new variable per iteration, and since the variables are moving outwards in scope, nothing vital is changed. There is no danger of loop dependencies, since each iteration still uses their own indexes for these hoisted variables.

After this work, the program is in a state to distribute the various loops. There is a mild amount of slowdown, since we allocate more memory to some of the variables, and we perform some initial calculations which otherwise performed per loop iteration. This is not a good stopping point, but it is necessary to continue.

0.3.3 Kernel Replacement

The Idea

After loop distribution, we began to convert the loops to CUDA kernels. These kernels are copied almost directly from the already existing CPU code. This allowed us to begin using the GPU, and CUDA.

The Work

For each of the distributed loops, we naïvely translated the code into a CUDA kernel. That is, we attempted to rewrite the functions so that they performed precisely the same task on the GPU as they did on the CPU. Figures 4 and 5, show the conversion of the function `void setPayoff` from `ProjCoreOrig.cpp` into the `__global__ void setPayoffKernel` in the file `ProjKernels.cu.h`. All distributed loops and matching functions were converted.

```
void setPayoff(const REAL strike, PrivGlobs& globs )
{
    for(unsigned i=0;i<globs.myX.size();++i)
    {
        REAL payoff = max(globs.myX[i]-strike, (REAL)0.0);
        for(unsigned j=0;j<globs.myY.size();++j)
            globs.myResult[i][j] = payoff;
    }
}
```

Figure 4: The distributed loop implementation of `void setPayoff`.

```

template<const unsigned T>
__global__ void setPayoffKernel(
    const unsigned outer,
    const unsigned numX,
    const unsigned numY,
    REAL* myX,
    REAL* myResult
)
{
    int i = blockIdx.x*T + threadIdx.x; // outer
    int j = blockIdx.y*T + threadIdx.y; // myX.size
    int k = blockIdx.z*T + threadIdx.z; // myY.size
    if (i < outer && j < numX && k < numY) {
        myResult[i * numX*numY + j * numY + k] = max(myX[i * numX +
            j]-0.001*i, (REAL)0.0);
    }
}

```

Figure 5: The equivalent CUDA kernel.

The next step was to load the appropriate variables to and from memory. figure 6 shows some of the variables that are now on device memory, to be accessed by the kernels. This is analogous to the `globs` structures from the sequential code.

```

// Arrays for rollback:
REAL* d_a, *d_b, *d_c, *d_y, *d_yy; // [outer][max(numX,numY)]
REAL *d_v, *d_u;
cudaMalloc((void**) &d_a, sizeof(REAL)*outer*numX*numY);
cudaMalloc((void**) &d_b, sizeof(REAL)*outer*numX*numY);
cudaMalloc((void**) &d_c, sizeof(REAL)*outer*numX*numY);
cudaMalloc((void**) &d_y, sizeof(REAL)*outer*numX*numY);
cudaMalloc((void**) &d_yy, sizeof(REAL)*outer*numX*numY);
cudaMalloc((void**) &d_u, sizeof(REAL)*outer*numX*numY);
cudaMalloc((void**) &d_v, sizeof(REAL)*outer*numY*numX);
REAL * timeline = (REAL*) malloc(sizeof(REAL)*numT);
cudaMemcpy(timeline,d_globs.myTimeline,sizeof(REAL)*numT,cudaMemcpyDeviceToHost);

```

Figure 6: Initializing device memory.

Correctness

It is no longer as easy to argue that this is a simple replacement. Transforming the distributed loops into kernels requires that the syntax has to change, and the semantics as well. The example in figure 5 shows that the loop iteration

has been replaced by `threadIdx`; at a fundamental level, we are moving from sequential loops to parallel threads. The equations are the same, with a slightly more complicated index (`threadIdx` instead of a simple `i`), and the variables are now all loaded into GPU memory.

There are two arguments for correctness. The first is that although the semantics and syntax have both changed, from distributed loops to threads, the process is a one-to-one mapping. Every kernel can be directly traced back to a loop, and every loop gave rise to a single kernel. As long as the input and output of the kernels and loops is the same, we have the same program.

The second is that we continuously validated the GPU kernels during the programming by copying intermediate CPU results to the GPU, calculating the GPU kernel results, and then comparing them to the CPU results. That is, each step was substituted piecewise to ensure that they functioned the same.

Although the program has now been distributed to the GPU, the runtime of the entire program took a hit. The problem was, that the first CUDA call of the program also loads the entire CUDA runtime library, which gives a major overhead compared to the actual running time on our datasets. Timing early versions of our code was highly dominated by the CUDA library load time. To combat this, we have added a `cudaFree(0);` line in the `ProjectMain.cu` before we start timing.

0.3.4 TRIDAG

The Idea

The Work

Correctness

0.4 Results

0.4.1 Comparisons

Figure 7: Results of the three different implementations.

Conclusion