

Performance Evaluation and Enhancements of a Flood Simulator Application for Heterogeneous
HPC Environments

Ryan Marshall

Department of Computer Science
Tennessee Technological University
Cookeville, TN, 38505, United States

Sheikh Ghafoor

Department of Computer Science
Tennessee Technological University
Cookeville, TN, 38505, United States

Mike Rogers

Department of Computer Science
Tennessee Technological University
Cookeville, TN, 38505, United States

Alfred Kalyanapu

Department of Civil and Environmental Engineering
Tennessee Technological University
Cookeville, TN, 38505, United States

and

Tigstu T. Dullo

Department of Civil and Environmental Engineering
Tennessee Technological University
Cookeville, TN, 38505, United States

Received: February 1, 2018

Revised: May 6, 2018

Accepted: June 4, 2018

Communicated by Susumu Matsumae

Abstract

This paper presents a practical implementation of a 2D flood simulation model using hybrid distributed-parallel technologies including MPI, OpenMP, CUDA, and evaluations of its performance under various configurations that utilize these technologies. The main objective of this research work was to improve the computational performance of the flood simulation in a hybrid architecture. Modern desktops and small cluster systems owned by domain researchers are able to perform these simulations efficiently due to multicore and GPU computing devices, but lack the expertise needed to fully utilize software libraries designed to take advantage of the latest hardware. By leveraging knowledge of our experimentation environment, we were able to incorporate MPI and multiple GPU devices to improve performance over a single-process OpenMP version up to 18x, depending on the size of the input data. We discuss some observations that have significant effects on overall performance, including process-to-device mapping,

communication strategies and data partitioning, and present some experimental results. The limitations of this work are discussed, and we propose some ideas to relieve or overcome such limitations in future work.

1 Introduction

Simulations are invaluable for forecasting and warning, however, simulations such as flood inundation are often computationally very expensive and takes long time even when using powerful computers. For flood mitigation purposes, civil authorities may decide to break a dam to save a very densely populated city. In this case, a flood simulator can be used to evaluate the damage the dam break would cause in the areas downstream. An accurate and efficient flood simulator running on a high performance computing (HPC) system can be used support the decision making process, especially if the situation is time-sensitive.

For more than a half of a century, computers have been used simulate floods [9]. These implementations have, traditionally, applied a one-dimensional (1D) dynamic wave approach. Examples include the Hydrologic Engineering Center River Analysis System (HEC-RAS) and MIKE 11 [27, 16]. The 1D approach is relatively simple and thus easy to implement, but they cannot simulate the lateral diffusion and introduce inaccuracies due to cross-section discretization and, therefore, can give inaccurate results, especially in urban flood-prone areas [3, 25]. 2D models [10, 21, 23] are superior to 1D models because they use higher-order topographic representation and preferential flood pathways in the simulations [3, 4, 8, 14]. In fact, the National Research Council in 2009 recommended that the Federal Emergency Management Agency (FEMA) promote using 2D hydraulic models for floodplain topography [17].

A disadvantage to 2D models is that compared to 1D models, they are computationally complex [10, 23]. Practitioners may find that implementing their 2D model that adheres to their specific requirements, such as accuracy constraints, data input requirements, and data output requirements, is difficult. Inefficient solutions can affect computational time, which can violate the overall time constraints defined by the practitioners. For example, the average modeling time should be less than an hour for flood risk and warning applications [17]. Furthermore, the necessary computational time of a simulation increases with the scale and resolution of the problem domain and complexity of the numerical model. Solving such problems and adhering to the practitioners time constraints require using parallel-computation platforms. Because this problem is well-suited for throughput processing, much of the computational work can be offloaded to a general purpose graphics processing unit (GPGPU) device. Future platforms for simulation, forecasting and warning systems are trending towards a higher degree of architectural diversity, including the use of pervasive grids and computation over mobile mesh networks [5, 24, 11]; adaptation of 2D flood simulators to support heterogeneity is a logical development choice for long-term projects. An efficient implementation of 2D models on heterogeneous platforms require in-depth knowledge of hardware, software, libraries, and their complex interactions. Practitioners may not have the in-depth knowledge of parallel platforms that is required, and spending precious time learning such complexities may result in spending less time designing more accurate simulations. The main contributions of this work are as follows:

- The improvement in efficiency of two existing 2D flood simulator applications (OpenMP-based and CUDA-based), achieved primarily by a redesign of the data handling operations to support higher concurrency, leaving the underlying numerical model largely untouched.
- A use case of a 2D iterative stencil-based application that supports hybrid distributed parallel technologies, including MPI+OpenMP and MPI+CUDA, with further modifications to support greater combinations of compute devices in HPC environments.
- Experimental results that focus on the performance effects of several test cases under different combinations of compute devices and different process-to-device mappings. Through experimentation, we identify several characteristics of the software organization that can be developed into a generalized solution for iterative stencil-based computations over n -dimensional gridded datasets.

The rest of this paper is organized as follows. Section 2 describes related background research. Section 3 describes the implementation of the 2D model. Sections 4 and 5 describe the experimental setup and present the experimental results. Section 6 describes the optimizations made to the simulation. Finally, Section 7 concludes the paper and describes future work.

2 Background

A variety of related research exists on the computational performance improvements of the 2D model via parallelization, the 2D numerical model and the 2D data model. In general, related works are narrow in scope and geared toward homogeneous environments.

2.1 Parallelizing the 2D Model

Recent research efforts have improved the computational performance of 2D models by implementing thread-based parallelization or incorporating GPUs. In [17], a desktop with an Intel Core 2 Duo processor required over 9 hours to compute 15 minutes of simulation time over a 624 x 1136 grid. Although the simulation application was redesigned to use CUDA on a Tesla C1060 GPU, and reported a speedup by a factor of 88, this application was hand optimized and not very portable because it employed a somewhat unique combination of hardware devices and software requirements. In [10], the authors report a speedup factor of 116 when running a simulation of their JFLOW 2D diffusion wave flood model versus a serialized version of the same program. RiverFLO-2D, an OpenMP implementation of the 2D numerical model, was shown to have a speedup of, approximately, a factor of 4.5 [15]. The authors of [6] developed a simulation application using C++, OpenGL and CUDA, that achieved 98% GPU utilization, and completed a 66 minute dam break simulation of the Malpasset Dam in 27 seconds. The application can produce a visualization of the simulation in near real time. An alternative to solving the underlying numerical model used for many flood simulations is given in [7]. The authors presented a method to improve surface effects over dry terrain and increased the accuracy of the underlying kernel function. Additionally, they show the SPH approach offers promising results for simulations that have complex boundary conditions.

Building on these previous works, we made modifications to a existing implementation of a flood simulation application with the following goals in mind:

- achieve a significant performance gain over the standalone OpenMP and CUDA versions,
- embrace heterogeneous and hybrid computing environments,
- provide a structural codebase that can be reused and built upon in future works.

The computation and memory access patterns for this application lend themselves to certain problem classes found in other domains, such as thermodynamics, biological systems and cellular automata. In general, these problems have a base layer containing static data, with one or more overlays of data that may change over time (specific data for our implementation is described in Section 2.3). The data is computed for the current time step based on adjacent cell values in one or more layers from the previous time step.

2.2 The 2D Numerical Model

The numerical algorithm used in this study is based on an existing version of the model, first described in [22]. The model uses a first-order accurate upwind finite difference scheme and solves the non-linear hyperbolic shallow water equations. These equations, also called the St. Venant Equations, are derived by integrating the horizontal momentum and continuity equations over depth. The equations are often referred to as depth-averaged or depth-integrated shallow water equations. The non-linear shallow water equations are presented as:

$$\frac{\delta h}{\delta t} + \frac{\delta(hu)}{\delta x} + \frac{\delta(hv)}{\delta y} = 0, \quad (1)$$

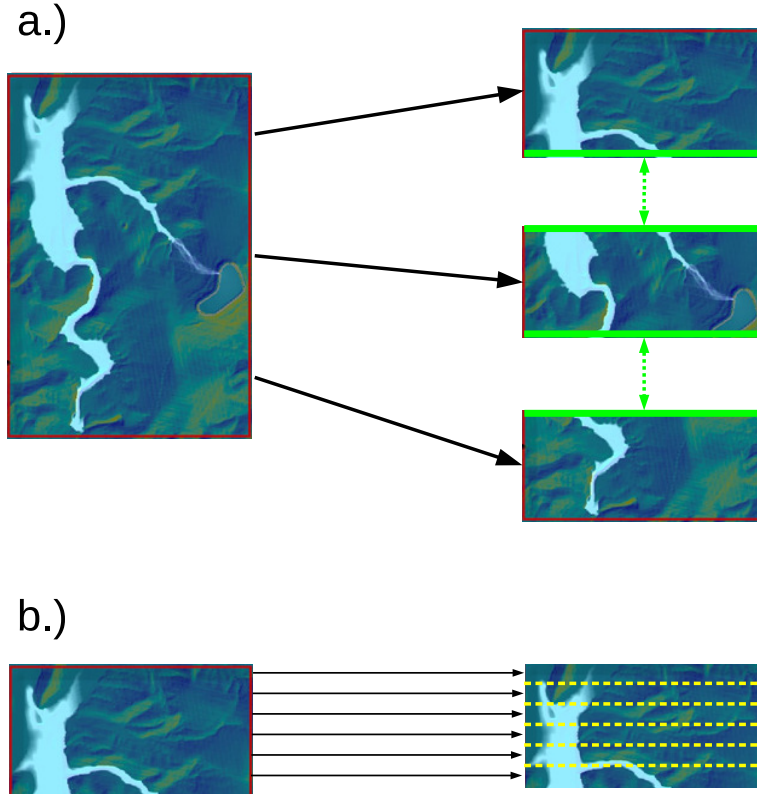


Figure 1: Partitioning of a domain by a.) MPI processes, where edge data communication occurs between process neighbors, and each process by b.) threads in shared memory space.

$$\frac{\delta u}{\delta t} + u \frac{\delta u}{\delta x} + v \frac{\delta u}{\delta y} + g \frac{\delta H}{\delta x} + g S_{fx} = 0, \quad (2)$$

$$\frac{\delta v}{\delta t} + u \frac{\delta v}{\delta x} + v \frac{\delta v}{\delta y} + g \frac{\delta H}{\delta y} + g S_{fy} = 0, \quad (3)$$

where the variable t is time in seconds, h is the water depth, H is the water surface elevation, u is the velocity in the x -direction, v is the velocity in the y -direction, t is the time, g is the acceleration due to gravity, S_{fx} is the friction slope in the x -direction, and S_{fy} is the friction slope in the y -direction. This model uses the upwind finite difference numerical scheme with the governing equations 1-3, which yields non-oscillatory solutions through numerical diffusion [12, 22]. The model also uses a staggered grid computational stencil to define the domain with the water depth (h) at the center of the cell and u and v velocities on the cell edges. The model propagates the simulation by constraining the time step using the Courant condition. Because these equations can be applied to every cell in the grid, the computation of a cell value depends on neighboring values and the state of the entire grid at time t depends on the state of the same grid at time $t - 1$, this model can be implemented as a cellular automata problem; we discuss the implementation details in Section 3.

2.3 Data Model

The input data requirements for the numerical algorithm include topography, surface roughness, flow hydrograph and its geographic source location (e.g., gauging station). The data model represents terrain with the Digital Elevation Model (DEM). DEM files are available in many different formats from different sources. The formats specified by USGS have been used in the present work [26].

The numerical algorithm computes the solution on a uniform grid, to take advantage of the use of downloadable DEM data from sources such as The National Map [2] or Google Maps API [1].

The commonly used Manning's n coefficient [20] models the surface roughness data. A single n value in 2D modeling applications [18] can represent the entire domain, but our model incorporates a spatial variation of Manning's n values using an appropriate raster dataset (e.g., National Land Cover Dataset from United States Geological Survey).

The flow hydrograph is an input dataset that can be developed from a hydrological model, dam break model, or empirical observations. We used hydrographs in our simulations that describe discharge from one or more source locations in cubic feet per second over simulation time in hours.

3 Implementation

To maintain consistency with the numerical model, we planned our implementation to create a set of 2D grids, H , U , V , to store (for each cell in the grid) depth, velocity in the x -direction, velocity in the y -direction, respectively. Additionally, we need a second copy of this set to store the values from the previous time step, so the full set includes H_a , U_a , V_a , H_b , U_b , V_b , and D for the terrain elevation data. To reduce the number of expensive copy operations, we employed a common swapping technique to swap the pointers between the a and b sets so the state data at the end of the current time step effectively becomes the data for the previous time step. The following limitations apply to our implementation:

- For most cases, we defer to the MPI runtime environment to determine process distribution over compute nodes. Because we use a heterogeneous system, there are some situations where we will need to manually map a process to a specific node in order to use a compute device not found on other nodes.
- We assume all data partitions are mutable- at time step t , each partition can potentially contain at least one cell whose value is not equal to the corresponding cell's value at time step $t - 1$.
- We use block row partitioning for MPI and block cyclic threading for OpenMP with the understanding that more efficient schemes possibly exist, but a full comparison and analysis is beyond the scope of this work.

3.1 Computational Model

Equations 1-3 provide the basis for the computational model. We used the following equations to calculate S_{fx} and S_{fy} , based on Manning's equation:

$$S_{fx} = n^2(u_{ij})\sqrt{\frac{u_{ij}^2 + \bar{v}_{ij}^2}{h_{ij} + h_{(i+1)j}}} \quad (4)$$

$$S_{fy} = n^2(v_{ij})\sqrt{\frac{v_{ij}^2 + \bar{u}_{ij}^2}{h_{ij} + h_{i(j+1)}}} \quad (5)$$

$$\bar{u}_{ij} = \frac{u_{ij} + u_{i(j-1)} + u_{(i-1)(j-1)} + u_{(i-1)j}}{4} \quad (6)$$

$$\bar{v}_{ij} = \frac{v_{ij} + v_{i(j+1)} + v_{(i+1)(j+1)} + v_{(i+1)j}}{4} \quad (7)$$

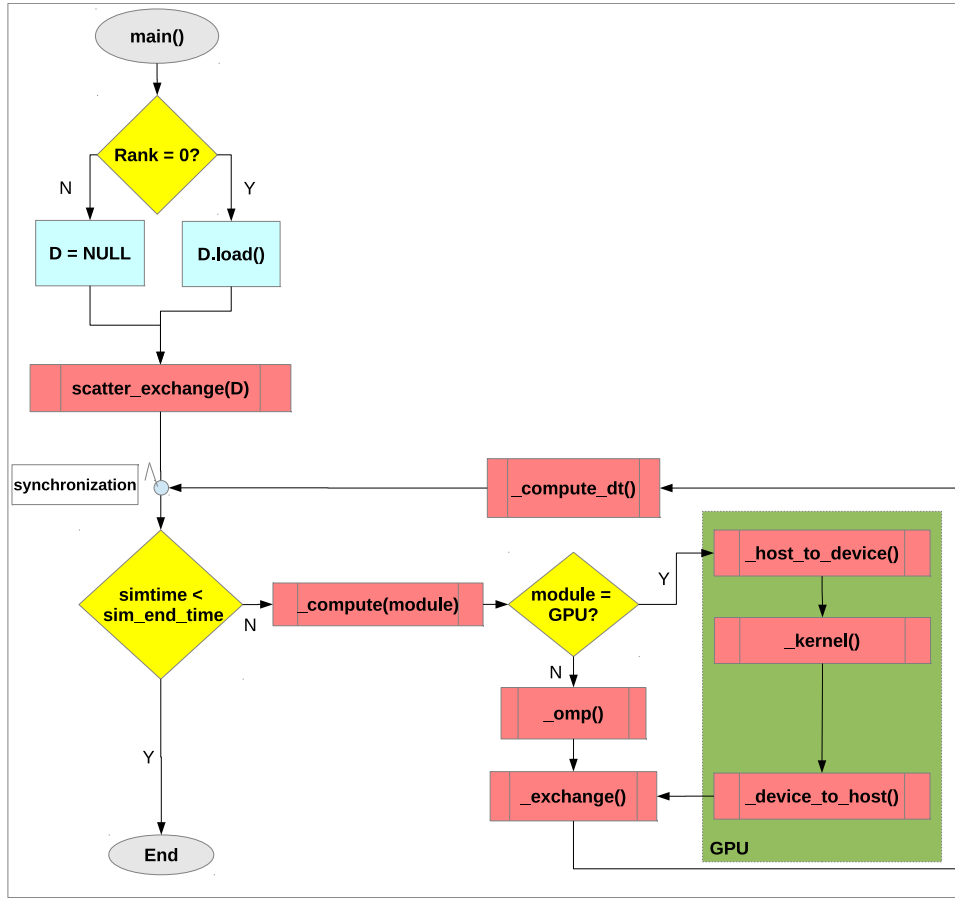


Figure 2: Execution control flow for our MPI+CUDA implementation. The function *scatter_exchange()* loads the input DEM using the root process (Process 0), performs a MPI scatter, and exchanges edge data. The *_omp()* function or *_kernel()* performs the core computations, depending on the module selected by the user.

In Equations 4 and 5, n is the Manning coefficient. For purposes of this study, we used a constant n of 0.035.

Temporal dependencies for each iteration include the interpolated flow that is derived from the input hydrograph using the following interpolation function:

$$h(r_t, f_t, s_t) = f_t + (s_t - r_{(t-1)}) \frac{f_t - f_{(t-1)}}{r_t - r_{(t-1)}} \tag{8}$$

We applied this function for the time step t , starting from a initial simulation time where $t = 1$ and s_t is the accumulated simulation time at time step t . Spatial dependencies for a grid cell include specific cell neighbors of the corresponding cell from the previous time step. By examining Equations 4-7, we identified the specific dependencies and treat this problem as an iterative stencil-based computation. When the domain is partitioned, we introduce a layer of complexity; if part of the stencil is on the edge of a distributed memory partition, the cell data must be communicated to the neighboring partition, as shown in Figure 1.

3.2 Implementing MPI+OpenMP

Figure 2 illustrates the workflow in general. The simulation begins by reading the configuration and input files from the root process (Process 0), then calculating the partition sizes based on the requested number of processes. Next, the root process determines which partitions will contain source cells, and performs the translation of source locations at global scope into local coordinates. Then, the root process distributes the D grid plus the localized configuration data over each partition. Each process will then initialize the H , U , V , grids based on the localized dimension data, and set the initial state. The b set contains the initial state at $t = 0$, and the a set will compute the state for time $t = 1$.

The simulation loop begins with an initial value *simtime* for the simulation time, and control will remain in the loop until the accumulated value of *simtime* equals or slightly exceeds the simulation end time given in the configuration file. The initial time step is set by the user and stored in the value *dt*, which can be configured by the user to be fixed or variable. After each iteration, *dt* is added to *simtime*, and in the case of variable *dt*, a new value is calculated by *_compute_dt()* for the next time step.

Variable *dt* can save time by increasing the *dt* value during periods where the water flows at a lower rate (where the flood extent undergoes small changes) and decreasing *dt* when water flow is high. We compute *dt* based on the following equations:

$$\Delta_x = \frac{dx}{(U_{max} + \sqrt{g(hx_{max} + \epsilon)})}, \quad (9)$$

$$\Delta_y = \frac{dy}{(V_{max} + \sqrt{g(hy_{max} + \epsilon)})}, \quad (10)$$

$$dt = [\text{Min}(\Delta_x, \Delta_y)] \quad (11)$$

Equation 9 finds U_{max} , the maximum value in the U grid, and hx_{max} , the value in the H grid located at the same index where U_{max} was found, and uses these values to compute Δ_x . Similarly, Equation 10 computes Δ_y by finding V_{max} and hy_{max} , and Equation 11 will set *dt* to the smaller of the two.

The compute function is called, then edge data is exchanged among the partitions as shown in Figure 1a, and the partition data is transferred from the host to the device. In our implementation, the *_omp()* function is a wrapper for a OpenMP *parallelfor* construct, which splits work among the requested number of threads as illustrated in Figure 1b. Synchronization of processes and threads occur after *dt* for the next time step is calculated.

3.3 Implementing MPI+CUDA

The MPI+CUDA version operates in a similar manner, excluding the OpenMP *parallelfor* and calling a kernel instead. Multiple processes can be allocated to the same device, if the dataset is sufficiently small. This strategy can be especially useful when the host-device link can support multiple communication channels.

The compute function invokes the CUDA API to execute the computational kernel that processes the data that was loaded on the device during the initial host-to-device transfer. After each iteration, the CUDA kernel will check to determine if a print operation is needed. If so, it will initiate a device-to-host transfer of the entire partition, and the synchronization point will be reached.

4 Experimental Setup

The experimental setup was designed to measure key performance metrics for specific test cases that we believe will illustrate general performance improvements. In the following subsections, we present these metrics and test cases, and describe the hardware and runtime configurations used in our experiments.

4.1 Key Performance Metrics

We used megacells per second and I/O seconds to measure the performance of our simulation runs.

4.1.1 Megacells per second

(Mc/s or M) is the metric we use to compare the computational efficiency across all implementations. Mc/s is defined as the number of cells (in millions) divided by the computational time required to process all cells. Higher Mc/s values correspond to higher performances, since more cells are computed in the same time reference window. This measure allows capturing the average speed during the computation and generating early performance projections early in the simulation run. We compute Mc/s according to Equation 12, where I is the number of iterations, rc is the total number of cells in the grid, and t is the runtime in seconds.

$$M = \frac{Irc}{t(10^{-6})} \quad (12)$$

4.1.2 I/O seconds

This metric is the time that is spent specifically on data reading/writing to disk and MPI communications. While disk I/O is a constant overhead, MPI messages may vary in size and frequency, depending on the runtime configurations.

4.2 Simulation Test Cases

We will use three different simulation datasets of varying size and complexity.

4.2.1 Taum Sauk

The first test case is based on a dam break that occurred to the Taum Sauk Dam in Reynolds County, Missouri in December 2005. It is a twin reservoir system, designed to produce electricity during peak periods, by discharging water from the upper reservoir on to the lower reservoir. Water is pumped back to the upper reservoir during off-peak periods. The upper reservoir built on the Proffit Mountain is approximately 760 feet above the floodplain of the East Fork Black River, with a storage capacity of 1.5 billion gallons. The upper reservoir failed on December 14th, 2005, at a 680-foot wide section on the northwest side of the Proffit Mountain, releasing the entire storage into East Fork Black River floodplain, through Johnsons Shut-Ins State Park, and into a lower storage reservoir within 25 minutes. In this study, this dam break scenario is simulated and a comparative analysis of the computational times reported by all the various implementations is performed.

The input data can be represented by a grid consisting of 1136 rows and 624 columns. The total simulation time is one hour. Using dynamically generated time steps that are 0.013 seconds on average, the simulation takes 290,453 iterations to complete

4.2.2 Cairo

The second test case simulates a flood event that occurred in 2011, during the weeks of massive flooding on the Mississippi River watershed. The levee system along the Birds Point-New Madrid Floodway was in danger of failure, and threatened to flood the town of Cairo, IL. Beginning on May

Table 1: MPI+OpenMP Configurations

| Run | OMP Threads | MPI Procs | Node(s) |
|-----|-------------|-----------|---------|
| R1 | 40 | 1 | N3 |
| R2 | 40 | 2 | N3, N4 |
| R3 | 40 | 4 | N1-N4 |
| R4 | 20 | 8 | N1-N4 |

2, 2011, a 2-mile section of the levee was detonated by the U.S. Army Corps of Engineers, causing over 100,000 acres of farmland within the watershed to become inundated [19].

The base topographical layer for our Cairo simulations is represented by an array containing a total of 1795 x 1555 cells, which we populate with elevation data to create the terrain. The simulation takes 80,384 iterations to compute 14,400 seconds of simulation time.

4.2.3 Large Random

The final test case is a large-scale simulation over randomly generated terrain. Random generation of the topography was performed using a variation of the Diamond-Square method [13], a inflow location representing the origin of the flood was chosen at random. The purpose of this test case is to observe the performance of the flood simulation when using input data requiring memory allocation near the capacity of certain GPU devices. By using a grid of 16387 x 16387 cells, the memory footprint for all 7 layers is near 7.5GB when using a 4-byte numerical type.

4.3 Hardware Specifications

For our simulation runs, we used a heterogeneous HPC cluster consisting of 4 nodes, each having Intel Xeon E5-2680 processors in 2 sockets, 10 cores each (with hyper-threading), for a maximum of 160 dedicated CPU processes across all 4 nodes. The heterogeneity is introduced with the GPGPU arrangement of each node, where node N1 is equipped with a single Nvidia Tesla K40M (k40) and node N2 with two K20M (k20) GPUs. The GPUs on nodes N3 and N4 will not be used for these simulations.

4.4 Runtime Configurations

Table 1 lists a set of simulation runs we will perform using the Taum Sauk and Cairo datasets. Run R1 is the baseline configuration, using one MPI process on a single node and 40 OpenMP threads. Runs R2 and R3 distribute the workload across two and four nodes, respectively, each using 40 OpenMP threads. Run R4 also uses 4 nodes, but uses two processes on each node. To maintain a maximum of 40 threads per node, we allow 20 OpenMP threads for each process. Using more than 40 OpenMP threads per node on our test system results in performance degradation because it exceeds the number of available CPU cores, so some threads will always be forced to wait until an active thread finishes its work and a core becomes available.

Table 2 lists a set of configurations where Runs R5-R24 will use the Taum Sauk and Cairo datasets and Runs R25-R27 will use the Large Random dataset. Runs R5-R11 divide the work between a number of processes on the same node; space for each process is allocated on the device, so multiple kernels allocated to separate work groups can run concurrently. Communication over PCIe supports full duplex transfers, but true concurrency is not guaranteed. A sufficiently large number of processes can create a communication bottleneck that overshadows performance gains made by effective utilization of the computational cores. Runs R12-R18 are similar to R5-R11, but they are performed on node N2, which is equipped with two k20 GPUs. Runs R19-R24 utilize two nodes and all GPU devices available across the nodes. Since there are three devices being used, we test the number of processes in increasing multiples of three. Runs R25-R27 utilize one or two

Table 2: MPI+OpenMP+GPU Configurations

| Run | MPI Procs | Devices | Node(s) |
|-----|-----------|---------------|---------|
| R5 | 2 | k40 | N1 |
| R6 | 4 | k40 | N1 |
| R7 | 6 | k40 | N1 |
| R8 | 8 | k40 | N1 |
| R9 | 12 | k40 | N1 |
| R10 | 16 | k40 | N1 |
| R11 | 18 | k40 | N1 |
| R12 | 2 | k20 x 2 | N2 |
| R13 | 4 | k20 x 2 | N2 |
| R14 | 6 | k20 x 2 | N2 |
| R15 | 8 | k20 x 2 | N2 |
| R16 | 12 | k20 x 2 | N2 |
| R17 | 16 | k20 x 2 | N2 |
| R18 | 18 | k20 x 2 | N2 |
| R19 | 3 | k20, k20, k40 | N1, N2 |
| R20 | 6 | k20, k20, k40 | N1, N2 |
| R21 | 12 | k20, k20, k40 | N1, N2 |
| R22 | 18 | k20, k20, k40 | N1, N2 |
| R23 | 24 | k20, k20, k40 | N1, N2 |
| R24 | 30 | k20, k20, k40 | N1, N2 |
| R25 | 2 | k20, k40 | N1, N2 |
| R26 | 8 | k20 x 2 | N2 |
| R27 | 8 | k20, k40 | N1, N2 |

nodes, each with a different selection of process distribution and devices used. The objective is to test the process and device distribution when the workload is near the limitations of the devices.

5 Results and Analysis

The results are organized by implementation- MPI+OpenMP, MPI+CUDA (same device type), MPI+CUDA (heterogeneous) and MPI+CUDA (heterogeneous) for extreme datasets.

5.1 MPI+OpenMP simulations of Taum Sauk and Cairo

Figure 3 shows the results in Mc/s of runs R1-R4 for both the Taum Sauk and Cairo simulations. Performance from 1 to 4 processes improves in both cases in a linear fashion, but the performance for the Cairo simulation is about half the performance of Taum Sauk. The row size and overall grid size for Taum Sauk is smaller, so the negative impact of communication overhead is expected to be seen at a lower number of processes than the Cairo simulation, which we see in the results for 8 processes. The performance improvement for Taum Sauk is much smaller when using 8 processes, while the

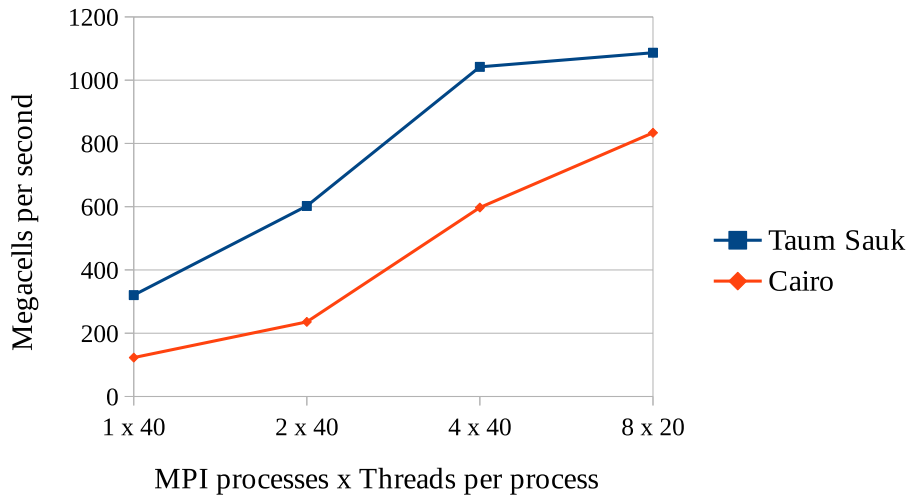


Figure 3: Comparison of MPI+OpenMP results for two simulations using an increasing number of processes and process mappings (MPI Processes x OMP threads).

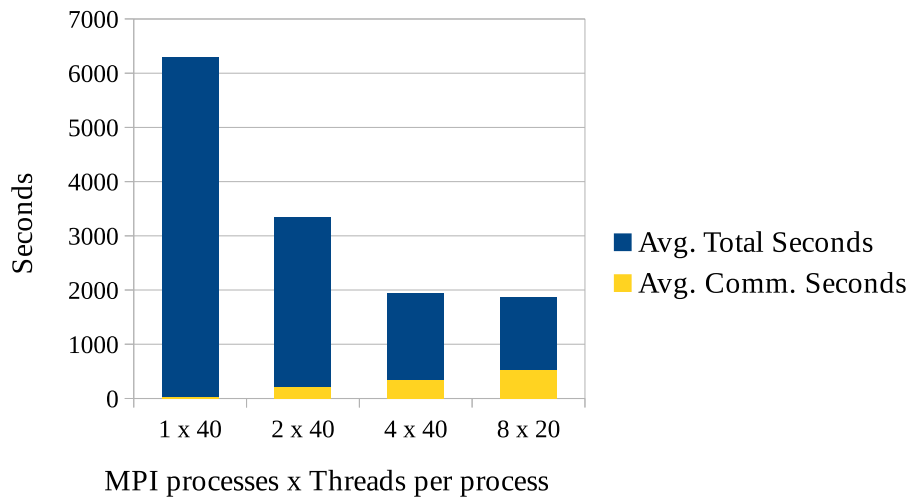


Figure 4: Comparison of MPI+OpenMP communication time versus total time for the Taum Sauk simulation using an increasing number of processes and process mappings (MPI Processes x OMP threads).

improvement for Cairo remains on a linear path. As the number of processes increase beyond 8, the performance of the Cairo simulation should continue to improve while the performance for the Taum Sauk simulation should approach a constant value or possibly degrade, since there is no added performance benefit when distributing increasingly smaller workloads.

The main reason for the lack of expected increase in the Taum Sauk performance involves MPI communication overhead as the partition sizes become smaller. We observe two very different patterns in Figures 4 and 5. In the Taum Sauk simulations, the communication time is significant under 4 and 8 processes, while it is less significant in the Cairo simulations, and nearly constant from 2 to 8 processes. The communication time for Taum Sauk increases with the number of processes until it is nearly half the total run time using 8 processes. The number of messages passed in the Taum Sauk simulations is the same as the number of messages passed in Cairo, but the size of the messages

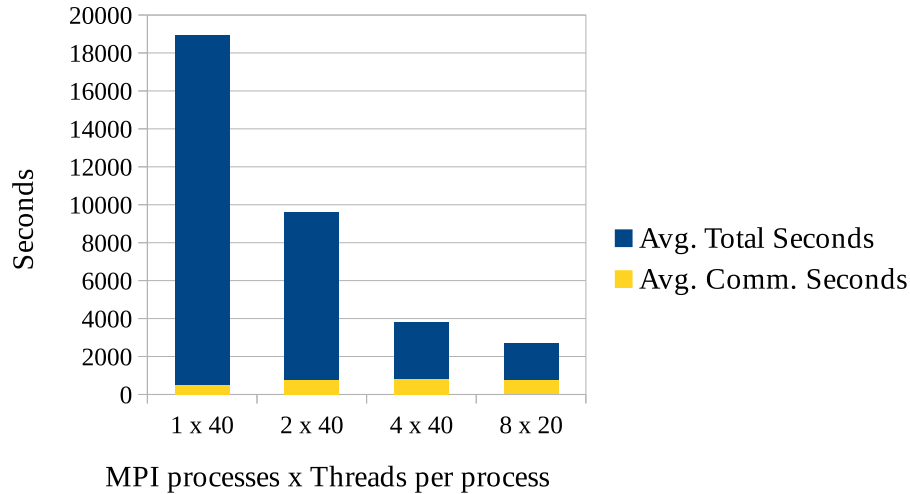


Figure 5: Comparison of MPI+OpenMP communication time versus total time for the Cairo simulation using an increasing number of processes and process mappings (MPI Processes x OMP threads).

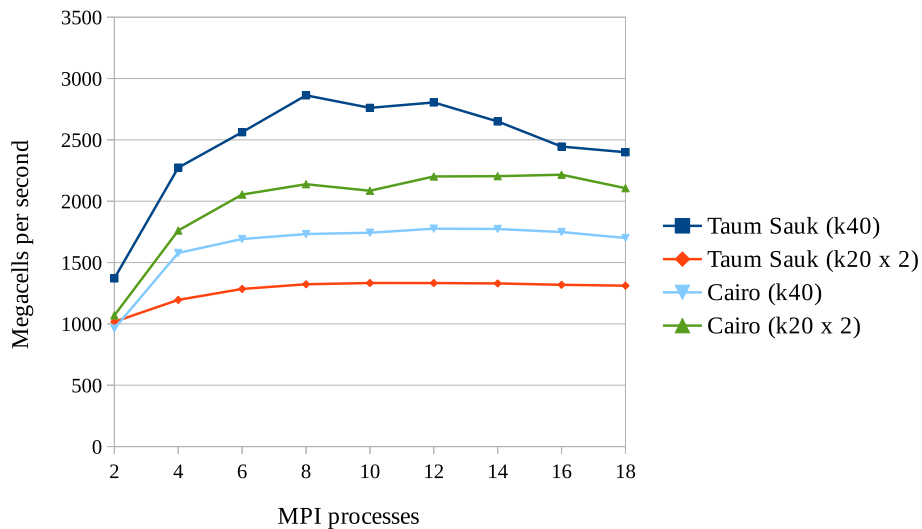


Figure 6: Comparison of MPI+CUDA runs of Taum Sauk and Cairo simulations using processes allocated to a single GPU versus processes balanced across 2 GPUs.

in Cairo is over twice the size of the Taum Sauk messages, since the row size in the Taum Sauk grid is 624 and the rows for the Cairo grid have 1555 elements. In Figure 5, the communication time for Cairo shows a less dramatic increase; we can also expect at some point that the communication overhead will be detrimental to this case, but likely at a higher number of processes than Taum Sauk.

5.2 MPI+CUDA simulations of Taum Sauk and Cairo

Figure 6 shows the impact of the process-to-device mapping of simulations using two different datasets and two generalized mappings, based on runs R5-R15. Performance is improved in all

cases while increasing to 8 processes. The most significant improvement is seen when using a single k40. Above 8 processes, we observe performance degradation caused by a bottleneck of host-device transfers over PCIe. Once the number of processes allocated to a GPU device becomes sufficiently large, kernels will be forced to wait for the transfers of other kernels to complete before their own transfer can be serviced.

The results for Runs R5-R18 are shown in Figure 6. The performance of Taum Sauk is dramatically better on a single k40 device than on the two k20 devices. Cairo performs slightly better on the k20 devices than on a single k40. The optimally efficient number of processes appears to be 8 for most cases.

5.3 MPI+CUDA heterogeneous simulations with three devices

Runs R19-R24 investigate the performance of Taum Sauk and Cairo using three available GPU devices across two nodes under various process-to-device mappings. In Section 3, we noted the MPI environment will govern the process to node allocation for most of our simulation runs. For this set of experiments, we will be manually assigning processes to nodes to have more control over the number of processes that can use a given device.

By default, the cluster environment we used will allocate processes in order for each host or node in the list. If the number of processes is greater than the number of hosts, the system will cycle back to the first host on the list and continue this sequence until all processes are allocated. Figure 7 illustrates the communication patterns that result from the combination of our partitioning strategy and the default "Cyclical" process allocation when using six processes over two hosts n0 and n1, where n0 has a single k40 GPU (n0d0) and n1 has two k20 GPUs (n0d0, n0d1). In our partitioning strategy, communications of edge data will occur between processes of rank r and rank $r + 1$, except for the process with the highest rank, and rank r and rank $r - 1$ except for the process with the rank of 0. A practical illustration of this partitioning is given in Figure 1, where edge data (highlighted in green) is exchanged after each time step.

We can manually assign processes to hosts in a way that complements our partitioning strategy, as shown in the "Adjacent" mapping in Figure 7. Under this mapping, we expect to see higher performance in many cases, because more intranode MPI communications will occur than would occur under the default mapping. We assume intranode communications in our test environment are at least as fast as internode communications, because more messages can be passed over the system bus without needing to travel over the full networking stack to another physical node. The difference may be negligible in some cases, but we expect to identify other cases where it would be advantageous to manually define which processes are assigned to each node.

Figure 8 shows the performance of Runs R19-R24 for both Taum Sauk and Cairo using both the cyclical and adjacent mapping strategies. For Cairo, the differences are relatively small and favor adjacent in most cases, while Taum Sauk results show a greater disparity as the number of processes increase to 30. Figure 9 gives some insight, as we see the communication times for Cairo having less variance than Taum Sauk, which shows greater increases in communication overhead from 12 to 30 processes. This is consistent with the results for MPI+OpenMP given in Figures 4 and 5 that also reveal an increasing communication overhead on the smaller grid.

5.4 MPI+CUDA heterogeneous simulations of a large random grid

Experimentation with the large grid shows different process-to-device mappings have a significant impact on performance, as shown in Figure 10. A heterogeneous HPC environment is expected to have different GPU devices available, some may be more powerful, or otherwise better suited to the problem instance than others. Because of its intentionally large size, attempts to run these simulations on a single k20 GPU will cause the CUDA runtime to generate out-of-memory errors.

The k40 has a higher capacity, so it can successfully run this dataset on its own. In Runs R25-R27, we first distribute the data into 2 processes and run on 2 different GPUs. Then we split into 8 processes and run on the pair of k20 devices, then a single k40 plus a single k20. Performance is improved dramatically from Runs R25 to R27, although the number of processes remain the same

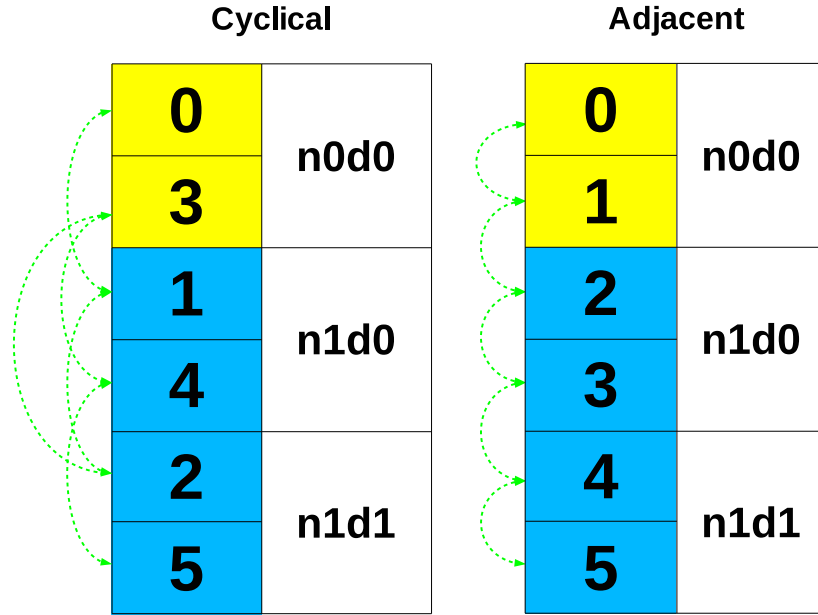


Figure 7: A cyclical distribution (left) and adjacent distribution (right) of six processes demonstrates the difference in communication patterns based on how processes are mapped. The cyclical distribution does not utilize any intranode links while the adjacent distribution utilizes four.

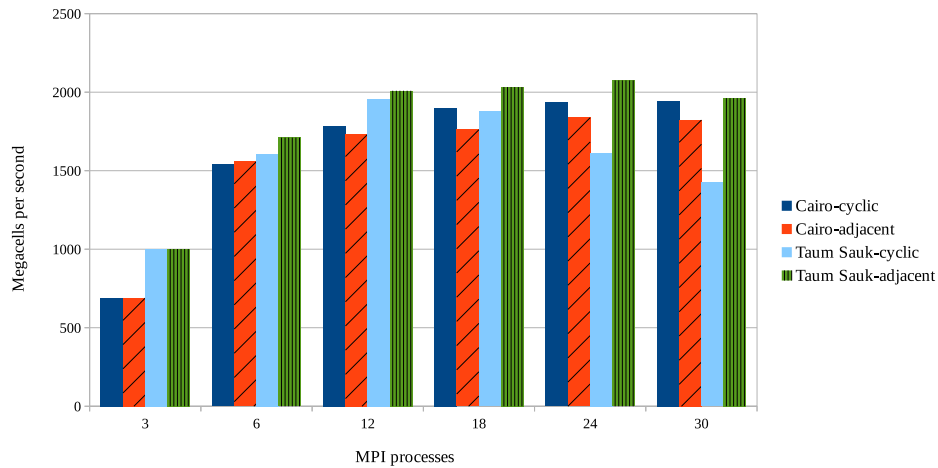


Figure 8: Comparison of MPI+CUDA runs using a k40 and two k20 GPUs across two physical nodes using cyclic versus adjacent process distribution.

and the 8 processes in Run R25 require no ethernet communications between the processes. The k40 is slightly faster than the k20, with a higher cache availability. Additionally, the k40 has higher throughput speeds, around 2Gbps versus 1.6Gbps for the k20.

6 Further Performance Enhancements

After examination of the initial set of results, there are a couple observations of note that we can use to enhance the performance in certain cases. This simulator accepts an initial state with one or

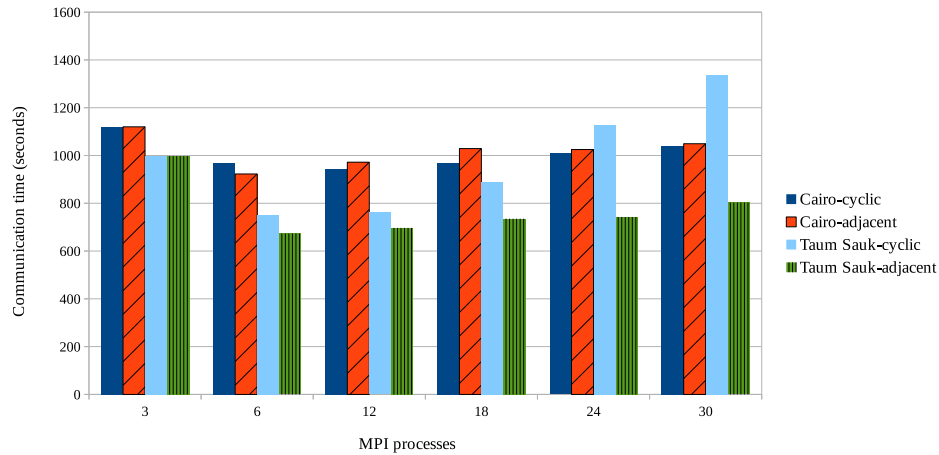


Figure 9: Comparison of communication times for MPI+CUDA runs using a k40 and two k20 GPUs across two physical nodes using cyclic versus adjacent process distribution.

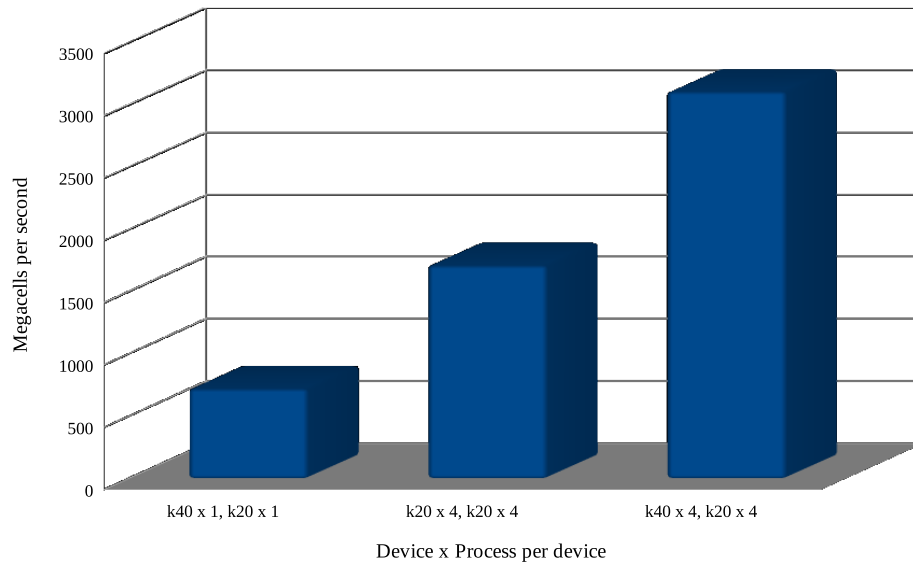


Figure 10: Comparison of MPI+CUDA simulation runs using the Large Random dataset under different process-to-device mappings.

more source cells, then propagates the flood extent into cells reachable at a Chebyshev distance k . For purposes here, we assume $k = 1$. Then the base case of the computational model applied to a single cell at time t affects its Moore neighborhood at time $t + 1$. We know if the h -value of the cell is zero and every neighbor is also zero, then the h -value (and by extension, u -value and v -value) of the cell at time $t + 1$ after applying the computational model will be zero and therefore, a redundant computation. A dry cell, i.e. h -value is in the range $[0, \epsilon)$, at time t will always be dry at time $t + 1$ if it is surrounded by dry cells. Applied to a larger scale, we possibly omit a large contiguous space equal to the size of a partition. Since the implementation allows for the exchange of edge cells after each iteration, we will use the upper and lower edges of a partition as our neighborhood. In cases where the flood extent does not expand into one or more additional partitions during any part of the simulation, we can use a simple method to determine if a partition is dry by checking these edge cell values of a known dry partition after the exchange is performed at the end of each time step. If

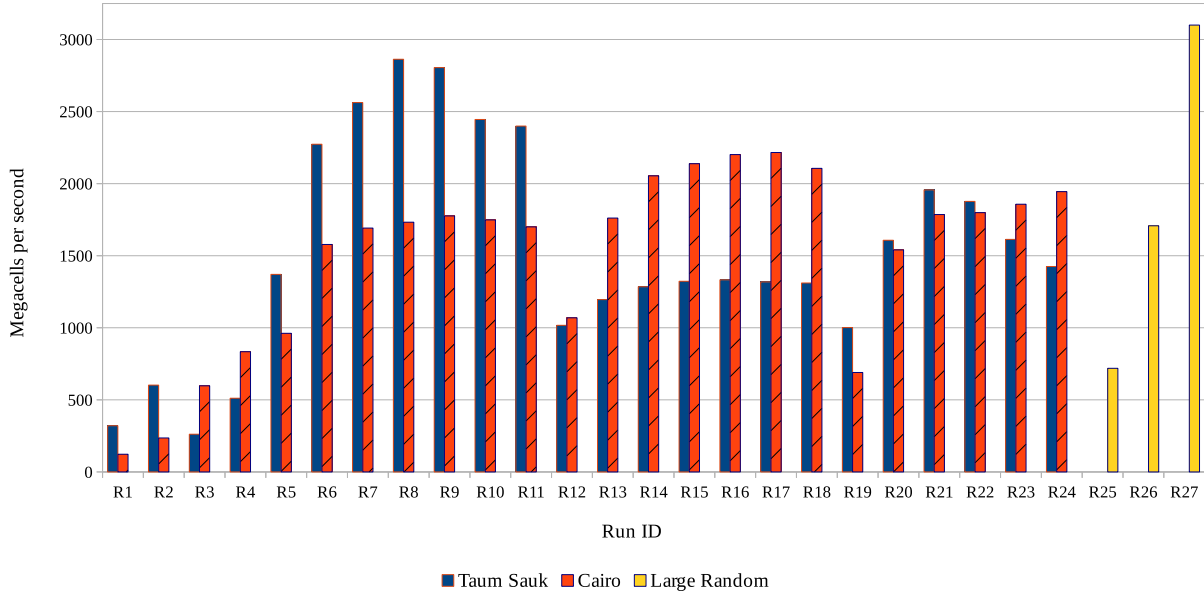


Figure 11: Performance in Mc/s for each of the runs listed in Tables 1 and 2. Runs R1-R4 used MPI+OpenMP, R5-R11 used MPI+CUDA on a single GPU, R12-R18 used two GPUs and R19-R24 used three GPUs. Runs R1-R24 used the Taum Sauk and Cairo datasets, while Runs R25-R27 used only the Large Random dataset.

any of the halo cells are non-zero, we mark the partition as wet, otherwise we can skip computation for the entire partition.

In general, print operations are performed at fixed intervals after a set of iterations are complete. With large datasets, this could take 30 seconds or more, and processes must wait until printing is complete before launching the next iteration. We can reduce the overall run time by launching a thread to perform file I/O operations asynchronously on data generated during time step t once computation begins for time step $t + 1$. For purposes of experimentation, we have done so on the MPI+CUDA version.

6.1 Processing partitions of wet cells only

At any given time step, the data domain may have partitions that are dry. A dry partition at time step t may become a wet partition at time $t + 1$ if any updated edge values from neighboring partitions are non-zero at the end of time step t . In Figure 12, the third partition from the top is dry after the first 6 minutes of simulation time, but contains a portion of the flood extent at 18 minutes. At a simulation time occurring between 6 and 18 minutes, the flood extent reaches the lower boundary of the second partition, and edge data is transmitted to the third partition from the top. Water may recede completely from a partition, so a lower bounds for the number of wet partitions at time step $t + 1$ is not known until time step t is complete.

In Figure 13, we repeated Runs R1-R4 with the wet partition modification. As expected, the single process run shows similar performance to R1, where the multiple process runs show significant improvement in each case. The greatest improvement is at R3, where we used 4 processes, but we also see at R4 a slight decrease in performance with the modified run where there was an improvement in the original runs from R3 to R4. As the number of partitions increase, data size per partition gets smaller and the number of communications increase. Additionally, the number of checks for dry partitions increase and in general, the amount of simulation time a given partition remains dry will decrease. Figure 14 shows a similar pattern, however we continue to see increased performance on the modified R4 run.

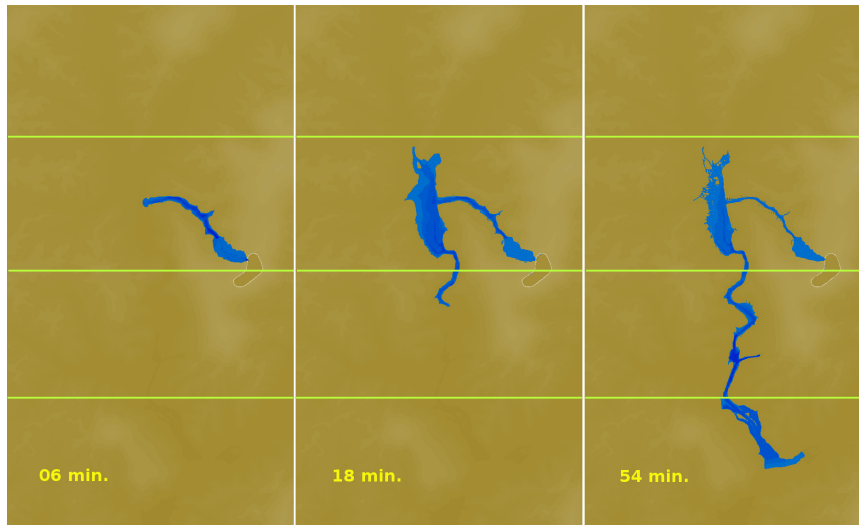


Figure 12: Rasterized outputs for Taum Sauk simulation at 6, 18 and 54 minutes, partitioned for 4 processes.

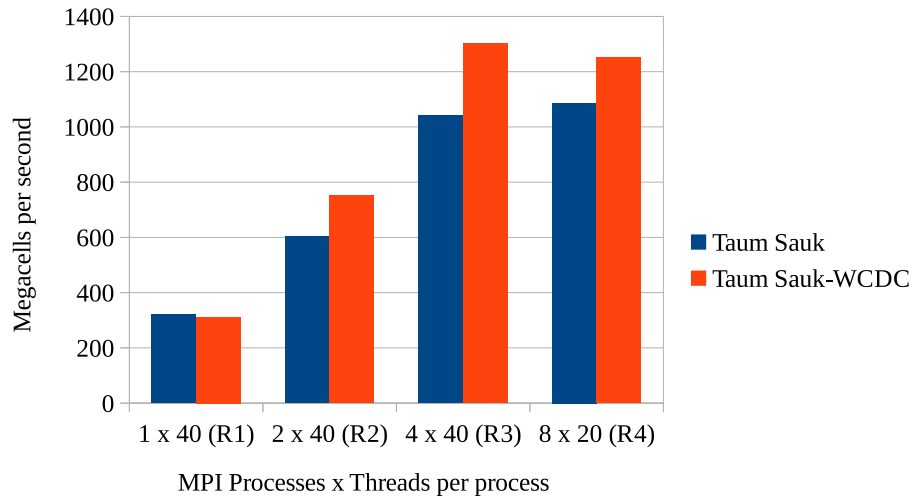


Figure 13: Results for Taum Sauk using MPI+OpenMP under original configuration (Runs R1-R4) versus processing wet partitions only (WCDC).

Figure 12 shows when divided into 4 processes, the flood extent does not include the topmost partition. From experimentation, we are able to determine that wet cells never reach this partition over the entire 60 minutes of simulation time. Furthermore, only a single partition is wet for the first 14 minutes, two partitions are wet for the next 14 minutes, then three partitions are wet for the final 32 minutes.

6.2 MPI+CUDA: Asynchronous I/O

Asynchronous I/O can be performed when MPI+CUDA is enabled, which overlaps computation with output by allowing the host process to perform output operations on data received from a previous device-to-host transfer while the device is computing results for the next device-to-host transfer request. In Figure 15, a dedicated thread is forked to perform disk writes after every n th

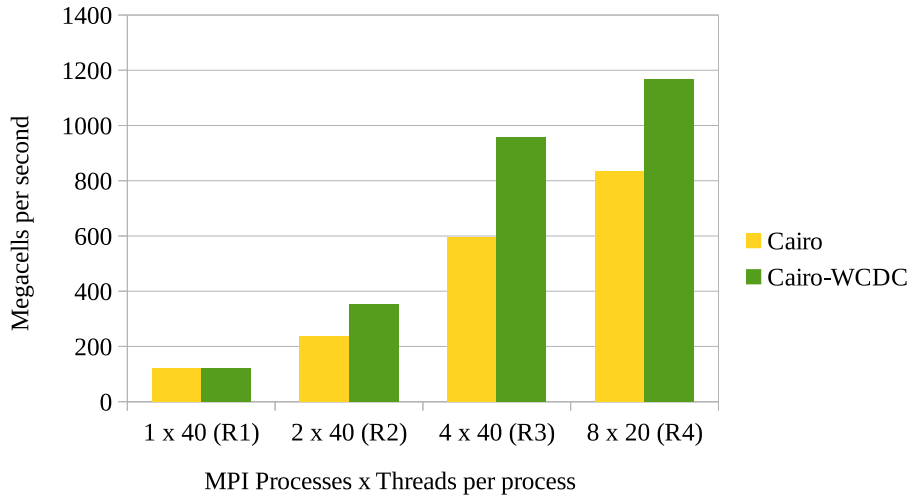


Figure 14: Results for Cairo using MPI+OpenMP under original configuration (Runs R1-R4) versus processing wet partitions only (WCDC).

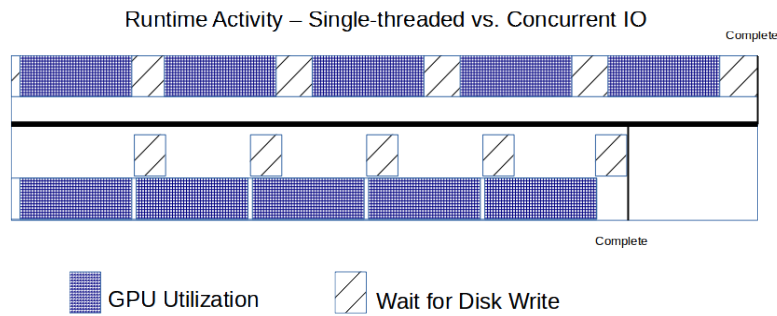
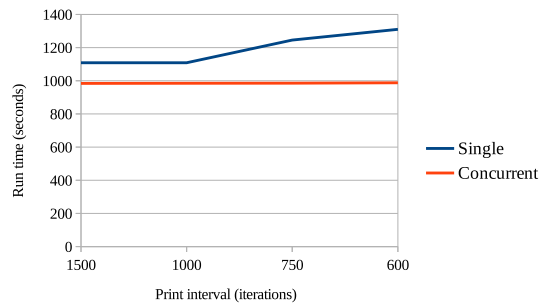


Figure 15: Example of time saved by using asynchronous I/O. The plot (top) shows, as I/O requests increase in frequency along the x -axis, the single thread version takes longer, while the concurrent version remains constant.

iteration while the main thread concurrently returns control to the device to resume execution of the compute kernel. With asynchronous I/O, decreasing n results in no change in the overall run time, though a lower limit can possibly exist for n where overall run time will begin to increase.

7 Conclusion and Future Work

We have shown in all cases that decomposing the problem and distributing the data across multiple processes will result in higher performance than running the same simulation with a single process. While not intended to be a direct comparison, the results reported in [6] for a domain size similar to Cairo range from approximately 300-1600 Mc/s while the results for Cairo on a single GPU range from 900-1800 Mc/s. Noting that each physical node is equipped with identical CPUs, Figure 11 gives a summary of the performance in Mc/s of all runs, where R1 represents the performance of the original OpenMP code, showing around a 9x improvement at R8 for Taum Sauk and around a 18x improvement for Cairo at R17. For the OpenMP implementation in Runs R1-R4, we demonstrated how improvement follows a linear path as the number of processes increase (although there is a limitation to these performance gains, based on the size of the simulation domain and communication patterns). In Runs R5-R18, we illustrate how knowledge of the underlying hardware can be a significant advantage, as it is possible in some cases to give a single device the same workload as two less powerful devices and achieve better results.

Runs R19-R24 highlighted how knowledge of process mapping on a system and software-defined partitioning can be used to increase performance if they have complementary strategies. With Runs R25-R27 we found that knowledge of the underlying hardware devices is essential for developing software meant for large scale grid simulations, since data transfer operations will fail when the data is too large to fit in device memory.

We have shown it can be beneficial to apply a wet/dry classification to each partition to avoid unnecessary work, and also to explore efficient methods of file I/O that are now available due to the execution environment. Future exploration includes developing a set of generalized data structures that support a number of different partitioning strategies, irregular computations as presented in [7], runtime configuration parameters that match a selected partitioning strategy with a complementary mapping strategy, and a metacomputing API that is simple and user-friendly, to allow domain scientists access to more efficient hardware without the expertise that is required today.

References

- [1] The google geocoding API - google maps API web services google developers.
- [2] The National Map: Topographic Maps for the 21st Century. Technical report, Reston, VA, 2002.
- [3] P D Bates and A P J De Roo. A simple raster-based model for flood inundation simulation. *Journal of Hydrology*, 236(12):54–77, 2000.
- [4] Paul D Bates, Malcolm G Anderson, Laura Baird, Des E Walling, and David Simm. Modelling floodplain flows using a two-dimensional finite element model. *Earth Surface Processes and Landforms*, 17(6):575–588, 1992.
- [5] Carlo Bertolli, Gabriele Mencagli, and Marco Vanneschi. Adaptivity in risk and emergency management applications on pervasive grids. *I-SPAN 2009 - The 10th International Symposium on Pervasive Systems, Algorithms, and Networks*, pages 550–555, 2009.
- [6] André R Brodtkorb, Martin L Sætra, and Mustafa Altinakar. Efficient Shallow Water Simulations on GPUs : Implementation , Visualization , Verification , and Validation. (0314), 2010.
- [7] Michal Chladek and Roman Durikovic. Particle-based shallow water simulation for irregular and sparse simulation domains. *Comput. Graph.*, 53(PB):170–176, December 2015.
- [8] National Research Council. *Mapping the Zone: Improving Flood Map Accuracy*. National Academies Press, Washington, D.C., 2009.

- [9] N.H. H Crawford, R K Linsley, Stanford University. Dept. of Civil Engineering, and F.K. Linsley Jr. *Digital simulation in hydrology: Stanford watershed model IV*. Technical report (Stanford University. Dept. of Civil Engineering). Dept. of Civil Engineering, Stanford University, Stanford, CA, 1966.
- [10] A Crossley, R Lamb, and S Waller. Fast solution of the Shallow Water Equations using GPU technology. pages 1–6, 2010.
- [11] Romano Fantacci, Marco Vanneschi, Carlo Bertolli, Gabriele Mencagli, and Daniele Tarchi. Next generation grids and wireless communication networks: towards a novel integrated approach. *Wireless Communications and Mobile Computing*, 9:445–467, 2009.
- [12] Joel H Ferziger and Milovan Perić. *Computational methods for fluid dynamics*, volume 3. Springer Berlin, 1996.
- [13] Alain Fournier, Don Fussell, and Loren Carpenter. Computer rendering of stochastic models. *Commun. ACM*, 25(6):371–384, June 1982.
- [14] P García-Navarro, P Brufau, J Burguete, and J Murillo. The shallow water equations : An example of hyperbolic system. *Monografías de la Real Academia de Ciencias de Zaragoza*, pages 89–119, 2008.
- [15] Noemi Gonzalez-Ramirez. Simulating Flood Propagation in Urban Areas using a Two-Dimensional Numerical Model. 2010.
- [16] Danish Hydraulic Institute. Mike 11 reference manual, appendix a. scientific background. *Danish Hydraulic Institute*, 2001.
- [17] Alfred J. Kalyanapu, Siddharth Shankar, Eric R. Pardyjak, David R. Judi, and Steven J. Burian. Assessment of GPU computational enhancement to a 2D flood model. *Environmental Modelling & Software*, 26(8):1009–1016, aug 2011.
- [18] Alfred Jayakar Kalyanapu, Steven J Burian, and Timothy N. McPherson. Effect of land use-based surface roughness on hydrologic model output. *Journal of Spatial Hydrology*, 9(2):51–71, 2009.
- [19] T A Koenig and R R Holmes. Documenting the stages and streamflows associated with the 2011 activation of the New Madrid Floodway, Missouri. 2013.
- [20] Robert Manning, John Purser Griffith, T F Pigot, and Leveson Francis Vernon-Harcourt. *On the flow of water in open channels and pipes*. 1890.
- [21] Kate Marks and Paul Bates. Integration of high-resolution topographic data with floodplain flow models. *Hydrological Processes*, 14(11-12):2109–2122, 2000.
- [22] Suhas V Patankar. *Numerical heat transfer and fluid flow*. Taylor & Francis, 1980.
- [23] P G Samuels. Cross-section location in 1-D models. In *Proc., Int. Conf. on River Flood Hydraulics*, pages 339–350. Wiley Chichester, UK, 1990.
- [24] Siqing Shan, Li Wang, Ling Li, and Yong Chen. An emergency response decision support system framework for application in e-government. *Information Technology and Management*, 13:411–427, 2012.
- [25] V Tayefi, S N Lane, R J Hardy, and D Yu. A comparison of one- and two-dimensional approaches to modelling flood inundation over complex upland floodplains. *Hydrological Processes*, 21(23):3190–3202, 2007.
- [26] USGS. National Mapping Program Technical Instructions: Standards for Digital Elevation Models. Technical report, U.S. Geological Survey National Mapping Division, Reston, Virginia, 1998.

- [27] Gary W. Brunner. Hec-ras river analysis system. hydraulic reference manual. version 1.0. page 143, 07 1995.