

# Accelerating multi-dimensional population balance model simulations with a highly scalable framework on GPUs.

Chaitanya Sampat<sup>a</sup>, Yukteshwar Baranwal<sup>a</sup>, Rohit Ramachandran<sup>a,\*</sup>

<sup>a</sup>Chemical and Biochemical Engineering, Rutgers University, Piscataway, NJ, USA - 08854

---

## Abstract

Population Balance models (PBMs) are widely used to simulate and optimize the distributed particulate process. Prediction of time evolution of the distribution of the particulate properties using these PBMs, is very important in understanding key dynamics. However, with an increase in the number of components and spatial configurations the process, the complexity of a PBM increases. This leads to multi-dimensional matrix calculations potentially requiring significant computational power. Solving such a system of equations with a traditional central processing unit (CPU) would be computationally intractable. This study focuses on the development of a scalable algorithm to parallelize the nested loops inside the PBM via a GPU framework. The developed PBM was unique since it adapts to the size of the problem and uses the GPU cores accordingly. This algorithm was parallelized specifically for NVIDIA GPUs as it was written in CUDA C/C++. The major bottleneck to such algorithms is communication time between the CPU and the GPU which in our studies contributed to less than 1% of the total run time leading to a maximum speedup of about 12 over the serial code being achieved on the GPU. The speed up reported for the GPU PBM was about 2x the PBM code on multi-core configuration on a desktop computer. The speed improvements are also reported for various CPU & GPU architectures and configurations.

---

\*Corresponding author

Email address: rohitrr@soemail.rutgers.com (Rohit Ramachandran)

*Keywords:* Population Balance Model, GPU, Parallel Computing, Granulation, MPI, CUDA

---

## 1. Introduction

Over the past decade, Population Balance Models (PBMs) have been used widely to predict dynamics of distributed processes (Ramkrishna and Singh, 2014). The population balance equation is essentially a number conservation of particles (entities) which regulate the behavior of the entire system. PBMs consist of several external (i.e spatial) and internal coordinates and with an increase in grids of these coordinates it can lead to a more accurate model. With The increase in the number of grid of these coordinates, it leads to an increase in calculations for each time step, leading to a higher simulation times. The calculations increase by a factor of  $n^4$  where n being the number of entity grids. PBMs describe the evolution of entities into different states and into newer entities. Deriving accurate equations to determine these rate of change of internal coordinates requires precise empirical correlations from experimental data. Another way to increase the accuracy of these models is to introduce a first-principle based kernel into these models (Barrasso and Ramachandran, 2015). An accurate model which incorporates higher number of grids as well as includes a mechanistic kernel in its calculations is expected to be sluggish to simulate and could take upto an hour (Barrasso et al., 2015) to complete. Such models and their solving techniques are not viable to be used in real time system control. Thus, there is a need to improve the time it takes to simulate the model.

The advancement of computers and its components in recent years have led to a great increase in computational resources leading to faster simulations. The recent central processing unit (CPU) now contain various cores thus making it possible to run multiple processes in parallel. In order to take advantage of a highly parallel framework, large number of cores are required which may not be possible in a personal desktop and a supercomputer cluster would be needed.

Another **computer component** that can be used to run a highly parallel code is the computer's graphic processing unit (GPU) (Prakash et al., 2013). ~~These~~ GPUs contain thousands of compute cores that can be used ~~run~~ tasks in parallel. Thus, a desktop equipped with a GPU could compute the same results as a CPU code on supercomputers in lesser amount of time as seen in Section 2.1. With the launch of Compute Unified Device Architecture (CUDA), NVIDIA ~~made~~ it easier to use GPUs for general parallel programming in an approach usually termed as general purpose computing on GPUs (~~GPGPUs~~).

Various chemical industries such as detergent, food, pharmaceutical, fertilizers, catalyst encounter particulate processing daily. These processes constitute to about 50% of the world's chemical production (Seville et al., 2012). In the pharmaceutical industry, particulate processes are ~~widely~~ used to increase the size of the granules, improve flow-ability, increase yield strength etc. One of major processes is granulation, in which fine pharmaceutical powder blends are converted to larger granules using a liquid or a dry binder (Chaturbedi et al., 2017). These larger granules help in better flow ability and strength to these mixtures aiding further processing.

The main objective of this work is to develop a **highly scalable framework** ~~for~~ to simulate PBMs on GPUs. This framework **would be** independent of the type of kernels as well as number of spatial compartments as well as the number of solid bins used. The framework **would be** universal and would be able to run of different NVIDIA GPUs without the need for any changes. The number **of threds** that the PBM will execute depends upon the size of the problem as well as the number of cores available on the GPU. This code was developed in NVIDIA CUDA C/C++. A similar code was also developed on C++ to be run on the CPU which has limited scalability due to number of CPU cores available on desktop computer. This work also enables the use of desktop computers to obtain numerical solutions to computationally intensive tasks rather than relying on supercomputers ~~for quicker results~~.

## 2. Background and previous works

### 2.1. Population balance modeling and wet granulation

Population balances equations have been successful in predicting physical phenomena occurring in granulation such as aggregation, breakage and consolidation. These models predict how groups of distinct entities behave on a bulk scale during granulation. A general representation of the model is:

$$\begin{aligned} & \frac{\partial}{\partial t} F(\mathbf{v}, \mathbf{x}, t) + \frac{\partial}{\partial \mathbf{v}} [F(\mathbf{v}, \mathbf{x}, t) \frac{d\mathbf{v}}{dt}(\mathbf{v}, \mathbf{x}, t)] + \frac{\partial}{\partial \mathbf{x}} [F(\mathbf{v}, \mathbf{x}, t) \frac{d\mathbf{x}}{dt}(\mathbf{v}, \mathbf{x}, t)] \\ &= \mathfrak{R}_{formation}(\mathbf{v}, \mathbf{x}, t) + \mathfrak{R}_{depletion}(\mathbf{v}, \mathbf{x}, t) + \dot{F}_{in}(\mathbf{v}, \mathbf{x}, t) - \dot{F}_{out}(\mathbf{v}, \mathbf{x}, t) \quad (1) \end{aligned}$$

where  $\mathbf{v}$  is a vector of internal coordinates.  $\mathbf{v}$  is commonly used to describe the solid, liquid, and gas content of each type of particle. The vector  $\mathbf{x}$  represents external coordinates, usually spatial.  $F$  represents the number of particles present inside the system,  $\dot{F}_{in}$  and  $\dot{F}_{out}$  is the rate of particles coming in and going out the system respectively.  $\mathfrak{R}_{formation}$  and  $\mathfrak{R}_{depletion}$  are the rate of formation and depletion due various phenomena.  $\mathfrak{R}_{formation}$  described in Equation 5 is the most computationally intensive component of the PBM. This calculation consists of a double integral for each entity being tracked. In a system tracking 2 particulates, it would require 4 nested loops to compute these integrals. Figure 5 also indicates that this calculation takes 25% of the total run time. Thus, parallelizing this calculation is crucial to the effectiveness of the parallel algorithm.

Wet granulation is the process of engineering granules from pharmaceutical powder blends with the addition of liquid or solid binders. This process is usually carried out to obtain granules with a certain PSD, bulk densities and other physical properties (Barrasso and Ramachandran, 2015). There are about three rate processes that occur due the addition of a liquid binder to the powder mixture are wetting and nucleation, consolidation and aggregation, and breakage and attrition (Sen et al., 2014). One way to perform wet granulation of

pharmaceutical powders is to use a high shear granulator. In a high shear granulator, as the liquid binder is distributed within the granulator, liquid bridges are formed between particles and are subjected to the impeller. The impeller rotation along with particle-particle and particle-wall interactions results in the aforementioned granulation mechanisms.

## *2.2. Parallel Computing*

Computing with its current infrastructure has become an important tool in science. Parallel computing is one of the more extensively used type of computing used by scientists to perform simulations. It is the process of splitting of larger calculations into many smaller processes executed concurrently (Almasi and Gottlieb, 1989). This type of execution helps achieve large speed gains over simulations run on a single core in a serial manner. The computational task can be decomposed by various means to help simulate the system in a reasonable amount of time. The simulation problem can be decomposed either at the task level or at the data level. Task parallelism involves each process to behave distinctively from another as they would each be performing different operations. These multiple operations could be performed on a single data set or on multiple data sets, known as multiple instruction single data (MISD) and multiple instruction multiple data (MIMD) respectively. On the other hand, data parallelism involves the distribution of data across various processes which usually perform same set of operations on the data (Solihin, 2015). This type of parallelism is known as single instruction multiple data (SIMD). MIMD and SIMD can also be combined in certain systems, thus decreasing the simulation times further.

## *2.3. GPU based parallel computing*

Traditionally, large parallel jobs needed to be run on supercomputers which had thousands of cores, but these require special components making them expensive. Graphic processing unit (GPU) were initially used for vector calculations to support graphics inside a computer system. But, lately GPU manufacturers have started to promote them general computing as well. This form

of computing has been gaining popularity among scientists to accelerate simulations (Kandrot and Sanders, 2011). These GPUs comprise of a massively parallel architecture with hundreds to thousands of computational cores which can have thousands of active threads running simultaneously (Keckler et al., 2011). This means that GPUs have large computing potential the GPU computing can be exploited using parallel programming languages such as OpenCL and CUDA.

CUDA is an application programming interface (API) developed by NVIDIA (NVIDIA Corporation, 2012) that enables users to program parallel code for execution on the GPU. This framework is an extension implemented on top of C/C++ or Fortran. Parallel code for the GPU is written as kernels, which theoretically are similar to functions or methods in traditional programming languages. As several parts of the code need to be executed only once during a simulation, only few sections of the code can be written in terms of kernel while the remaining has to be executed in serial on the CPU of the system. The `nvcc` compiler from the CUDA toolkit prioritizes the compilation of these kernels before passing the serial section of the code to the native C/C++ compiler inside the system. There are three main parallel abstractions that exist in CUDA are grids, blocks and threads (Santos et al., 2013). Each CUDA kernel is executed in a serial manner during the execution of the program unless specified, where the kernels can be run in parallel using CUDA streams. Each kernel executes as a grid which in turn consists of various blocks which are constituted by various threads. This thread-block-grid hierarchy helps obtain fine grained data level and thread level parallelism. An illustration of this hierarchy is observed in Figure 1.

Another important aspect related to GPU parallelization is the data communication between the threads. The GPU consists of various memory modules with different access limitations as shown in Figure 1. The threads inside each block can communicate with each other using the shared memory. This memory is local to the block where these threads exist i.e. they are not accessible by threads from other blocks. In addition to the shared memory each thread has

its own local memory where local/temporary variables for each kernel can be saved to them. The threads from different communicate with each other using the global memory which is visible to all blocks inside the GPU at the cost of higher communication times. Accessing data from the local memory is the fastest for a thread and its slows down when the thread needs to retrieve data from the shared memory inside a block and accessing data is the slowest from the global memory of the GPU.

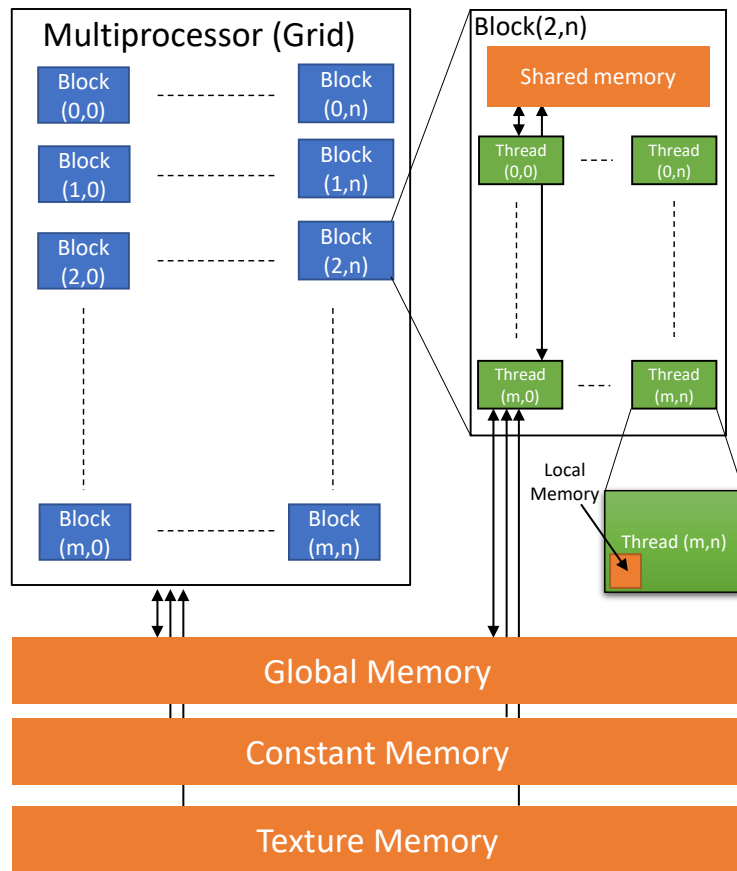


Figure 1: The parallel structure matrix inside the GPU and the various memories associated with each structure.

#### *2.4. Previous parallelized PBM works*

PBMs have been computationally intensive notably ones with larger number of internal coordinates and higher dimensions. Thus, several researchers have made attempts to increase the speed of these simulations. Gunawan et al. (2008) developed a parallelization technique using a high-resolution finite volume solution of the PBM. The studies carried out by Gunawan et al. (2008) for their parallel PBM algorithm achieved good parallel efficiency upto 100 cores. This study was limited by the number of grids used for the PBM, which meant the problem size was not computationally very heavy. This study also suggested that an algorithm with a shared memory model could help improve simulation speeds further. A hybrid memory model which uses both shared and local memory was implemented by Bettencourt et al. (2017) to obtain speed improvements of about 98% from the serial code. This implementation took into account both Message Passing Interface (MPI) as well as Open Multi-processing (OMP). A similar PBM parallelization approach was also undertaken in Sampat et al. (2018) and an approximate speedup of 13 was obtained. The reduction in speed was attributed to the use of dynamic arrays used in their PBM framework to accommodate the hybrid nature of the model being used.

Algorithms to parallelize the PBM codes on GPU have been studied briefly by Prakash et al. (2013) using the inbuilt MATLAB's parallel computing toolbox (PCT). This study was able to achieve good speedup values but could have been higher if the code had been implemented in native programming languages such as C or FORTRAN. Since MATLAB is a high level language it internally converts the written code to native programming languages before it is sent to the processor leading to excess computation which can be avoided (MathWorks<sup>TM</sup> Documentation, 2017b). Other works that have used GPU acceleration to improve computation times for their population balance simulations include those from various other chemical engineering processes such as crystallization (Szilágyi and Nagy, 2016), combustion (Shi et al., 2012), multi-phase flow (Santos et al., 2013), coagulation dynamics (Xu et al., 2015). Several of these works fail to address the problem size as a factor in determining num-



ber of GPU cores used for the simulation. These algorithms are restricted to a maximum number of cores that are used and with an increase in the problem size, the algorithm would not adapt to increase the number of cores used. This could potentially lead to computation power that may go unused in high-end GPUs that have the capacity to perform several teraflops of double-precision calculations.

### 3. Method and implementation

#### 3.1. PBM implementation

The overall population balance equation with a lumped liquid and gas coordinates can be represented as:

$$\begin{aligned} \frac{d}{dt}F(s_i, x, t) &= \mathfrak{R}_{agg}(s_i, x, t) + \mathfrak{R}_{break}(s_i, x, t) \\ &+ \dot{F}_{in}(s_i, x, t) - \dot{F}_{out}(s_i, x, t) \end{aligned} \quad (2)$$

where,  $F(s_i, x)$  represents the number of solid particles of type  $i$  being studied in each spatial compartment  $x$  of the granulator. The rate of aggregation  $\mathfrak{R}_{agg}(s_i, x)$  and the rate of breakage  $\mathfrak{R}_{break}(s_i, x)$  determines the rate at which particles density changes within different size classes. The rate of particles entering,  $\dot{F}_{in}(s_i, x)$  and exiting,  $\dot{F}_{out}(s_i, x)$ , the spatial compartment due to particle transfer also affects their number in each size class. The rate of change of internal liquid volume in each particle can be calculated as:

$$\begin{aligned} \frac{d}{dt}F(s_i, x)l(s_i, x) &= \mathfrak{R}_{liq,agg}(s_i, x) + \mathfrak{R}_{liq,break}(s_i, x) + \dot{F}_{in}(s_i, x)l_{in}(s_i, x) \\ &- \dot{F}_{out}(s_i, x)l_{out}(s_i, x) + F(s_i, x)\dot{l}_{add}(s_i, x) \end{aligned} \quad (3)$$

where,  $l(s_i, x)$  is the internal liquid volume in each particle with  $s_i$  as the solid volume for solid type 1 in the spatial compartment  $x$ .  $\mathfrak{R}_{liq,agg}(s_i, x)$  and  $\mathfrak{R}_{liq,break}(s_i, x)$  are the rates at which liquid is transferred between size classes due to aggregation and breakage respectively.  $l_{in}(s_i, x)$  and  $l_{out}(s_i, x)$  are the

internal liquid volumes of the particles entering and exiting the spatial compartment.  $\dot{l}_{add}(s_i, x)$  is the rate of volume of liquid acquired by each particle in the compartment at every time step due to external liquid addition. Similarly, the rate of change of gas volume is calculated using the following equation:

$$\begin{aligned} \frac{d}{dt}F(s_i, x)g(s_i, x) = & \Re_{gas,agg}(s_i, x) + \Re_{gas,break}(s_i, x) + \dot{F}_{in}(s_i, x)g_{in}(s_i, x) \\ & - \dot{F}_{out}(s_i, x)g_{out}(s_i, x) + F(s_i, x)\dot{g}_{cons}(s_i, x) \end{aligned} \quad (4)$$

where,  $g(s_i, x)$  is the gas volume of each particle with solid volumes of  $s_i$  in the spatial compartment  $x$ .  $\Re_{gas,agg}(s_i, x)$  and  $\Re_{gas,break}(s_i, x)$  are the rates of gas transferred between size classes due to aggregation and breakage respectively.  $g_{in}(s_i, x)$  and  $g_{out}(s_i, x)$  are the gas volume of the particles entering and leaving the spatial compartment respectively.  $\dot{g}_{cons}(s_i, x)$  represents rate of the volume of gas particles formed due to process of consolidation occurring inside the system. The rate of aggregation,  $\Re_{agg}(s_i, x)$  in Equation 2 is calculated as (Chaturbedi et al., 2017):

$$\begin{aligned} \Re_{agg}(s_i, x) = & \frac{1}{2} \int_0^{s_i} \int_0^{s'_i} \beta(s'_i, s_i - s'_i, x) F(s'_i, x) F(s_i - s'_i, x) ds_i ds'_i \\ & - F(s_i, x) \int_0^{s_{max_i} - s_i} \beta(s_i, s'_i, x) F(s'_i, x) ds'_i \end{aligned} \quad (5)$$

where,  $\beta(s_i, s'_i, x)$  is the aggregation kernel and is expressed as a function of collision frequency ( $C$ ) and collision efficiency ( $\psi$ ). Further information on the model can be found in Appendix A.1.

Similarly, the breakage rate can be expressed as follows:

$$\begin{aligned} \Re_{break}(s_i, x) = & \int_0^{s_{max_i}} K_{break}(s'_1, s'_2, x) F(s'_1, s'_2, x) ds'_i \\ & - K_{break}(s_i, x) F(s_i, x) ds_i \end{aligned} \quad (6)$$

where,  $K_{break}(s_i, x)$  is the breakage kernel. The formulation for the breakage kernel is discussed in more detail in Appendix A.2.

The rate of increase of liquid volume of inside a particle,  $\dot{l}_{add}(s_i, x)$  is expressed as:

$$\dot{l}_{add}(s_i, x) = \frac{\sum_i s_i \times \dot{m}_{spray}}{m_{solid}(x)} \quad (7)$$

where,  $\sum_i(s_i)$  is the total solid volume of the particle;  $\dot{m}_{spray}$  is the rate of external liquid addition and  $m_{solid}$  is the total amount of solid present in the compartment.

Particle transfer rate,  $\dot{F}_{out}(s_i, x)$  in Equation 2 is calculated as:

$$\dot{F}_{out}(s_i, x) = \dot{F}(s_i, x) \frac{\nu_{compartment}(x) \times dt}{d_{compartment}} \quad (8)$$

where,  $\nu_{compartment}(x)$  and  $d_{compartment}$  are respectively the average velocity of particles in compartment  $x$  and the distance between the mid-points of two adjacent compartment, which is the distance particles have to travel to move to the next spatial compartment.  $dt$  is the time-step.

A finite difference method was used to solve the developed system of ordinary differential equations (ODEs) (Barrasso and Ramachandran, 2015). Euler integration was used as the numerical integration technique for its speed improvements and its minimal impact on accuracy (Barrasso et al., 2013). To avoid numerical instability due to the explicit nature of the Euler integration, Courant-Friedrichs-Lewis (CFL) condition must be satisfied (Courant et al., 1967). For our PBM model, time-step was calculated at each iteration such that, the number of particles leaving a particular bin at any time was less than the number of particles present at that time (Ramachandran and Barton, 2010).

### 3.2. MPI implementation

The message passing interface (MPI) parallel implementation of the PBM was focused towards equal distribution of the task and memory. The implementation used in this work differs from the hybrid implementation used by (Bettencourt et al., 2017) and Sampat et al. (2018) as only MPI was used to

parallelize the code. It was pointed by Sampat et al. (2018) that open multi-processing(OMP) does not provide significant speed improvements due to limitations with usage of dynamic vectors which are essential for such a system. Thus, the OMP implementation was avoided which also meant that lesser number of cores would be required. The focus of this study to localize the computation power rather than depend on supercomputers/clusters. A pseudo code has been presented in Algorithm 1 to illustrate the distribution of tasks. For each time step, a MPI process is responsible for a certain section of the problem, usually a spatial chunk inside the geometry (also referred to as compartment).

Simulations for this study were performed on a computer with an Intel Core i7-7700K processor clocked at 4.2GHz and 32 GB of RAM. For maximum performance while data reading and writing a SSD was used. GCC version 7.4 with openMPI 2.0 was used to compile the parallel C++ code.

---

**Algorithm 1** CPU-based Parallel Population Balance Model

---

```

1: procedure PBM( $N_{Comp}, N_{MPI}$ )  $\triangleright N_{Comp}$  is the number of compartments
2:   Divide  $N_{Comp}$  in  $N_{MPI}$ 
3:   while  $t < t_{final}$  do
4:     for  $\forall n_{Comp}$  in 1 MPI process do
5:       Calculate  $\mathcal{R}_{aggregation}$  for solid bins  $s_1, s_2$ 
6:       Calculate  $\mathcal{R}_{breakage}$  for solid bins  $s_1, s_2$ 
7:       Calculate  $n_{particles}$  using Euler's method
8:     end for
9:     Collect  $n_{particles}$  from  $N_{MPI}$   $\triangleright$  Master process collects all data
10:    Calculate  $timestep$  using CFL condition
11:     $t_{new} = t + timestep$ 
12:  end while
13: end procedure

```

---

### 3.3. GPU implementation

NVIDIA's CUDA toolkit extends the C language such that user defined functions called kernels can be created to be run on the GPU. These kernels can be executed several number of times in parallel using large number of threads. A thread is a sequence of programmed instructions that can be managed by the computer's scheduler. A kernel depending upon the dimensions of the data can

execute instructions in 1-D, 2-D or 3-D thread blocks. The kernel can also launch multiple thread blocks at once, thus increasing the number of parallel process executions known as grid. Similar to a thread block, a grid can be up to 3-D depending upon the data under study. The code execution was split between the CPU (also called host) and GPU (also called device). Time sensitive calculations and calculations that require collection of data from the GPU were handled on the CPU with a single core, while the more computationally intensive tasks were distributed on to the GPU using kernels. Similar to the parallelization using MPI in CPUs, the geometry was split into multiple compartments. These compartments determined the number of blocks inside each GPU kernel. The number of solids used helped formed the threads in each of these blocks. The work flow of the execution can be found in Figure 2. The arrows in Figure 2 indicate the interchange of data between the CPU memory to the GPU memory.

---

**Algorithm 2** GPU-based Parallel Population Balance Model

---

```

1: procedure PBM( $N_{Comp}$ )  $\triangleright N_{Comp}$  is the number of compartments
2:   Copy initial variables from CPU memory to GPU memory
3:   GPU initial calculation kernel call from CPU  $\triangleright$  Performed on GPU
4:   Divide  $N_{Comp}$  in  $N_{blocks}$ 
5:   Copy back initial values to the CPU RAM
6:   while  $t < t_{final}$  do
7:     Copy time-dependent process variables from CPU to GPU RAM
8:     GPU aggregation rate kernel call from CPU
9:     Calculate  $\mathcal{R}_{agg}$  for solid bins  $s_1, s_2$   $\triangleright$  Performed on GPU
10:    GPU aggregation rate kernel call from CPU
11:    Calculate  $\mathcal{R}_{breakage}$  for solid bins  $s_1, s_2$   $\triangleright$  Performed on GPU
12:    Calculate  $n_{particles}$  using Euler's method
13:    Copy back process rate data back to the CPU RAM
14:    Calculate  $timestep$  using  $CFL$  condition  $\triangleright$  Performed on CPU
15:     $t_{new} = t + timestep$ 
16:  end while
17:  Clear GPU memory
18: end procedure

```

---

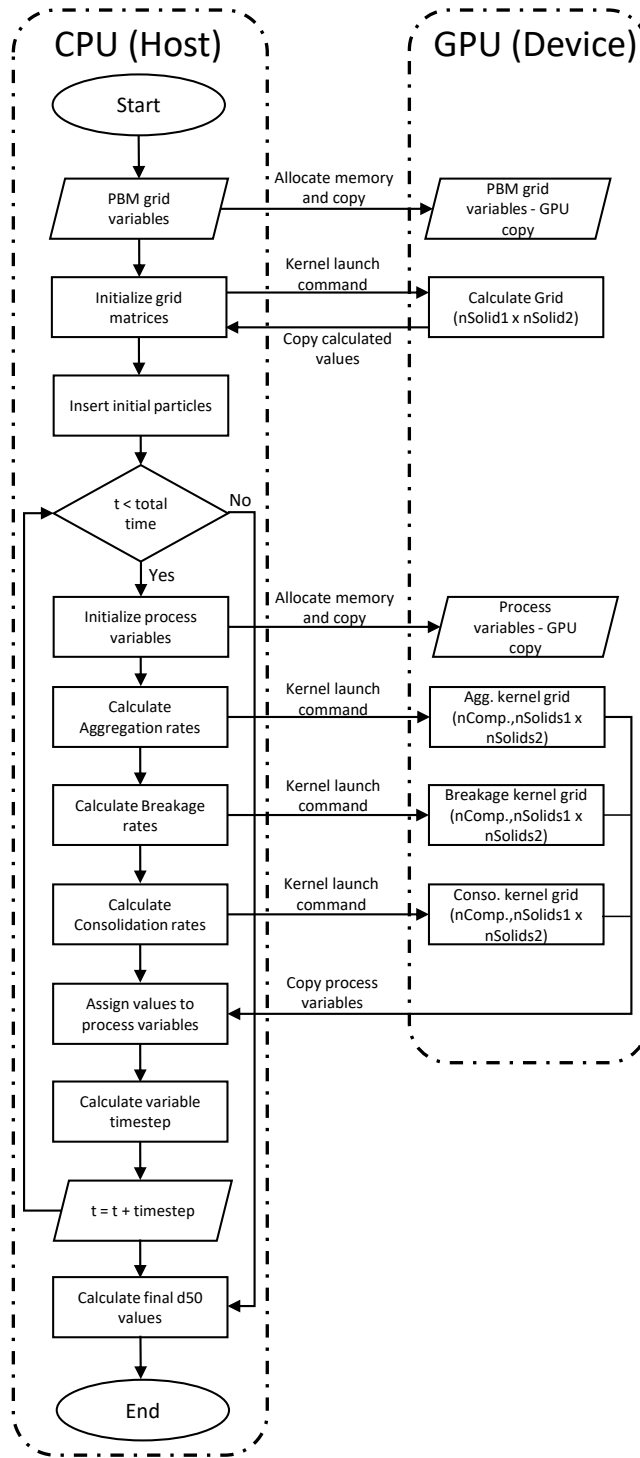


Figure 2: Workflow of the GPU code indicating data transfers and execution timeline of the code.

## 4. Results and discussions

### 4.1. Performance metrics

Parallel efficiency of an algorithm can be tested either by strong scaling, where the problem size remains the same and number of processing elements are increased or by weak scaling, where the number of processing elements remain the same and the problem size is increased. In this study, the number of processing units were limited due to the architecture of the GPU and the CUDA C++ code developed did not utilize more than one GPU during execution. Thus, a weak scaling approach was preferred in such a scenario. The parallel performance of a code is usually measured in terms of on ratio of time taken to solve the run the simulations on one core to the time taken to run the simulation on N cores. It is depicted in Equation 9, where  $t_1$  is the time taken to the run the problem on one core where as  $t_N$  is the time taken to run the problem on N cores.

$$Speedup = \frac{t_1}{t_N} \quad (9)$$

The problem size was varied by increasing the number of compartments inside the PBM. This in turn increased the total number of calculations performed without increasing the amount of work that needed to be performed by each processing unit (core). The number of cores required during the simulation on the GPU was determined by the product of the number of compartments and the number of solid bins for each solid type used. Thus, the number of GPU cores used in this study varied from 256 cores for 1 compartment to 8192 for 32 compartments.

Parallelization of code can lead to deviation in calculations from the serial execution. This discrepancy in calculation can be attributed to the difference in the precision of the cores of the CPU or GPU being used. This change in precision can lead to a small error in one timestep, which can percolate and produce results that can drastically different from the serial code. To check

the validity of the presented algorithm a relative least sum of squares error was calculated as shown in Equation 10:

$$rSSE = \sum_{i=0}^{t_{end}} \frac{(d_{50_{serial_i}} - d_{50_{parallel_i}})^2}{N} \quad (10)$$

The  $rSSE$  was calculated for the algorithm considering the single core MPI solution as the base and determining the error values wrt to this simulation. The error for each of the case can be found in Figure ???. The error for the GPU simulations varied from 0.45% to about 6% which is an acceptable range. Thus, making the solutions from the algorithm quicker with a very small loss in accuracy. There was an increase in the error with the increase in number of compartments, which could be due to the increase in the number of data points compared.

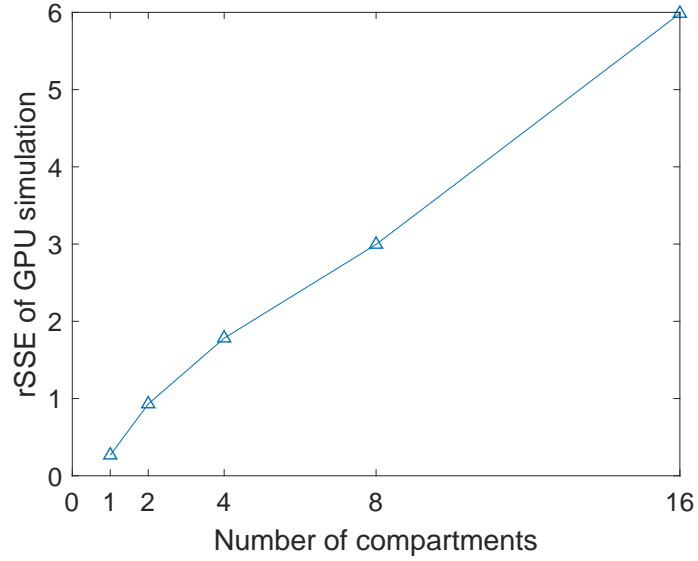


Figure 3: Relative sum of square error observed for the GPU simulations compared to the serial simulation



#### *4.2. Algorithm performance on a desktop GPU*

The desktop configuration used for these studies comprised of a Intel *i7* – 7700K CPU clocked at 4.5 GHz with 32 GB DDR4 RAM and a NVIDIA Quadro P4000 GPU. The NVIDIA Quadro GPU used had 1792 CUDA cores with 8 GB of GDDR5 RAM. CUDA version 9.0 paired with GCC 7.3 was used to the run the desktop GPU simulations with Ubuntu 18.04 operating system (OS).

The number of solid bins for the 2 different types of solids used were 16, this meant that there was a maximum of 65,536 calculations that needed to be performed for each time step for each compartment. While, the number of calculations increased to over 2 million per time step when the number of compartments was increased to 32. Each GPU consisted of several streaming multi-processors(SM) which help distribute the problem to the 1792 cores inside the GPU. Once the calculations pass from the CPU to the GPU, the SMs take over and allocate work to the GPU in blocks of 32 threads each. SMs divide these calculations in blocks of 32 threads and send it to the cores for calculations which accounts for some overheard time during the simulation. This overhead is present for each timestep, which can be compensated by the number of calculations running in parallel on the GPU.

The PBM simulation was run for 90 seconds which included 45 seconds of mixing and 45 seconds of liquid addition. The algorithm performance was tested by weak scaling the problem by changing the number of compartments from 1 and doubling them in each simulation until the number of compartments reached 32. Figure 4 shows the time taken these simulation. It can be observed that the amount of time taken for the simulation remains almost constant till the number of compartments reaches 8, followed by an increase in the time taken as the number of compartment are increased further to 32. According to parallelization procedure used the cores utilized to run the code is directly proportional to the number of compartments and the number of solid bins present in the problem. The constant time is a result of the problem size being smaller than the Quadro P4000's 1792 CUDA core, i.e. the algorithm was not able to utilize all the CUDA cores till the compartment number was 8. Since the algorithm uses about

256 cores to simulate each compartment, the cores would not suffice once the compartment number reaches 16 and the SMs would have to wait to distribute the calculation to cores once initially allocated calculations are finished. This wait time leads to the increase in the time of simulations as seen in the case for 16 and 32 compartments.

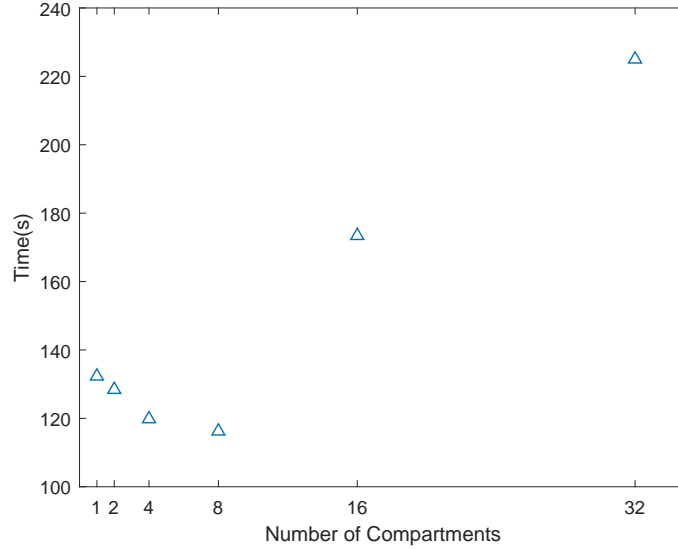


Figure 4: Time taken to complete 90 seconds of PBM simulation with varying number of compartments on NVIDIA Quadro P4000 GPU

The above argument was supported by the profile of the code that was obtained using NVIDIA's inbuilt code profiler *nvprof*. Code profiling is an important step in algorithm development. The profiler results can be varied based on the options chosen to obtain the parameters being studied. In this case, API calls and GPU activity were exported to understand the performance bottlenecks and those sections of the code were rectified to improve the speed of the algorithm. The profiler was executed for each case and it was observed that aggregation kernel calculations took the most time for execution followed by the breakage kernel. Consolidation kernel and other calculations comprised

of less than 1% of execution time. The other parameter studied was the number of times each API was called by the code and the time spent. API calls included the synchronization of threads working inside the GPU device, memory allocation for arrays, etc. Each thread inside the GPU operates independently, thus all threads may not be at the same section of the code at a given moment of time, thus some threads may finish calculations before others. The time taken to synchronize these threads required the most amount of time during the execution of the code. When such an API is called by the code further execution of the code is paused until all the threads of the GPU are in the same line of code. This accounted for 99% of the total other API call time. This indicated that there were not many places where the code could have been optimized further since synchronization statements were only added before calculations where complete array of data was required. If further reduction in these statements was undertaken, it would lead to data loss and possibly incorrect final calculated particle size distribution. A comparison of times taken by each process in the simulations is shown in Figure 5. A similar distribution of times was observed for all simulations on the GPU.

#### *4.3. Performance on GPUs compared to CPUs*

The NVIDIA Quadro P4000 GPU used had its cores at a base clock speed of 1202 MHz while the CPU cores had a base clock of 4000 MHz. The algorithm used to parallelize on the GPU did not permit the use of only one core of the GPU for simulation. Thus, a single MPI core CPU simulation was used as the baseline for all comparisons. Theoretically, it would take longer on a single core of the GPU to run a similar simulation than on a single GPU core

The CPU version of the parallel PBM was run on the desktop with the aforementioned configuration. This meant that the number of MPI cores available for the simulations was limited to 4. Weak scaling of the problem by changing the number of compartments was performed for this study. Figure 6 shows a comparison of the times taken by the simulation to run on 1, 2 & 4 MPI cores. The times indicate that with the increase in the number of cores the model took

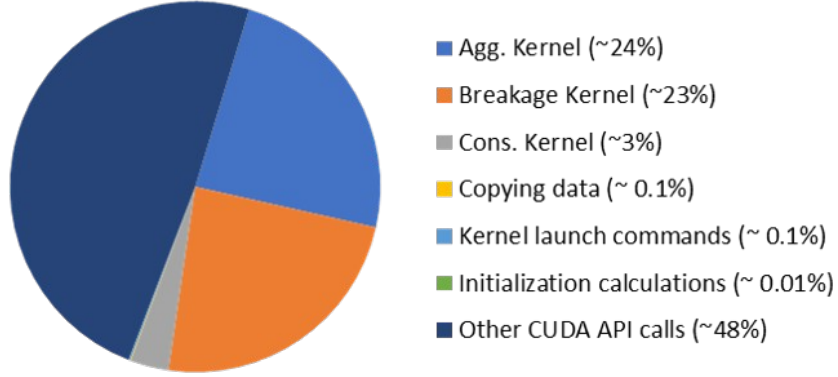


Figure 5: Distribution of times taken by different processes inside the GPU-parallelized PBM

less time to complete calculations for the same number of compartments. It can also be seen that for the same number of MPI cores used in a soft scaling the amount of time increases with increase in the number of cores. This increase can be attributed to the increase in the number of calculations with addition of new compartments. There is a plateau in the times when 2 and 4 MPI cores were used for 8 and 16 compartments respectively. When a further analysis of the rates was for each compartment was undertaken it was observed that till the particles did not reach the last few compartments of the granulator no aggregation or breakage occurred in those spatial sections thus reducing the compute time and leading to similar simulation times.

Speedup is important to understand the scalability and parallel performance of a code. Speedup for a code is directly proportional to the number of cores used for a simulation. In Figure 7, the speedup increases with the increase in the number of MPI cores. This increase in speedup can be attributed to the increase in the computation power. One unusual trend observed in the case of 8, 16 and 32 number of compartments for both 2 MPI and 4 MPI core simulations,

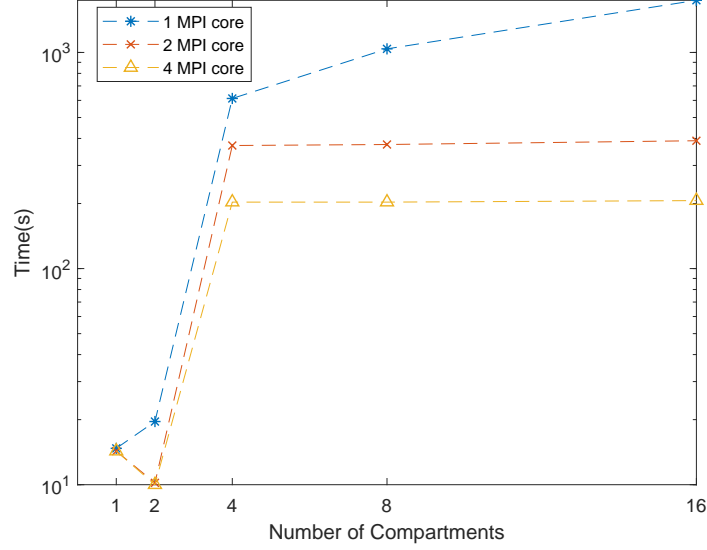


Figure 6: Time taken to complete 90 seconds of PBM simulation with varying number of compartments on desktop CPU with varying number of MPI cores

the speedup is higher than 2 and 4 respectively. This phenomena is known as super linear speedup which occurs when the speedup is greater than the number of cores used. In rare cases like these speedup increases due to increase in cache memory and random access memory(RAM) available (Benzi and Damodaran, 2009). The simulations with the GPU parallel code showed an overall increase in the speedup as the number of compartments as seen in Figure 7. The speedup was low for compartment numbers 1 and 2 since the amount of time spent in communication in between the CPU and the GPU as well as the time taken by the thread synchronization in the GPU to had a larger contribution to the simulation time. The increase in number of compartments diminished this communication time effect as amount of calculations is significantly higher. The highest speedup achieved for a GPU simulation was about 12.3, which means it took 12 times less time than a serial CPU computation.

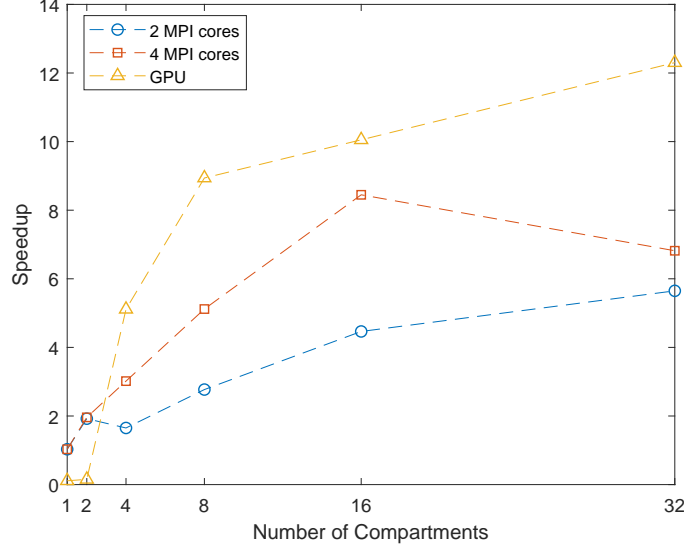


Figure 7: Comparing speedup for CPU and GPU simulations to respective serial simulations

#### 4.4. Server level GPU code performance

The high performance computing (HPC) device used to run the parallel PBM GPU code was present at Rutgers at the School of Engineering (SoE). The SoE HPC cluster was equipped with a NVIDIA Kepler K20 GPU. This GPU contains 2496 CUDA cores which have a base clock of 706 MHz and only 5 GB of GDDR5 of memory. The Kepler series GPUs were a couple of generations older than the Pascal generation Quadro P400 used in desktop simulation studies. The clock speed and memory of the desktop GPU were higher than the one present on the HPC.

Time taken to complete the 90 second PBM simulation on the HPC's GPU are shown in Figure 8a. The time taken to run the PBM initially remains constant upto 16 compartments, but a large increase in the time is observed for the simulation with 32 compartments. This increase could be attributed to the saturation of CUDA cores of the GPU and that SMs had to wait for the previous calculations to complete before the threads were assigned the remnant of calculations. This trend is similar to the trend observed in Figure 7 for the

desktop GPU. A serial simulation was performed on the HPC and used as a baseline for speedup calculations. Speedup from the server GPUs are shown in Figure 8b. The increase in the speed of the simulation for these studies is lower than the desktop studies which could be directly connected to clock speeds of the CUDA cores. The server GPU cores were clocked at a lower frequency which meant the rate of calculations would decrease. One other reason for reduced speedup could be the older architecture of Kepler GPU which are slower in floating point calculations (NVIDIA Corporation, 2016).

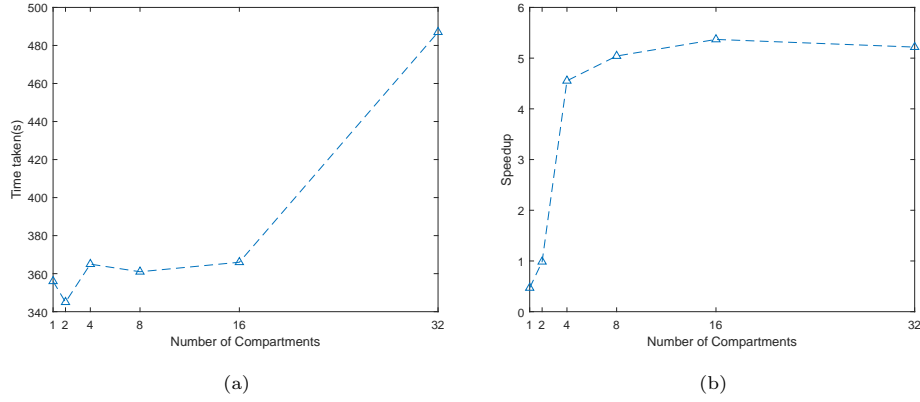


Figure 8: (a)Time taken to run 90s PBM simulation on HPC's Kepler K20 GPU (b)Speedup achieved for the GPU simulation over the serial simulation on the HPC device

## 5. Conclusions

In the presented study, a PBM was developed using NVIDIA's CUDA C/C++ language to run in parallel on a GPU. The time of the simulations on the GPU were compared to CPUs. In cases with larger problem size, it was observed that GPU simulations were faster than CPU simulations and there was minimal loss in accuracy. It can be observed that GPUs are more efficient when the complexity of the problem is high in terms of compartments, which is a commonly observed in general application of PBMs. The adaptive structure of the algorithm enabled the simulation to use varying number of GPU cores to parallelize the PBM simulations. The GPU architecture also plays a major role

in the simulation time. This work also highlighted that a desktop PC could be sufficient for a computationally intensive simulation instead of a utilizing a supercomputer or cluster. This work can be extended in the future by testing it on newer GPUs from NVIDIA like the Volta and Turing platforms, which are more optimized for float point calculations than the Pascal platform GPU used in this study. The presented algorithm can also be improved further by eliminating loops inside the kernels using dynamic parallelization supported by newer versions of CUDA.

## Software and Data

Source code and input scripts for reproduction of the experiments can be found at:

<https://github.com/csampat/pbmOnGPUs>

## Appendix A. PBM model

### *Appendix A.1. Aggregation Kernel*

The aggregation kernel used in this work was formulated as in Barrasso et al. (2015):

$$\beta(s_i, s'_i, x) = C(s_i, s'_i, x)\psi(s_i, s'_i, x) \quad (\text{A.1})$$

The collision frequency of the solid particles was evaluated from the existing DEM data from (Sampat et al., 2018). To facilitate this study, it was assumed that the collision frequency was independent of the liquid particles present in the system.

The collision efficiency  $\psi$  was estimated based on Stokes , which states that a collision is successful when the Stokes number  $St_v$  associated with the collision is lesser than the critical Stokes number  $St_v^*$  for the particles. These number are calculated as follows:



$$St_v = \frac{8\tilde{m}U}{3\pi\tilde{d}^2\mu} \quad (\text{A.2})$$

$$St_v^* = \left(1 + \frac{1}{e}\right) \log\left(\frac{h}{h_a}\right) \quad (\text{A.3})$$

Here,  $\tilde{m}$  &  $\tilde{d}$  represent the harmonic mean of the masses and diameters of the particles respectively.  $U$  is the collision velocity,  $\mu$  is the viscosity of the system and  $e$  is the coefficient of restitution. The thickness of the liquid on the surface of the particle  $h$  and the height of surface asperities  $h_a$  were obtained from (Barrasso et al., 2015).  $U_{critical}$  is defined as the ratio of the critical Stokes number to the Stokes number associated with the collision. The collision frequency  $\psi$  is defined as:

$$\Psi = \int_0^{U_{critical}} p(U) dU \quad (\text{A.4})$$

where it is assumed that the collision velocities follow a log normal distribution:

$$p(U) = \frac{1}{U\sqrt{2\pi}\sigma} \exp\left[-\frac{(\ln U - \mu)^2}{2\sigma^2}\right] \quad (\text{A.5})$$

#### Appendix A.2. Breakage Kernel

The breakage kernel  $K_{break}(s_i, x)$  is formulated as:

$$K_{break}(s_i, x) = C_{impact} \int_{U_{break}}^{\infty} p(U) dU \quad (\text{A.6})$$

Similar to the aggregation kernel,  $C_{impact}$  is defined as rate at which the particles impact with the geometry in the DEM simulation. Critical velocity for breakage to occur is defined as:

$$U_{break} = \frac{2St_d^{*def}}{\rho_{s_i}} \frac{9}{8} \frac{(1-\epsilon)^2}{\epsilon} \frac{9\mu}{16d_{p_i}} \quad (\text{A.7})$$

where  $2St_d^{*def}$  is defined as critical Stokes deformation number (Iveson et al., 2001) and  $d_{p_i}$  is diameter of the solid particle  $s_i$ .

## References

- Almasi, G.S., Gottlieb, A., 1989. Highly Parallel Computing. Benjamin-Cummings Publishing Co., Inc., Redwood City, CA, USA.
- Barrasso, D., Eppinger, T., Pereira, F.E., Aglave, R., Debus, K., Bermingham, S.K., Ramachandran, R., 2015. A multi-scale, mechanistic model of a wet granulation process using a novel bi-directional PBM–DEM coupling algorithm. *Chemical Engineering Science* 123, 500–513.
- Barrasso, D., Ramachandran, R., 2015. Multi-scale modeling of granulation processes: bi-directional coupling of PBM with DEM via collision frequencies. *Chemical Engineering Research and Design* 93, 304–317.
- Barrasso, D., Walia, S., Ramachandran, R., 2013. Multi-component population balance modeling of continuous granulation processes: a parametric study and comparison with experimental trends. *Powder technology* 241, 85–97.
- Benzi, J., Damodaran, M., 2009. Parallel three dimensional direct simulation monte carlo for simulating micro flows, in: *Parallel computational fluid dynamics 2007: implementations and experiences on large scale and grid computing*. Springer Science & Business Media. volume 67, pp. 91–98.
- Bettencourt, F.E., Chaturbedi, A., Ramachandran, R., 2017. Parallelization methods for efficient simulation of high dimensional population balance models of granulation. *Computers and Chemical Engineering* 107, 158–170.
- Chaturbedi, A., Bandi, C.K., Reddy, D., Pandey, P., Narang, A., Bindra, D., Tao, L., Zhao, J., Li, J., Hussain, M., Ramachandran, R., 2017. Compartment based population balance model development of a high shear wet granulation process via dry and wet binder addition. *Chemical Engineering Research and Design* 123, 187–200.
- Courant, R., Friedrichs, K., Lewy, H., 1967. On the partial difference equations of mathematical physics. *IBM journal of Research and Development* 11, 215–234.

- Gunawan, R., Fusman, I., Braatz, R.D., 2008. Parallel high-resolution finite volume simulation of particulate processes. *AIChE journal* 54, 1449–1458.
- Iveson, S.M., Litster, J.D., Hapgood, K., Ennis, B.J., 2001. Nucleation, growth and breakage phenomena in agitated wet granulation processes: a review. *Powder technology* 117, 3–39.
- Kandrot, E., Sanders, J., 2011. *Cuda By Example: An Introduction To General-Purpose Gpu Programming*. Addison-Wesley Professional. URL: <https://books.google.com/books?id=6mwanQAACAAJ>.
- Keckler, S.W., Dally, W.J., Khailany, B., Garland, M., Glasco, D., 2011. Gpus and the future of parallel computing. *IEEE Micro* 31, 7–17.
- NVIDIA Corporation, 2012. *NVIDIA CUDA C Programming Guide*. version 4.2 ed. NVIDIA Corporation, 2701 San Tomas Expressway, Santa Clara, CA 95050. URL: [https://developer.download.nvidia.com/compute/DevZone/docs/html/C/doc/CUDA\\_C\\_Programming\\_Guide.pdf](https://developer.download.nvidia.com/compute/DevZone/docs/html/C/doc/CUDA_C_Programming_Guide.pdf).
- NVIDIA Corporation, 2016. *NVIDIA Tesla P100: The most advanced datacenter accelerator ever built*. Technical Report.
- Prakash, A.V., Chaudhury, A., Ramachandran, R., 2013. Parallel simulation of population balance model-based particulate processes using multicore CPUs and GPUs. *Modelling and Simulation in Engineering* 2013, 2.
- Ramachandran, R., Barton, P.I., 2010. Effective parameter estimation within a multi-dimensional population balance model framework. *Chemical Engineering Science* 65, 4884–4893.
- Ramkrishna, D., Singh, M.R., 2014. Population balance modeling: current status and future prospects. *Annual review of chemical and biomolecular engineering* 5, 123–146.
- Sampat, C., Bettencourt, F., Baranwal, Y., Paraskevagos, I., Chaturbedi, A., Karkala, S., Jha, S., Ramachandran, R., Ierapetritou, M., 2018. A parallel

- unidirectional coupled dem-pbm model for the efficient simulation of computationally intensive particulate process systems. *Computers and Chemical Engineering* 119, 128 – 142.
- Santos, F.P., Senocak, I., Favero, J.L., Lage, P.L., 2013. Solution of the population balance equation using parallel adaptive cubature on gpus. *Computers and Chemical Engineering* 55, 61–70.
- Sen, M., Barrasso, D., Singh, R., Ramachandran, R., 2014. A multi-scale hybrid CFD-DEM-PBM description of a fluid-bed granulation process. *Processes* 2, 89–111.
- Seville, J., Tüzün, U., Clift, R., 2012. *Processing of particulate solids*. volume 9. Springer Science & Business Media.
- Shi, Y., Green, W.H., Wong, H.W., Oluwole, O.O., 2012. Accelerating multi-dimensional combustion simulations using gpu and hybrid explicit/implicit ode integration. *Combustion and Flame* 159, 2388 – 2397.
- Solihin, Y., 2015. *Fundamentals of parallel multicore architecture*. Chapman and Hall/CRC, New York, USA.
- Szilágyi, B., Nagy, Z.K., 2016. Graphical processing unit (gpu) acceleration for numerical solution of population balance models using high resolution finite volume algorithm. *Computers and Chemical Engineering* 91, 167 – 181.
- MathWorks<sup>TM</sup> Documentation, 2017b. Parallel computing toolbox - MATLAB<sup>®</sup>. <https://www.mathworks.com/products/parallel-computing.html>.
- Xu, Z., Zhao, H., Zheng, C., 2015. Accelerating population balance-monte carlo simulation for coagulation dynamics from the markov jump model, stochastic algorithm and gpu parallel computing. *Journal of Computational Physics* 281, 844 – 863.