

Splay Thread Cooperation on Ray Tracing as a Load Balancing Technique in Speculative Parallelism and GPGPU

Suma Shivaraju¹ and Gopalan Pudur²

¹Department of Computer Scienc, Bharathiar University, India

²Department of Computer Applications, National Institute of Technology, India

Abstract: *The introduction of the speculative parallelism into any models can improve the performance and provides significant benefits and increases the Instruction level parallelism (ILP) and Thread level Parallelism (TLP). A General Purpose Graphics Processing Unit (GPGPU) is the future computing technology working with both Graphics Processing Unit (CPU) and GPU to solve many real-world problems not only the graphics problems but also the general purpose applications. As GPU uses data parallelism tasks, the dynamic memory creation and the splay trees which are self adjusting allows for the increase in throughput and load balancing. The frequently used nodes near to the root are an advantage for finding locality of threads as well as for caching and garbage collection. A technique used to render and to study complex scenes into images and to render color, intensity of pixels, distance between pixels is referred as Ray tracing. Multithreading is a promising technique which increases the performance of the computer systems by increasing the instruction level parallelism and thread level speculation. In this paper a new technique is proposed for workload balancing on the Graphics processors and CPU that can be implemented on the graphics processors along with the CPU which provides the optimal result with the speculation techniques and Lorentz Transformation, which is used to determine color and brightness of the ray which are refracted or reflected and also the relative distance between the thread spawning which results in time dilation and contraction. A Graphics Processing Unit Ocelot (GPUOCELOT) is a compilation framework, a simulator used for the execution of the programs which has resulted in the increase in the performance of the instructions which uses the amortized cost.*

Keywords: *Load balancing, graphics processors, splay trees, optimization, instruction level parallelism, thread level speculation, amortized cost, speculative multithreading, ray tracing, lorentz transformation.*

Received September 26, 2014; accepted February 10, 2015

1. Introduction

Increase in the clock rate, exploiting the Instruction Level Parallelism (ILP) through pipelining, out-of-order execution and multithreading are the techniques which results in reducing the dependencies and increase in the efficiency and performance of the microprocessors. Speculation techniques are implemented in the modern micro architectures to improve performance and to utilize the computer resources as well as the scheduling of the instructions based on the optimization techniques [1, 6]. The general purpose computing on graphics processing units handles the computations executed on the graphics processor or in kernels. The graphics processing units are computing on the data parallelism with fast access to memory and high throughput on parallel tasks. Since CPU's performs well on task parallelism and GPU's performs well on data parallelism. The memory systems are designed to stream data when the pattern can be accessed linearly and that can be prefetched [5, 14]. The combined technology of speculative execution on graphics processors still increases the performance and efficiency of the execution of the programs.

Jaikrishnan Menon, Marc de Kruijf, Karthikeyan Sankaralingam proposed iGraphics Processing Unit (iGPU) architecture executing the idempotent regions Load balancing is a method for distributing the workloads across the processors [9]. The sharing of the workload across CPU and GPU needs to be balanced as the architectural designs of the CPU and GPU differs. Christian Lauterback, Qi Mo, Dinesh Manocha proposed explicit balancing on the GPU's and dynamic work Distribution [8]. In this model the hybrid architecture of CPU and GPU together holds threads of processing units and CPU does the task parallelism and GPU does the data parallelism [3, 13]. Distributing work across GPU and CPU processor's is an intellectual task so that the performance, scalability and tolerance factors of the processor is reduced instead of increase in the processor performance in ILP and Thread Level Parallelism (TLP). Many Researchers' worked on the parallel ray tracing on the hybrid architectures with thread level speculation. Architectures are proposed to increase the performance through dynamic load balancing, reducing communication time, and reducing memory latencies, predicting value reuse, speculative computation reuse

techniques between the CPU and GPU. But the overall speedup of computation, performance considerably is the linear speedup and the overheads incurred, granularity of the code presented a major focus of issue. GPU supports fine grained parallelism over the data but the CPU works on the Coarse grained parallelism as the tasks. Transferring of the data and control values from CPU to GPU memory is challenging. Though with the speculative parallelization the state of the system is not altered, if the speculation fails, the normal execution of the task is carried out by the nonspeculative thread to recover the system state to safe. To efficiently distribute the workload across the graphics processors a new technique is proposed which is a hybrid of CPU and GPU. A self adjusting splay tree provides load balancing on the GPU.

2. Splay Trees

Splay trees are the data structures where the binary search trees are self adjusting. Whenever we access a node of the tree, whether for retrieval, insertion or deletion a newly accessed node becomes the root of the modified tree. In splay trees the nodes which are frequently accessed move towards the root whereas the infrequent ones moves far away from the root. The different splaying steps which use the bottom-up approach are similar to the rebalancing operations on Adelson-Velskii and Landis (AVL) and Red-Black trees. The cases for the splaying are ZIG-ZIG and ZIG-ZAG. In Zig-Zig, the splaying a node X which is a inside grandchild and it is a double rotation. In Zig-Zag, the splaying a node X which is a outside grandchild Where the grandchild is pulled to become the root of the sub tree [12].

2.1. Speculative Work Load Balancing

In this paper an algorithm is proposed where splaying steps happens for the threads in a block. The architecture of the speculatively parallelized graphics processors where the computations which are massively parallel in nature can be executed on both CPU and GPU, where the nonspeculative thread starts executing tasks on the CPU. The speculative threads execute tasks on the GPU. Another technique is proposed in this paper to balance the work load of the threads on the graphic processors.

The locality of reference to the code blocks on the graphic processors is scheduled using an algorithm where each thread of the block of a graphic processor is represented in the form of the nodes on a binary search tree. The splaying steps are applied on the tree of nodes such that the threads which are active and processing the computations are rotated near to the root node of the tree, which are always accessed frequently. The different types of threads are running to balance the work load i.e., some threads are in

- a. Active state.
- b. Passive state.
- c. Squashed state.
- d. Misspeculation state.

Some threads are at the initialization state. Some threads are waiting for some event.

The threads in the block are checked for the status. The different states the nodes are determined as nodes which are active, passive, nodes under initialization and misspeculation nodes.

2.2. Splaying on GPU's

On the graphic processor, the splaying operation is carried out, where root node is the nonspeculative thread and other child nodes are speculative threads. We apply splaying steps on these threads and determine whether the node is in what state i.e., threads on the GPU. Initially the current thread accessed is its states are determined based on the priority.

If the thread is active, its priority is high and that thread is right rotated and made as the child node of the nonspeculative thread. Successive access to a node newly increases the efficiency and locality of referencing to the nodes becomes easier. If the status of the node is determined, based on the status of the nodes the active nodes are moved towards the node and the nodes which are misspeculated/passive are moved away from the root node. Only active nodes become the children of the nonspeculative thread such that the load can be balanced between active nodes. Other nodes are squashed or eliminated from the binary search tree.

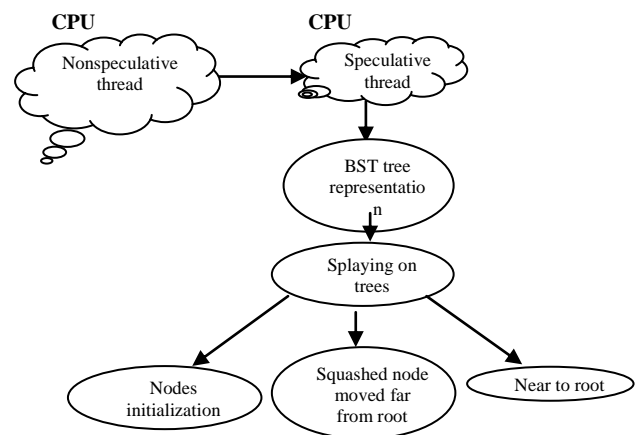


Figure 1. Diagrammatic representation of splaying operation on Threads on GPU.

Figure 1 describes the steps such that the state of the thread is determined based its locality of reference after splaying.

2.3. Design Methodology

The general algorithm for the splaying data structure is

Algorithm 1: The general algorithm for the splaying data structure.

Algorithm SplaySpec ()

This algorithm works on the CPU.

Step 1. Initialize the structure which holds elements
No, left, right and the size of the nodes on the tree representation where each node represents the threads on the GPU.

Step 2. Function to create treeofthreads ()

Step 3. If tree is empty,

Step 4. Create the threads of tree.

Step 5. Print the tree

Step 6. Print the each node of the thread of trees.

Algorithm 2: The Algorithm for the splaying operations.

```
#include<stdio.h>
```

```
#include<cuda.h>
```

```
# define NUM_BLOCKS 20;
```

```
Step 1. Structure for the tree.
```

```
Typedef struct splay_node splay;
```

```
Struct splay node {
```

```
    Splay *r,*l;
```

```
    Int info;
```

```
}
```

```
_device_int * ptr [NUM_BLOCKS];
```

```
Step 2. Memory Allocation.
```

```
_global void allocate memory ()
```

```
{
```

```
    If (threadID.x==0)
```

```
    ptr [blockIdx.x]=(int *) malloc(blockDim.x*4);
```

```
    -syncthreads ();
```

```
    If (ptr[blockIdx.x]==NULL)
```

```
    Return;
```

```
    ptr[blockIdx.x][threadIdx]=0
```

```
}
```

```
_global void memalloc()
```

```
{
```

```
    Int ptr1=ptr[blockIdx.x];
```

```
    If ( ptr 1!=NULL)
```

```
    Ptr1(threadIDX.X) +=threadIDX.X
```

```
}
```

```
Step 3. Function splay tree manipulation.
```

```
Tree *splay (int i ,tree *s)
```

```
{
```

```
    Tree n,*l,*r,*y;
```

```
    If (t==NULL) return t;
```

```
    For (;;
```

```
    {
```

```
        If (i<t->info)
```

```
        { if (2t+1)==NULL) break;
```

```
        If ( i < (2t+1).info)
```

```
        { z=2t+1
```

```
        2t+1=2t+2
```

```
        2t+2=z
```

```
        t=z
```

```
        }
```

```
        If( 2t+1)==NULL ) break;
```

```
    }
```

```
    r.2t+1=t
```

```
    r=t
```

```
    t=2t+1
```

```
    } else if (i > t.info)
```

```
    { if ((2t+2)==NULL) break ;
```

```
    If (i >2t+2.info)
```

```
    { z=2t+2
```

```
    2t+2=2t+1
```

```
    2t+1=t
```

```
    t=z
```

```
    }
```

```
    If ((2t+2)==NULL) break;
```

```
    l.2t+2=t
```

```
    l=t
```

```
    t=2t+2
```

```
    Else
```

```
    {break ;}
```

```
    l.2t+2=2t+1
```

```
    2t+1=2t+2
```

```
    2t+2=n.2t+1
```

```
    2t+1=m.2t+2
```

```
    Return t
```

```
}
```

Initially, the memory is allocated dynamically using the functions `allocatememory ()` and `memalloc ()` functions [10]. The function `splay ()` uses this array's and does the splaying by swapping the left child nodes to the right to balance the tree. Since the implementation of linked lists concept in GPU is little difficult as the memory allocation becomes difficult, a proposal to use the dynamic array is implemented for the splaying operation. Dynamic arrays are the growable arrays consisting of variable sized data structure where the addition and deletion happens and supported by the many of the modern programming languages.

The dynamic arrays are constructed using the geometric expansion of doubling in the size of the fixed array while it is used for reversed use such that n elements take $O(n)$ time such that each insertion takes amortized constant time.

These types of dynamic arrays are supported by the GPU's also. Compute Unified Device Architecture (CUDA) supports dynamic allocation of memory for a thread or a block. The performance issue relates to the dynamic array are that it includes Locality of reference and data cache utilization and it is random access also. The dynamic arrays supports for the faster indexing also [10]. Since the Amortized cost definitely guaranteed for the worst case analysis and performance than the speculation of the states in the program structure.

2.4. Splaying Transformations

The Figure 2 below shows the transformations on the splaying rotations of the threads in zigzag and zig-zig. Consider the threads on the GPU, each of the blocks of the threads holds 8 small threads inside numbered from 0 to 7 and each thread is identified by a threadID. The computations are allocated to the block of threads instead of individual threads on the block. Each block's threads are referred by the ThreadIDs and Blocks are represented by the BlockID's. BlockID=0 BlockID=1 BlockID=20 to 7 ThreadID's.

Suppose a computation is assigned to the BlockID=0. There are 7 threads whose synchronization is carried out using splaying zig-zag or zig-zig

rotations as these threads are represented in the binary search tree.

The Zig-Zig rotation is carried out for checking the Activeness, passiveness, initialization of nodes and misspeculated nodes. Consider BlockID=0 which consists of ThreadID's from 0 to 7.

The Zig-Zig operation on this BlockID is that the thread 0, 1, 3, 5 are active, 6 is inactive, 7 is squashed and 2, 4 are misspeculated threads.

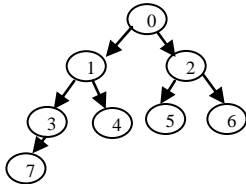


Figure 2. The binary search tree of threads of a single block.

The Figure 3 describes the tree of threads constructed from the 8 threads of the single block.

If t is a node considered for splaying operation.

Parent(t)=(t-1)/2 if t !=0.

Left child(t) =2t+1 if 2t+1 <=n.

Right child=2t+2 if 2t+2 <=n.

Left Sibling (t) = t-1 if t is even.

Right Sibling (t) =t+1 if t is odd and t+1 <=n [4, 15].

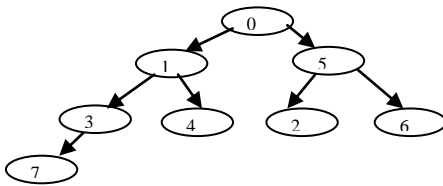


Figure 3. The tree after the splaying operation.

Since it is a self adjusting tree, the nodes which are active are moved near to the root and other nodes are moved to the end. Table 1 represent the identification of the parent and child threads on the GPU such that after splaying operation nodes are self adjusted.

Table1. Illustrates the parent and child node relationship.

	0	1	2	3	4	5	6	7	8
Parent p	-	0	0	1	1	2	2	3	3
Left Child LC	1	3	5	7	-	-	-	-	-
Right Child RC	2	4	6	8	-	-	-	-	-
Left Sibling LS	-	-	1	-	3	-	5	-	7
Right Sibling RS	-	2	-	4	-	6	-	8	-

To transfer the data to the speculative threads the table is referred to check for the parent and child threads on the GPU.

At the next level the same operation can be extended to the block level also since the blocks are numbered as the B0, B1, B2 the splaying operation can bring the blocks which are active to the nearer to the root by self adjustment that saves the memory as well as theCollection also.

3. Ray Tracing Using Splaying

In computer graphics, Ray Tracing is a technique for generating an image by tracing the path of the light through pixels in an image plane and simulating the effects of its encounters with virtual objects. In Ray Tracing problem, we send a ray from the eye/camera through each pixel on the virtual screen to compute the color of the pixel. Some light could also have been reflected or refracted by this object. The other light ray is blocked by another object. If the multiple reflections or refractions we trace recursively the reflected or refracted rays until they do not hit any object. Finally, the energy contributions of all rays used to get the color of the screen pixel.

Figure 4 describes the ray tracing technique. Ray tracing can be efficiently solved by the technique of speculative parallelization and multithreading. Speculative parallelism technique is a promising technique used in current technologies to enhance the instruction level parallelism and thread level speculation. Speculative multithreading uses speculative Architectural threads where it has two threads. A non speculative thread and a set of speculative threads. A non speculative thread starts executing the program and at the spawning point where the program or loop and instructions can be parallelized. In the Speculative Architectural Thread paradigm, the technique is to execute the different parts of the program in parallel in different sections. The execution of the program results in correct, if the values are computed by the speculative threads and the code executed by them need not be reexecuted by the main thread. Considering the Ray tracing techniques, the execution of the loop or instruction where the ray is being traced or send is traced by the nonspeculative thread. Once the ray is reflected or refracted, the point of reflection or refraction is termed as spawning point or spawning pair [15].

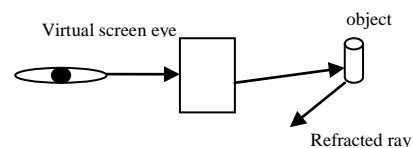


Figure 4. Ray tracing technique.

Feng *et al.* [3] have developed the augmented design for the indicating the dynamic data structures with speculative parallelism, and also have designed scheduling policies for cross-iteration dependences as well as irregular control flows.

Reinhard and Jansen [11] have discussed about data driven and demand driven tasks for good load balancing and spreading communication evenly in the network. Kobayashi *et al.* [7] proposed hierarchical multiprocessor system for dynamic load balancing with the static one.

Hence from the spawning point or spawning pair, the non speculative thread spawns the speculative thread, the speculative thread starts tracing the ray which is either reflected or refracted. All the rays are traced by the speculative threads and the energy is calculated and passed onto the nonspeculative thread. Each time a ray is refracted or reflected, a speculative thread is spawned and a ray suggests that a nonspeculative thread executes the entire program. $T(n)$ is the time taken for the sequential execution.

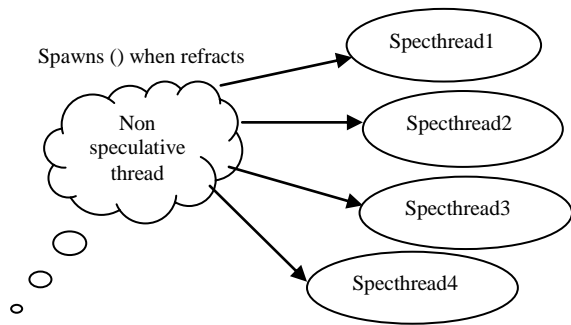


Figure 5. Process of spawning threads.

When speculative threads are tracing the path of a ray, then non speculative thread initializes all the live in variables and initial values to the consumer threads i.e., speculative threads. Figure 5 describes the process of spawning threads. After the proper live in values are propagated to the consumer threads, the consumer threads determines the energy of the rays or the color of the rays. If the value determined by the speculative thread if it is correct, the state of the thread is stored and the values are moved to the producer thread. If the speculation is wrong, the thread is squashed and the nonspeculative thread starts executing normally.

When the speculative thread is spawned by the nonspeculative thread the tracing of the refracted ray is carried out by speculative thread. Then the ray of light passes as it is not hitting the object. The eye/camera is kept at a fixed frame and the ray of the light is moving either in refracted/reflected. i.e., it is in the moving frame. The relativity between the fixed frame and the object moving frame related through the transformations called the lorentz transformation [10].

3.1. Lorentz Transformation

It is a linear transformation. This transformation measures the relative distances between the space and the time, the ordering of the events, elapsed time and preserves the space time interval between the two events [10]. In our paper we apply Lorentz transformation to determine the elapsed time, relativity between the nonspeculative and speculative threads. The time dilation when the speculative threads are spawned.

The eye/camera where the ray of light is spawned is referred as fixed frame. From the fixed frame the coordinates are taken as fixed reference values as $(x1,y1)$.

Once the ray is passed through the virtual screen, the ray is traced and if it touches any object the ray gets either refracted or reflected. i.e., the ray is in the moving frame. When an observer from the fixed frame is observing the tracing of the ray from a fixed frame, the moving frame, the time dilation happens. The coordinate of the moving frame are treated as $(x2, y2)$. $\Delta t1$ is the time taken for the ray to move from fixed frame. Using this logic we are determining the spawning time of the thread and the thread distance between the speculative thread and the non speculative thread.

When the light ray touches the object the ray gets refracted, the light ray is in the moving frame. The moving frame is moving with velocity v in the x and y directions with respect to the reference frame (fixed frame).

$$\begin{aligned} X2 &= X1 - vt/\sqrt{1-v^2/c^2} \\ Y2 &= Y1 - vt/\sqrt{1-v^2/c^2} \end{aligned} \quad (1)$$

To locate the ray from the fixed frame the reverse transformation happens

$$\begin{aligned} X2 &= X1 + vt/\sqrt{1-v^2/c^2} \\ Y2 &= Y1 + vt/\sqrt{1-v^2/c^2} \end{aligned} \quad (2)$$

The time determined as

$$t' = t - vt/c^2/\sqrt{1-v^2/c^2} \quad (3)$$

We determine the values of β and α
Where $\alpha = v/c$ $\beta = 1/\sqrt{1-v^2/c^2}$

The distance conserved under a coordinate rotation is

$$\begin{aligned} x' &= x \cos \theta + y \sin \theta \quad y' = -x \sin \theta + y \cos \theta \\ x'^2 + y'^2 &= (x \cos \theta + y \sin \theta)^2 + (-x \sin \theta + y \cos \theta)^2 \\ x'^2 + y'^2 &= x^2 \cos^2 \theta + 2xy \cos \theta \sin \theta + y^2 \sin^2 \theta + x^2 \sin^2 \theta \\ &\quad + 2xy \cos \theta \sin \theta + y^2 \cos^2 \theta \\ x'^2 + y'^2 &= x^2 (\cos^2 \theta + \sin^2 \theta) + y^2 (\sin^2 \theta + \cos^2 \theta) \\ x'^2 + y'^2 &= x^2 + y^2 \end{aligned} \quad (4)$$

If the ray of light is refracted/reflected in the reference frame the ray gets reflects or refracted.

The coordinate values are $(t1, x1, y1)$ $(t2, x2, y2)$.

The speed of the ray of light in the reference frame of the system are determined as

$$\begin{aligned} X' &= x2 - x1/t2 - t1 = \Delta x / \Delta t \\ Y' &= y2 - y1/t2 - t1 = \Delta y / \Delta t \end{aligned} \quad (5)$$

The relativity between the different speeds are

$$\begin{aligned} Ux &= \Delta x / \Delta t \\ &= ((\Delta x' + v \Delta t') / \sqrt{1-v^2/c^2}) / ((\Delta t' + v \Delta x' / c^2) / \sqrt{1-v^2/c^2}) \\ &= \Delta x' / \Delta t' + v / 1 + v \Delta x' / \Delta t' c^2 \end{aligned} \quad (6)$$

- *Speculative Time Dilation*: The time gap between the spawning of the speculative thread and the spawning pair if it happens in the same coordinate system is referred as correct time of spawning. The correct time is determined as

$$\begin{aligned} \Delta t &= t_2 - t_1 \\ &= (t_2 - v x_2 / c^2 \sqrt{1 - v^2/c^2}) - (t_1 - v x_1 / c^2 \sqrt{1 - v^2/c^2}) \\ &= t_2 - v x_2 / c^2 \sqrt{1 - v^2/c^2} - t_1 + v x_1 / c^2 \sqrt{1 - v^2/c^2} \\ &= t_2 - t_1 + \sqrt{1 - v^2/c^2} (v x_1 / c^2 - v x_2 / c^2) \\ &= \Delta t / \sqrt{1 - v^2/c^2} \end{aligned} \tag{7}$$

• *Speculative Distance Contraction:* The distance between the spawned speculative thread and the nonspeculative thread is determined as the contraction.

The correct distance between the spawned speculative thread and nonspeculative thread $L = x_2 - x_1$

$$\begin{aligned} T &= (x_2 - v t_2 / \sqrt{1 - v^2/c^2}) - (x_1 - v t_1 / \sqrt{1 - v^2/c^2}) \\ &= x_2 - x_1 / \sqrt{1 - v^2/c^2} \\ &= L / \sqrt{1 - v^2/c^2} \end{aligned} \tag{8}$$

3.2. A Regular Sturm-Liouville's Equation

The Sturm-Liouville explains a finite dimensional vector space. We consider the inner product of two vectors $X = (X_1, X_2, X_3)$ $Y = (Y_1, Y_2, Y_3)$ which inner product is $X \cdot Y = (X, Y) = X_1 Y_1 + X_2 Y_2 + X_3 Y_3$. The two nonzero vectors X and Y are said Orthogonal if $X \cdot Y = 0$. A set of nonzero vectors is said to be orthogonal if any two distinct vectors from this set are orthogonal. The Sturm-Liouville consists of the general class of boundary value problems with sets of solutions that are mutually orthogonal.

A Sturm-Liouville problem is a boundary value problem on a closed finite interval $[a, b]$ of the form

$$[p(x) y']' + [q(x) + r(x)] y = 0, a < x < b \tag{9}$$

$$c_1 y(a) + c_2 y'(a) = 0 \tag{10}$$

$$d_1 y(b) + d_2 y'(b) = 0 \tag{11}$$

Where c_1, c_2 and at least d_1, d_2 are nonzero and is a parameter.

The nonzero solutions of a Sturm-Liouville problem are called the Eigen functions of the problem and the values which consist of the nonzero solutions are referred as Eigen values.

The Eigen values and Eigen functions of the sturm-liouville problem is defined as $Y'' + \lambda y = 0$ $y(0) = y(\pi) = 0$.

This equation is similar to the Equation (8) with $p(x) = 1, g(x) = 0, r(x) = 1$.

The boundary conditions $a = 0, b = \pi$ with $c_1 = d_1 = 1, c_2 = d_2 = 0$.

Considering the three cases for the solution of the sturm-liouville solution is:

- *Case 1.* When $\lambda < 0$, so $\lambda = -\alpha^2$ where $\alpha > 0$. the equation becomes $y'' - \alpha^2 y = 0$ the solution is $y = c_1 \sinh \alpha x + c_2 \cosh \alpha x$. When $y = 0, c_2 = 0$ when $y(\pi) = 0$ the equation will be $0 = c_1 \sinh \alpha \pi$ as $\sinh x \neq 0$ so there are no nonzero solution.

- *Case 2.* When $\lambda = 0$, the solution of the equation is $y = c_1 x + c_2$, the boundary conditions $c_1, c_2 = 0$ as there exists no nonzero solution.
- *Case 3.* When $\lambda > 0, \lambda = \alpha^2$ and $\alpha > 0$ so the equation becomes $y'' + \alpha^2 y = 0$. The solution is $y = c_1 \cos \alpha x + c_2 \sin \alpha x$. When $y(0) = 0$ then $0 = c_1 \cos 0 + c_2 \sin 0$ so $y = c_2 \sin \alpha x$.

The other boundary condition is $c_2 \neq 0$ then we get $\sin \alpha \pi = 0$ it has the eigen values as $y_1 = \sin x, y_2 = \sin 2x, y_3 = \sin 3x, \dots$

According to the 3 cases, considering the values of λ , when $\lambda < 0$ and when $\lambda = 0$ there exists no solution but when $\lambda > 0$, there exists solution to the problem and eigen values also exists.

We consider the sphere equation that provides quadratic equation as $X^2 - b + \sqrt{b^2 - 4ac}/2a$ $X = -b \pm \sqrt{b^2 - 4ac}/2a$

It is similar to representing and replacing $b^2 - 4ac$ by λ and considering the cases for finding the roots of the equation. According to the sturm-liouville equation, there exists no solution when $\lambda < 0$ and $\lambda = 0$. Because when $\lambda < 0, \lambda = 0$ the Eigen values remains at nonzero solution. So the ray intersects or passes through the object when $\lambda > 0$ as the boundary values for the solution exists or the solution itself exists [2].

The ray is represented as an origin and direction i.e., vector. Origin $O = [X_0, Y_0]$ Direction $D = [X_f, Y_f]$.

The ray consists of points $R(K) = O + X_f * K$.

The ray object intersection is determined by the equation $f(O + X_f * K) = 0$.

The ray sphere intersections with center and radius is determined as $(X_1 - A)^2 + (Y_1 - B)^2 = R^2$ where A, B, C are the centres of the Sphere, R is radius and X_1, Y_1 are the points on the sphere.

The parametric equation for ray are $X = X_0 + X_f * K$ $Y = Y_0 + Y_f * K$ where X_0, Y_0 is the origin of the ray X_f, Y_f are the camera rays direction. To find the intersection, the ray equation into sphere equation is $(X_0 + X_f * K - A)^2 + (Y_0 + Y_f * K - B)^2 = R^2$ which is equal to $(X_f + Y_f)^2 * K^2 + [2[X_f * (X_1 - A) + Y_f * (Y_1 - B)]] * K + [(X_f - A)^2 + (Y_f - B)^2 - R^2] = 0$


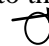
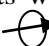
The quadratic equation is also in the form of $MX^2 + NX + C = 0$. Where

$$M = [(X_f^2 + Y_f^2)]$$

$$N = [2[X_f * (X_1 - A) + Y_f * (Y_1 - B)]]$$

$$C = [(X_1 - A)^2 + (Y_1 - B)^2 - R^2]$$

Considering the sphere equation above $b^2 - 4ac$ is discriminant. Based on the values of the discriminant,

- *Case 1.* The solution $b^2 - 4ac < 0$ it leads to imaginary value, the ray and sphere do not intersect in real plane. 
- *Case 2.* The solution $b^2 - 4ac = 0$ leads to the boundary conditions with nonzero solutions. 
- *Case 3.* If $b^2 - 4ac > 0$ exists in real roots where ray and sphere intersect with each other. 

The intersection point of a ray with the plane surface is determined as $AX+BY+D=0$.

The ray equation with the origin and the distance d is represented into the plane equation as $A(X_0+Xf*K)+B(Y_0+Yf*K)+D=0$.

$$AX_0+AXf*K+BY_0+BYf*K+D=0 \quad [2].$$

In this paper, we substitute the discriminant with the value λ . considering the 3 cases of the λ from equations the rays are cast.

Since the discriminant itself is replaced with λ , the computing one square root, 2 multiplications and one subtraction is reduced for each and every computation for computing the color of each pixel in the object which tremendously reduces the computation comparatively.

Considering the values of λ is implemented for tracing the rays. The speculative time dilation and contraction is also considered for optimizing the performance.

3.3. Speculative Parallel Ray Tracing

A nonspeculative thread is computing the execution of the ray tracing program where a ray is cast from the eye/camera through the scene which is a image file.

If the ray hits an object, the color of the object is assigned and it gets reflected/refracted. These rays are called secondary rays. The secondary rays are traced by speculative threads parallelly.

When the secondary rays are traced by speculative threads, the color of each pixel is not computed but background color is defined. Hence when the secondary rays are tracing the rays there is no need to compute the background colour.

In a speculative parallel ray tracer the algorithm is divided among number of threads on GPU. Since the numbers of processors on GPU are significantly more, can compute fine grained parallelism to compute the color of the each of the pixel on the image or object simultaneously. The data is distributed from the CPU to the GPU threads. To keep track of the timing events the start and end of the spawning of the threads on non speculative threads. Compute Unified Device Architecture (CUDAEVENT) event to keep track of the spawned threads to calculate time contraction between start and end of the event.

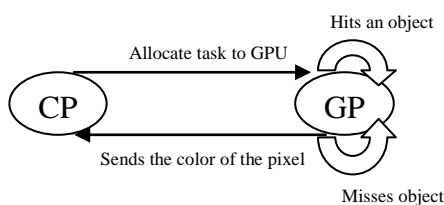


Figure 6. Allocation of task by CPU.

The steps shown in Figure 6 are as follows. The primary rays shot by the camera or the viewport are stored in the buffered array of the CUDA program and

a tree of splay trees are formed and splaying operation is carried out when the rays which are in the moving frame of the rays which hits an object are determined and after splaying operation the rays are traced by the threads on GPU. The rays which hits an object may either gets reflected or refracted such that the speculative threads spawn secondary rays are traced by them. Each time splaying operation is carried for the rays that hit an object. Each of the ray traces the path of the rays and determines the color of the pixels.

For defining the image, Open Graphics Library (OPENGL) programming and for ray tracing program, the CUDA is being used. The program is run on Graphics Processing Unit ocelot (GPUOCELOT), a parallel programming simulator for CUDA programming is used for the results and analysis.

The steps for the speculated parallel ray tracer are provided below in the Figure 7.

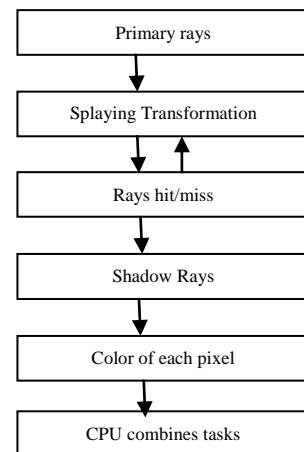


Figure 7. Steps for Speculated Parallel ray tracer.

During the execution of the ray tracing program, the load balancing of the threads is carried out through the technique of splaying operation on GPU's which is discussed in the section 1.1 of this paper. The algorithm of the speculative parallel ray tracer is provided.

Algorithm 3: The algorithm for the speculative parallel ray tracer.

Step1. The image file is created through OPENGL

Step 2. For each pixel of the object

{ color =0;

For (row=0;row< nrows;row+=blocksize

For (col=0;col<ncols; col+=blocksize

{

Step 3. Construct the primary rays.

Step 4. Call splay transform function.

Step5. The nonspeculative thread spawns speculative threads.

Step 6. The rays intersects the object called object intersection.

Step 7. The rays from the camera either hit or miss the objects.

Step 8. The CUDAEVENT records the speculative time dilation and contraction.

Step 9. If the rays hit an object, secondary rays are formed and they either reflect or refract.

Step 10. The GPU processors collect the color of each pixel and group them, finally the CPU determines the colors of the pixels.

}

Algorithm 4: Code for intersection using strum-liouville Equation.

Step 1. Scene is designed using OPENGL with objects, the camera and the light sources are set up.

Step 2. Ray tracer function.

```

__global__ void raytracing (start, direct)
{CUDA_EVENT_T begin, end;
  Float M,N,C,D;
  CUDA_EVENT_CREATE (begin)
  CUDA_EVENT_RECORD (begin,0)
  M= Xf2+Yf2
  N=2(Xf[XI-A]+Yf[yI-B])
  C=(XI-A)2+(YI-B)2-(r*r)
  D=N*N-4*M*C
  If (D < 0)
  Return false
  Dsqr =sqrt(D)
  Dsqr=λ
  If (λ>0)
  S0=-(N-λ)/2.0
  Else
  S1=- (N+λ)/2.0
  If (λ==0 || λ<0)
  Return false
  CUDA_EVENT_SYNCHRONIZE (end)
  Float executed time
  CUDA_EVENT_ELAPSED_TIME (and executed time)
}
    
```

4. Experimental Setup

The algorithm is run on the Inspiron5050 using the simulator GPUOCELOT simulator. The efficiency of the algorithm is determined. The amortized complexity of the algorithm is determined as follows.

4.1. Analysis

Amortized Cost: Many of the online algorithm works whose time complexities are determined by the Amortized Time complexities. Amortization is the process of determining the actual cost involved for insertions, deletions and rotation operations. The time complexities remain in O (logn) for the average time complexities. In the worst time analysis remains to be amortized. The average time complexity of the splaying operation follows the Amortized time complexity of O (logn) instead of O(n) and the cost incurred for the insertions and deletions remains the constant or amortized. So this technique is efficiently balances the load without incurring extra time for any primitive operations on the tree. Consider an empty tree and start S splay operations on the tree, the total time of running is O(S logn) but the average time for the operation is O (logn).

The total amortized cost for splaying operation is the sum of all the amortized cost for splay steps.

$$\begin{aligned}
 \text{Amortized cost} &= \sum \text{cost of splay steps.} \\
 &= \sum 3(\text{ra}(\text{root})-\text{ra}(\text{S}))+1 \\
 &= \sum 3\log n+1 \\
 &= O(\log n).
 \end{aligned}$$

ra(root) is the rank of the root node is logn.
 ra(S) rank of the node after splaying.

The performance efficiency of the splay operations on the amortized cost is always O(1) as it is self balancing, self optimized, no storage of data reducing memory requirement and it is very efficient with the uniform accessing of the data.

5. Conclusions

The splay thread cooperation for load balancing results in good performance comparatively on the graphics processors. The amortized time complexities reduce the time of execution task comparatively. The GPU consist of thousands of threads the application of the above technique comparatively reduces the memory access and balancing the work of the active nodes can be realized through the graph.

6. Results

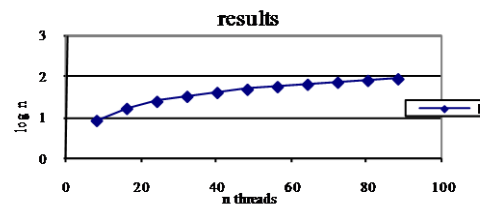


Figure 8. Graph of nthreads vs logn.

The above graph Figure 8 describes as the n number of threads increases for splaying on the GPU's the x axis consists of the n threads and y axis consists of log n resulting in linear speedup Figure 9,10,11,12,13. shows the different stages of Ray Traced image. Figures 14 and 15 is a snapshot of the speculative parallel ray tracing algorithm performance on processors.

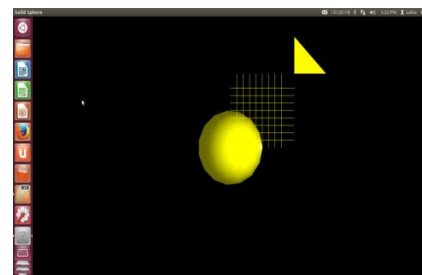


Figure 9. Initial scene for Ray tracing.

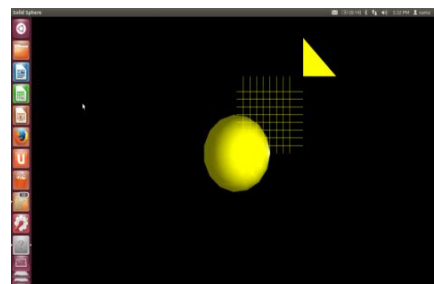


Figure 10. Scene with transformations.

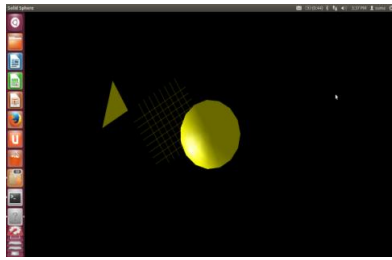


Figure 11. Reflections and splaying.

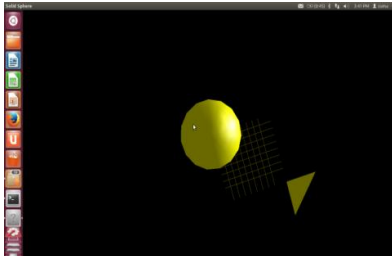


Figure 12. Ray tracing image.

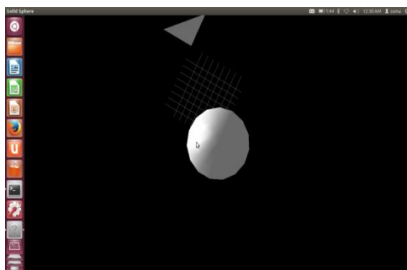


Figure 13. Ray traced image.



Figure 14. Performance on processors.



Figure 15. Performance on processors.

References

- [1] Aguila J. and Campero K., “An Explicit Parallelism Study Based On Thread Level Speculation,” *CLEI Electronic Journal*, vol. 17, no. 2, pp. 1-12, 2014.
- [2] Asmer N., *Partial Differential Equations with Fourier series and Boundary Value Problems*, Pearson, 2015.
- [3] Feng M., Gupta R., and Bhuyan L., “Speculative Parallelization on GPGPU’s,” in *Proceeding of the 17th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, Orleans, pp. 293-294, 2012.
- [4] Horowitz E., Sahni S., and Freed A., *Fundamentals of Data Structures in C*, Computer Science Press, 1992.
- [5] Han M., Wang Z., and Yuan J., “Mining Closed and Multi-Supports-Based Sequential Pattern in High Dimensional Dataset,” *The International Arab Journal of Information Technology*, vol. 12, no. 4, pp. 360-369, 2015.
- [6] Kaeli D. and Yew P., *Speculative Execution in High Performance Computer Architectures*, Chapman and Hall/CRC, 2005.
- [7] Kobayashi H., Nishimura S., Kubota H., Nakamura T., and Shigei Y., “Load Balancing Strategies a Parallel Ray-Tracing System Based on Constant Subdivision,” *The Visual Computer*, vol. 4, no. 4, pp. 197-209, 1988.
- [8] Lauterback C., Mo Q, and Manocha D., “Work Distribution Methods on GPUs,” Technical Report TR009-16, 2009.
- [9] Menon J., Kruijff M., and Sankaralingam K., “IGPU: Exception Support and Speculative Execution on GPUs,” in *Proceeding of 39th International Symposium on Computer Architecture*, Portland, pp. 72-83, 2012.
- [10] NVIDIA Corporation, *CUDA C programming guide PG-02829001-v5.0*, NVIDIA Corporation, 2012.
- [11] Reinhard E. and Jansen F., “Rendering Large Scenes Using Parallel Ray Tracing,” *Parallel Computing*, vol. 23, no. 7, pp. 873-885, 1997.
- [12] Sleator D. and Tarjan R., “Self-Adjusting Binary Search Trees,” *Journal of the Association for Computing Machinery*, vol. 32, no. 3, pp. 652-686, 1985.
- [13] Tian C., Feng M., and Gupta R., “Supporting Speculative Parallelism in Presence of Dynamic Data Structures,” in *Proceeding of the 31st ACM SIGPLAN Conference on Programming Language Design and Implementation*, Toronto, pp. 62-73, 2010.
- [14] Uht A., Morano D., Khalafi A., and Kaeli D., “A Scalable Processor with High IPC,” *journal of instruction level parallelism*, vol. 5, pp. 1-28, 2003.

- [15] Wood D., *Data Structures, Algorithms, and Performance*, Addison-Wesley, 1993.



Gopalan Pudur, PhD (IISCBangalore) M.E. (NIT Trichy) N.P.Gopalan received M.Sc degree with university rank in Mathematics from Madras University in 1978, PhD in Applied Mathematics from Indian Institute of Science,

Bangalore in 1983. He received M.E degree in computer Science and Engineering from NIT, Tiruchirappalli and currently serving as the Professor of Computer Applications in it. He has authored books on Web Technology, TCP/IP, Unix Programming, Data Mining object oriented Programming and Signals and Systems. His areas of interest includes Algorithms, Combinatorics, Data Mining, and Distributed Computing.



Suma Shivaraju, currently pursuing PhD from Bharti University Coimbatore. Her area of interest consists of Parallel and Distributed Computing, Data Structures Analysis and Design of algorithms, Automata Theory and Compiler

Design.