

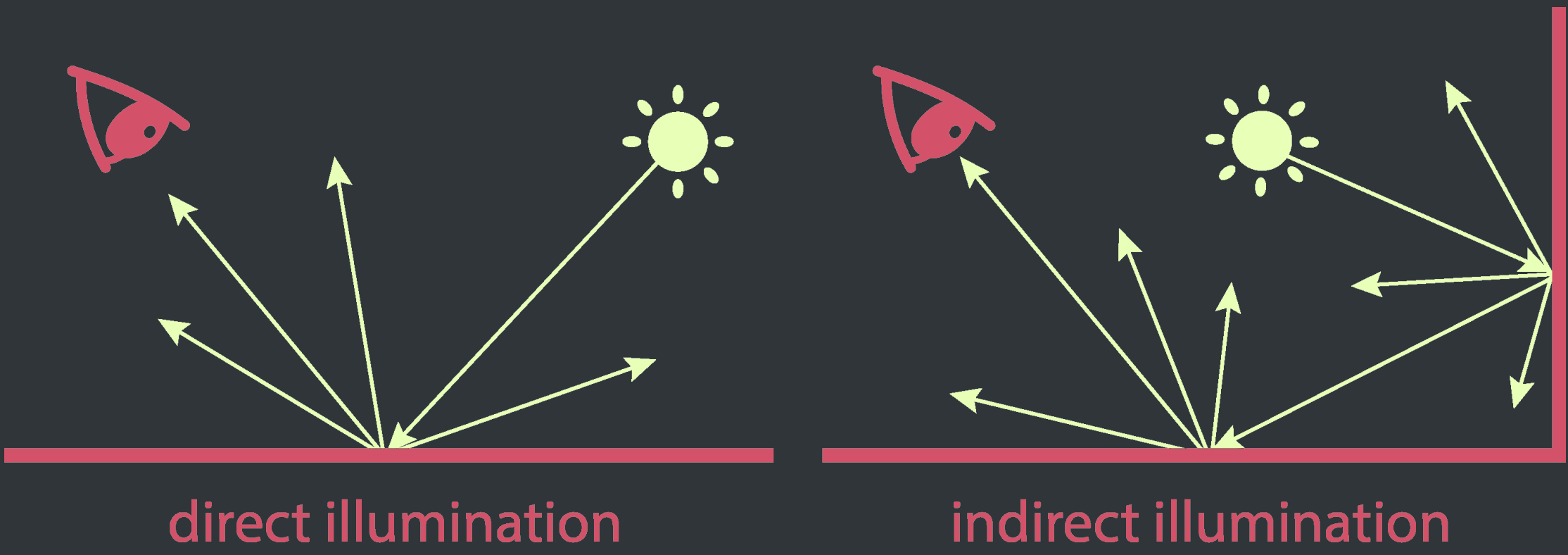
Performant GPU Pathtracing

Introduction Raytracing

Raytracing is a common method used in computer graphics to render three-dimensional images. Commercial renderers focus primarily on image fidelity and secondarily on performance (speed), but for professionals like filmmakers, architects, and game developers, the extended time involved in rendering a complex scene limits the number of iterations and refinements possible during the creative process. This project tests the effectiveness of developing a unidirectional pathtracer (a type of raytracer) that prioritizes performance over fidelity. The hypothesis is that with a focus on GPGPU (General-Purpose Graphical Processing Unit) computing and a novel implementation of algorithms, render times can be reduced without significantly impacting image quality. This could greatly improve iteration speed for large projects.

Pathtracing

Pathtracing uses Monte-Carlo integration (a method for approximating an integral) to simulate global illumination (indirect-illumination). It would take infinitely long to truly calculate the lighting for a scene, but by using Monte-Carlo integration, the global illumination for a point can be approximated to a high degree of accuracy in very little time. The pathtracer shoots out rays for each pixel on the screen. It then goes along each ray until it hits something in the scene. At the point of intersection, it accumulates colour and lighting data from the collider and shoots more (pseudo) random rays (samples of the Monte-Carlo integration) off into the scene again. It repeats this process until it either doesn't collide with anything or reaches the maximum number of bounces. The averaged result of these samples is the returned colour for the original pixel with approximated global illumination.



k-d Tree Traversal

k-d trees are the multidimensional analogue of Binary Search Trees, where a hyperplane is used to split elements of a node (this process happens recursively until a termination function is satisfied). In short, a k-d tree is a method of spatially splitting a scene into nodes. In a scene without spatial splitting, each ray would have to iterate through each object in the scene and separately do a collision test. This can be extremely costly in large scenes with millions of triangles (the base unit of geometry in a scene). In large scenes, each pixel will do hundreds of samples and each sample will have ten or so bounces, resulting in billions of intersection tests. With spatial splits, you can traverse through the tree and only do collision tests with triangles that are likely to intersect the ray. During runtime, if a ray is colliding with the scene, it can recurse down the tree with each node that it enters, removing an exponential amount of intersection tests. This dramatically improves render times, going from an $O(n^2)$ algorithm to $O(\log n)$. In reality, no scene can be perfectly divided so that there are no overlapping triangles, so issues arise with traversal where multiple nodes can contain triangles that could potentially intersect with the ray. Solutions are discussed below.

Procedure

Research was conducted and a Parallel GPU Pathtracer was developed to meet the criteria of the study. (See concatenated source code, Addendum A, display binder.) The final product was tested comparing its performance to that of commercial pathtracers. Using the same hardware, each test scene was rendered across three different pathtracers - Cycles CPU, Cycles GPU, and my own. The tests were repeated with a varying number of samples to compare performance with increasing fidelity. The results were recorded with their time and image. As computers are deterministic machines, multiple tests for each scenario were not necessary.

Algorithms & Implementation Groundwork and Strategy

Work began by laying down the core structural features of the pathtracer: Scene Loading; Custom Scene Format; OBJ Loading; Scene Serialization; OSX and Windows Front ends; OpenCL workload distribution utilities; a server serving the UI (as a website); and a remote debugger via websocket server. Then to meet the criteria of the study, it was most valuable to focus on the acceleration structure of the pathtracer. The two algorithms chosen were: the $O(n \log^2 n)$ Surface Area Heuristic (SAH) k-d tree construction algorithm; and a persistent short stack k-d tree traversal algorithm.

SAH k-d Tree Construction

SAH k-d tree construction is one of the best-known ways to find the cost of a plane split. The SAH of a split is the potential cost of traversal and intersection for a given split. The SAH itself provides no way of finding the minimal costs. Luckily, the minima can only be found on minimum and maximum bounds of an object within the voxel being split. This is because the change in cost is linear between each edge bound of an object. It is then easy to sweep across all possible minima with only $O(n \log^2 n)$ complexity. While there are faster algorithms, the increased implementation time for them would not have been worth it for the minimal speed increases. Recently the practice of using a favouring function (λ) to increase the chance of a split where one of the sides has no objects has become quite common showing consistent speed increases. My modified implementations of these algorithms can be found below:

Favouring Function

$$\lambda(P_L, P_R, N_L, N_R) = \begin{cases} 80\% & (N_L \equiv 0 \vee N_R \equiv 0) \wedge (P_L \neq 1 \wedge P_R \neq 1) \\ 100\% & \text{otherwise} \end{cases}$$

Cost Function

$$C(P_L, P_R, N_L, N_R) = \lambda(P_L, P_R, N_L, N_R) (K_T + K_I(P_L N_L + P_R N_R))$$

Algorithm .1: Surface Area Heuristic.

```
1 (Cov, Side) SAH(p, V, N_L, N_R, Np)
2 {
3   Voxel V_L, V_R;
4   voxel_split(p, V, &V_L, &V_R);
5   P_L = SAH(p, V_L, N_L, Np);
6   P_R = SAH(p, V_R, N_R, Np);
7   C_L = C(P_L, P_R, N_L + Np, Np);
8   C_R = C(P_R, P_R, N_R + Np, Np);
9   return min((C_L, LEFT), (C_R, RIGHT));
10 }
```

Algorithm .3: Classify.

```
1 (T_L, T_R) classify(tree, T, p, N_L, N_R, Np)
2 {
3   T_L = null; T_R = null;
4   T_L = nullloc(N_L + Np);
5   T_R = nullloc(N_R + Np);
6   T_L = tree; T_R = tree;
7   for(i < T)
8   {
9     isLeft = isRight = false;
10    for(j = 0; j < 3; j++)
11    {
12      b_j = T[j];
13      if(b_j < p_j)
14        isLeft = true;
15      if(b_j > p_j)
16        isRight = true;
17    }
18    if(!isLeft & !isRight) // planar
19    {
20      if(p_j == RIGHT)
21        T_L = T;
22      else
23        T_R = T;
24      T_L = T;
25    }
26    if(isLeft)
27      T_L = T;
28    if(isRight)
29      T_R = T;
30    T_L = T;
31    T_R = T;
32    return (T_L, T_R);
33  }
```

Algorithm .4: GenerateTree.

```
1 Node gen_node(tree, V, T, depth)
2 {
3   Node n;
4   (N_L, N_R, Np, p, side) = find_plane(tree, V, T);
5   Np = p;
6   if(depth == tree_max_depth || p > K(T))
7   {
8     n = T;
9     return n;
10  }
11  }
12  }
13  }
14  }
15  }
16  }
17  }
18  }
19  }
20  }
```

Algorithm .2: Find Split Plane $O(n \log n)$.

```
1 (N_L, N_R, Np, p, side) find_plane(tree, V, T)
2 {
3   best_cost = inf; best_p = null; side = null; E = null;
4   best_N_L = best_N_R = best_Np = 0;
5   for(i = 0; i < tree->sk; i++)
6   {
7     j = 0;
8     E = nullloc(2T);
9     for(j < T)
10    {
11      Voxel B;
12      B = voxel_gen_from_tri(i);
13      B = voxel_clip(B, V);
14      if(voxel_is_planar(B, i))
15      {
16        E_L = {r, B.data, k, PLANAR};
17      }
18      else
19      {
20        E_L = {r, B.data, k, START};
21        E_R = {r, B.data, k, END};
22      }
23    }
24    sort(E); // sort the events by b
25    N_L = N_R = 0;
26    Np = |E|;
27    for(i = 0; i < |E|)
28    {
29      p = E[i];
30      P_start = P_end = P_planar = 0;
31      while(i < |E| & E[i].p_j == p_j & E[i].type == END)
32        i++;
33      while(i < |E| & E[i].p_j == p_j & E[i].type == PLANAR)
34        P_planar++;
35      while(i < |E| & E[i].p_j == p_j & E[i].type == START)
36        P_start++;
37      Np = P_start; Np += P_planar; Np += P_end;
38      sub_data = SAH(k, p_j, V, N_L, N_R, Np);
39      if(sub_data.cost < best_cost)
40      {
41        best_cost = sub_data.cost;
42        best_p = p;
43        best_side = sub_data.side;
44        best_N_L = N_L; best_N_R = N_R; best_Np = Np;
45      }
46      Np += P_planar; Np += P_start; Np += P_end;
47    }
48    return (best_N_L, best_N_R, best_Np, best_p, best_side);
49  }
```

Persistent Threading

Persistent Threading utilizes the improved performance of warp synchronous execution (a warp is Nvidia's unit for a group of processors that run using SIMT, also known as a Wavefront on AMD hardware). It enforces that only a single task (thread) is assigned to each core of the GPU. This allows for improved SIMT performance between cores of a warp/wavefront as it enforces equal distribution of work tasks (especially when there is a non-trivial workload). It gets around the limited workloads of warp synchronous programming by implementing a global work queue that each thread pulls off of. Persistent threading also alleviates the unbalanced workload distribution that can happen with algorithms that don't fit the SIMT model well (such as tree traversal). The GPU scheduler can sometimes distribute work incorrectly where certain warps will be doing all of the heavy lifting for an operation, but because each core only has one thread, persistent threading can bypass the schedulers as the workload must be evenly distributed. This can improve the performance of these algorithms as it takes advantage of the entire GPU.



Results



The CPU pathtracer performed the worst by a fairly large margin. My pathtracer was roughly twelve times faster in the Norway scene than Cycles CPU. Results primarily come down to the fact that the test GPU was faster than the test CPU (this is very common in computers). This validates the focus on parallelism. One reason for the speed difference could, potentially, be the complexity of the shaders that Cycles uses. They use PBR-shader pipelines, which are significantly more complicated than my simple diffuse/emission shader set-up. That being said shading itself is not a performance-draining task and wouldn't account for the extent of the performance gap. More research will have to be conducted to determine why there is such a marked difference.

Future Directions

Physically Based Rendering: Farther off in the distance, implementing a proper Physically Based Rendering (PBR) pipeline would be ideal as there would be greatly improved graphical fidelity. PBR leverages physically accurate lighting models as opposed to Phong lighting, which is a very high-level abstraction of how a surface interacts with light. It results in more realistic images.

Conclusion

The hypothesis proved true. Removing extraneous features and designing a pathtracer for GPGPU computing can boost render times. The program, while limited feature-wise, has use even in its current state. Where someone wants a quick high detail pre-visualization of a more complicated render, this pathtracer would render a much more accurate image than a cheap rasterized Phong imitation. With minimal work on this pathtracer, it could be possible to generate better images that would be comparable to even the best of commercial pathtracers.



Bibliography

- Alia, T., & Laine, S. (2009). Understanding the Efficiency of Ray Traversal on GPUs. Proceedings of the SIGGRAPH/Eurographics Conference on High Performance Graphics 2009, 145-149. ACM. doi:10.1145/1587927.1587992
- Castano, I. (2007). Irradiance Caching – Part I. Retrieved from <http://www.thewyrtan.net/news/2007/07/irradiance-caching-part-1/>
- Castano, I. (2014). Irradiance Caching – Continued. Retrieved from <http://www.ludicri.com/castano-blog/2014/07/irradiance-caching-continued/>
- Chen, S., Choi, M., Voun, B., & Kim, S. (2018). Comparison of BVH and KD-tree for the GPU acceleration on real mobile devices. Retrieved from https://www.researchgate.net/publication/330219705_Comparison_of_BVH_and_KD-tree_for_the_GPU_acceleration_on_real_mobile_devices
- Cormen, T. H., Leiserson, C. E., Rivest, R. L., & Stein, C. (1999). Introduction to Algorithms. Cambridge, MA: MIT Press.
- Cullum, Brian. (2017). Anti-aliasing and Monte Carlo Path Tracing, CSE 557 Course Lecture Notes, University of Washington. Retrieved from courses.cs.washington.edu/courses/cse557/17.1/assets/lectures/aa-and-mcmt-5pp.pdf
- Driscoll, V. (2009). Better sampling. Retrieved from <http://www.raytrid.org/2009/01/07/better-sampling/>
- Frisol, V., Vostyakov, K., Kharlamov, A., & Galaktionov, V. (2013). Implementing Irradiance caching in a GPU realistic render. Transactions on Computational Science XIX, 17-32. Springer. doi:10.1007/978-3-642-37974-2_2
- Gottsch, K., Stuart, J. A., & Owens, D. J. (2012). A study of persistent threads vs. CUDA programming for GPU workloads. 2012 Innovative Parallel Computing (InPar). 1414-1421. IEEE. doi:10.1109/InPar.2012.6335956
- Harris, M. (2013). Optimizing cuda. SCOT' High Performance Computing with CUDA. NVIDIA. Retrieved from <http://gpgpu.conf.nvidia.com/2013/07/27/CUDA-5-Optimization-IEEE.pdf>
- Horn, D. R., Supeman, J., Houston, M., & Hanrahan, P. (2007). Interactive k-D Tree GPU Raytracing. Proceedings of the 2007 Symposium on Interactive 3D Graphics and Games. 167-174. ACM. doi:10.1145/1273004.1273074
- Kernighan, B., & Ritchie, D. M. (2017). The C Programming Language. Upper Saddle River, NJ: Prentice Hall.
- Lapere, S. (2016, September 20). Optimised BVH building, faster traversal and intersection kernels and HDR environment lighting. GPU Path Tracing Tutorial 4. Retrieved from <http://www.raytraced.com/2016/09/20/optimised-bvh-building-faster-traversal-and-intersection-kernels-and-hdr-environment-lighting/>
- Lindholm, E., Booth, K., Oberman, S., & Montyorn, J. (2008). NVIDIA Tesla unified Graphics and Computing Architecture. IEEE Micro (2008), 28 (2), 39-55. doi:10.1109/MICR.2008.31
- MacDonald, D. L., & Nickolls, K. S. (1990). Heuristics for ray tracing using a dedicated Graphics. The Visual Computer (1990), 6 (3), 153-166. Springer-Verlag. <https://doi.org/10.1007/BF01901006>
- Mikkelsen, P. T. (2009). Tesla CUDA Ray Tracing Tutorial. Retrieved from <https://docs.google.com/alexandradk/22x78>
- Nvidia Corporation (2016). NVIDIA Tesla P100 Whitepaper. Retrieved from <https://images.nvidia.com/content/pdf/tesla/whitepaper/pascal-architecture-whitepaper.pdf>
- Pharr, M., & Fernando, R. (Eds.). (2005). GPU Gems 2: Programming Techniques for High-Performance Graphics and General-Purpose Computing. Boston, MA: Addison-Wesley Professional.
- Pharr, M., Jakob, W., & Humphreys, G. (2017). Physically Based Rendering: From Theory to Implementation. Cambridge, MA: Elsevier.
- Popov, S., Günther, J., Seidel, H. P., & Susulak, P. (2007, September). Stacked kd-tree traversal for high performance GPU ray tracing. Computer Graphics Forum, 26 (3), 415-424. The Eurographics Association and Blackwell Publishing. doi:10.1111/j.1467-8659.2007.01064.x
- Prunier, J. C. et al. (2016). Scratchapixel. Retrieved from <https://www.scratchapixel.com/>
- Scarpino, M. (2012). OpenCL in Action: How to Accelerate Graphics and Computations. Shelter Island, NY: Manning Publications.
- Spencer, S. et al. (1993). Examining Radiosity. 1993 ACM SIGGRAPH Educator's Slide Set. Retrieved from <http://education.sigraph.org/archive/slide-sets/>
- Vivo, P. G., & Lowe, J. (2008, October 24). The Book of Shaders. Retrieved from <http://thebookofshaders.com/>
- Wald, I., & Havar, V. (2010). On building fast kd-trees for Ray Tracing, and on doing that in O (N log N) time. 2006 IEEE Symposium on Interactive Ray Tracing, 61-69. IEEE. doi:10.1109/IRT.2006.280216
- Ward, C. (2007). Irradiance caching algorithm, SIGGRAPH '07 Special Interest Group on Computer Graphics and Interactive Techniques Conference, Article 3. New York, NY, ACM. doi:10.1145/1281500.1281619
- Ward, C. (2007). Implementation of irradiance caching in real time. SIGGRAPH '07 Computer Graphics and Interactive Techniques Conference, Article 4. ACM. doi:10.1145/1281500.1281620
- Zatseukin, M., & Havar, V. (2010). "Ray tracing on a GPU with CUDA-comparative study of three algorithms." (2010). Retrieved from https://www.researchgate.net/publication/22574023_Ray_tracing_on_a_GPU_with_CUDA-comparative_study_of_three_algorithms