

Surface Area Heuristic

The Surface Area Heuristic (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 $O(n \log n)$ complexity. While there are faster algorithms, the time that it would have taken to implement them would not have been worth it for the minimal speed increases that they would have given. 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.

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

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

Algorithm 1: Surface Area Heuristic.

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

Algorithm 2: Classify.

```

1  ( $T_R, T_L$ ) classify ( $tree, T, p, N_L, N_R, N_P$ )
2  {
3       $T_L = \text{null}; T_R = \text{null};$ 
4       $T_L = \text{malloc}(N_L + N_P);$ 
5       $T_R = \text{malloc}(N_R + N_P);$ 
6       $i_{TL} = i_{TR} = 0$ 
7      for ( $t < T$ )
8      {
9           $is\_left = is\_right = \text{false}$ 
10         for ( $j = 0; j < 3; j++$ )
11         {
12              $t_b = t_{j,p_k};$ 
13
14             if ( $t_b < p_b$ )
15                  $is\_left = \text{true};$ 
16
17             if ( $t_b > p_b$ )
18                  $is\_right = \text{true};$ 
19         }
20
21         if ( $\neg is\_left \wedge \neg is\_right$ ) // planar
22         {
23             if ( $p_{side} \equiv \text{RIGHT}$ )
24                  $T_R[i_{TR}++] = t;$ 
25             else
26                  $T_L[i_{TL}++] = t;$ 
27         }
28         if ( $is\_left$ )
29              $T_L[i_{TL}++] = t;$ 
30         if ( $is\_right$ )
31              $T_R[i_{TR}++] = t;$ 
32     }
33     return ( $T_R, T_L$ );
34 }
35 
```

Algorithm 4:

```

1  Node gen_node ( $tree, V, T, depth$ )
2  {
3      Node  $n;$ 
4
5      ( $N_L, N_R, N_P, p, side$ ) = find_plane ( $tree, V, T$ );
6       $n_p = p;$ 
7
8      if ( $depth \equiv tree_{max.depth} \vee p_{cost} > K_t |T|$ )
9      {
10          $n_T = T;$ 
11         return  $n;$ 
12     }
13
14     Voxel  $V_L, V_R;$ 
15     voxel_split ( $p, V, \&V_L, \&V_R$ );
16     ( $T_R, T_L$ ) = classify ( $tree, T, p, N_L, N_R, N_P$ );
17      $n_{left} = \text{gen\_node} (tree, V_L, T_L, depth+1);$ 
18      $n_{right} = \text{gen\_node} (tree, V_R, T_R, depth+1);$ 
19     return  $n;$ 
20 }
21 
```

Algorithm 3: Find Split Plane. $O(n \log n)$

```

1  ( $N_L, N_R, N_P, p, side$ ) find_plane ( $tree, V, T$ )
2  {
3       $best\_cost = \infty; best\_p = \text{null}; side = \text{null}; E = \text{null};$ 
4       $best\_N_L = best\_N_R = best\_N_P = 0;$ 
5
6       $j = 0;$ 
7      for ( $k = 0; k < tree \rightarrow k; k++$ )
8      {
9           $E = \text{malloc}(2|T|);$ 
10         for ( $t \text{ in } T$ )
11         {
12             Voxel  $B;$ 
13             voxel_gen_from_tri (& $B, t$ );
14             voxel_clip (& $B, V$ );
15             if (voxel_is_planar ( $B, k$ ))
16             {
17                  $E[j++] = \{t, B_{min,k}, k, \text{PLANAR}\};$ 
18             }
19             else
20             {
21                  $E[j++] = \{t, B_{min,k}, k, \text{START}\};$ 
22                  $E[j++] = \{t, B_{max,k}, k, \text{END}\};$ 
23             }
24         }
25         sort ( $E$ ); // explain later
26          $N_L = N_P = 0;$ 
27          $N_P = |E|;$ 
28         for ( $i = 0; i < |E|$ )
29         {
30              $p = E[i];$ 
31              $P_{START} = P_{END} = P_{PLANAR} = 0;$ 
32
33             while ( $i < |E| \wedge E[i]_b \equiv p_b \wedge E_{type} \equiv \text{END}$ )
34                  $\{P_{END}++; i++;\}$ 
35
36             while ( $i < |E| \wedge E[i]_b \equiv p_b \wedge E_{type} \equiv \text{PLANAR}$ )
37                  $\{P_{PLANAR}++; i++;\}$ 
38
39             while ( $i < |E| \wedge E[i]_b \equiv p_b \wedge E_{type} \equiv \text{START}$ )
40                  $\{P_{START}++; i++;\}$ 
41
42              $N_P = P_{PLANAR}; N_R \leftarrow P_{PLANAR}; N_L \leftarrow P_{END};$ 
43
44              $sah\_data = \text{SAH}(k, p_b, V, N_L, N_R, N_P);$ 
45
46             if ( $sah\_data_{cost} < best\_cost$ )
47             {
48                  $best\_cost = sah\_data_{cost};$ 
49                  $best\_p = p;$ 
50                  $best\_side = sah\_data_{side};$ 
51                  $best\_N_L = N_L; best\_N_R = N_R; best\_N_P = N_P;$ 
52             }
53              $N_L += P_{PLANAR}; N_L += P_{START}; N_P = 0;$ 
54         }
55     }
56     return ( $best\_N_L, best\_N_R, best\_N_P, best\_p, best\_side$ );
57 }
58 
```

Persistent Threading

Persistent threading is a relatively new algorithm that takes advantage of a rather old strategy. It takes advantage of 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). Persistent threading 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. Another issue that persistent threading alleviates is unbalanced workload distribution that can happen with algorithms that don't fit the SIMD 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 greatly improve the performance of these algorithms as it takes advantage of the entire GPU.

Short stack k-d tree traversal.

This is one of the fastest methods of k-d tree traversal. It takes advantage of both the Persistent Threading algorithm and the Short Stack k-d tree traversal algorithm. The short stack k-d tree traversal algorithm is a method for traversal that uses the best parts of k-d restart (with pushdown modification) and traversal algorithms with stacks by using a bounded stack. It will keep descending the tree (and pushing down the root) until the traversal algorithm finds a split where both nodes intersect the ray. The algorithm will then disable push down and pushes the farther node to the stack. It will then continue traversing down the closest node. If no intersection is found along the closest path it will then pop the first split off the stack (the next closest to the origin) and traverse. If the stack becomes full, and then doesn't hit anything while it is being emptied, the algorithm will resort to use k-d restart with the aforementioned pushdown modification where it will then restart right before the first split.

Algorithm 5: Persistent Short Stack k-d Tree Traversal.

```

1 void update_state(tree_buffer, index, *type, *node, *leaf)
2 {
3     *type = *(uint8*)(tree_buffer+index); // first value of both node and leaf is the type
4     if(*type == LEAF)
5     {
6         kd_serialized.leaf sleaf = *(kd.leaf*)(tree_buffer + index); //need to use serialized types for this
7         leaf->type = *type;
8         leaf->tri_n = sleaf.tri_n; //number of triangles
9         leaf->tri_offset = index+sizeof(kd_serialized.leaf);
10    }
11    else //NODE
12    {
13        *node = *(kd_node*)(tree_buffer + index);
14    }
15 }
16
17 (index,t,u,v) traverse(ray_buffer, indices, vertices, tree_buffer)
18 {
19     blocksize_x = STREAM_PROCESSORS.PER_SIMT_GROUP;
20     blocksize_y = SIMT_GROUPS.PER_STREAM_MULTIPROCESSOR;
21
22     x = SM.ID % blocksize_x; //ld within the SIMT GROUP
23     y = SM.ID / blocksize_x; //ld of the SIMT GROUP within the Stream Multiprocessor
24
25     //NOTE: shared memory is called local memory in OpenCL
26     shared volatile next_ray_array[blocksize_y]; //shared across all processors in the multiprocessor
27     shared volatile ray_count_array[blocksize_y];
28
29     //NOTE: In the implementation, the warp_counter is initialised on the cpu and copied to the GPU.
30     global volatile warp_counter; //global memory is shared across the entire device.
31
32     next_ray_array[y] = 0;
33     ray_count_array[y] = 0;
34
35     (node_ptr, min, max) stack[STACK_SIZE];
36     stack_pointer = 0;
37
38     ray r;
39     t_hu = 0;
40     hit_info = {0};
41     tringle_index = 0;
42     t_min = t_max = 0;
43     scene_min = 0; scene_max = 0;
44     kdtree_node root, node;
45     kdtree_leaf leaf;
46     current_type = false;
47     pushdown = false;
48     ray_index = 0;
49
50     while (true)
51     {
52         shared volatile int* local_pool_ray_count = ray_count_array+y; //get this SIMT groups ray count
53         shared volatile int* local_pool_next_ray = next_ray_array+y; //get this SIMT groups next ray
54
55         if(x==0 && *local_pool_ray_count <= 0)
56         {
57             *local_pool_next_ray = atomic.add(warp_counter, BATCH_SIZE); //retrieve and increment
58             *local_pool_ray_count = BATCH_SIZE;
59         }
60
61         ray_index = *local_pool_next_ray + x;
62         if(ray_index >= ray_buffer)
63             break;
64
65         if(x == 0)
66         {
67             *local_pool_next_ray += 32;
68             *local_pool_ray_count -= 32;
69         }
70
71         r = ray_buffer[ray_index];
72         t_hu = 0;
73
74         if(!collides_voxel(SCENE.V, r, &scene_min, &scene_max))
75         {
76             scene_max = 0;
77         }
78
79         stack_pointer = 0;
80         root = *tree_buffer;
81
82         while (t_max < scene_max)
83         {
84             if(stack_pointer == 0)
85             {
86                 node = root;
87                 current_type = NODE;
88                 t_min = t_max;
89                 t_max = scene_max;
90                 pushdown = true;
91             }
92             else //pop a node off the stack
93             {
94                 stack_pointer--;
95                 svalue = stack[stack_pointer];
96                 t_min = svalue_min;

```

```

98     t_max = svalue_max;
99     update_state(tree_buffer, svalue_node_pointer, &current_type, &node, &leaf);
100     pushdown = false;
101 }
102
103 while (current_type != LEAF)
104 {
105     t_split = (node_t - r_origin_k) / r_dir_k;
106
107     left_is_close = (r_origin_k < node_b) & (r_origin_k <= node_b & r_dir_k <= 0);
108
109     first = left_is_close ? node_left : node_right;
110     second = left_is_close ? node_right : node_left;
111
112     if (t_split > t_max & t_split <= 0)
113     {
114         update_state(tree_buffer, first, &current_type, &node, &leaf);
115     }
116     else if (t_split < t_min)
117     {
118         update_state(tree_buffer, second, &current_type, &node, &leaf);
119     }
120     else
121     {
122         stack[stack_pointer++] = {second, t_split, t_max}; //push
123         update_state(tree_buffer, first, &current_type, &node, &leaf);
124         t_max = t_split;
125         pushdown = false;
126     }
127
128     if (pushdown)
129     {
130         root = node;
131     }
132 }
133
134 for (i = 0; i < leaf_num_triangles; i++)
135 {
136     vec3 tri[3];
137     offset = *(uint*)(kd_tree + leaf_triangle_start + (i * sizeof(uint)));
138
139     for (j = 0; j < 3; j++) //read triangle indices
140         tri[j] = read_texture(vertices, read_texture(indices, offset+j), 0);
141
142     hit_coords = {0};
143
144     if (collides_triangle(tri, &hit_coords, r))
145     {
146         if (hit_coords_t <= 0)
147             continue;
148         if (hit_coords_t < t_hu)
149         {
150             t_hu = hit_coords_t;
151             hit_info = hit_coords;
152             tri_index = offset;
153         }
154     }
155 }
156
157 result = {0};
158 if (t_hu != 0)
159 {
160     result = {tri_index, t_hu, hit_info, hit_info};
161 }
162
163 return result;
164 }
165 }

```