



SIGGRAPH 2024
DENVER+ 28 JUL — 1 AUG

THE PREMIER CONFERENCE
& EXHIBITION ON
COMPUTER GRAPHICS &
INTERACTIVE TECHNIQUES

ACHIEVING SCALABLE PERFORMANCES FOR LARGE SCALE GAME COMPONENTS WITH CBTS



© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

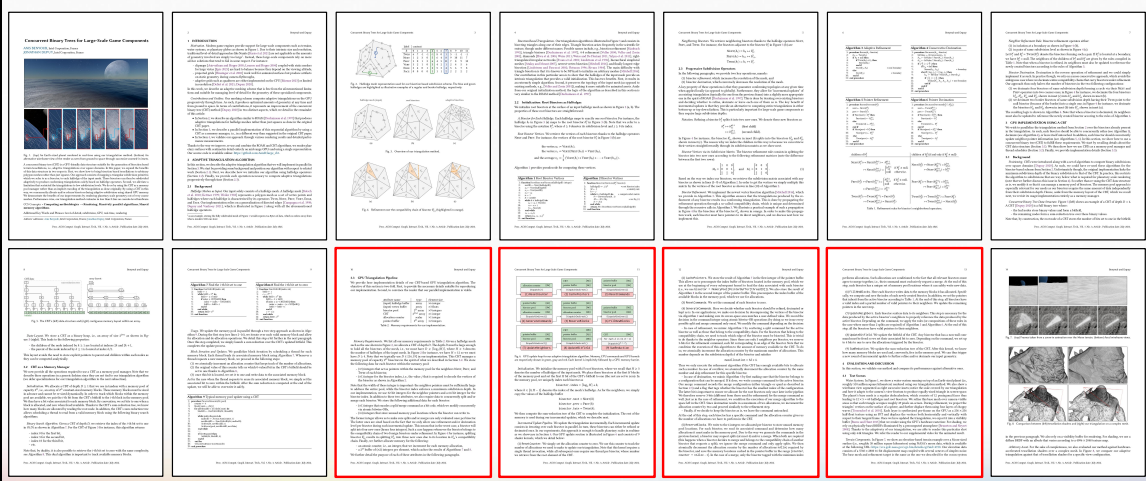
"Concurrent Binary Trees for Large-Scale Game Components" HPG24



This presentation is a follow-up of our HPG paper that we presented just a couple of days ago.

The paper is called "Concurrent Binary Trees for Large-Scale Game Components"...

...and most of the content of the paper is shown here, including...



implementation details

4 pages of implementation details shown in red!

The reason I'm showing these is that it turns out that since we wrote the paper 6 months ago, Anis re-worked the implementation to make it even faster.

The image displays a grid of 12 technical slides from a SIGGRAPH 2024 presentation. The slides contain text, diagrams, and code snippets. A red box highlights a section of slide 3, which is titled "3.1 GPU Transpilation Pipeline" and includes a diagram of a transpilation pipeline and a table of metrics.

Slide 1: Overview of the research, including a title slide and a list of authors.

Slide 2: Introduction to the research, including a list of authors and a summary of the work.

Slide 3: Overview of the research, including a list of authors and a summary of the work. A red box highlights a section titled "3.1 GPU Transpilation Pipeline" which includes a diagram of a transpilation pipeline and a table of metrics.

Slide 4: Overview of the research, including a list of authors and a summary of the work.

Slide 5: Overview of the research, including a list of authors and a summary of the work.

Slide 6: Overview of the research, including a list of authors and a summary of the work.

Slide 7: Overview of the research, including a list of authors and a summary of the work.

Slide 8: Overview of the research, including a list of authors and a summary of the work.

Slide 9: Overview of the research, including a list of authors and a summary of the work.

Slide 10: Overview of the research, including a list of authors and a summary of the work.

Slide 11: Overview of the research, including a list of authors and a summary of the work.

Slide 12: Overview of the research, including a list of authors and a summary of the work.

(simplified) implementation details

So what we describe in the paper is now more or less deprecated already!

The image displays a grid of 14 presentation slides from a SIGGRAPH 2024 presentation. The slides are arranged in three rows and four columns. The top row contains slides 1 through 4. The middle row contains slides 5 through 8, which are highlighted with a red border. The bottom row contains slides 9 through 12. The slides contain various technical content, including diagrams, tables, and text. The bottom row of slides shows a sequence of images illustrating a rendering process. The overall layout is professional and technical, typical of a SIGGRAPH presentation.

(simplified) implementation details
more details in this talk!

And so in this presentation Anis will share all the details of his new implementation.

Before he does so, I will quickly recap what the paper is about to bring everyone up to speed.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



The paper introduces a new algorithm to deal with large-scale game components.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



A large-scale game component is typically what makes the virtual world of your game look "big".

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



terrains



Red Dead Redemption 2
Rockstar

Taking a few video games as example, that would be terrains, ...

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



terrains



Red Dead Redemption 2
Rockstar

oceans



Assassin's Creed Origins
Ubisoft

..., oceans, ...

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



terrains



Red Dead Redemption 2
Rockstar

oceans



Assassin's Creed Origins
Ubisoft

planets



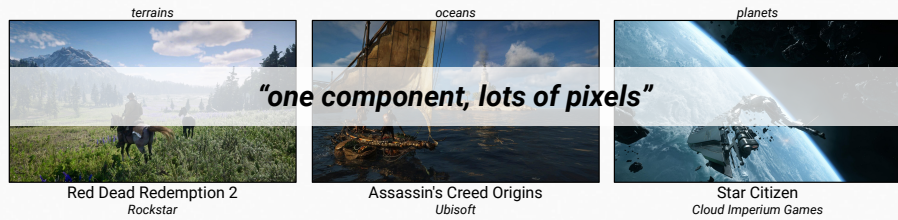
Star Citizen
Cloud Imperium Games

... or entire planets.

In each of these screenshots, pretty much everything except the characters in the foreground is rendered using a dedicated system.

This dedicated system is what we refer to a "large-scale game component".

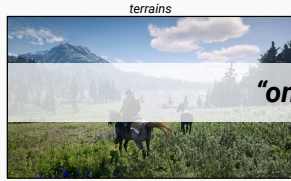
"Concurrent Binary Trees for Large-Scale Game Components" HPG24



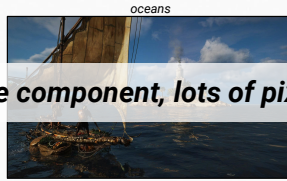
Most of the time, such components occupy a lot of pixels so it's important to have a set of efficient algorithms to render them as fast as possible.

The goal of our paper is to contribute to this set of efficient algorithms, ...

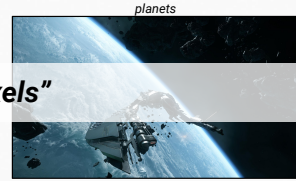
"Concurrent Binary Trees for Large-Scale Game Components" HPG24



Red Dead Redemption 2
Rockstar



Assassin's Creed Origins
Ubisoft



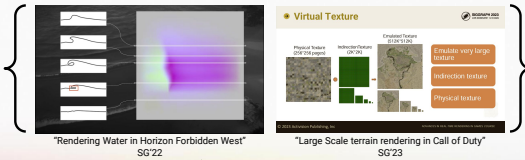
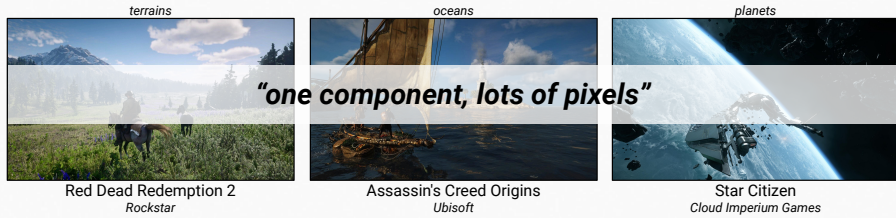
Star Citizen
Cloud Imperium Games

"one component, lots of pixels"

2 sub-problems:

...which can typically be classified into two categories.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24

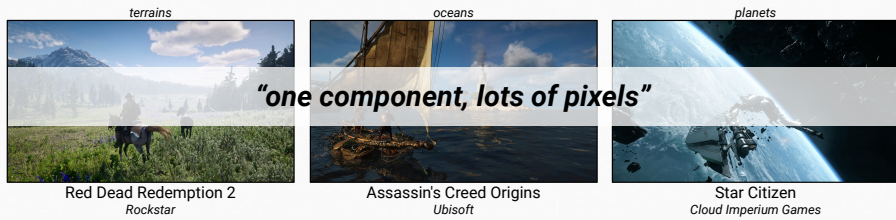


2 sub-problems:

1) data generation

First there's the data-generation category, which addresses how to generate textures, sprites, instances, etc.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



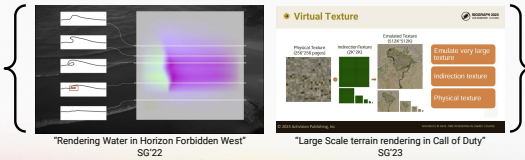
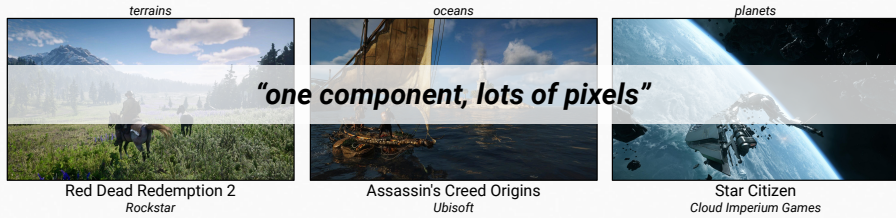
2 sub-problems:

1) data generation

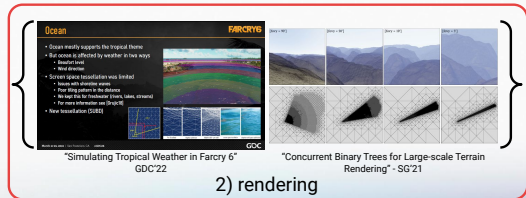
2) rendering

Second the triangulation / rendering category, which focuses on how to render this data.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



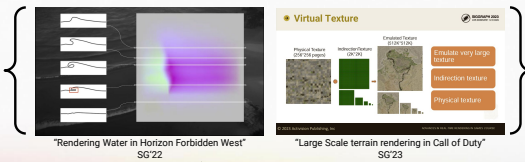
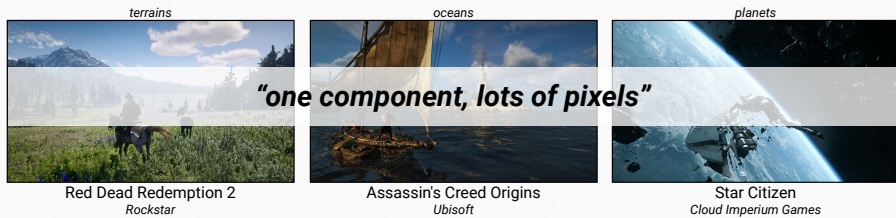
2 sub-problems: 1) data generation



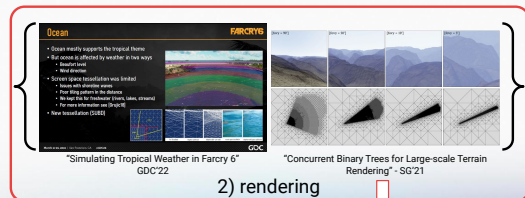
our contribution 2) rendering

Our paper contributes only to the latter (so we won't be discussing data-generation here) with a triangulation method capable of handling very large environments that can be explored at any different scales.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



2 sub-problems: 1) data generation



2) rendering
our contribution: "CBT-V2"

And actually, our contribution is an improvement over something called "Concurrent Binary Trees" (CBTs), which we presented in the same course 3 years ago.

I like to refer to our improvement as "CBT version 2", or simply CBT-V2. I will explain what CBTs are in just a minute.



© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

Before I do just that, here is a look at what CBT-V2 can render.

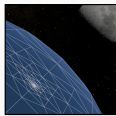
Let's have a look at a video that captures the result of our method, which runs at 250+FPS on a PS5 level hardware at full HD resolution.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



To provide a better sense of scale of what our CBTV2 produces, we are going to have a look at how dense the triangulation of this shot is.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



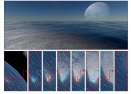
Here is an alternative view of the same shot and we are going to zoom into it until we reach the resolution of the mesh

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



"Concurrent Binary Trees for Large-Scale Game Components" HPG24

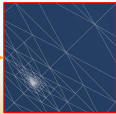
Concurrent Binary Trees for Large-Scale Game Components



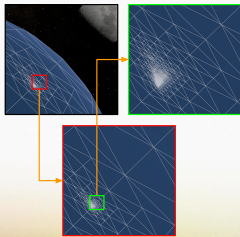
By the end of this presentation, you will be able to:

- Understand the challenges of rendering large-scale game components.
- Describe the concurrent binary tree data structure and how it is used for rendering.
- Explain how the concurrent binary tree data structure is used to render large-scale game components.
- Describe the performance benefits of the concurrent binary tree data structure.

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

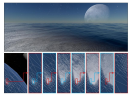


"Concurrent Binary Trees for Large-Scale Game Components" HPG24

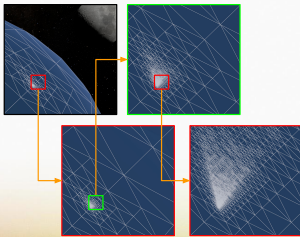


"Concurrent Binary Trees for Large-Scale Game Components" HPG24

Examining Binary Trees for Large-Scale Game Components

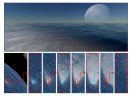


By the end of this presentation, you will be able to understand the importance of binary trees in game development, how they are used to manage large-scale game components, and how they can be optimized for performance. You will also learn about the challenges of implementing binary trees in a game engine and how to overcome them.



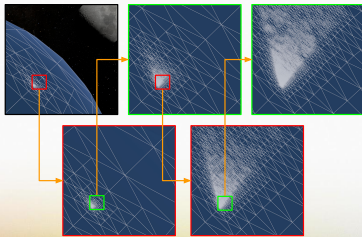
"Concurrent Binary Trees for Large-Scale Game Components" HPG24

Examining Binary Trees for Large-Scale Game Components



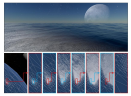
By the end of this presentation, you will be able to:

- Understand the challenges of rendering large-scale game components.
- Analyze the performance of different rendering techniques.
- Apply concurrent binary trees to optimize rendering.
- Evaluate the results of the optimization process.

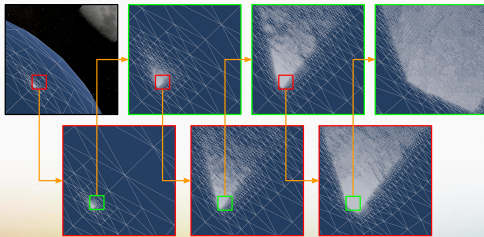


"Concurrent Binary Trees for Large-Scale Game Components" HPG24

Concurrent Binary Trees for Large-Scale Game Components

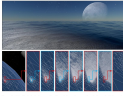


By the end of this presentation, you will be able to understand the concurrent binary trees for large-scale game components rendering. You will learn how to design a concurrent binary tree structure for large-scale game components rendering. You will learn how to implement a concurrent binary tree structure for large-scale game components rendering. You will learn how to optimize a concurrent binary tree structure for large-scale game components rendering. You will learn how to evaluate a concurrent binary tree structure for large-scale game components rendering. You will learn how to compare a concurrent binary tree structure for large-scale game components rendering. You will learn how to conclude a concurrent binary tree structure for large-scale game components rendering.

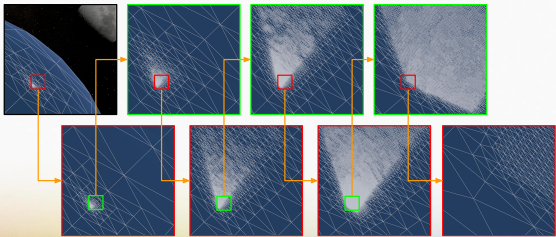


"Concurrent Binary Trees for Large-Scale Game Components" HPG24

Concurrent Binary Trees for Large-Scale Game Components
by [unreadable]
[unreadable]



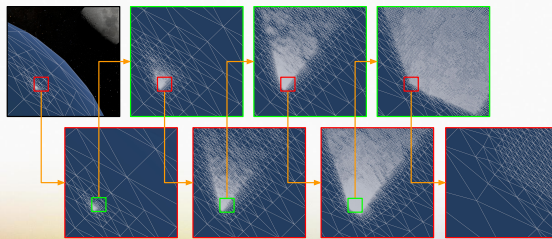
[unreadable text]



"Concurrent Binary Trees for Large-Scale Game Components" HPG24



→ full-resolution mesh requires exabytes* of data...
(* 1 **exabyte** = 1,000,000 **terabytes**)



In terms of numbers, it turns out the full resolution mesh of our Earth model, i.e., without any form of level-of-detail as we do here, would require exabytes of data.

An exabyte is a million terabytes, so it would not even fit in a large SSD.

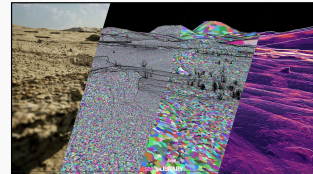
"Concurrent Binary Trees for Large-Scale Game Components" HPG24



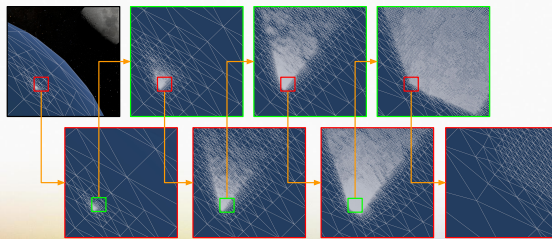
→ full-resolution mesh requires exabytes* of data...

(*) 1 **exabyte** = 1,000,000 **terabytes**

... so we can't use this:



"A Deep Dive into Nanite Virtualized" - SG21
Geometry Clusters + DAG

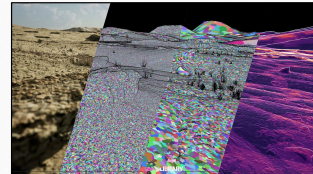


This prevents the use of LOD systems like Nanite, which requires the full resolution mesh as input.

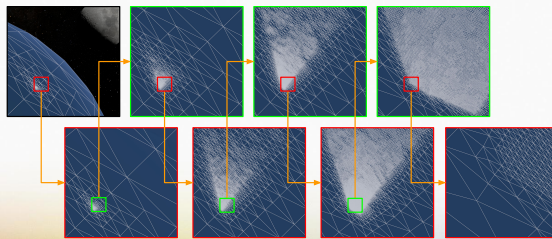
"Concurrent Binary Trees for Large-Scale Game Components" HPG24



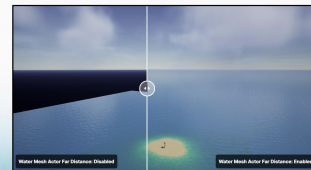
→ full-resolution mesh requires exabytes* of data...
 (*) 1 **exabyte** = 1,000,000 **terabytes**
 ... so we can't use this:



"A Deep Dive into Nanite Virtualized" - SG21
 Geometry Clusters + DAG



alternative in UE 5.3:



Unreal Engine 5.3 Water Documentation

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

Note that in UE5.3, you could still handle very large meshes using their very own large-scale game component, which consists in coupling two representations: one for close scale, and the other for far away scales.

Unfortunately, hybrid representation are hard to use especially with free-flight cameras because it becomes really tricky to set the location of the transition between both representations.

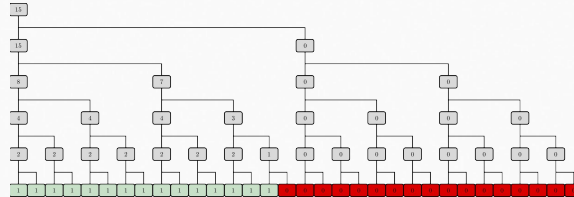
A nice advantage of our method is that it relies on a single representation so you don't need to worry about these issues.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



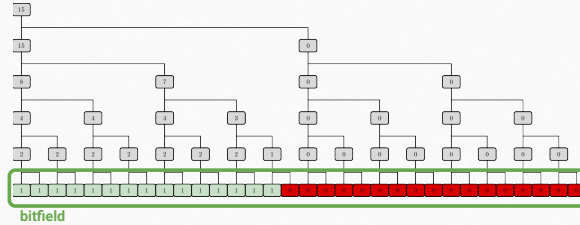
Right, now I am going to quickly explain what a concurrent binary tree (CBT) is and how our method works with them, then Anis will dive into the details.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



A concurrent binary tree of CBT is a full binary tree (so each node has exactly two children except for the leaves) with two main parts.

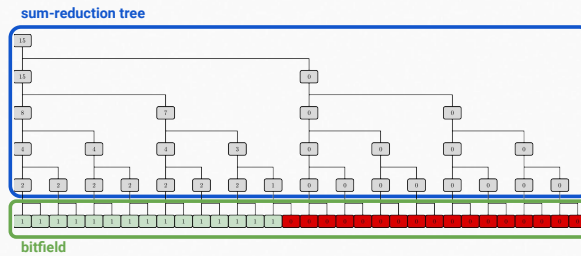
"Concurrent Binary Trees for Large-Scale Game Components" HPG24



The first part is a bitfield located at the bottom of the tree.

So the leaf nodes only store binary values.

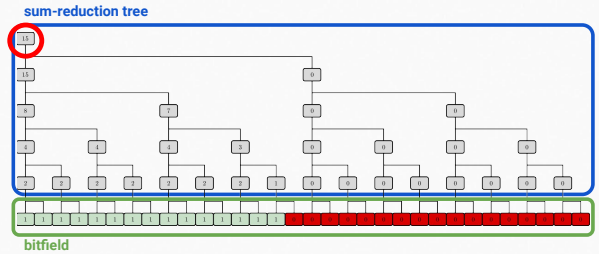
"Concurrent Binary Trees for Large-Scale Game Components" HPG24



The second part is a sum-reduction tree of this bitfield.

So all remaining nodes store the number of green bits, i.e, bits set to one, in its corresponding subtree.

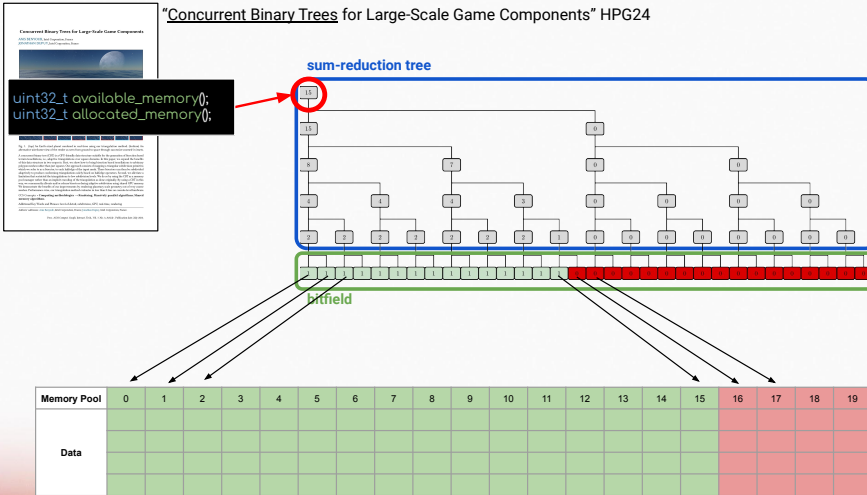
"Concurrent Binary Trees for Large-Scale Game Components" HPG24



This means that the root node gives the total number of green bits.

In addition, the sum-reduction tree makes it possible to iterate over the green bits, even if the green bits aren't located sequentially in the bitfield.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



```
uint32_t available_memory();
uint32_t allocated_memory();
```

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

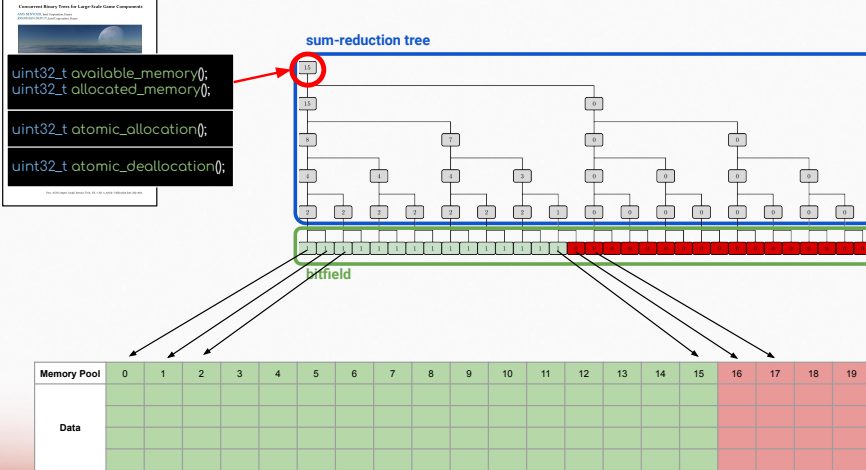
These properties makes it possible to use a CBT as memory pool manager that tracks allocated and available memory.

A memory pool is simply an array of whatever data you want to store and we set its capacity to that of the bitfield.

We then track allocated entries using green bits and available memory with red ones.

This way the root of the CBT gives us how much memory is available / allocated.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

We can then implement a simple allocation and de-allocation operator as follows.

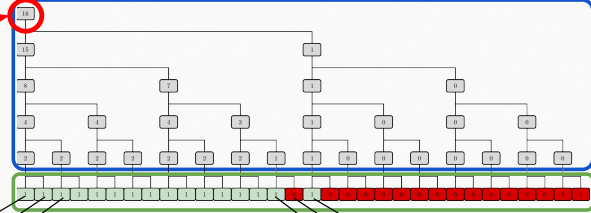
"Concurrent Binary Trees for Large-Scale Game Components" HPG24



```
uint32_t available_memory();
uint32_t allocated_memory();

uint32_t atomic_allocation();
uint32_t atomic_deallocation();
```

sum-reduction tree

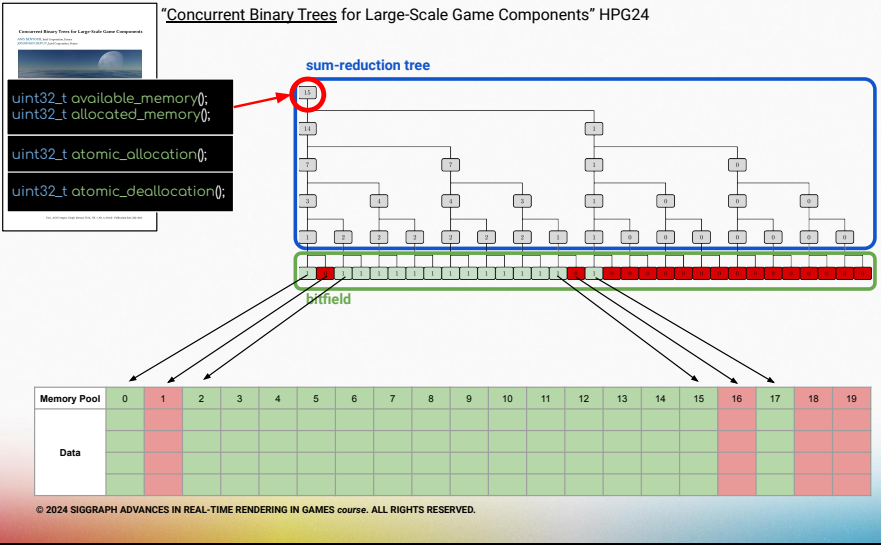


Memory Pool	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
Data																				

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

For allocation we set a bit to one and update the sum-reduction tree.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24



For de-allocation we set a bit to zero and again update the sum-reduction tree.

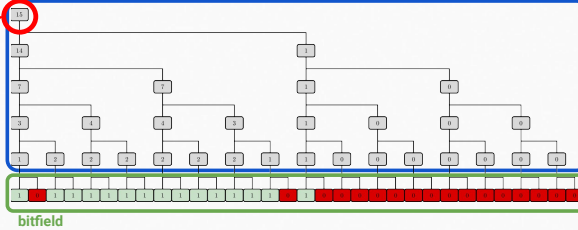
"Concurrent Binary Trees for Large-Scale Game Components" HPG24

```

uint32_t available_memory();
uint32_t allocated_memory();

uint32_t atomic_allocation();
uint32_t atomic_deallocation();
    
```

sum-reduction tree



Memory Pool	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
Data																				

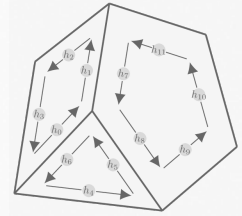
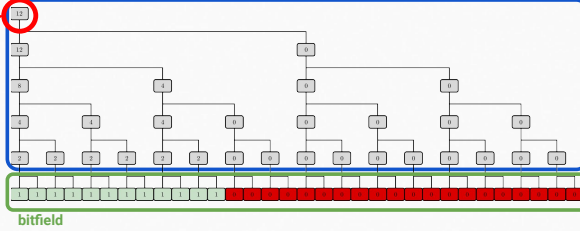
Now here is what we do for our triangulations.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24

```
uint32_t available_memory();
uint32_t allocated_memory();

uint32_t atomic_allocation();
uint32_t atomic_deallocation();
```

sum-reduction tree



Memory Pool	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
Data																				

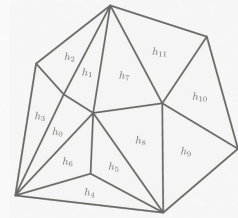
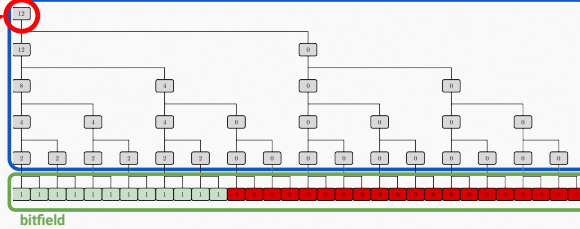
Our method takes a halfedge mesh as input...

"Concurrent Binary Trees for Large-Scale Game Components" HPG24

```
uint32_t available_memory();
uint32_t allocated_memory();

uint32_t atomic_allocation();
uint32_t atomic_deallocation();
```

sum-reduction tree



Memory Pool	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	
Data																					

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

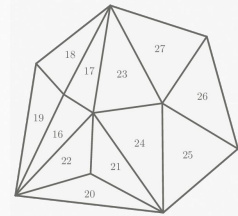
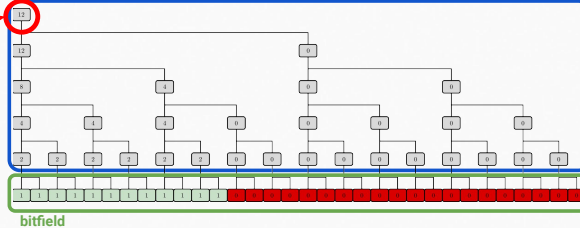
... and produces a triangle for each halfedge (feel free to go back and forth between this slide and the previous one).

"Concurrent Binary Trees for Large-Scale Game Components" HPG24

```
uint32_t available_memory();
uint32_t allocated_memory();

uint32_t atomic_allocation();
uint32_t atomic_deallocation();
```

sum-reduction tree



Memory Pool	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	
Data	16	17	18	19	20	21	22	23	24	25	26	27									

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

We then compress the 3 vertices of each triangle into a single integer value that we call a heapID and store it in a dedicated entry of the memory pool.

In this example we have 12 triangles so we require 12 slots in the memory pool.

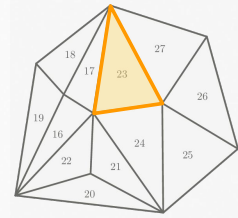
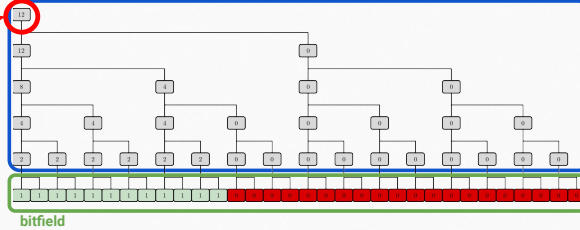
"Concurrent Binary Trees for Large-Scale Game Components" HPG24

```

uint32_t available_memory();
uint32_t allocated_memory();

uint32_t atomic_allocation();
uint32_t atomic_deallocation();
    
```

sum-reduction tree



Memory Pool	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	
Data	16	17	18	19	20	21	22	23	24	25	26	27									

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

As an example, the triangle 23 is stored in slot 7.

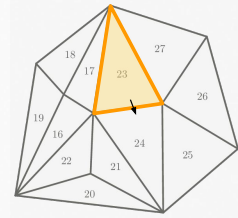
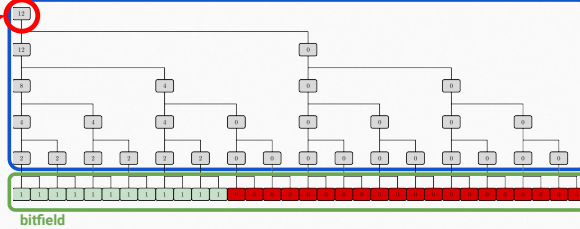
In addition to the heapID, we also store neighborhood information with pointers.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24

```
uint32_t available_memory();
uint32_t allocated_memory();

uint32_t atomic_allocation();
uint32_t atomic_deallocation();
```

sum-reduction tree



Memory Pool	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	
Data								23	24												
								8													

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

Continuing with triangle 23, the neighbors are triangle 24, which is located at slot 8...

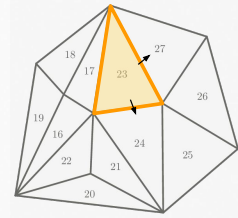
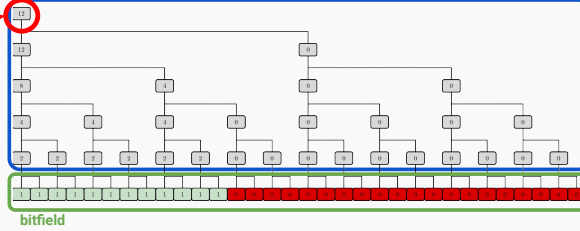
"Concurrent Binary Trees for Large-Scale Game Components" HPG24

```

uint32_t available_memory();
uint32_t allocated_memory();

uint32_t atomic_allocation();
uint32_t atomic_deallocation();
    
```

sum-reduction tree



Memory Pool	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	
Data	16	17	18	19	20	21	22	23	24	25	26	27									
							8														
							11														

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

... triangle 27, which is located at slot 11 ...

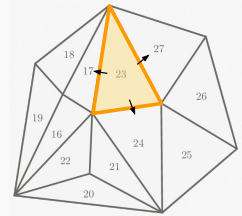
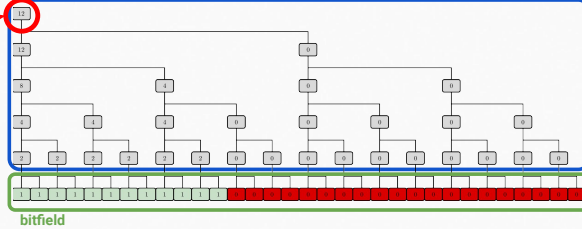
"Concurrent Binary Trees for Large-Scale Game Components" HPG24

```

uint32_t available_memory();
uint32_t allocated_memory();

uint32_t atomic_allocation();
uint32_t atomic_deallocation();
    
```

sum-reduction tree



Memory Pool	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	
Data	16	17	18	19	20	21	22	23	24	25	26	27									
								8													
								11													
								1													

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

... and triangle 17, which is located at slot 1.

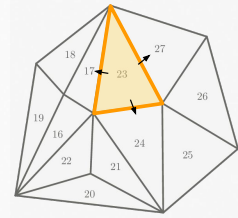
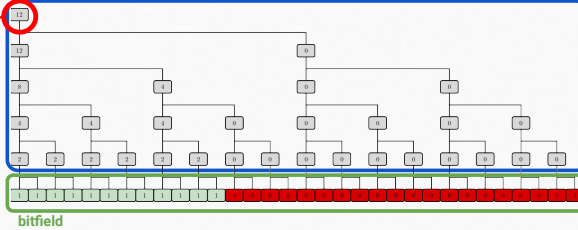
"Concurrent Binary Trees for Large-Scale Game Components" HPG24

```

uint32_t available_memory();
uint32_t allocated_memory();

uint32_t atomic_allocation();
uint32_t atomic_deallocation();
    
```

sum-reduction tree



Memory Pool	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
	16	17	18	19	20	21	22	23	24	25	26	27								
Data	1	2	3	0	5	6	4	8	9	10	11	7								
	3	0	1	2	6	4	5	11	7	8	9	10								
	6	7	null	null	null	8	0	1	14	null	null	null								

And naturally we do this for each triangle.

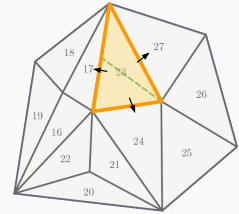
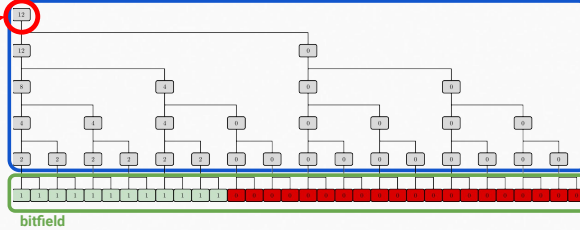
"Concurrent Binary Trees for Large-Scale Game Components" HPG24

```

uint32_t available_memory();
uint32_t allocated_memory();

uint32_t atomic_allocation();
uint32_t atomic_deallocation();
    
```

sum-reduction tree



Memory Pool	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	
	16	17	18	19	20	21	22	23	24	25	26	27									
Data	1	2	3	0	5	6	4	8	9	10	11	7									
	3	0	1	2	6	4	5	11	7	8	9	10									
	6	7	null	null	null	8	0	1	14	null	null	null									

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

The reason we store this information is because we implement a bisection scheme that can split triangles into two new ones.

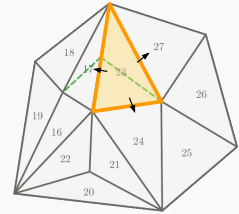
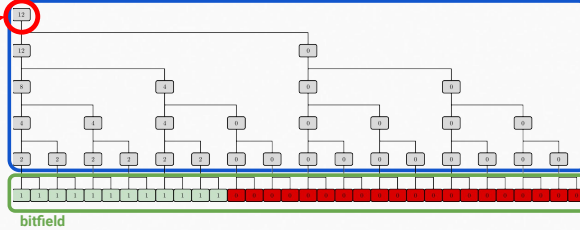
Naively bisecting a triangle would produce a T-junction, which would result in cracks in the final surface.

"Concurrent Binary Trees for Large-Scale Game Components" HPG24

```
uint32_t available_memory();
uint32_t allocated_memory();

uint32_t atomic_allocation();
uint32_t atomic_deallocation();
```

sum-reduction tree



Memory Pool	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
	16	17	18	19	20	21	22	23	24	25	26	27								
Data	1	2	3	0	5	6	4	8	9	10	11	7								
	3	0	1	2	6	4	5	11	7	8	9	10								
	6	7	null	null	null	8	0	1	14	null	null	null								

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

But thanks to the neighborhood information, we can propagate bisections across multiple triangles to guarantee crack-free surfaces.

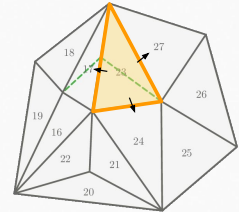
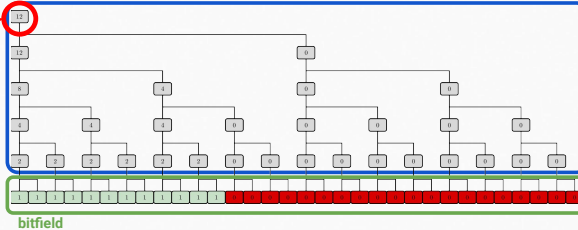
"Concurrent Binary Trees for Large-Scale Game Components" HPG24

Concurrent Binary Trees for Large-Scale Game Components

```
uint32_t available_memory();
uint32_t allocated_memory();

uint32_t atomic_allocation();
uint32_t atomic_deallocation();
```

sum-reduction tree



Note: our demo uses a 128k memory pool (requires 7 MB of memory)

Memory Pool	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	
	16	17	18	19	20	21	22	23	24	25	26	27									
Data	1	2	3	0	5	6	4	8	9	10	11	7									
	3	0	1	2	6	4	5	11	7	8	9	10									
	6	7	null	null	null	8	0	1	14	null	null	null									

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

Last thing for me: the memory pool we use for our demo is 128k wide, which requires 7 MB of memory in total.

That concludes my overview of the paper.

The key takeaway here is that CBTs provide a way to allocate, release, and iterate over all the elements of a memory pool.

And as Anis will show, all this can be done efficiently on the GPU.

Your turn Anis!



The slide features a dark background with a rainbow gradient at the top left. On the left side, there are two images: the top one shows a highway sign with 'EXIT 12' and 'LEFT' above it, and a green sign with 'Overview' and 'Implementation details' with an arrow pointing right; the bottom one shows a car on a road with the text 'This presentation' overlaid. On the right side, the SIGGRAPH 2024 logo is displayed, followed by the text 'SIGGRAPH 2024 DENVER+ 28 JUL - 1 AUG'. Below this, the main title 'Parallel Update Routine' is shown in white, with '(Simplified implementation details)' in green underneath. At the bottom left, there is a small copyright notice: '© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.' and a small rainbow logo at the bottom right.

SIGGRAPH 2024
DENVER+ 28 JUL - 1 AUG

Parallel Update Routine

(Simplified implementation details)

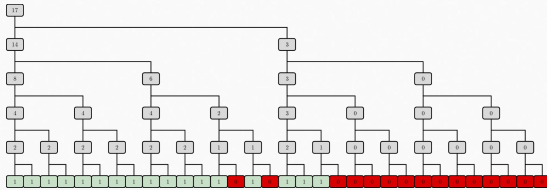
This presentation

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

We're going to dig into some of the implementation details that allowed us to reach reasonable performance numbers with this method.

That said, we'll not go into code details due to the limited time we have.

The full source code of the demo is released and is available for you to explore and play with it!



```
// Store in shared memory, atomic friendly
groupshared uint32_t gs_cbt[      ];
```

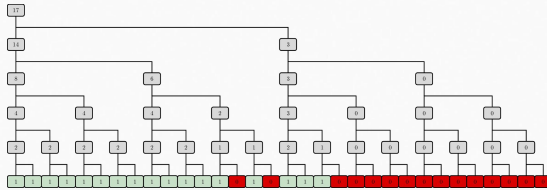
First, We'll precise the implementation of the CBT.

We'd like to store the CBT in the group shared memory to make the access more efficient when reading and writing.

As an example, we'll use the 128k elements CBT in this presentation, but as we mentioned before we can go higher (or lower) depending on the needs of the application

The CBT needs to be readable and writable from worker threads. We could use `uint32_t` to benefit from the atomic intrinsics, but there is an issue with that:

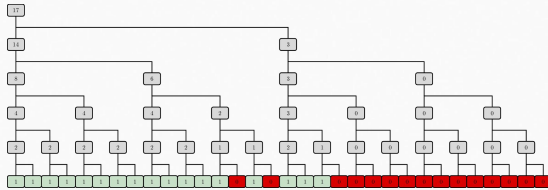
Parallel Update Routine (CBT Implementation 128K)



```
// Store in shared memory, atomic friendly
const uint32_t cbt_num_nodes = 2 * (128 * 1024); // 262144
groupshared uint32_t gs_cbt[cbt_num_nodes]; // 1 MB
```

The number of nodes of the CBT is twice the size of the bitfield, which would, if naively stored would be around one 1MB

Parallel Update Routine (CBT Implementation 128K)

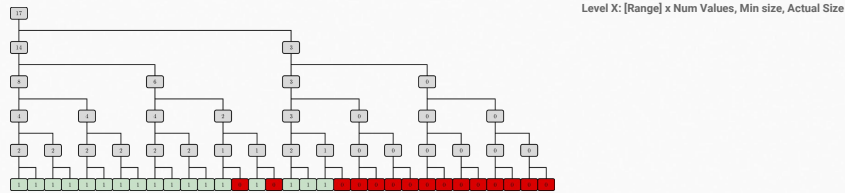


```
// Store in shared memory, atomic friendly
const uint32_t cbt_num_nodes = 2 * (128 * 1024); // 262144
groupshared uint32_t gs_cbt[cbt_num_nodes]; // 1 MB
```

Limited to 32 KB on DX12 !!!

However, the group shared memory storage is limited to 32KB on dx12 so we need a better representation!

Parallel Update Routine (CBT Implementation 128K)

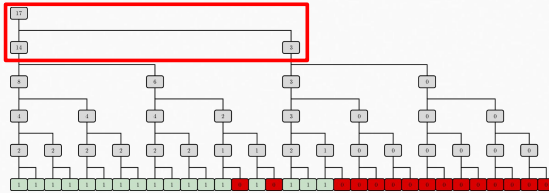


Each level of the tree is defined by 4 things:

- The number of nodes within the level
- The range of values each node can represent
- The minimal size of each node to represent that range
- The final size used to represent each node

So what we'll do is decompose the tree into multiple parts with different constraints

Parallel Update Routine (CBT Implementation 128K)



Level X: [Range] x Num Values, Min size, Actual Size
Level 0: [0, 131072] x 1, Min 18 bits (rounded up to 32 bits for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to 32 bits for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to 32 bits for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to 32 bits for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to 32 bits for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to 32 bits for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to 32 bits for alignment and atomic operations)

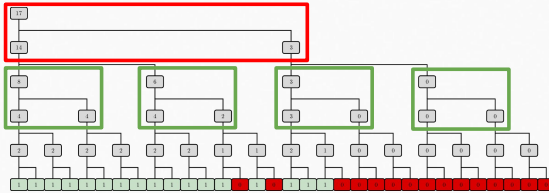
- Atomic-friendly Sub-tree (uint32_t aligned, group shared memory)

The first part in red that is an atomic-friendly subtree

During modification, all threads can write safely to any node of the tree.

Each node is represented by uint32_t to benefit from the atomic intrinsics. This is stored in the group shared memory;

Parallel Update Routine (CBT Implementation 128K)



Level X: [Range] x Num Values, Min size, Actual Size
Level 0: [0, 131072] x 1, Min 18 bits (rounded up to 32 bits for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to 32 bits for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to 32 bits for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to 32 bits for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to 32 bits for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to 32 bits for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to 32 bits for alignment and atomic operations)

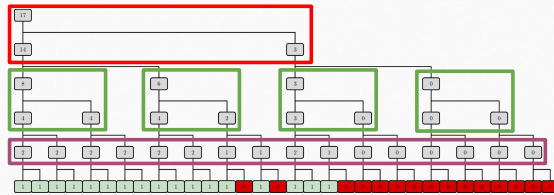
Level 7: [0, 1024] x 128, Min 11 bits (rounded up to 16 bits for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to 16 bits for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to 16 bits for alignment)
Level 10: [0, 128] x 1024, Min 8 bits

- Atomic-friendly Sub-tree (uint32_t aligned, group shared memory)
- Thread Sub-trees (Min size, aligned on powers of 2, group shared memory)

in green, we have several subtrees, each subtree will only be modified by one thread at a time.

The size of each node is rounded to the closest power of two to represent the data underneath. This is also stored in the group shared memory;

Parallel Update Routine (CBT Implementation 128K)



- Atomic-friendly Sub-tree (uint32_t aligned, group shared memory)
- Thread Sub-trees (Min size, aligned on powers of 2, group shared memory)
- Virtual Tree (not explicitly represented)

Level X: [Range] x Num Values, Min size, Actual Size

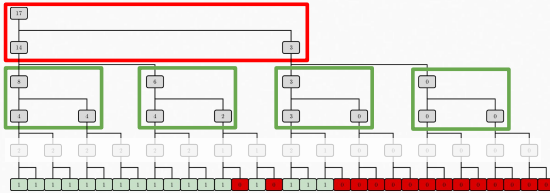
- Level 0: [0, 131072] x 1, Min 18 bits (rounded up to 32 bits for alignment and atomic operations)
- Level 1: [0, 65536] x 2, Min 17 bits (rounded up to 32 bits for alignment and atomic operations)
- Level 2: [0, 32768] x 4, Min 16 bits (bumped to 32 bits for atomic operations)
- Level 3: [0, 16384] x 8, Min 15 bits (rounded up to 32 bits for alignment and atomic operations)
- Level 4: [0, 8192] x 16, Min 14 bits (rounded up to 32 bits for alignment and atomic operations)
- Level 5: [0, 4096] x 32, Min 13 bits (rounded up to 32 bits for alignment and atomic operations)
- Level 6: [0, 2048] x 64, Min 12 bits (rounded up to 32 bits for alignment and atomic operations)

- Level 7: [0, 1024] x 128, Min 11 bits (rounded up to 16 bits for alignment)
- Level 8: [0, 512] x 256, Min 10 bits (rounded up to 16 bits for alignment)
- Level 9: [0, 256] x 512, Min 9 bits (rounded up to 16 bits for alignment)
- Level 10: [0, 128] x 1024, Min 8 bits

- Level 11: [0, 64] x 2048, Min 7 bits (Virtual)
- Level 12: [0, 32] x 4096, Min 6 bits (Virtual)
- Level 13: [0, 16] x 16384, Min 5 bits (Virtual)
- Level 14: [0, 8] x 32768, Min 4 bits (Virtual)
- Level 15: [0, 4] x 65536, Min 3 bits (Virtual)
- Level 16: [0, 2] x 131072, Min 2 bits (Virtual)

Then there are what we call virtual levels. These are not represented explicitly,

Parallel Update Routine (CBT Implementation 128K)



- Atomic-friendly Sub-tree (uint32_t aligned, group shared memory)
- Thread Sub-trees (Min size, aligned on powers of 2, group shared memory)
- Virtual Tree (not explicitly represented)

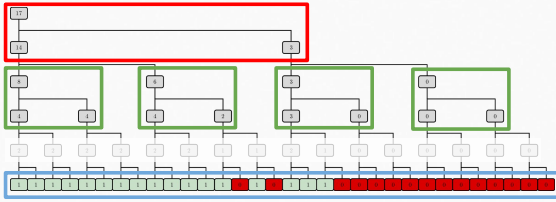
Level X: [Range] x Num Values, Min size, Actual Size
Level 0: [0, 131072] x 1, Min 18 bits (rounded up to 32 bits for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to 32 bits for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to 32 bits for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to 32 bits for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to 32 bits for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to 32 bits for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to 32 bits for alignment and atomic operations)

Level 7: [0, 1024] x 128, Min 11 bits (rounded up to 16 bits for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to 16 bits for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to 16 bits for alignment)
Level 10: [0, 128] x 1024, Min 8 bits

Level 11: [0, 64] x 2048, Min 7 bits (Virtual)
 Level 12: [0, 32] x 4096, Min 6 bits (Virtual)
 Level 13: [0, 16] x 16384, Min 5 bits (Virtual)
 Level 14: [0, 8] x 32768, Min 4 bits (Virtual)
 Level 15: [0, 4] x 65536, Min 3 bits (Virtual)
 Level 16: [0, 2] x 131072, Min 2 bits (Virtual)

but can be deduced from the last level

Parallel Update Routine (CBT Implementation 128K)



- Atomic-friendly Sub-tree (uint32_t aligned, group shared memory)
- Thread Sub-trees (Min size, aligned on powers of 2, group shared memory)
- Virtual Tree (not explicitly represented)
- Raw Bitfield Buffer (uint64_t, StructuredBuffer)

Level X: [Range] x Num Values, Min size, Actual Size
Level 0: [0, 131072] x 1, Min 18 bits (rounded up to 32 bits for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to 32 bits for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to 32 bits for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to 32 bits for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to 32 bits for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to 32 bits for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to 32 bits for alignment and atomic operations)

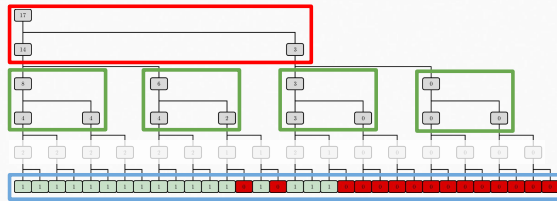
Level 7: [0, 1024] x 128, Min 11 bits (rounded up to 16 bits for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to 16 bits for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to 16 bits for alignment)
Level 10: [0, 128] x 1024, Min 8 bits

Level 11: [0, 64] x 2048, Min 7 bits (Virtual)
 Level 12: [0, 32] x 4096, Min 6 bits (Virtual)
 Level 13: [0, 16] x 16384, Min 5 bits (Virtual)
 Level 14: [0, 8] x 32768, Min 4 bits (Virtual)
 Level 15: [0, 4] x 65536, Min 3 bits (Virtual)
 Level 16: [0, 2] x 131072, Min 2 bits (Virtual)

Level 17: Raw 64 bits representation

Which is the rawbitfield. This bitfield it is represented using uint64_t to benefit from some intrinsic functions and is stored in a structured buffer outside of the shared memory.

Parallel Update Routine (CBT Implementation 128K)



Level X: [Range] x Num Values, Min size, Actual Size
Level 0: [0, 131072] x 1, Min 18 bits (rounded up to 32 bits for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to 32 bits for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to 32 bits for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to 32 bits for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to 32 bits for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to 32 bits for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to 32 bits for alignment and atomic operations)
Level 7: [0, 1024] x 128, Min 11 bits (rounded up to 16 bits for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to 16 bits for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to 16 bits for alignment)
Level 10: [0, 128] x 1024, Min 8 bits
Level 17: Raw 64 bits representation

- Atomic-friendly Sub-tree (uint32_t aligned, group shared memory)
- Thread Sub-trees (Min size, aligned on powers of 2, group shared memory)
- Virtual Tree (not explicitly represented)
- Raw Bitfield Buffer (uint64_t, StructuredBuffer)

And this the memory footprint of a 128k bit CBT. In the implementation, you'll find the layout for the other sizes of the CBT

Parallel Update Routine (CBT Implementation 128K)



```
// Graphics Buffers that hold the CBT
```

Level 0: [0, 131072] x 1, Min 18 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to **32 bits** for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to **32 bits** for alignment and atomic operations)

Level 7: [0, 1024] x 128, Min 11 bits (rounded up to **16 bits** for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to **16 bits** for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to **16 bits** for alignment)
Level 10: [0, 128] x 1024, Min **8 bits**

Level 17: Raw 64 bits representation

Now let's look in practice what that maps to:

Parallel Update Routine (CBT Implementation 128K)

```
// Graphics Buffers that hold the CBT
```

```
RWStructuredBuffer<uint32_t> TreeBufferRW.register(CBT_BUFFER0_BINDING_SLOT); // RED + GREEN
```

```
RWStructuredBuffer<uint64_t> BitfieldBufferRW.register(CBT_BUFFER1_BINDING_SLOT); // BLUE
```

Level 0: [0, 131072] x 1, Min 18 bits (rounded up to 32 bits for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to 32 bits for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to 32 bits for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to 32 bits for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to 32 bits for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to 32 bits for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to 32 bits for alignment and atomic operations)

Level 7: [0, 1024] x 128, Min 11 bits (rounded up to 16 bits for alignment)

Level 8: [0, 512] x 256, Min 10 bits (rounded up to 16 bits for alignment)

Level 9: [0, 256] x 512, Min 9 bits (rounded up to 16 bits for alignment)

Level 10: [0, 128] x 1024, Min 8 bits

Level 17: Raw 64 bits representation

We have two structured buffers:

- The first one stores `uint32_t` and contains the red and green parts of the tree
- The second stores `uint64_t` and contains the bitfield (in blue)

Parallel Update Routine (CBT Implementation 128K)



```
// Graphics Buffers that hold the CBT
RWStructuredBuffer<uint32_t> _TreeBufferRW.register(CBT_BUFFER0_BINDING_SLOT); // RED + GREEN
RWStructuredBuffer<uint64_t> _BitfieldBufferRW.register(CBT_BUFFER1_BINDING_SLOT); // BLUE
RWStructuredBuffer<uint32_t> _AllocationBufferRW.register(CBT_BUFFER2_BINDING_SLOT);
```

Level 0: [0, 131072] x 1, Min 18 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to **32 bits** for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to **32 bits** for alignment and atomic operations)

Level 7: [0, 1024] x 128, Min 11 bits (rounded up to **16 bits** for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to **16 bits** for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to **16 bits** for alignment)
Level 10: [0, 128] x 1024, Min **8 bits**

Level 17: Raw 64 bits representation

In addition to that we need an allocation buffer which is a uint32_t buffer, we'll see how it is used in a second

Parallel Update Routine (CBT Implementation 128K)



```
// Graphics Buffers that hold the CBT
RWStructuredBuffer<uint32_t> _TreeBufferRW.register(CBT_BUFFER0_BINDING_SLOT); // RED + GREEN
RWStructuredBuffer<uint64_t> _BitfieldBufferRW.register(CBT_BUFFER1_BINDING_SLOT); // BLUE
RWStructuredBuffer<uint32_t> _AllocationBufferRW.register(CBT_BUFFER2_BINDING_SLOT);

// _TreeBufferRW in the shared memory
const uint32_t tree_num_slots = (1 * 32 + 2 * 32 + ... + 1024 * 8) / 32;
```

Level 0: [0, 131072] x 1, Min 18 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to **32 bits** for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to **32 bits** for alignment and atomic operations)

Level 7: [0, 1024] x 128, Min 11 bits (rounded up to **16 bits** for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to **16 bits** for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to **16 bits** for alignment)
Level 10: [0, 128] x 1024, Min **8 bits**

Level 17: Raw 64 bits representation

Using the memory footprint table on the right, we can deduce the size of the tree

Parallel Update Routine (CBT Implementation 128K)



```
// Graphics Buffers that hold the CBT
RWStructuredBuffer<uint32_t> _TreeBufferRW.register(CBT_BUFFER0_BINDING_SLOT); // RED + GREEN
RWStructuredBuffer<uint64_t> _BitfieldBufferRW.register(CBT_BUFFER1_BINDING_SLOT); // BLUE
RWStructuredBuffer<uint32_t> _AllocationBufferRW.register(CBT_BUFFER2_BINDING_SLOT);

// _TreeBufferRW in the shared memory
const uint32_t tree_num_slots = (1 * 32 + 2 * 32 + ... + 1024 * 8) / 32;
groupshared uint32_t gs_cbtTree[tree_num_slots]; // ~3KB
```

Level 0: [0, 131072] x 1, Min 18 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to **32 bits** for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to **32 bits** for alignment and atomic operations)

Level 7: [0, 1024] x 128, Min 11 bits (rounded up to **16 bits** for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to **16 bits** for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to **16 bits** for alignment)
Level 10: [0, 128] x 1024, Min **8 bits**

Level 17: Raw 64 bits representation

And that defines the group shared memory space our CBT occupies. In this case the CBT is about 3 KB, which we can consider reasonable.

Parallel Update Routine (CBT Implementation 128K)



```
// Graphics Buffers that hold the CBT
RWStructuredBuffer<uint32_t> _TreeBufferRW.register(CBT_BUFFER0_BINDING_SLOT); // RED + GREEN
RWStructuredBuffer<uint64_t> _BitfieldBufferRW.register(CBT_BUFFER1_BINDING_SLOT); // BLUE
RWStructuredBuffer<uint32_t> _AllocationBufferRW.register(CBT_BUFFER2_BINDING_SLOT);

// _TreeBufferRW in the shared memory
const uint32_t tree_num_slots = (1 * 32 + 2 * 32 + ... + 1024 * 8) / 32;
groupshared uint32_t gs_cbtTree[tree_num_slots]; // ~3KB
```

Level 0: [0, 131072] x 1, Min 18 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to **32 bits** for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to **32 bits** for alignment and atomic operations)

Level 7: [0, 1024] x 128, Min 11 bits (rounded up to **16 bits** for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to **16 bits** for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to **16 bits** for alignment)
Level 10: [0, 128] x 1024, Min **8 bits**

Level 17: Raw 64 bits representation

Book & Allocate & Cancel

The way the memory manager is used every frame is the following.

- First we're gonna book the worst case memory we need
- Then we're gonna allocate a bunch of bits, that will depend on the subdivision case we're processing,
- We'll finally return to the memory manager memory space what we didn't consume

Parallel Update Routine (CBT Implementation 128K)

```
// Graphics Buffers that hold the CBT
RWStructuredBuffer<uint32_t> _TreeBufferRW; register(CBT_BUFFER0_BINDING_SLOT); // RED + GREEN
RWStructuredBuffer<uint64_t> _BitfieldBufferRW; register(CBT_BUFFER1_BINDING_SLOT); // BLUE
RWStructuredBuffer<uint32_t> _AllocationBufferRW; register(CBT_BUFFER2_BINDING_SLOT);

// _TreeBufferRW in the shared memory
const uint32_t tree_num_slots = (1 * 32 + 2 * 32 + ... + 1024 * 8) / 32;
groupshared uint32_t gs_cbtTree[tree_num_slots]; // ~3KB

// Memory booking
void cbt_rev_frame0();
void book_memory_space(uint32_t numSlots);
void cancel_memory_booking(uint32_t numSlots);
```

Level 0: [0, 131072] x 1, Min 18 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to **32 bits** for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to **32 bits** for alignment and atomic operations)

Level 7: [0, 1024] x 128, Min 11 bits (rounded up to **16 bits** for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to **16 bits** for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to **16 bits** for alignment)
Level 10: [0, 128] x 1024, Min **8 bits**

Level 17: Raw 64 bits representation

Book & Allocate & Cancel

The booking and cancelling is done using these functions

Parallel Update Routine (CBT Implementation 128K)

```
// Graphics Buffers that hold the CBT
RWStructuredBuffer<uint32_t> _TreeBufferRW; register(CBT_BUFFER0_BINDING_SLOT); // RED + GREEN
RWStructuredBuffer<uint64_t> _BitfieldBufferRW; register(CBT_BUFFER1_BINDING_SLOT); // BLUE
RWStructuredBuffer<uint32_t> _AllocationBufferRW; register(CBT_BUFFER2_BINDING_SLOT);

// _TreeBufferRW in the shared memory
const uint32_t tree_num_slots = (1 * 32 + 2 * 32 + ... + 1024 * 8) / 32;
groupshared uint32_t gs_cbtTree[tree_num_slots]; // ~3KB

// Memory booking
void cbt_new_frame();
void book_memory_space(uint32_t numSlots);
void cancel_memory_booking(uint32_t numSlots);

// Allocation
uint32_t allocate_next_available_slot(uint32_t bisectorID);
```

Level 0: [0, 131072] x 1, Min 18 bits (rounded up to 32 bits for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to 32 bits for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to 32 bits for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to 32 bits for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to 32 bits for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to 32 bits for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to 32 bits for alignment and atomic operations)

Level 7: [0, 1024] x 128, Min 11 bits (rounded up to 16 bits for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to 16 bits for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to 16 bits for alignment)
Level 10: [0, 128] x 1024, Min 8 bits

Level 17: Raw 64 bits representation

Book & Allocate & Cancel

And this function that allows us to allocate the next available memory slot

Parallel Update Routine (CBT Implementation 128K)

```
// Graphics Buffers that hold the CBT
RWStructuredBuffer<uint32_t> _TreeBufferRW; register(CBT_BUFFER0_BINDING_SLOT); // RED + GREEN
RWStructuredBuffer<uint64_t> _BitfieldBufferRW; register(CBT_BUFFER1_BINDING_SLOT); // BLUE
RWStructuredBuffer<uint32_t> _AllocationBufferRW; register(CBT_BUFFER2_BINDING_SLOT);

// _TreeBufferRW in the shared memory
const uint32_t tree_num_slots = (1 * 32 + 2 * 32 + ... + 1024 * 8) / 32;
groupshared uint32_t gs_cbtTree[tree_num_slots]; // ~3KB

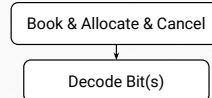
// Memory booking
void cbt_new_frame();
void book_memory_space(uint32_t numSlots);
void cancel_memory_booking(uint32_t numSlots);

// Allocation
uint32_t allocate_next_available_slot(uint32_t bisectorID);
```

Level 0: [0, 131072] x 1, Min 18 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to **32 bits** for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to **32 bits** for alignment and atomic operations)

Level 7: [0, 1024] x 128, Min 11 bits (rounded up to **16 bits** for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to **16 bits** for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to **16 bits** for alignment)
Level 10: [0, 128] x 1024, Min **8 bits**

Level 17: Raw 64 bits representation



Then once all the allocations have been done, we're gonna decode the location of the bits that we've allocated

Parallel Update Routine (CBT Implementation 128K)

```
// Graphics Buffers that hold the CBT
RWStructuredBuffer<uint32_t> _TreeBufferRW; register(CBT_BUFFER0_BINDING_SLOT); // RED + GREEN
RWStructuredBuffer<uint64_t> _BitfieldBufferRW; register(CBT_BUFFER1_BINDING_SLOT); // BLUE
RWStructuredBuffer<uint32_t> _AllocationBufferRW; register(CBT_BUFFER2_BINDING_SLOT);

// _TreeBufferRW in the shared memory
const uint32_t tree_num_slots = (1 * 32 + 2 * 32 + ... + 1024 * 8) / 32;
groupshared uint32_t gs_cbtTree[tree_num_slots]; // ~3KB

// Memory booking
void cbt_new_frame();
void book_memory_space(uint32_t numSlots);
void cancel_memory_booking(uint32_t numSlots);

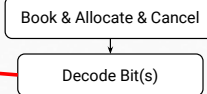
// Allocation
uint32_t allocate_next_available_slot(uint32_t bisectorID);

// Load and export
void load_buffer_to_shared_memory(uint32_t groupIndex);
```

Level 0: [0, 131072] x 1, Min 18 bits (rounded up to 32 bits for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to 32 bits for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to 32 bits for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to 32 bits for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to 32 bits for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to 32 bits for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to 32 bits for alignment and atomic operations)

Level 7: [0, 1024] x 128, Min 11 bits (rounded up to 16 bits for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to 16 bits for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to 16 bits for alignment)
Level 10: [0, 128] x 1024, Min 8 bits

Level 17: Raw 64 bits representation



To do so, we first have to load the CBT into shared memory for our searches, we do it using this function, this will copy per workgroup the tree buffer into the group shader memory

Parallel Update Routine (CBT Implementation 128K)

```
// Graphics Buffers that hold the CBT
RWStructuredBuffer<uint32_t> _TreeBufferRW; register(CBT_BUFFER0_BINDING_SLOT); // RED + GREEN
RWStructuredBuffer<uint64_t> _BitfieldBufferRW; register(CBT_BUFFER1_BINDING_SLOT); // BLUE
RWStructuredBuffer<uint32_t> _AllocationBufferRW; register(CBT_BUFFER2_BINDING_SLOT);

// _TreeBufferRW in the shared memory
const uint32_t tree_num_slots = (1 * 32 + 2 * 32 + ... + 1024 * 8) / 32;
groupshared uint32_t gs_cbtTree[tree_num_slots]; // ~3KB

// Memory booking
void cbt_new_frame();
void book_memory_space(uint32_t numSlots);
void cancel_memory_booking(uint32_t numSlots);

// Allocation
uint32_t allocate_next_available_slot(uint32_t bisectorID);

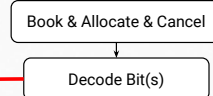
// Load and export
void load_buffer_to_shared_memory(uint32_t groupIndex);

// Find i-th bit set to zero
uint32_t decode_bit_complement(uint32_t index);
```

Level 0: [0, 131072] x 1, Min 18 bits (rounded up to 32 bits for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to 32 bits for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to 32 bits for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to 32 bits for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to 32 bits for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to 32 bits for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to 32 bits for alignment and atomic operations)

Level 7: [0, 1024] x 128, Min 11 bits (rounded up to 16 bits for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to 16 bits for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to 16 bits for alignment)
Level 10: [0, 128] x 1024, Min 8 bits

Level 17: Raw 64 bits representation



Then we actually operate the tree descents taking advantage of the very specific memory layout of the CBT which helps us to accelerate significantly the operation.

Parallel Update Routine (CBT Implementation 128K)

```
// Graphics Buffers that hold the CBT
RWStructuredBuffer<uint32_t> _TreeBufferRW; register(CBT_BUFFER0_BINDING_SLOT); // RED + GREEN
RWStructuredBuffer<uint64_t> _BitfieldBufferRW; register(CBT_BUFFER1_BINDING_SLOT); // BLUE
RWStructuredBuffer<uint32_t> _AllocationBufferRW; register(CBT_BUFFER2_BINDING_SLOT);

// _TreeBufferRW in the shared memory
const uint32_t tree_num_slots = (1 * 32 + 2 * 32 + ... + 1024 * 8) / 32;
groupshared uint32_t gs_cbtTree[tree_num_slots]; // ~3KB

// Memory booking
void cbt_new_frame();
void book_memory_space(uint32_t numSlots);
void cancel_memory_booking(uint32_t numSlots);

// Allocation
uint32_t allocate_next_available_slot(uint32_t bisectorID);

// Load and export
void load_buffer_to_shared_memory(uint32_t groupIndex);

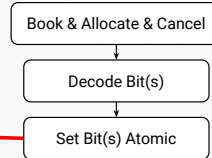
// Find i-th bit set to zero
uint32_t decode_bit_complement(uint32_t index);

// Atomic operations
void set_bit_atomic(uint32_t bitID, bool state);
```

Level 0: [0, 131072] x 1, Min 18 bits (rounded up to 32 bits for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to 32 bits for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to 32 bits for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to 32 bits for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to 32 bits for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to 32 bits for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to 32 bits for alignment and atomic operations)

Level 7: [0, 1024] x 128, Min 11 bits (rounded up to 16 bits for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to 16 bits for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to 16 bits for alignment)
Level 10: [0, 128] x 1024, Min 8 bits

Level 17: Raw 64 bits representation



Then we'll set atomically some bit to 1 (for the allocations that we've done) and 0 for the bit's we've deallocated.

It is done using this function that will operate directly on the bitfield buffer.

Parallel Update Routine (CBT Implementation 128K)

```
// Graphics Buffers that hold the CBT
RWStructuredBuffer<uint32_t> _TreeBufferRW; register(CBT_BUFFER0_BINDING_SLOT); // RED + GREEN
RWStructuredBuffer<uint64_t> _BitfieldBufferRW; register(CBT_BUFFER1_BINDING_SLOT); // BLUE
RWStructuredBuffer<uint32_t> _AllocationBufferRW; register(CBT_BUFFER2_BINDING_SLOT);

// _TreeBufferRW in the shared memory
const uint32_t tree_num_slots = (1 * 32 + 2 * 32 + ... + 1024 * 8) / 32;
groupshared uint32_t gs_cbtTree[tree_num_slots]; // ~3KB

// Memory booking
void cbt_new_frame();
void book_memory_space(uint32_t numSlots);
void cancel_memory_booking(uint32_t numSlots);

// Allocation
uint32_t allocate_next_available_slot(uint32_t bisectorID);

// Load and export
void load_buffer_to_shared_memory(uint32_t groupIndex);

// Find i-th bit set to zero
uint32_t decode_bit_complement(uint32_t index);

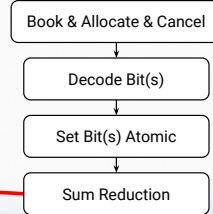
// Atomic operations
void set_bit_atomic(uint32_t bitID, bool state);

// Reduction (Detailed later)
void reduce(...);
```

Level 0: [0, 131072] x 1, Min 18 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to **32 bits** for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to **32 bits** for alignment and atomic operations)

Level 7: [0, 1024] x 128, Min 11 bits (rounded up to **16 bits** for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to **16 bits** for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to **16 bits** for alignment)
Level 10: [0, 128] x 1024, Min **8 bits**

Level 17: Raw 64 bits representation



And finally we'll operate a sum reduction on the tree to be able to return to a valid state. We can then take advantage of the memory layout we defined to get it performant, but i'll cover it a tad later.

Parallel Update Routine (CBT Implementation 128K)

```
// Graphics Buffers that hold the CBT
RWStructuredBuffer<uint32_t> _TreeBufferRW; register(CBT_BUFFER0_BINDING_SLOT); // RED + GREEN
RWStructuredBuffer<uint64_t> _BitfieldBufferRW; register(CBT_BUFFER1_BINDING_SLOT); // BLUE
RWStructuredBuffer<uint32_t> _AllocationBufferRW; register(CBT_BUFFER2_BINDING_SLOT);

// _TreeBufferRW in the shared memory
const uint32_t tree_num_slots = (1 * 32 + 2 * 32 + ... + 1024 * 8) / 32;
groupshared uint32_t gs_cbtTree[tree_num_slots]; // ~3KB

// Memory booking
void cbt_new_frame();
void book_memory_space(uint32_t numSlots);
void cancel_memory_booking(uint32_t numSlots);

// Allocation
uint32_t allocate_next_available_slot(uint32_t bisectorID);

// Load and export
void load_buffer_to_shared_memory(uint32_t groupIndex);

// Find i-th bit set to zero
uint32_t decode_bit_complement(uint32_t index);

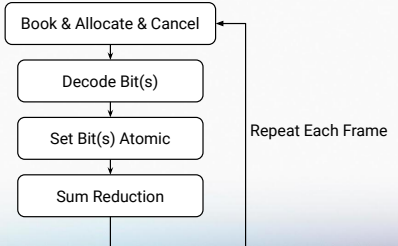
// Atomic operations
void set_bit_atomic(uint32_t bitID, bool state);

// Reduction (Detailed later)
void reduce(...);
```

Level 0: [0, 131072] x 1, Min 18 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 1: [0, 65536] x 2, Min 17 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 2: [0, 32768] x 4, Min 16 bits (bumped to **32 bits** for atomic operations)
Level 3: [0, 16384] x 8, Min 15 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 4: [0, 8192] x 16, Min 14 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 5: [0, 4096] x 32, Min 13 bits (rounded up to **32 bits** for alignment and atomic operations)
Level 6: [0, 2048] x 64, Min 12 bits (rounded up to **32 bits** for alignment and atomic operations)

Level 7: [0, 1024] x 128, Min 11 bits (rounded up to **16 bits** for alignment)
Level 8: [0, 512] x 256, Min 10 bits (rounded up to **16 bits** for alignment)
Level 9: [0, 256] x 512, Min 9 bits (rounded up to **16 bits** for alignment)
Level 10: [0, 128] x 1024, Min **8 bits**

Level 17: Raw 64 bits representation



And we'll repeat this operation every frame and that defines how the CBT works as a parallel memory manager

Parallel Update Routine (Memory Pool)

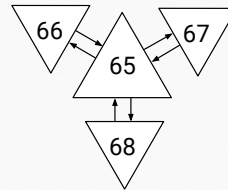
```
// Graphics Buffers that hold the Memory pool
```



On the other side, we need to specify what is the memory pool like and how it's used

Parallel Update Routine (Memory Pool)

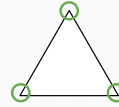
```
// Graphics Buffers that hold the Memory pool  
RWStructuredBuffer<uint54> _HeapIDBufferRW; register(PPOOL_BUFFER0_BINDING_SLOT);  
RWStructuredBuffer<uint3> _NeighborsBufferRW; register(PPOOL_BUFFER1_BINDING_SLOT);
```



As mentioned before, we need two buffers to store the HeapID and Neighbors

Parallel Update Routine (Memory Pool)

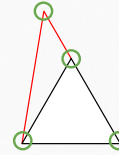
```
// Graphics Buffers that hold the Memory pool
RWStructuredBuffer<uint3> _HeapIDBufferRW.register(PPOOL_BUFFER0_BINDING_SLOT);
RWStructuredBuffer<uint3> _NeighborsBufferRW.register(PPOOL_BUFFER1_BINDING_SLOT);
RWStructuredBuffer<float3> _PositionRWSBufferRW.register(PPOOL_BUFFER2_BINDING_SLOT);
```



Then we also need to be able to store the camera relative position. Obviously, there are the three positions of the vertices of the triangles

Parallel Update Routine (Memory Pool)

```
// Graphics Buffers that hold the Memory pool
RWStructuredBuffer<uint3> _HeapIDBufferRW.register(PPOOL_BUFFER0_BINDING_SLOT);
RWStructuredBuffer<uint3> _NeighborsBufferRW.register(PPOOL_BUFFER1_BINDING_SLOT);
RWStructuredBuffer<float3> _PositionRWBufferRW.register(PPOOL_BUFFER2_BINDING_SLOT);
```



4 positions per bisector

But we actually need to store a 4th position which allows to rebuild the parent triangle and is required for split/merge decisions

Parallel Update Routine (Memory Pool)

```
// Graphics Buffers that hold the Memory pool
RWStructuredBuffer<uint> _HeapIDBufferRW; register(PPOOL_BUFFER0_BINDING_SLOT);
RWStructuredBuffer<uint3> _NeighborsBufferRW; register(PPOOL_BUFFER1_BINDING_SLOT);
RWStructuredBuffer<float3> _PositionRWSBufferRW; register(PPOOL_BUFFER2_BINDING_SLOT);
RWStructuredBuffer<UpdateData> _UpdateDataBuffer; register(PPOOL_BUFFER3_BINDING_SLOT);

struct UpdateData
{
};
```

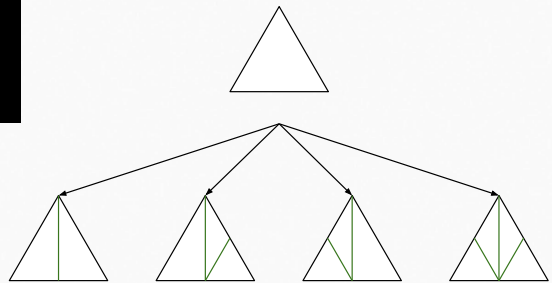


In addition to that, we need some temporary data, that we'll call those update data

Parallel Update Routine (Memory Pool)

```
// Graphics Buffers that hold the Memory pool
RWStructuredBuffer<uint> _HeapIDBufferRW; register(PPOOL_BUFFER0_BINDING_SLOT);
RWStructuredBuffer<uint> _NeighborsBufferRW; register(PPOOL_BUFFER1_BINDING_SLOT);
RWStructuredBuffer<float> _PositionRWSBufferRW; register(PPOOL_BUFFER2_BINDING_SLOT);
RWStructuredBuffer<UpdateData> _UpdateDataBuffer; register(PPOOL_BUFFER3_BINDING_SLOT);

struct UpdateData
{
};
```

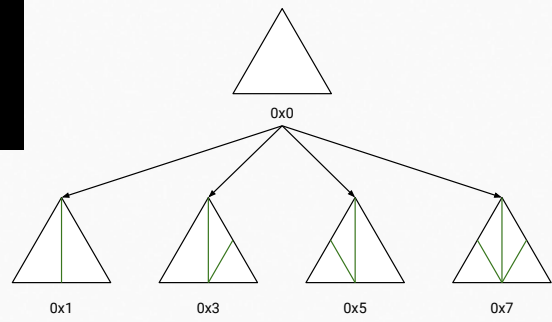


At each update loop, a bisector can be in one of 5 states, unchanged, single split, double splits (two cases) or triple split.

Parallel Update Routine (Memory Pool)

```
// Graphics Buffers that hold the Memory pool
RWStructuredBuffer<uint64_t> _HeapIDBufferRW; register(PPOOL_BUFFER0_BINDING_SLOT);
RWStructuredBuffer<uint3_t> _NeighborsBufferRW; register(PPOOL_BUFFER1_BINDING_SLOT);
RWStructuredBuffer<float3_t> _PositionRWSBufferRW; register(PPOOL_BUFFER2_BINDING_SLOT);
RWStructuredBuffer<UpdateData> _UpdateDataBuffer; register(PPOOL_BUFFER3_BINDING_SLOT);

struct UpdateData
{
    // Holds the subdivision pattern of a bisector, atomic friendly
    uint32_t subdivPattern;
};
```



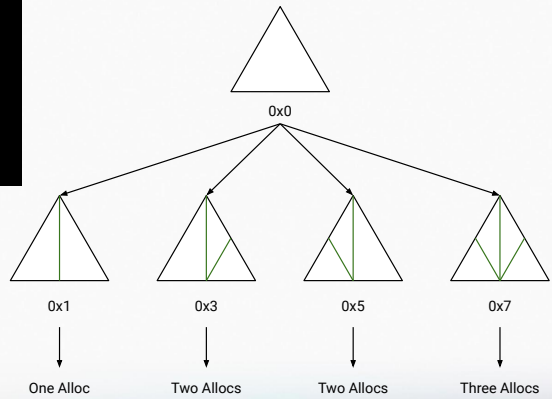
We encode that split state into three bits and store it in a `uint32_t` to use the atomic intrinsics

Parallel Update Routine (Memory Pool)

```
// Graphics Buffers that hold the Memory pool
RWStructuredBuffer<uint64> _HeapIDBufferRW: register(POOL_BUFFER0_BINDING_SLOT);
RWStructuredBuffer<uint3> _NeighborsBufferRW: register(POOL_BUFFER1_BINDING_SLOT);
RWStructuredBuffer<float3> _PositionRWSBufferRW: register(POOL_BUFFER2_BINDING_SLOT);
RWStructuredBuffer<UpdateData> _UpdateDataBuffer: register(POOL_BUFFER3_BINDING_SLOT);

struct UpdateData
{
    // Holds the subdivision pattern of a bisector, atomic friendly
    uint32_t subdivPattern;

    // Allocation Indices, we're always reusing the previous bisector
    uint3 allocations = {UINT32_MAX, UINT32_MAX, UINT32_MAX};
};
```



Depending on the split, we'll need up to 3 new slots to store the location of produced bisectors.

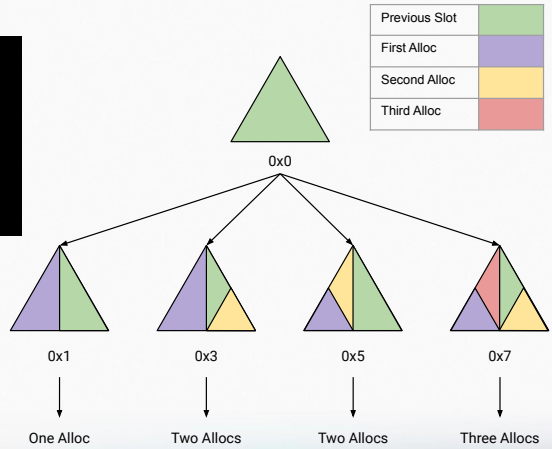
In this implementation, we always reuse the current bisector slot for performance and compactness reasons

Parallel Update Routine (Memory Pool)

```
// Graphics Buffers that hold the Memory pool
RWStructuredBuffer<uint64_t> _HeapIDBufferRW: register(POOL_BUFFER0_BINDING_SLOT);
RWStructuredBuffer<uint3_t> _NeighborsBufferRW: register(POOL_BUFFER1_BINDING_SLOT);
RWStructuredBuffer<float3> _PositionRWSBufferRW: register(POOL_BUFFER2_BINDING_SLOT);
RWStructuredBuffer<UpdateData> _UpdateDataBuffer: register(POOL_BUFFER3_BINDING_SLOT);

struct UpdateData
{
    // Holds the subdivision pattern of a bisector, atomic friendly
    uint32_t subdivPattern;

    // Allocation Indices, we're always reusing the previous bisector
    uint3 allocations = {UINT32_MAX, UINT32_MAX, UINT32_MAX};
};
```



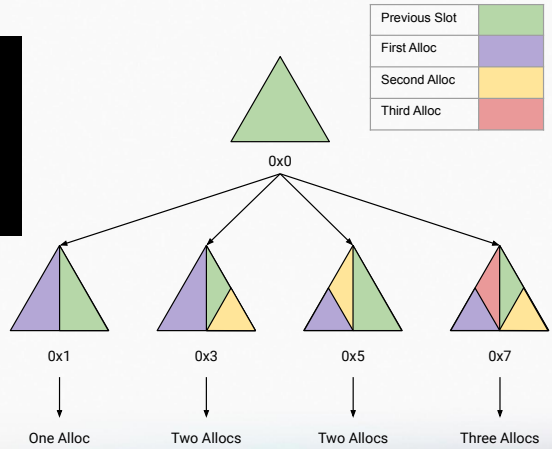
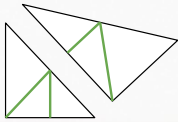
The way the indices are assigned has to be consistent

Parallel Update Routine (Memory Pool)

```
// Graphics Buffers that hold the Memory pool
RWStructuredBuffer<uint3> _HeapIDBufferRW: register(POOL_BUFFER0_BINDING_SLOT);
RWStructuredBuffer<uint3> _NeighborsBufferRW: register(POOL_BUFFER1_BINDING_SLOT);
RWStructuredBuffer<float3> _PositionRWSBufferRW: register(POOL_BUFFER2_BINDING_SLOT);
RWStructuredBuffer<UpdateData> _UpdateDataBuffer: register(POOL_BUFFER3_BINDING_SLOT);

struct UpdateData
{
    // Holds the subdivision pattern of a bisector, atomic friendly
    uint32_t subdivPattern;

    // Allocation Indices, we're always reusing the previous bisector
    uint3 allocations = {UINT32_MAX, UINT32_MAX, UINT32_MAX};
};
```



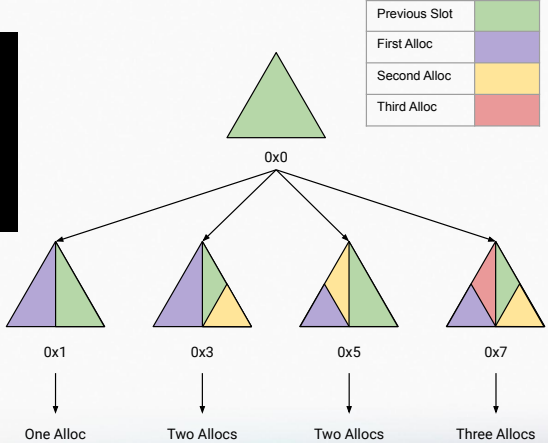
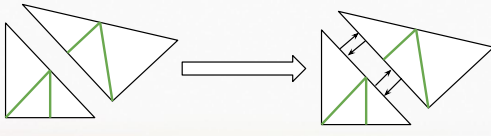
In that way, for any split combination at an interface between two bisectors

Parallel Update Routine (Memory Pool)

```
// Graphics Buffers that hold the Memory pool
RWStructuredBuffer<uint3> _HeapIDBufferRW: register(POOL_BUFFER0_BINDING_SLOT);
RWStructuredBuffer<uint3> _NeighborsBufferRW: register(POOL_BUFFER1_BINDING_SLOT);
RWStructuredBuffer<float3> _PositionRWSBufferRW: register(POOL_BUFFER2_BINDING_SLOT);
RWStructuredBuffer<UpdateData> _UpdateDataBuffer: register(POOL_BUFFER3_BINDING_SLOT);

struct UpdateData
{
    // Holds the subdivision pattern of a bisector, atomic friendly
    uint32_t subdivPattern;

    // Allocation Indices, we're always reusing the previous bisector
    uint3 allocations = {UINT32_MAX, UINT32_MAX, UINT32_MAX};
};
```



We're able to predict the neighbor data using simply the update and neighbor data of the current and neighbor bisectors.

This is important and it is the cornerstone of the sync free parallel implementation

Parallel Update Routine (Memory Pool)

```
// Graphics Buffers that hold the Memory pool
RWStructuredBuffer<uint64_t> _HeapIDBufferRW; register(PPOOL_BUFFER0_BINDING_SLOT);
RWStructuredBuffer<uint3> _NeighborsBufferRW; register(PPOOL_BUFFER1_BINDING_SLOT);
RWStructuredBuffer<float3> _PositionRWSBufferRW; register(PPOOL_BUFFER2_BINDING_SLOT);
RWStructuredBuffer<UpdateData> _UpdateDataBuffer; register(PPOOL_BUFFER3_BINDING_SLOT);

struct UpdateData
{
    // Holds the subdivision pattern of a bisector, atomic friendly
    uint32_t subdivPattern;

    // Allocation Indices, we're always reusing the previous bisector
    uint3 allocations = {UINT32_MAX, UINT32_MAX, UINT32_MAX};

    // Visibility and modification flags of a bisector
    uint32_t flags;

    // Used for patching neighbors after modifications
    uint2 propagation;
};
```

In addition to that we need an array to store additional flags and neighbors propagation data

Parallel Update Routine



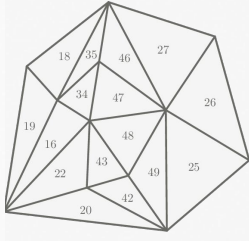
© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction
- Evaluate LEB
- Deformation

And with that, we've defined the CBT implementation and the Memory pool layout.

Now let's look at the update routine itself, each block on the right maps to one or multiple compute shader

Parallel Update Routine

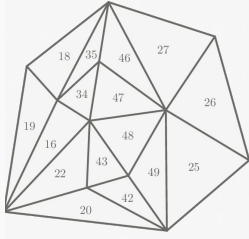


Index	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
Heap ID	18	34	18	19	20	42	22	46	48	25	26	27	35	43	47	49			

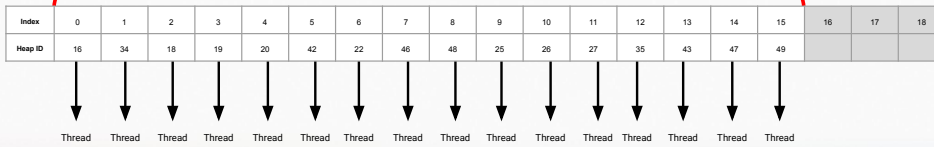
- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction
- Evaluate LEB
- Deformation

So for a given triangulation, there is a corresponding memory pool

Parallel Update Routine



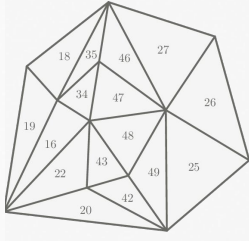
Indirect Dispatches



- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction
- Evaluate LEB
- Deformation

Most of the passes will be an indirect dispatch based on the a given set of bisectors.

Parallel Update Routine



Indirect Dispatches

Index	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
Heap ID	18	34	18	19	20	42	22	46	48	25	26	27	35	43	47	49			

Thread

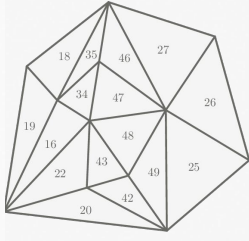
Thread Thread

Thread

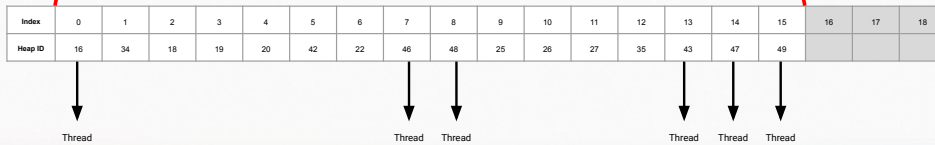
- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction
- Evaluate LEB
- Deformation

The set of bisectors will change based on the pass

Parallel Update Routine



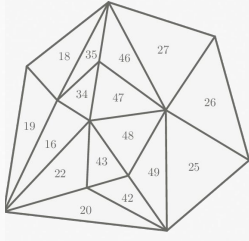
Indirect Dispatches



- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction
- Evaluate LEB
- Deformation

and the subset will be specified at each pass of the pipeline

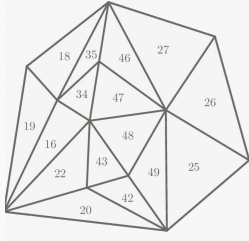
Parallel Update Routine



- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction
- Evaluate LEB
- Deformation

Alright, the first step of the pipeline is “Reset”

Parallel Update Routine (Reset)



CPU

```
m_cmdBuffer->dispatch(1, 1, 1); // One thread
```

GPU

```
// CBT function  
cbt_new_frame();  
  
// Bisector queues for indirect dispatches  
reset_queues();
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

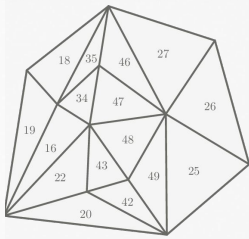
Evaluate LEB

Deformation

This pass is a single thread dispatch that:

- prepares the allocations for the frame
- resets the bisector queues for indirect dispatches

Parallel Update Routine (Classify)



CPU

```
m_cmdBuffer->dispatch_indirect(ACTIVE_BISECTORS); // HeapIDs: 16, 18, 19, ..., 48, 49
```

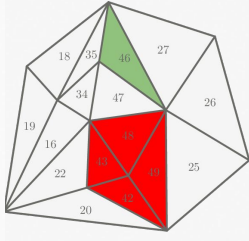
GPU

- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction
- Evaluate LEB
- Deformation

Index	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
Heap ID	16	34	18	19	20	42	22	46	48	25	28	27	35	43	47	49			
State																			

For the second pass, "Classify", each thread will run for one active bisector

Parallel Update Routine (Classify)



CPU

```
m_cmdBuffer->dispatch_indirect(ACTIVE_BISECTORS); // HeapIDs: 16, 18, 19, ..., 48, 49
```

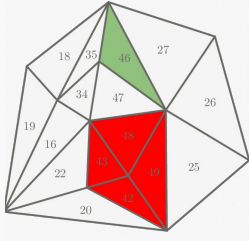
GPU

- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction
- Evaluate LEB
- Deformation

Index	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
Heap ID	18	34	18	19	20	42	22	46	48	25	26	27	35	43	47	49			
State	Keep	Keep	Keep	Keep	Keep	Merge	Keep	Split	Merge	Keep	Keep	Keep	Keep	Merge	Keep	Merge			

The idea is to classify each bisector in one of 3 states, Keep, Merge or Split. Depending on what we're trying to achieve the criteria for these operations can change, but for the demo here are the ones we used

Parallel Update Routine (Classify)



CPU

```
m_cmdBuffer->dispatch_indirect(ACTIVE_BISECTORS); // HeapIDs: 16, 18, 19, ..., 48, 49
```

GPU

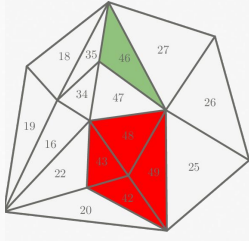
```
// Is the current bisector visible?  
if (NdotV < 0.0)  
    return MERGE;
```

- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction
- Evaluate LEB
- Deformation

Index	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
Heap ID	16	34	18	19	20	42	22	46	48	25	26	27	35	43	47	49			
State	Keep	Keep	Keep	Keep	Keep	Merge	Keep	Split	Merge	Keep	Keep	Keep	Keep	Merge	Keep	Merge			

First, we do a NdotV test and flag it for merge if it fails

Parallel Update Routine (Classify)



CPU

```
m_cmdBuffer->dispatch_indirect(ACTIVE_BISECTORS); // HeapIDs: 16, 18, 19, ..., 48, 49
```

GPU

```
// Is the current bisector visible?
if (NdotV < 0.0)
    return MERGE;

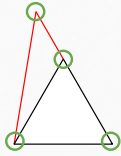
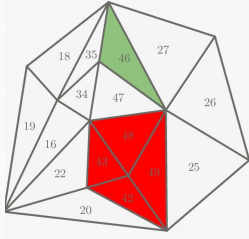
// Does the bisector intersect the camera(s) frustum(s)
if (not intersect_frustum(...))
    return MERGE;
```

- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction
- Evaluate LEB
- Deformation

Index	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
Heap ID	18	34	18	19	20	42	22	46	48	25	26	27	35	43	47	49			
State	Keep	Keep	Keep	Keep	Keep	Merge	Keep	Split	Merge	Keep	Keep	Keep	Keep	Merge	Keep	Merge			

We then intersect one or multiple frustums and flag it for merge if it does not intersect any of them

Parallel Update Routine (Classify)



CPU

```
m_cmdBuffer->dispatch_indirect(ACTIVE_BISECTORS); // HeapIDs: 16, 18, 19, ..., 48, 49
```

GPU

```
// Is the current bisector visible?
if (NdotV < 0.0)
    return MERGE;

// Does the bisector intersect the camera(s) frustum(s)
if (not intersect_frustum(...))
    return MERGE;

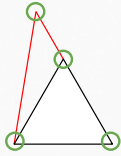
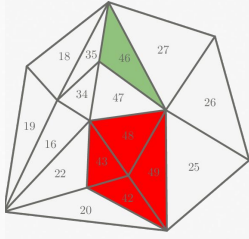
// Project the current triangle
float currentArea, parentArea;
project_bisector(...);
```

- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction
- Evaluate LEB
- Deformation

Index	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
Heap ID	16	34	18	19	20	42	22	46	48	25	26	27	35	43	47	49			
State	Keep	Keep	Keep	Keep	Keep	Merge	Keep	Split	Merge	Keep	Keep	Keep	Keep	Merge	Keep	Merge			

Then we'll project the 4 vertices and evaluate the projected area of the triangle and the parent triangle

Parallel Update Routine (Classify)



CPU

```
m_cmdBuffer->dispatch_indirect(ACTIVE_BISECTORS); // HeapIDs: 16, 18, 19, ..., 48, 49
```

GPU

```
// Is the current bisector visible?
if (NdotV < 0.0)
    return MERGE;

// Does the bisector intersect the camera(s) frustum(s)
if (not intersect_frustum(...))
    return MERGE;

// Project the current triangle
float currentArea, parentArea;
project_bisector(...);

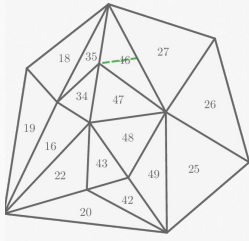
if (currentArea > TRIANGLE_SIZE)
    return SPLIT;
if (parentArea < TRIANGLE_SIZE)
    return MERGE;
return KEEP;
```

- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction
- Evaluate LEB
- Deformation

Index	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
Heap ID	16	34	18	19	20	42	22	46	48	25	26	27	35	43	47	49			
State	Keep	Keep	Keep	Keep	Keep	Merge	Keep	Split	Merge	Keep	Keep	Keep	Keep	Merge	Keep	Merge			

We compare those area to a triangle size that is defined by the application and flag the bisector for merge, split or keep depending on the result

Parallel Update Routine (Split)



CPU

```
m_cmdBuffer->dispatch_indirect(SPLIT_BISECTORS); // HeopIDs: 46
```

GPU

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

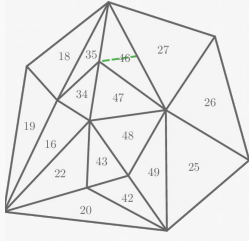
Evaluate LEB

Deformation

Then we move on to the “Split” pass, this will only run on the bisectors that have been tagged for as requiring a split.

In the example that would be only the bisector with the heapid 46.

Parallel Update Routine (Split)



CPU

```
m_cmdBuffer->dispatch_indirect(SPLIT_BISECTORS); // HeoplDs: 46
```

GPU

```
// Define how much memory we need to reserve for the worst case  
// The naive solution is  $2 * depth - 1$ , not viable, need to do better  
uint32_t maxMemory = eval_max_memory(L);  
  
// Try to book the maximum memory we would need  
if (not book_memory(maxMemory)) return;
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

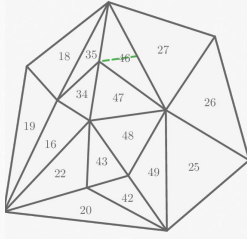
Evaluate LEB

Deformation

The first thing we need to do is figure out what is the worst case allocation size we need to guarantee the propagation of this split.

If we cannot guarantee the required memory space, the split cannot go through, that would break the LEB scheme

Parallel Update Routine (Split)



book 2 * depth - 1

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

CPU

```
m_cmdBuffer->dispatch_indirect(SPLIT_BISECTORS); // HeapIDs: 46
```

GPU

```
// Define how much memory we need to reserve for the worst case  
// The naive solution is 2 * depth - 1, not viable, need to do better  
uint32_t maxMemory = eval_max_memory(L);  
  
// Try to book the maximum memory we would need  
if (not book_memory(maxMemory)) return;
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

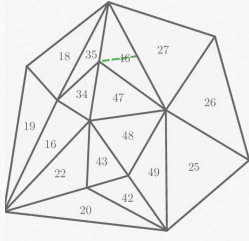
Reduction

Evaluate LEB

Deformation

The naive solution that is $2 * \text{depth} - 1$ ends up causing tessellation issues when moving at a high speed and this is not viable

Parallel Update Routine (Split)



CPU

```
m_cmdBuffer->dispatch_indirect(SPLIT_BISECTORS); // HeoplDs: 46
```

GPU

```
// Define how much memory we need to reserve for the worst case  
// The naive solution is 2 * depth - 1, not viable, need to do better  
uint32_t maxMemory = eval_max_memory(L);  
  
// Try to book the maximum memory we would need  
if (not book_memory(maxMemory)) return;
```



book 2 * depth - 1



book smart

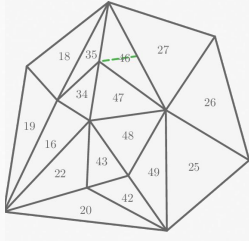
© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction
- Evaluate LEB
- Deformation

There is a better way of doing it that we'll not detail in the presentation but you can find the details in the source code we share

Once we have that estimation, we'll try to book that amount of memory from the CBT. Again we fail to guarantee that amount, this split won't go through as we cannot preserve the LEB scheme

Parallel Update Routine (Split)



CPU

```
m_cmdBuffer->dispatch_indirect(SPLIT_BISECTORS); // HeaplDs: 46
```

GPU

```
// Define how much memory we need to reserve for the worst case
// The naive solution is 2 * depth - 1, not viable, need to do better
uint32_t maxMemory = eval_max_memory(-);

// Try to book the maximum memory we would need
if (not book_memory(maxMemory)) return;

// Propagate the split until we are done
uint32_t usedMemory = 1;
while (not done)
{
    if (currentDepth == twinDepth || already_split(twinID))
    {
        InterlockedOr(twinUpdateData.subdivPattern, 0x1);
        usedMemory += 1;
        allocate_next_available_slot();
        done = true;
    }
    else
    {
        InterlockedOr(twinUpdateData.subdivPattern, rightSuddiv ? 0x3 : 0x5);
        usedMemory += 2;
        allocate_next_available_slot(); allocate_next_available_slot();
        move_to_twin();
    }
}
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

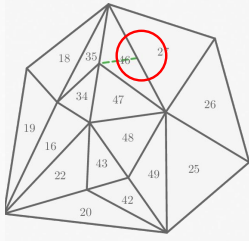
Reduction

Evaluate LEB

Deformation

If we manage to book the worst case, we de-recursify the propagation algorithm into a while loop that will split each triangle, allocate one or multiple slots, propagate the subdivision to the twins all the bits we need until we're done.

Parallel Update Routine (Split)



CPU

```
m_cmdBuffer->dispatch_indirect(SPLIT_BISECTORS); // HeoplDs: 46
```

GPU

```
// Define how much memory we need to reserve for the worst case
// The naive solution is 2 * depth - 1, not viable, need to do better
uint32_t maxMemory = eval_max_memory(-);

// Try to book the maximum memory we would need
if (not book_memory(maxMemory)) return;

// Propagate the split until we are done
uint32_t usedMemory = 1;
while (not done)
{
    if (currentDepth == twinDepth || already_split(twinID))
    {
        InterlockedOr(twinUpdateData.subdivPattern, 0x1);
        usedMemory += 1;
        allocate_next_available_slot();
        done = true;
    }
    else
    {
        InterlockedOr(twinUpdateData.subdivPattern, rightSuddiv ? 0x3 : 0x5);
        usedMemory += 2;
        allocate_next_available_slot(); allocate_next_available_slot();
        move_to_twin();
    }
}
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

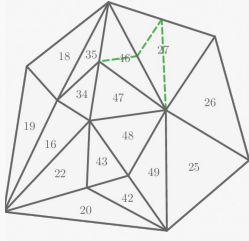
Reduction

Evaluate LEB

Deformation

So in this case we do the first allocation, but we have a T junction so we continue on our path

Parallel Update Routine (Split)



CPU

```
m_cmdBuffer->dispatch_indirect(SPLIT_BISECTORS); // HeoplDs: 46
```

GPU

```
// Define how much memory we need to reserve for the worst case
// The naive solution is 2 * depth - 1, not viable, need to do better
uint32_t maxMemory = eval_max_memory(-);

// Try to book the maximum memory we would need
if (not book_memory(maxMemory)) return;

// Propagate the split until we are done
uint32_t usedMemory = 1;
while (not done)
{
    if (currentDepth == twinDepth || already_split(twinID))
    {
        InterlockedOr(twinUpdateData.subdivPattern, 0x1);
        usedMemory += 1;
        allocate_next_available_slot();
        done = true;
    }
    else
    {
        InterlockedOr(twinUpdateData.subdivPattern, rightSuddiv ? 0x3 : 0x5);
        usedMemory += 2;
        allocate_next_available_slot(); allocate_next_available_slot();
        move_to_twin();
    }
}
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

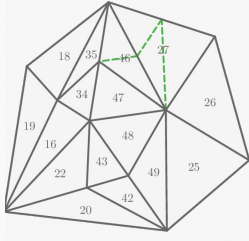
Reduction

Evaluate LEB

Deformation

Which leads us to doing two more allocations and we're good to go.

Parallel Update Routine (Split)



© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course.

CPU

```
m_cmdBuffer->dispatch_indirect(SPLIT_BISECTORS); // HeaplDs: 46
```

GPU

```
// Define how much memory we need to reserve for the worst case
// The naive solution is 2 * depth - 1, not viable, need to do better
uint32_t maxMemory = eval_max_memory(L);

// Try to book the maximum memory we would need
if (not book_memory(maxMemory)) return;

// Propagate the split until we are done
uint32_t usedMemory = 1;
while (not done)
{
    if (currentDepth == twinDepth || already_split(twinID))
    {
        InterlockedOr(twinUpdateData.subdivPattern, 0x1);
        usedMemory += 1;
        allocate_next_available_slot();
        done = true;
    }
    else
    {
        InterlockedOr(twinUpdateData.subdivPattern, rightSuddiv ? 0x3 : 0x5);
        usedMemory += 2;
        allocate_next_available_slot(); allocate_next_available_slot();
        move_to_twin();
    }
}

// Return the memory that we did not use
cancel_memory_booking(maxMemory - usedMemory);
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

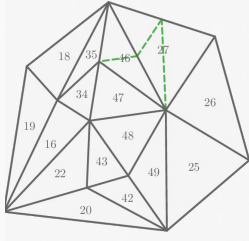
Reduction

Evaluate LEB

Deformation

Finally we're return to the CBT the memory space that we didn't allocate

Parallel Update Routine (Decode)



CPU

```
m_cmdBuffer->dispatch_indirect(SUBDIV_BISECTORS); // HeapIDs: 46, 27
```

GPU

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

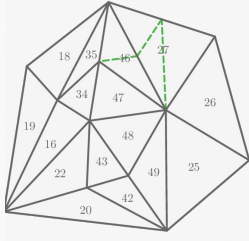
Deformation

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

The next step “Decode”, runs on the bisectors that have been flagged for a subdivision, (46 and 27 in this case),

This will associate each allocated slot with a free bit of the CBT

Parallel Update Routine (Decode)



CPU

```
m_cmdBuffer->dispatch_indirect(SUBDIV_BISECTORS); // HeapIDs: 46, 27
```

GPU

```
// We need to load the CBT into the shared memory  
load_buffer_to_shared_memory(...);  
  
// Depending on the subdivision pattern, the number of bits we need to allocate varies  
uint32_t numAllocations = countbits(updateData.subdivPattern);  
  
// Allocate the bits  
for (uint32_t allocIdx = 0; allocIdx < numAllocations; ++allocIdx)  
updateData.allocations[allocIdx] = decode_bit_complement(...);
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

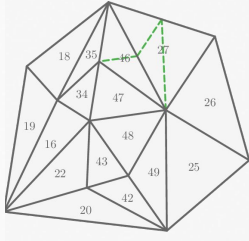
For this pass we need to load the CBT into the shared memory first

The number of bits set to 1 of the subdivision pattern gives us the number of decodes we need to do.

We do those decodes and store them in the update data

Once this is done, we move on to the next step

Parallel Update Routine (Bisect)



CPU

```
m_cmdBuffer->dispatch_indirect(SUBDIV_BISECTORS); // HeapIDs: 46, 27
```

GPU

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

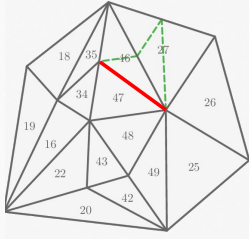
Reduction

Evaluate LEB

Deformation

The following pass, “Bisect”, that runs on the same set of bisectors than the decode pass, aka every bisector that will be splitted

Parallel Update Routine (Bisect)



CPU

```
m_cmdBuffer->dispatch_indirect(SUBDIV_BISECTORS); // HeapIDs: 46, 27
```

GPU

```
void update_bisectors(uint32_t ID0, uint32_t ID1, ...)  
{  
    // Depending on the pattern and the neighbors pattern  
    // Modify the neighbors, heapID and flag which subdivision needs to be propagated  
}  
  
// Depending on the subdivision pattern, we need to modify the allocated bisectors  
if (subdivPattern == 0x1)  
    update_bisectors(currentID, upData.allocation[0]);  
else if (subdivPattern == 0x3 || subdivPattern == 0x5)  
    update_bisectors(currentID, upData.allocation[0], upData.allocation[1]);  
else if (subdivPattern == 0x7)  
    update_bisectors(currentID, upData.allocation[0], upData.allocation[1], upData.allocation[2]);
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

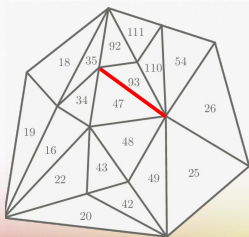
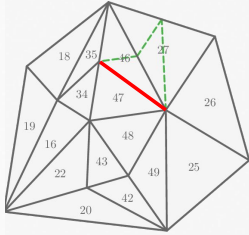
Evaluate LEB

Deformation

Based on the subdivision pattern of the current bisector and and the one of its neighbors, we'll modify the heap ID and adjust neighbor pointers.

This only partially updates the neighbors and we still have problematic interfaces (marked in red in the example) that will be resolved in the following pass

Parallel Update Routine (Bisect)



© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

CPU

```
m_cmdBuffer->dispatch_indirect(SUBDIV_BISECTORS); // HeapIDs: 46, 27
```

GPU

```
void update_bisectors(uint32_t ID0, uint32_t ID1, ...)  
{  
    // Depending on the pattern and the neighbors pattern  
    // Modify the neighbors, heapID and flag which subdivision needs to be propagated  
}  
  
// Depending on the subdivision pattern, we need to modify the allocated bisectors  
if (subdivPattern == 0x1)  
    update_bisectors(currentID, upData.allocation[0]);  
  
else if (subdivPattern == 0x3 || subdivPattern == 0x5)  
    update_bisectors(currentID, upData.allocation[0], upData.allocation[1]);  
  
else if (subdivPattern == 0x7)  
    update_bisectors(currentID, upData.allocation[0], upData.allocation[1], upData.allocation[2]);
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

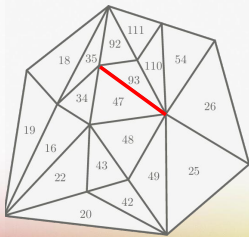
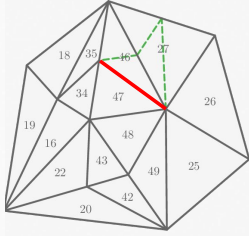
Reduction

Evaluate LEB

Deformation

This will produce the new state of the triangulation

Parallel Update Routine (Bisect)



© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

CPU

```
m_cmdBuffer->dispatch_indirect(SUBDIV_BISECTORS); // HeapIDs: 46, 27
```

GPU

```
void update_bisectors(uint32_t ID0, uint32_t ID1, ...)  
{  
    // Depending on the pattern and the neighbors pattern  
    // Modify the neighbors, heapID and flag which subdivision needs to be propagated  
}  
  
// Depending on the subdivision pattern, we need to modify the allocated bisectors  
if (subdivPattern == 0x1)  
    update_bisectors(currentID, upData.allocation[0]);  
  
else if (subdivPattern == 0x3 || subdivPattern == 0x5)  
    update_bisectors(currentID, upData.allocation[0], upData.allocation[1]);  
  
else if (subdivPattern == 0x7)  
    update_bisectors(currentID, upData.allocation[0], upData.allocation[1], upData.allocation[2]);  
  
// Atomic modification of the bitfield  
for (uint32_t allocIdx = 0; allocIdx < numAllocations; ++allocIdx)  
    set_bit_atomic(upData.allocation[allocIdx], true);
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

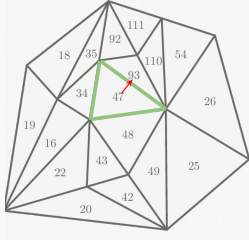
Reduction

Evaluate LEB

Deformation

And finally we'll update the bits of the CBT based on the number of allocations that were required for this split

Parallel Update Routine (Propagate Bisect)



CPU

```
m_cmdBuffer->dispatch_indirect(PROPGATE_SPLIT_BISECTORS); // HeaplDs: 93
```

GPU

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

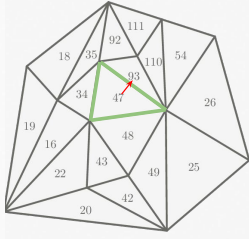
Reduction

Evaluate LEB

Deformation

The next step, "Propagate Bisect", will process these problematic interfaces by changing the neighbors when required

Parallel Update Routine (Propagate Bisect)



CPU

```
m_cmdBuffer->dispatch_indirect(PROPAGATE_SPLIT_BISECTORS); // HeapIDs: 93
```

GPU

```
// Did the neighbor split?  
if (was_split(neighborID))  
    // Adjust based on the subdivision pattern of the neighbor  
else  
    adjust_neighbors(neighborID, currentID);
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

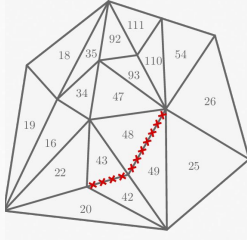
Evaluate LEB

Deformation

The patching routine has to take into account if the neighbor was split (and how it was split, once, twice or thrice) and adjust the neighbors accordingly

That covers it for the split part of the algorithm. Now we need to process the merges, It is roughly the same approach as the splits, but simpler as there are no propagation to account for.

Parallel Update Routine (Prepare Simplify)



CPU

```
m_cmdBuffer->dispatch_indirect(MERGE_BISECTORS); // HeapsIDs: 42, 43, 48, 49
```

GPU

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

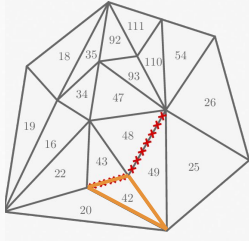
Reduction

Evaluate LEB

Deformation

The first pass is “Prepare simplify”, This dispatch will run on the 4 bisectors that have been flagged for simplification

Parallel Update Routine (Prepare Simplify)



CPU

```
m_cmdBuffer->dispatch_indirect(MERGE_BISECTORS); // HeapIDs: 42, 43, 48, 49
```

GPU

```
// The bisector with the smallest heapID of the 4 is in charge of doing the operation
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

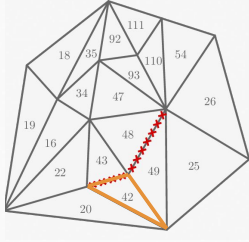
Deformation

To avoid synchronization and racing conditions, the bisector with the smallest heapID in the four is in charge of checking and registering for the simplification.

In the example, that would be the bisector with the heapID 42

(Actually in practice, we don't enqueue the bisectors with odd heap IDs into the merge queue as an optimization)

Parallel Update Routine (Prepare Simplify)



CPU

```
m_cmdBuffer->dispatch_indirect(MERGE_BISECTORS); // HeapIDs: 42, 43, 48, 49
```

GPU

```
// The bisector with the smallest heapID of the 4 is in charge of doing the operation  
if (heapID % 2 != 0 || pairDepth != currentDepth)  
    return;  
  
if (twinHeapID < heapID  
    || twinDepth != currentDepth  
    || twinPairDepth != currentDepth)  
    return;
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

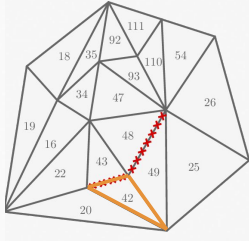
Reduction

Evaluate LEB

Deformation

We can easily identify which one it is using the depth of the bisectors and their heapIDs

Parallel Update Routine (Prepare Simplify)



CPU

```
m_cmdBuffer->dispatch_indirect(MERGE_BISECTORS); // HeapIDs: 42, 43, 48, 49
```

GPU

```
// The bisector with the smallest heapID of the 4 is in charge of doing the operation  
if (heapID % 2 != 0 || pairDepth != currentDepth)  
    return;  
  
if (twinHeapID < heapID  
    || twinDepth != currentDepth  
    || twinPairDepth != currentDepth)  
    return;  
  
// All four bisectors need to request a merge and be at the same depth  
if (!merge_request(currentID)  
    || !merge_request(pairID)  
    || !merge_request(twinID)  
    || !merge_request(twinPairID))  
    return;  
  
// This bisector need to process the simplification  
enqueue_for_simplify(currentID);
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

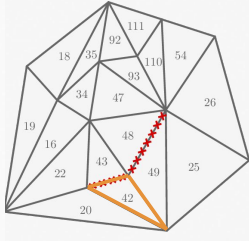
Reduction

Evaluate LEB

Deformation

That bisector is then in charge of making sure they all required a merge operation and if it is the case, it will enqueue for the next step

Parallel Update Routine (Simplify)



CPU

```
m_cmdBuffer->dispatch_indirect(SIMPLIFY_BISECTORS); // HeapIDs: 42
```

GPU

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

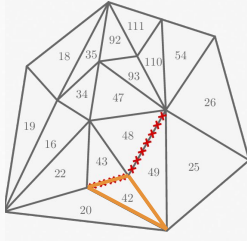
Reduction

Evaluate LEB

Deformation

Which is the “Simplify” step, this will only run for the bisector 42 in this case

Parallel Update Routine (Simplify)



CPU

```
m_cmdBuffer->dispatch_indirect(SIMPLIFY_BISECTORS); // HeapIDs: 42
```

GPU

```
// Update the current bisector and free it's pair  
update_bisector(currentID);  
free_bisector_slot(pairID);
```

```
// Update the twin bisector and free it's pair  
update_bisector(twinID);  
free_bisector_slot(twinPairID);
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

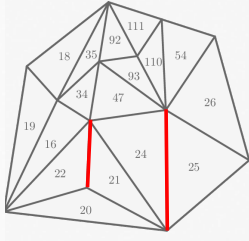
Reduction

Evaluate LEB

Deformation

We update the bisectors that will remain (the ones with the smallest heapIDs for each pair, 42 and 48 in this case) and free the one other ones (43 and 49).

Parallel Update Routine (Simplify)



CPU

```
m_cmdBuffer->dispatch_indirect(SIMPLIFY_BISECTORS); // HeapIDs: 42
```

GPU

```
// Update the current bisector and free it's pair  
update_bisector(currentID);  
free_bisector_slot(pairID);
```

```
// Update the twin bisector and free it's pair  
update_bisector(twinID);  
free_bisector_slot(twinPairID);
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

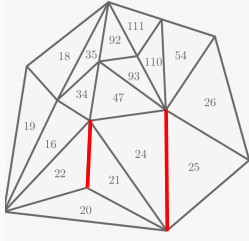
Reduction

Evaluate LEB

Deformation

The same way than for split, we update the heapID and only partially the neighbors that will be updated in the following pass. This generates the new triangulation.

Parallel Update Routine (Simplify)



CPU

```
m_cmdBuffer->dispatch_indirect(SIMPLIFY_BISECTORS); // HeapIDs: 42
```

GPU

```
// Update the current bisector and free it's pair  
update_bisector(currentID);  
free_bisector_slot(pairID);
```

```
// Update the twin bisector and free it's pair  
update_bisector(twinID);  
free_bisector_slot(twinPairID);
```

```
// Set the bits to zero  
set_bit_atomic(pairID, false);  
set_bit_atomic(twinPairID, false);
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

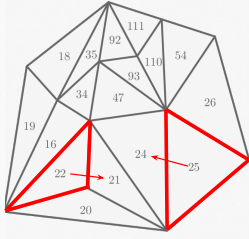
Reduction

Evaluate LEB

Deformation

Then we set the bits in the CBT to zero for the unused slots

Parallel Update Routine (Propagate Simplify)



CPU

```
m_cmdBuffer->dispatch_indirect(PROPGATE_SIMPLIFY_BISECTORS); // HeoplDs: 21, 24
```

GPU

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

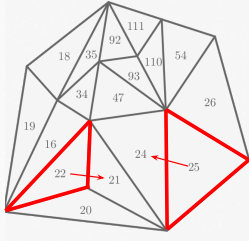
Reduction

Evaluate LEB

Deformation

In this last step of the modification “Propagate simplify”, we dispatch on the bisectors that generated an inconsistency in the neighbors

Parallel Update Routine (Propagate Simplify)



CPU

```
m_cmdBuffer->dispatch_indirect(PROPGATE_SIMPLIFY_BISECTORS); // HeapIDs: 21, 24
```

GPU

```
// Was the neighbor merged?  
if (was_merged(neighborID))  
{  
    if (was_deleted(neighborID))  
        adjust_neighbors(neighborPairID, currentID);  
    else  
        adjust_neighbors(neighborID, currentID);  
}  
else  
    adjust_neighbors(neighborID, currentID);
```

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

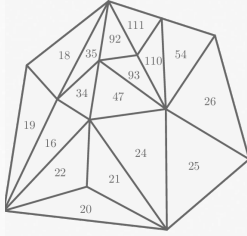
Reduction

Evaluate LEB

Deformation

And depending on if the neighbors were merged we or not, the routine needs to take account if the bisector was deleted, etc

Parallel Update Routine (Propagate Simplify)



CPU

```
m_cmdBuffer->dispatch_indirect(PROPGATE_SIMPLIFY_BISECTORS); // HeapIDs: 21, 24
```

GPU

```
// Was the neighbor merged?  
if (was_merged(neighborID))  
{  
    if (was_deleted(neighborID))  
        adjust_neighbors(neighborPairID, currentID);  
    else  
        adjust_neighbors(neighborID, currentID);  
}  
else  
    adjust_neighbors(neighborID, currentID);
```



© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

And with that we've processed our split and merge requests and have a coherent and LEB compatible triangulation

Parallel Update Routine (Reduction)



CPU Commands



Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

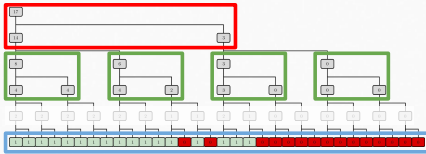
But we're not done yet, there are couple steps left before we get to the end of our update routine.

We did modify the bit field in the previous steps, but the tree now doesn't match anymore. The next one is operating a sum reduction on the CBT's tree

Parallel Update Routine (Reduction 128K)

CPU Commands

- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction**
- Evaluate LEB
- Deformation



Level 0: [0, 131072] x 1
Level 1: [0, 65536] x 2
Level 2: [0, 32768] x 4
Level 3: [0, 16384] x 8
Level 4: [0, 8192] x 16
Level 5: [0, 4096] x 32
Level 6: [0, 2048] x 64

Level 7: [0, 1024] x 128
Level 8: [0, 512] x 256
Level 9: [0, 256] x 512
Level 10: [0, 128] x 1024

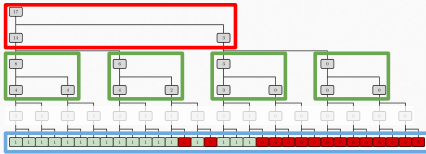
Level 17: Raw 64 bits representation

Let's again take the example of the 128k CBT

Parallel Update Routine (Reduction 128K)

CPU Commands

```
m_cmdBuffer->dispatch(1024 / 4 / 64); // REDUCE_PREPASS
```



Level 0: [0, 131072] x 1
Level 1: [0, 65536] x 2
Level 2: [0, 32768] x 4
Level 3: [0, 16384] x 8
Level 4: [0, 8192] x 16
Level 5: [0, 4096] x 32
Level 6: [0, 2048] x 64

Level 7: [0, 1024] x 128
Level 8: [0, 512] x 256
Level 9: [0, 256] x 512
Level 10: [0, 128] x 1024

Level 17: Raw 64 bits representation

Prepass

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

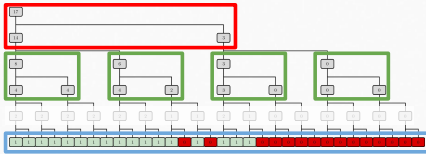
Deformation

First we do a reduction to move from our bitfield to the first explicitly stored level of the tree

Parallel Update Routine (Reduction 128K)

CPU Commands

```
m_cmdBuffer->dispatch(1024 / 4 / 64); // REDUCE_PREPASS  
m_cmdBuffer->dispatch(512 / 64); // REDUCE_FIRST_PASS
```



Level 0: [0, 131072] x 1
Level 1: [0, 65536] x 2
Level 2: [0, 32768] x 4
Level 3: [0, 16384] x 8
Level 4: [0, 8192] x 16
Level 5: [0, 4096] x 32
Level 6: [0, 2048] x 64

Level 7: [0, 1024] x 128
Level 8: [0, 512] x 256
Level 9: [0, 256] x 512
Level 10: [0, 128] x 1024

Level 17: Raw 64 bits representation

First Pass
Prepass

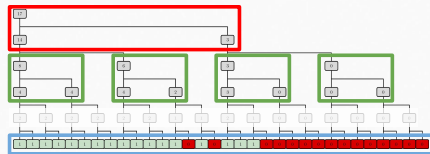
- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction
- Evaluate LEB
- Deformation

Then another pass within the pet-thread atomic part of the tree

Parallel Update Routine (Reduction 128K)

CPU Commands

```
m_cmdBuffer->dispatch(1024 / 4 / 64); // REDUCE_PREPASS  
m_cmdBuffer->dispatch(512 / 64); // REDUCE_FIRST_PASS  
m_cmdBuffer->dispatch(1); // REDUCE_SECOND_PASS
```



Level 0: [0, 131072] x 1 ← Second Pass
Level 1: [0, 65536] x 2
Level 2: [0, 32768] x 4
Level 3: [0, 16384] x 8
Level 4: [0, 8192] x 16
Level 5: [0, 4096] x 32
Level 6: [0, 2048] x 64

Level 7: [0, 1024] x 128 ← First Pass
Level 8: [0, 512] x 256
Level 9: [0, 256] x 512
Level 10: [0, 128] x 1024 ← Prepass

Level 17: Raw 64 bits representation

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

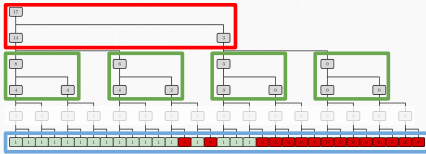
Finally one single workgroup will do the full sum reduction on the atomic friendly tree. With that we're able to do the have the full tree updated at a reasonable cost.

Parallel Update Routine (Reduction 128K)

CPU Commands

```
m_cmdBuffer->dispatch(1024 / 4 / 64); // REDUCE_PREPASS  
m_cmdBuffer->dispatch(512 / 64); // REDUCE_FIRST_PASS  
m_cmdBuffer->dispatch(1); // REDUCE_SECOND_PASS
```

CBT Size	Intel Arc 770 (ms)	AMD 6650 XT (ms)	Nvidia RTX 4090 (ms)
128k	0.055	0.01	0.008
256k	0.055	0.012	0.01
512k	0.062	0.013	0.01
1m	0.065	0.017	0.011



- Level 0: [0, 131072] x 1
 - Level 1: [0, 65536] x 2
 - Level 2: [0, 32768] x 4
 - Level 3: [0, 16384] x 8
 - Level 4: [0, 8192] x 16
 - Level 5: [0, 4096] x 32
 - Level 6: [0, 2048] x 64
 - Level 7: [0, 1024] x 128
 - Level 8: [0, 512] x 256
 - Level 9: [0, 256] x 512
 - Level 10: [0, 128] x 1024
 - Level 17: Raw 64 bits representation
- Second Pass
- First Pass
- Prepass

- Reset
- Classify
- Split
- Decode
- Bisect
- Propagate Bisect
- Prepare Simplify
- Simplify
- Propagate Simplify
- Reduction
- Evaluate LEB
- Deformation

To illustrate what i mean by reasonable, here you can see a table that recaps the reduction time for each CBT size that we provide in our demo and for 3 GPUs (Intel Arc 770, AMD 6650 XT and Nvidia 4090)

Parallel Update Routine (Evaluate LEB)



CPU

```
m_cmdBuffer->dispatch_indirect(MODIFIED_BISECTORS);
```

GPU

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

The next step is processing the modified bisectors to ensure that the positions are up to date

Parallel Update Routine (Evaluate LEB)

HeapID = 0 0 1 0 ... 0 1

CPU

m_cmdBuffer->dispatch_indirect(MODIFIED_BISECTORS);

GPU

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

As we mentioned before, the geometry information is encoded using the heapID and if we break it down to a binary representation, it would look something like this

Parallel Update Routine (Evaluate LEB)

HeapID = 0 0 1 0 ... 0 1

$$M_t = M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1$$

CPU

```
m_cmdBuffer->dispatch_indirect(MODIFIED_BISECTORS);
```

GPU

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

The final position is obtained using a chain of matrix multiplications, each bit of the heapID will correspond to one of two matrices

Parallel Update Routine (Evaluate LEB)

HeapID = 0 0 1 0 ... 0 1

$$M_i = M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1$$

$$M_0 = \begin{bmatrix} 0.0 & 1 & 0 \\ 0.5 & 0.0 & 0.5 \\ 1 & 0 & 0.0 \end{bmatrix} \quad M_1 = \begin{bmatrix} 0.0 & 0 & 1 \\ 0.5 & 0.0 & 0.5 \\ 0 & 1 & 0.0 \end{bmatrix}$$

i-th bit is 0 i-th bit is 1

CPU

```
m_cmdBuffer->dispatch_indirect(MODIFIED_BISECTORS);
```

GPU

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

These are the two matrices we use for the demo

Parallel Update Routine (Evaluate LEB)

HeapID = 0 0 1 0 ... 0 1

$M_i = M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1$

Up to 64 3x3 double matmul

$$M_0 = \begin{bmatrix} 0.0 & 1 & 0 \\ 0.5 & 0.0 & 0.5 \\ 1 & 0 & 0.0 \end{bmatrix} \quad M_1 = \begin{bmatrix} 0.0 & 0 & 1 \\ 0.5 & 0.0 & 0.5 \\ 0 & 1 & 0.0 \end{bmatrix}$$

i-th bit is 0 i-th bit is 1

CPU

```
m_cmdBuffer->dispatch_indirect(MODIFIED_BISECTORS);
```

GPU



Fig. 9. Impact of double floating point precision on planetary-scale tessellation.

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

There are two issues here:

- First, even if the output position are stored using simple precision floating point and camera relative, we have to do the multiplications in double space (otherwise we break the floating points due to precisions issues due to our method)
- Depending on the subdivision level, we can have up to 64 3x3 matrix multiplications to evaluate which is a lot

Parallel Update Routine (Evaluate LEB)

HeapID = 0 0 1 0 ... 0 1

$$M_t = M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1$$

Up to 64 3x3 double matmul

$$M_0 = \begin{bmatrix} 0.0 & 1 & 0 \\ 0.5 & 0.0 & 0.5 \\ 1 & 0 & 0.0 \end{bmatrix} \quad M_1 = \begin{bmatrix} 0.0 & 0 & 1 \\ 0.5 & 0.0 & 0.5 \\ 0 & 1 & 0.0 \end{bmatrix}$$

i-th bit is 0 i-th bit is 1

$$M_t = M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1$$

32 3x3 float matmul X2 + 1 3x3 double matmul

CPU

```
m_cmdBuffer->dispatch_indirect(MODIFIED_BISECTORS);
```

GPU

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

The first thing we can do is split this multiplication into two 32 3x3 float multiplications with a 3x3 double matmul at the middle which reduces the ALU consumption and better uses the FLOPS of the card, but is still not good enough

Parallel Update Routine (Evaluate LEB)

HeapID = 0 0 1 0 ... 0 1

$$M_t = M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1$$

Up to 64 3x3 double matmul

$$M_0 = \begin{bmatrix} 0.0 & 1 & 0 \\ 0.5 & 0.0 & 0.5 \\ 1 & 0 & 0.0 \end{bmatrix} \quad M_1 = \begin{bmatrix} 0.0 & 0 & 1 \\ 0.5 & 0.0 & 0.5 \\ 0 & 1 & 0.0 \end{bmatrix}$$

i-th bit is 0 i-th bit is 1

$$M_t = M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1$$

32 3x3 float matmul X2 + 1 3x3 double matmul

CPU

```
m_cmdBuffer->dispatch_indirect(MODIFIED_BISECTORS);
```

GPU

00000
00001
00010
00011

11111

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

If we decompose the heap ID into sequences of bits of equal size

Parallel Update Routine (Evaluate LEB)

HeapID = 0 0 1 0 ... 0 1

$$M_t = M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1$$

Up to 64 3x3 double matmul

$$M_0 = \begin{bmatrix} 0.0 & 1 & 0 \\ 0.5 & 0.0 & 0.5 \\ 1 & 0 & 0.0 \end{bmatrix} \quad M_1 = \begin{bmatrix} 0.0 & 0 & 1 \\ 0.5 & 0.0 & 0.5 \\ 0 & 1 & 0.0 \end{bmatrix}$$

i-th bit is 0 i-th bit is 1

$$M_t = M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1$$

32 3x3 float matmul X2 + 1 3x3 double matmul

CPU

m_cmdBuffer->dispatch_indirect(MODIFIED_BISECTORS);

GPU

- 00000 → $M_{c0} = M_0 * M_0 * M_0 * M_0 * M_0$
- 00001 → $M_{c1} = M_0 * M_0 * M_0 * M_0 * M_1$
- 00010 → $M_{c2} = M_0 * M_0 * M_0 * M_1 * M_0$
- 00011 → $M_{c3} = M_0 * M_0 * M_0 * M_1 * M_1$
- ...
- 11111 → $M_{c31} = M_1 * M_1 * M_1 * M_1 * M_1$

Matrix cache to speed up the evaluation (Depth 5)

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

The way we can accelerate this even more is by building what we call a matrix cache, it is a premultiplication of all combinations over a certain number of levels (in practice, we've measure that a depth of 5 is the best speedup/size compromise because we need to load this into shared memory as we'll be accessing it quite a bit).

Parallel Update Routine (Evaluate LEB)

HeapID = 0 0 1 0 ... 0 1

↓

$$M_t = \underbrace{M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1}_{\text{Up to 64 3x3 double matmul}}$$

$$M_0 = \begin{bmatrix} 0.0 & 1 & 0 \\ 0.5 & 0.0 & 0.5 \\ 1 & 0 & 0.0 \end{bmatrix} \quad M_1 = \begin{bmatrix} 0.0 & 0 & 1 \\ 0.5 & 0.0 & 0.5 \\ 0 & 1 & 0.0 \end{bmatrix}$$

i-th bit is 0 i-th bit is 1

$$M_t = \underbrace{M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1}_{32 \text{ 3x3 float matmul X2 + 1 3x3 double matmul}}$$

$$M_t = \underbrace{M_{c4} * M_{c12} * \dots * M_{c22} * M_{c18} * M_{c18} * \dots * M_{c1}}_{6 \text{ 3x3 float matmul X2 + 1 3x3 double matmul}}$$

CPU

m_cmdBuffer->dispatch_indirect(MODIFIED_BISECTORS);

GPU

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

00000 → $M_{c0} = M_0 * M_0 * M_0 * M_0 * M_0$

00001 → $M_{c1} = M_0 * M_0 * M_0 * M_0 * M_1$

00010 → $M_{c2} = M_0 * M_0 * M_0 * M_1 * M_0$

00011 → $M_{c3} = M_0 * M_0 * M_0 * M_1 * M_1$

...

11111 → $M_{c31} = M_1 * M_1 * M_1 * M_1 * M_1$

Matrix cache to speed up the evaluation (Depth 5)

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

By doing this, we reduce the number of matrices to up to two sets of 6 3x3 mat muls and one double 3x3 matmul which allows us to get very reasonable execution times as you'll see in a second.

This shader executes only for modified bisectors which is a fraction of the total bisectors

Parallel Update Routine (Evaluate LEB)

HeapID = 0 0 1 0 ... 0 1

$$M_i = M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1$$

Up to 64 3x3 double matmul

$$M_0 = \begin{bmatrix} 0.0 & 1 & 0 \\ 0.5 & 0.0 & 0.5 \\ 1 & 0 & 0.0 \end{bmatrix} \quad M_1 = \begin{bmatrix} 0.0 & 0 & 1 \\ 0.5 & 0.0 & 0.5 \\ 0 & 1 & 0.0 \end{bmatrix}$$

i-th bit is 0 i-th bit is 1

$$M_i = M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1$$

32 3x3 float matmul X2 + 1 3x3 double matmul

$$M_i = M_{c4} * M_{c12} * \dots * M_{c22} * M_{c18} * M_{c18} * \dots * M_{c1}$$

6 3x3 float matmul X2 + 1 3x3 double matmul

CPU

```
m_cmdBuffer->dispatch_indirect(MODIFIED_BISECTORS);
```

GPU

```
// Load the matrix cache
load_leb_matrix_cache_to_shared_memory(L);
```

- 00000 → $M_{c0} = M_0 * M_0 * M_0 * M_0 * M_0$
- 00001 → $M_{c1} = M_0 * M_0 * M_0 * M_0 * M_1$
- 00010 → $M_{c2} = M_0 * M_0 * M_0 * M_1 * M_0$
- 00011 → $M_{c3} = M_0 * M_0 * M_0 * M_1 * M_1$
- ...
- 11111 → $M_{c31} = M_1 * M_1 * M_1 * M_1 * M_1$

Matrix cache to speed up the evaluation (Depth 5)

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

In practice, we would load the matrix cache into shared memory

Parallel Update Routine (Evaluate LEB)

HeapID = 0 0 1 0 ... 0 1

$$M_t = M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1$$

Up to 64 3x3 double matmul

$$M_0 = \begin{bmatrix} 0.0 & 1 & 0 \\ 0.5 & 0.0 & 0.5 \\ 1 & 0 & 0.0 \end{bmatrix} \quad M_1 = \begin{bmatrix} 0.0 & 0 & 1 \\ 0.5 & 0.0 & 0.5 \\ 0 & 1 & 0.0 \end{bmatrix}$$

i-th bit is 0 i-th bit is 1

$$M_t = M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1$$

32 3x3 float matmul X2 + 1 3x3 double matmul

$$M_t = M_{c4} * M_{c12} * \dots * M_{c22} * M_{c18} * M_{c18} * \dots * M_{c1}$$

6 3x3 float matmul X2 + 1 3x3 double matmul

CPU

```
m_cmdBuffer->dispatch_indirect(MODIFIED_BISECTORS);
```

GPU

```
// Load the matrix cache
load_leb_matrix_cache_to_shared_memory(L);

// Flatten the heapIDs
double3 bisectorPos[4];
evaluate_positions(heapID);
```

$$M_{c0} = M_0 * M_0 * M_0 * M_0 * M_0$$

$$M_{c1} = M_0 * M_0 * M_0 * M_0 * M_1$$

$$M_{c2} = M_0 * M_0 * M_0 * M_1 * M_0$$

$$M_{c3} = M_0 * M_0 * M_0 * M_1 * M_1$$

$$\dots$$

$$M_{c31} = M_1 * M_1 * M_1 * M_1 * M_1$$

Matrix cache to speed up the evaluation (Depth 5)

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

We would flatten the heapID into 4 positions (3 for the current bisector and 1 for the parent)

Parallel Update Routine (Evaluate LEB)

HeapID = 0 0 1 0 ... 0 1

$$M_t = M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1$$

Up to 64 3x3 double matmul

$$M_0 = \begin{bmatrix} 0.0 & 1 & 0 \\ 0.5 & 0.0 & 0.5 \\ 1 & 0 & 0.0 \end{bmatrix} \quad M_1 = \begin{bmatrix} 0.0 & 0 & 1 \\ 0.5 & 0.0 & 0.5 \\ 0 & 1 & 0.0 \end{bmatrix}$$

i-th bit is 0 i-th bit is 1

$$M_t = M_0 * M_1 * \dots * M_0 * M_1 * M_0 * \dots * M_1$$

32 3x3 float matmul X2 + 1 3x3 double matmul

$$M_t = M_{c4} * M_{c12} * \dots * M_{c22} * M_{c18} * M_{c18} * \dots * M_{c1}$$

6 3x3 float matmul X2 + 1 3x3 double matmul

CPU

```
m_cmdBuffer->dispatch_indirect(MODIFIED_BISECTORS);
```

GPU

```
// Load the matrix cache
load_leb_matrix_cache_to_shared_memory(_);

// Flatten the heapIDs
double3 bisectorPos[4];
evaluate_positions(heapID);

// Apply planet coord remapping, camera relative positioning
_PositionRWSBufferRW[4 * currentID] = to_rws_planet_coord(bisectorPos[0]);
_PositionRWSBufferRW[4 * currentID + 1] = to_rws_planet_coord(bisectorPos[1]);
_PositionRWSBufferRW[4 * currentID + 2] = to_rws_planet_coord(bisectorPos[2]);
_PositionRWSBufferRW[4 * currentID + 3] = to_rws_planet_coord(bisectorPos[3]);
```

- 00000 → $M_{c0} = M_0 * M_0 * M_0 * M_0 * M_0$
- 00001 → $M_{c1} = M_0 * M_0 * M_0 * M_0 * M_1$
- 00010 → $M_{c2} = M_0 * M_0 * M_0 * M_1 * M_0$
- 00011 → $M_{c3} = M_0 * M_0 * M_0 * M_1 * M_1$
- ...
- 11111 → $M_{c31} = M_1 * M_1 * M_1 * M_1 * M_1$

Matrix cache to speed up the evaluation (Depth 5)

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

Reduction

Evaluate LEB

Deformation

And then we would apply a transformation to map this to a spherical planet and convert it to relative world space

Parallel Update Routine (Deformation)

Earth

Moon

Displacement

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING

Reset

Classify

Split

Decode

Bisect

Propagate Bisect

Prepare Simplify

Simplify

Propagate Simplify

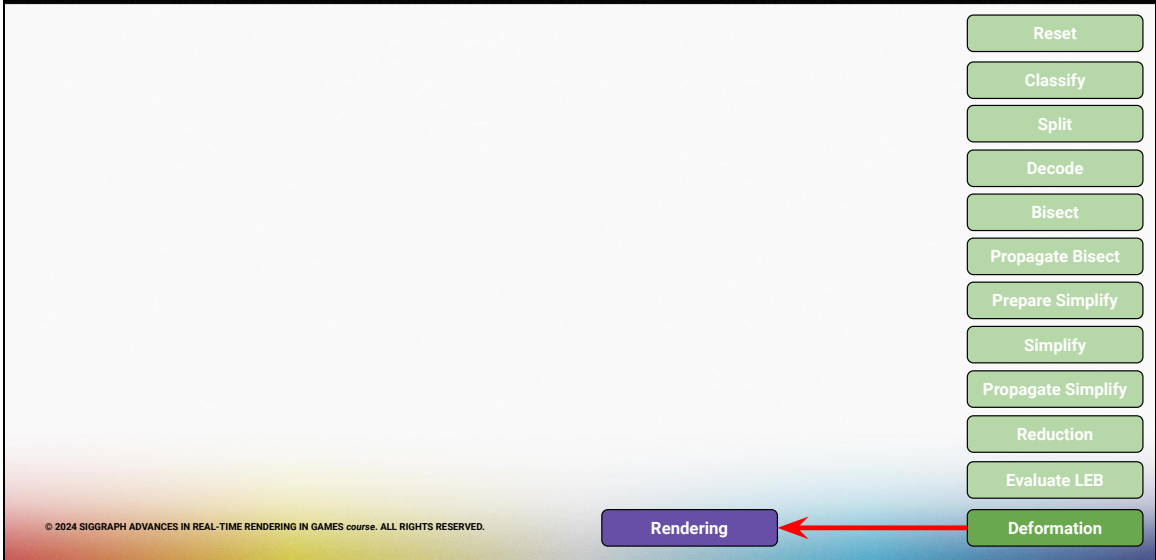
Reduction

Evaluate LEB

Deformation

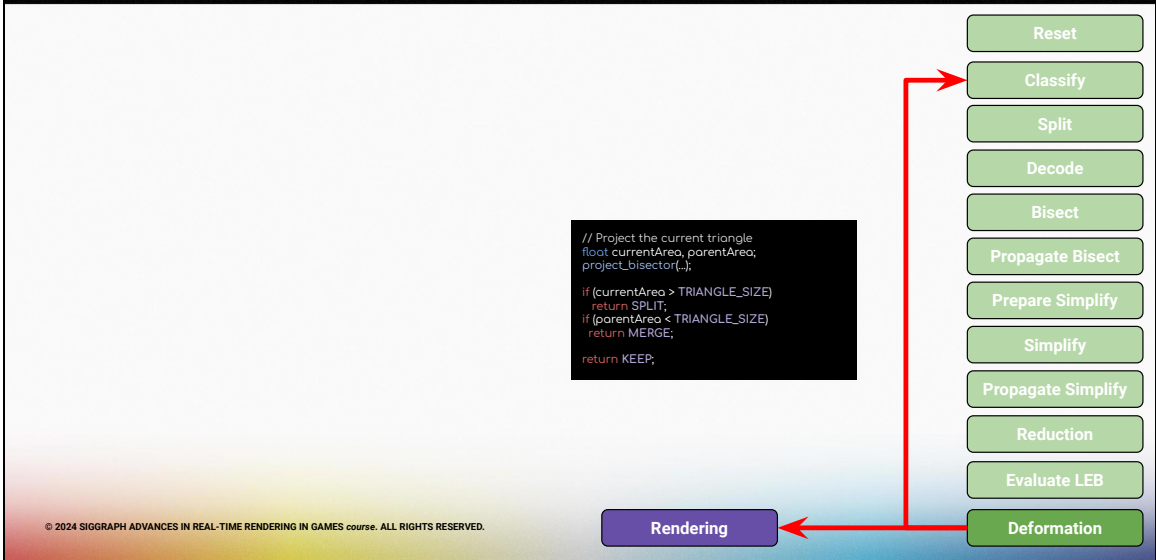
Finally to all existing positions, we would apply the relevant deformation, for the earth it would be the result of a bunch of FFTs and the moon we use the displacement maps that the NASA provides on it's website, plus a bunch of noise functions for high frequency details

Parallel Update Routine (Deformation)



These produced positions would then be used for the rendering rasterized or ray traced

Parallel Update Routine (Deformation)

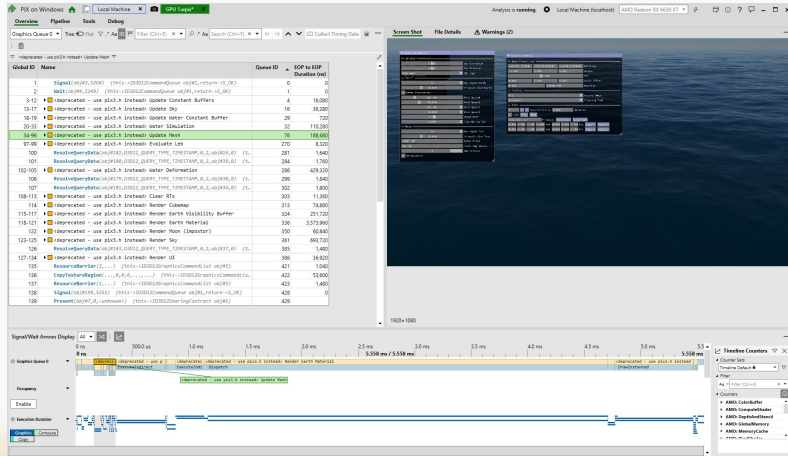


But also for the following frame to operate the classification



And with that, we are ready to render our mesh to screen

Performance



The last thing we would like to cover is the performance numbers of the method, this is a screenshot of a capture of the demo running at the surface of the earth

Performance



Scene	Earth Close Up	Earth Close Up	Earth Close Up	Moon Close Up	Moon Close Up	Moon Close Up
GPU	Intel Arc 770	AMD 6650 XT	Nvidia RTX 4090	Intel Arc 770	AMD 6650 XT	Nvidia RTX 4090
Update Proportion	High	High	High	Low	Low	Low
Active Primitives	~64k	~64k	~64k	~95k	~95k	~95k
Reset (ms)	0.014	0.003	0.004	0.013	0.003	0.005
Classify (ms)	0.078	0.058	0.012	0.07	0.072	0.013
Split (ms)	0.042	0.019	0.014	0.023	0.017	0.012
Allocate (ms)	0.024	0.026	0.013	0.031	0.021	0.011
Copy (ms)	0.011	0.017	0.006	0.012	0.017	0.006
Bisect (ms)	0.016	0.005	0.007	0.012	0.006	0.005
Propagate Bisect (ms)	0.016	0.007	0.009	0.02	0.01	0.01
Prepare Simplify (ms)	0.02	0.008	0.007	0.025	0.008	0.007
Simplify (ms)	0.02	0.009	0.011	0.012	0.006	0.012
Propagate Simplify (ms)	0.017	0.006	0.011	0.017	0.006	0.01
Reduce (ms)	0.058	0.01	0.008	0.058	0.01	0.008
Indexation (ms)	0.03	0.014	0.013	0.03	0.014	0.013
Evaluate LEB (ms)	0.025	0.008	0.01	0.019	0.008	0.01
Update Pass (ms)	0.371	0.19	0.125	0.342	0.198	0.122

Rendered at 1920x1080 - CBT 128K

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

This table recaps the timings for each of the passes that we've explained. We did the profiling on three GPUs:

- The Arc 770
- the AMD 6650 XT which is roughly equivalent to a PS5
- A higher end GPU, the Nvidia 4090

Performance



Scene	Earth Close Up	Earth Close Up	Earth Close Up	Moon Close Up	Moon Close Up	Moon Close Up
GPU	Intel Arc 770	AMD 6650 XT	Nvidia RTX 4090	Intel Arc 770	AMD 6650 XT	Nvidia RTX 4090
Update Proportion	High	High	High	Low	Low	Low
Active Primitives	~64k	~64k	~64k	~95k	~95k	~95k
Reset (ms)	0.014	0.003	0.004	0.013	0.003	0.005
Classify (ms)	0.078	0.058	0.012	0.07	0.072	0.013
Split (ms)	0.042	0.019	0.014	0.023	0.017	0.012
Allocate (ms)	0.024	0.026	0.013	0.031	0.021	0.011
Copy (ms)	0.011	0.017	0.006	0.012	0.017	0.006
Bisect (ms)	0.016	0.005	0.007	0.012	0.006	0.005
Propagate Bisect (ms)	0.016	0.007	0.009	0.02	0.01	0.01
Prepare Simplify (ms)	0.02	0.008	0.007	0.025	0.008	0.007
Simplify (ms)	0.02	0.009	0.011	0.012	0.006	0.012
Propagate Simplify (ms)	0.017	0.006	0.011	0.017	0.006	0.01
Reduce (ms)	0.058	0.01	0.008	0.058	0.01	0.008
Indexation (ms)	0.03	0.014	0.013	0.03	0.014	0.013
Evaluate LEB (ms)	0.025	0.008	0.01	0.019	0.008	0.01
Update Pass (ms)	0.371	0.19	0.125	0.342	0.198	0.122

Rendered at 1920x1080 - CBT 128K

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

The important bit there is that we're able to stay at reasonable cost at all time and less than 0.2ms on a hardware equivalent to a PS5 which is more than reasonable given the achievements of the method



The rest of the pipeline is dependent on the actual deformation that we'd use and the general structure of the rendering pipeline.
For the demo we're using a visibility buffer followed by a material pass.

Performance



Scene	Earth Close Up	Moon Close Up
GPU	AMD 6650 XT	AMD 6650 XT
Update Proportion	High	Low
Active Primitives	~64k	~95k
Reset (ms)	0.003	0.003
Classify (ms)	0.058	0.072
Split (ms)	0.019	0.017
Allocate (ms)	0.026	0.021
Copy (ms)	0.017	0.017
Bisect (ms)	0.005	0.006
Propagate Bisect (ms)	0.007	0.01
Prepare Simplify (ms)	0.008	0.008
Simplify (ms)	0.009	0.006
Propagate Simplify (ms)	0.006	0.006
Reduce (ms)	0.01	0.01
Indexation (ms)	0.014	0.014
Evaluate LEB (ms)	0.008	0.008
Update Pass (ms)	0.19	0.198
Deformation (ms)	0.43	1
Render Vis Buffer (ms)	0.25	0.27
Render Material (ms)	3.5	3.2
Frame Total	4.18	4.47

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

These are numbers we got for the AMD 6650 XT, but these are really dependent on how you generate your surface data really.

Performance



Scene	Earth Close Up	Moon Close Up
GPU	AMD 6650 XT	AMD 6650 XT
Update Proportion	High	Low
Active Primitives	~64k	~95k
Reset (ms)	0.003	0.003
Classify (ms)	0.058	0.072
Split (ms)	0.019	0.017
Allocate (ms)	0.026	0.021
Copy (ms)	0.017	0.017
Bisect (ms)	0.005	0.006
Propagate Bisect (ms)	0.007	0.01
Prepare Simplify (ms)	0.008	0.008
Simplify (ms)	0.009	0.006
Propagate Simplify (ms)	0.006	0.006
Reduce (ms)	0.01	0.01
Indexation (ms)	0.014	0.014
Evaluate LEB (ms)	0.008	0.008
Update Pass (ms)	0.19	0.198
Deformation (ms)	0.43	1
Render Vis Buffer (ms)	0.25	0.27
Render Material (ms)	3.5	3.2
Frame Total	4.18	4.47
Build RTAS (ms)	3.5	3.6
Trace RTAS (ms)	3.6	2.9

© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

In addition to that i'd like to share some number ray tracing related. The method is compatible with real-time hardware ray tracing and given that the number of primitives remains quite "low", we're able to get something decent even by doing a full rebuild every frame of the tris and blas



https://github.com/AnisB/large_cbt

Concurrent Binary Trees for Large-Scale Game Components

ANIS BENYOUB, Intel Corporation, France
JONATHAN DUPUY, Intel Corporation, France

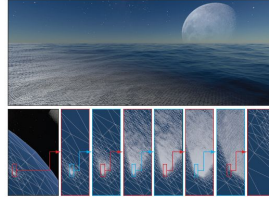


Fig. 3. (Top) An Earth-based planet rendered in real-time using our triangulation method. (Bottom) An alternative white-noise view of the render as seen from ground to space through successive zoomed-in views.

A concurrent binary tree (CBT) is a CPU-friendly data-structure suitable for the generation of heuristic based terrain tessellations, i.e., adaptive triangulations over square domains. In this paper, we expose the benefits of this data-structure in two aspects. First, we show how to bring heuristic based tessellations to arbitrary poly-topologies rather than just squares. Our approach consists of mapping a triangular mesh into a quadtree which we split into a hierarchy to match the edge of the target mesh. These hierarchies can then be automatically adapted to produce conforming triangulations solely based on halving operations. Second, we illustrate a technique that reduces the triangulation to low-resolution meshes. We do so by using the CBT as a memory pool manager rather than an on-the-fly encoding of the triangulation as done originally. By using a CBT in this way, we automatically alternate and/or reduce the number of triangles subdivisions using shared CPU memory. We demonstrate the benefits of our improvements by rendering planetary scale geometry out of very coarse meshes. Performance-wise, our triangulation method enables to use three times as much mesh hardware.

CCS Concepts • Computing methodologies • Rendering; Massively parallel algorithms; Shared memory algorithms.

Additional Key Words and Phrases: level-of-detail, subdivision, CPU, real-time, rendering

Authors' address: Anis Benyoub, Intel Corporation, France, Jonathan Dupuy, Intel Corporation, France.

Proc. ACM Comput. Graph. Interact. Tech., Vol. 7, No. 1, Article . Publication date: July 2024.

And that's it. The full code of the demo that we've show is open-source and can be found in the github that is linked here. Also, don't hesitate to check the paper for more detailed explanations about the method, in the meantime, we are available for questions