



SIGGRAPH 2024
DENVER+ 28 JUL — 1 AUG

THE PREMIER CONFERENCE
& EXHIBITION ON
COMPUTER GRAPHICS &
INTERACTIVE TECHNIQUES

SEAMLESS RENDERING ON MOBILE

THE MAGIC OF ADAPTIVE LOD PIPELINE

SHUN CAO (TENCENT GAMES)



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

Hello everyone, I'm Shun from Tencent Games. My topic is seamless rendering on mobile.



We hope this rendering pipeline can simplify developers' work of customizing and tailoring mesh resources for different platforms, especially mobile, and reduce errors caused by manual operations. Of course, we also aim to ensure game graphics quality and interactive experience. It means refined rendering in areas that attract players' attention and no visual jumps or lags due to dynamic resource loading during gameplay.

SEAMLESS CLUSTER RENDERING ON MOBILE



Talking alone might be abstract, so let's demonstrate with a test scenario. This is a rendering scene with 80 million triangles achievable on mobile. We only used high-precision assets in production, but during rendering, we control the granularity of clusters based on projection area and distance from the camera at the cluster level, ensuring optimal GPU resources for suitable effects.

CURRENT TECHNOLOGY LANDSCAPE



DESKTOP AND CONSOLE

- Mesh Production Pipeline
 - Mesh Mipmaps (Generate LODs for Mesh)
 - **Mesh Clusters (Hierarchical Cluster Levels)**
- Mesh Render Pipeline
 - Instance / **Cluster Culling**
 - DeferredLighting (Gbuffer)
 - **DeferredTexture (VisibilityBuffer)**
- Hardware Support
 - TaskShader/MeshShader
 - **Bindless**
 - **Atomic64**
 - Fast bandwidth
 - **SSD or not**
 - Wave Operations

MOBILE (MAIN-STREAMING DEVICES)

- Mesh Production Pipeline
 - Mesh Mipmaps (Generate LODs for Mesh)
- Mesh Render Pipeline
 - Instance Culling
 - DeferredLighting (Gbuffer)
- Hardware Support
 - VertexShader/PixelShader
 - ComputeShader
 - DrawIndirect
 - **Bandwidth limitation**
 - **Tile based**

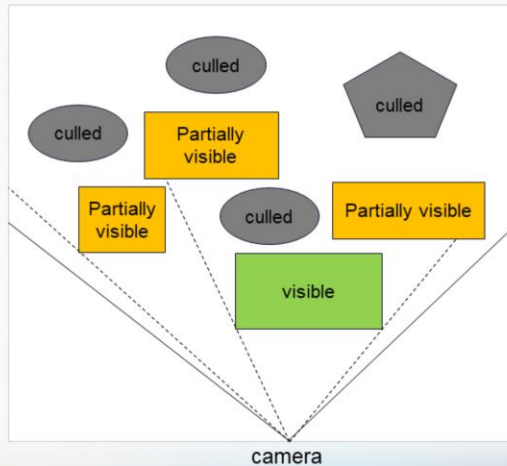
© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

However, compared to Desktop and console, mobile mainstream devices lack many hardware and software features necessary for seamless rendering, like rendering pipeline limitations, support for mesh shader and bindless features. Bandwidth and IO impacts on performance are also concerns. In addition, the mobile platform has its own tiled base GPU architecture.



So, how do we achieve seamless cluster rendering under these limitations? We aim for simplicity and efficiency, simplifying parts heavily reliant on hardware and bandwidth as much as possible.

- Cluster Based Rendering
 - More fine-grained culling
 - Uniform and small size triangles
 - Streaming friendly
- VisibilityBuffer
 - Small triangle
 - Quad Utilization Efficiency
 - Lighter than GBuffer



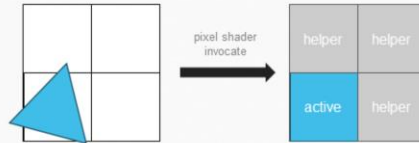
© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

Before diving into technical details, let's quickly review the background knowledge. First, Cluster-based Rendering: Occlusion culling is common in games, solving overdraw issues caused by direct object occlusion. But large objects often get fully rendered even if only part is visible. By dividing original meshes into small clusters and culling at the cluster level, GPUs can skip many invalid triangles. With fewer triangles per cluster, we can reduce vertex index precision, like using 8 bits. We can also load clusters on demand based on bounding boxes, improving GPU memory utilization.

- Cluster Based Rendering
 - More fine-grained culling
 - Uniform and small size triangles
 - Streaming friendly
- VisibilityBuffer
 - Small triangle
 - Quad Utilization Efficiency
 - Lighter than Gbuffer

Triangle index	Instance index
22 bits	10 bits

1 pixel triangle draw



Shading invocations

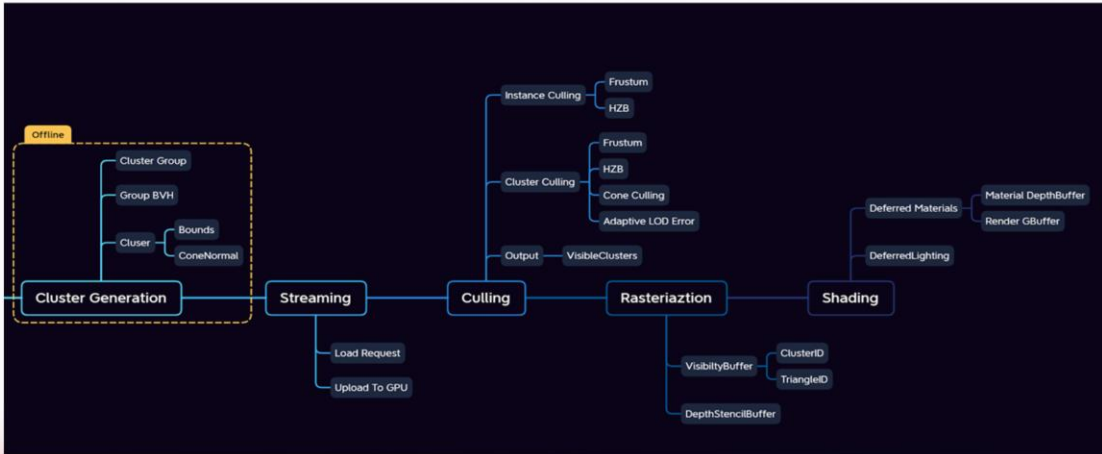
	Material	Lighting
Forward	4x	
Deferred	4x	1x
Visibility	1x	1x

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

The second technology is Visbuffer. When rendering fine triangles, Visbuffer provides higher Quad Utilization than forward and deferred pipelines, with just 32 bits of extra overhead. John Hable will talk about vis buffer further in this 2024 advances.

- Cluster Based Rendering
 - More fine-grained culling
 - Improved Culling for Tiled and Clustered , Michal Drobot (https://advances.realtimerendering.com/s2017/2017_Sig_Improved_Culling_final.pdf)
 - Uniform and small size triangles
 - The Journey to Nanite - Brian Karis, Epic Games (https://www.youtube.com/watch?v=NRnj_InpORU)
 - Streaming friendly
- VisibilityBuffer
 - 4K Rendering Breakthrough: The Filtered and Culled Visibility Buffer. Wolfgang Engel. (<https://www.gdcvault.com/play/1023792/4K-Rendering-Breakthrough-The-Filtered>)
 - Small triangle
 - Visibility Buffer Rendering with Material Graphs,John Hable (<http://filmicworlds.com/blog/visibility-buffer-rendering-with-material-graphs/>)
 - Quad Utilization Efficiency
 - Lighter than Gbuffer
- More details
 - GPU-Driven Rendering Pipelines, Sebastian Aaltonen (https://advances.realtimerendering.com/s2015/aaltonenhaar_siggraph2015_combined_final_footer_220dpi.pdf)
 - References
 - Variable Rate Shading with Visibility Buffer Rendering,John Hable,Jun 30th

Thanks to Drobot's presentation in Advances, as well as Brian Karis' talk, and Wolfgang Engel's vis buffer original deck, also Sebastian's SIGGRAPH 2015 talk . we can get more detail about those technologies . And John Hable will talk about vis buffer further in this 2024 advances

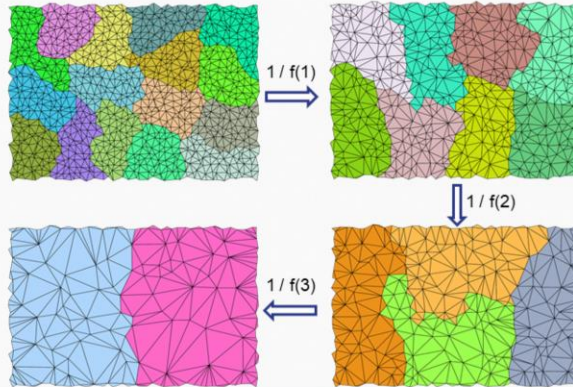


© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

After introducing the background technologies, let's move on to the rendering pipeline architecture. The pipeline has two parts: offline and runtime. The offline stage imports and generates custom cluster data, while the runtime part includes streaming, culling, rasterization, and shading.

REDESIGN EXPRESSION FORMAT OF MESH

- Mesh to Clusters
 - Fine-grained culling
 - Cluster boundingbox
 - Cluster normal cone
 - Cluster Error
 - 8 bits triangle index
- Hierarchical Clusters
 - Cluster-level LODs
 - Seamless LODs
 - On-demand Loading

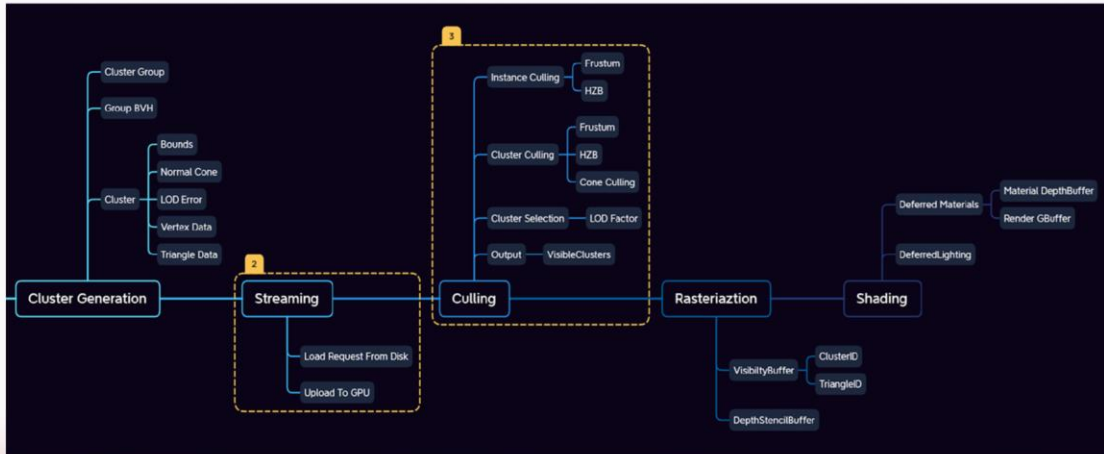


$$f(\text{level}) = \text{floor}(\text{pow}(\text{base}, \text{level}))$$
$$\text{base} = 2, \dots, \text{max_base}$$

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

In the offline part, before anything else, we redesign the mesh storage structure. Instead of traditional multi-level LODs or simple clusters, we split meshes into many clusters and continuously recombine them into new, coarser-level clusters, iterating similarly. It's like UE5's Nanite, but we add a merge coefficient function to ensure controllability of assets and manually generated mesh LODs. This makes multi-level cluster data smaller, supporting data structures for culling and streaming. For each cluster, we store its bounding box and normal cone for efficient occlusion and backface culling. Each cluster has up to 128 triangles, so we can store indices with just 8 bits, saving space. For non-leaf nodes, we record errors from merging.

THE RENDER PIPELINE

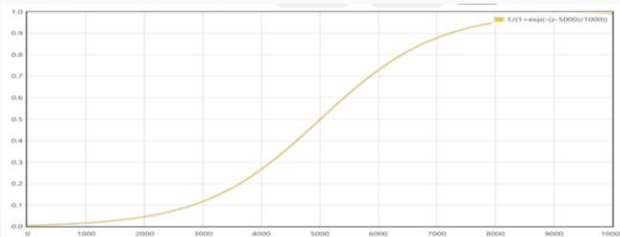


© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

Next is the streaming and culling stages. Streaming starts with the roughest clusters from each instance, then loads on demand based on culling pass output.

- GPUScene
 - InstanceBuffer
 - ClusterBuffer
 - OutVisibleClusters
- Culling Pass
 - Instance Level Culling
 - Cluster Culling
 - Adaptive Selection
- Cluster Streaming
 - Feedback to CPU
 - On-demand Loading by frames

- Target
 - $device_factor = device_power * device_level$
 - $decay_factor = 1 / (1 + \exp(distance_to_view/decay_distance));$
 - $threshold = projected_area * device_factor * decay_factor;$
- Selection
 - $parentErr > threshold \ \&\& \ clusterErr < threshold$

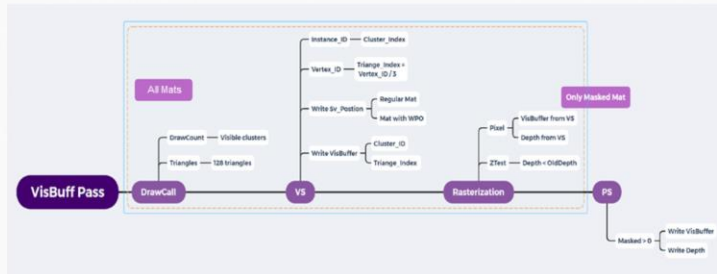


© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

In culling, we let the GPU handle cluster control and selection. Objects are stored in an instance buffer, clusters in a cluster buffer, and IDs of successfully culled clusters in a visible cluster buffer. The cluster buffer loads dynamically based on GPU selection. Culling starts with fast instance-level culling using HZB, then applies various techniques like frustum culling, occlusion culling, and normal cone-based backface culling. We use LOD error similar to UE5's Nanite for cluster LOD culling, but we design a distance curve for LodFactor selection, ensuring enough rendering precision nearby and rougher triangles far away, mimicking manual mesh LOD factor control.

RENDER VISIBLE CLUSTERS TO VISIBILITYBUFFER

- Binning Pass
 - InVisibleClusters
 - OutBinningClusters
 - Regular Geometry
 - Deformed Geometry
 - Special Material Geometry
- Raster Pass
 - Simplest Indirect Draw
 - Only hardware rasterization
 - No extra depth buffer rewrite
 - Light-weight buffer and bandwidth
 - 32-bits VisBuffer
 - Cluster index and triangle index



VisibilityBuffer(32bits)

cluster index	triangle index
25	7

During rasterization, to reduce power consumption and GPU usage, we store intermediate rendering results in a 32-bit visbuffer. Since all triangles are in clusters, ideally, one draw call renders them all to the visbuffer. For GPU parallelism, we categorize objects: type 1 records cluster and triangle IDs in the visbuffer in the vertex shader, skipping pixel shader computation. Type 2, like skinned meshes, recalculates vertex positions with bone info in the vertex shader. Type 3, like masked vegetation, filters visbuffer generation based on masked textures in the pixel shader. We avoid soft rasterization due to extra scene depth passes, lack of atomic64 support, and higher bandwidth from 64-bit visbuffers. Since we render by clusters, not instances, 7 bits for triangle indices suffice, leaving space for clusters.

- GPUScene
 - ClusterBuffer/VisibleClusterBuffer
 - VisBuffer
 - OutMatDepthTexture
- Material Depth Pass
 - Only VertexShader
 - Draw MaterialID to Depth
- Material Shading Pass
 - Generate Tiles for Materials
 - 64x64 pixels / Tile
 - Draw Tiles for Materials
 - DepthTest: Equal

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

Since bindless is not yet supported on mobile devices, to achieve material shading, we need to render each material separately with a drawcall. Since it's rare for a single material to cover the entire screen, we divide the screen into many small tiles, and each material has its own list of tiles. This way, during rendering, we don't need to render each material across the entire screen. For each tile's rendering, we extract the materialID from the cluster information in the visbuffer and compare it with the current shading materialID. Only when they are equal will the material shading proceed.

- GPUScene
 - InstanceBuffer
 - BoneTransformBuffer
- Cluster data
 - Main-bone index
 - Max-Influence bone in cluster triangles
 - Normal Cone
 - All Triangles
 - Conservative main-bone space boundingbox
- Culling Pass
 - Apply main-bone transform matrix to conservative boundingbox
 - Apply instance transform to conservative boundingbox
 - Culling with transformed conservative boundingbox
- Raster Pass
 - Transform vertex position with influence bones

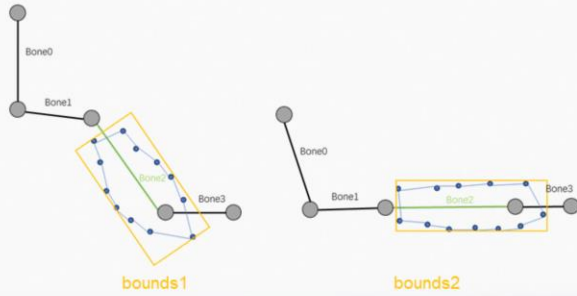


© 2024 SIGGRAPH ADVANCES IN REAL-TIME RENDERING IN GAMES course. ALL RIGHTS RESERVED.

Rendering deformable meshes is inevitable in games. Take skinned meshes as an example. During animation, each cluster's bounds keep changing, making pre-calculated bounds for LOD and culling inaccurate. Real-time vertex-based bounds calculation is too time-consuming. So, we need special handling for skinned meshes. The approach is dynamic bound box calculation and culling based on clusters in the main bone space.

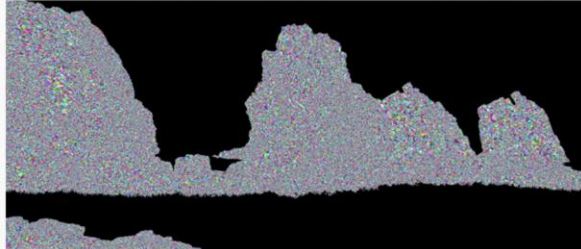
- GPUScene
 - InstanceBuffer
 - BoneTransformBuffer
- Cluster data
 - Main-bone index
 - Max-Influence bone in cluster triangles
 - Normal Cone
 - All Triangles
 - Conservative main-bone space boundingbox
- Culling Pass
 - Apply main-bone transform matrix to conservative boundingbox
 - Apply instance transform to conservative boundingbox
 - Culling with transformed conservative boundingbox
- Raster Pass
 - Transform vertex position with influence bones

main-bone = bone2
 boundingbox = max(bounds1,bounds2,...)



Offline, we compute each cluster's main bone, normal cone, and a conservative bounding box in the main bone space. The main bone is the one with the highest vertex weight in the cluster. The normal cone is the average normal of all cluster triangles, with a cone angle covering them all. The bound box covers the cluster's maximum range in the main bone space across all animations. Runtime, the CPU sends bone transform data to the GPU. During cluster culling, it reads the bone transform, transforms the conservative bound box to mesh space, then proceeds with normal cluster culling. Skinning rasterization follows normal skeletal animation calculations.

- Lightmaps ?
 - Pixel level triangles
 - How to set Lightmap resolution ?
 - Is it easy to unwrap lightmap uv?
 - Memory / Storage controllable ?
- Hierarchical Voxelization Global Illumination



Atlas specification	width x height	mipmaps	format	Memory(MB)
Low Atlas	512 x 512	9	R8G8B8A8	> 1 MB
Medium Atlas	1024 x 1024	10	R8G8B8A8	> 4 MB
High Atlas	2048 x 2048	11	R8G8B8A8	> 16 MB

For Example, 100 instance will be from 100MB to 1600MB memory usage on mobile

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

Seamless rendering, due to its inherent characteristics, employs special implementation methods and optimization techniques in lighting calculations. For instance, when calculating global illumination, we consider whether lightmaps can be utilized and whether their precision can cover pixel-level triangles. These factors significantly impact memory and storage. Hence, we've implemented a fully dynamic GI algorithm.

- No Lightmap
- Hierarchical Voxelization Global Illumination
 - BrickGI
 - Low Memory Usage
 - GPU Friendly Algorithm
 - Voxelization for Clusters
- Gathering lighting data from Bricks
 - Tracing from BrickGroup => Bricks => Voxels
 - Get lighting data from BrickVisFacesAtlas

```
// Voxel cell size is 0.5 meters
BrickVoxelSize=0.5
// a virtual voxel group
VoxelGroupSize=2x2x2
// voxel group cell size is 2 x 0.5 = 1 meters
VoxelGroupCellSize=1
// 8x8x8 voxels in one brick
BrickSize=8x8x8
// 4 meters per brick, 8 cells multiply 0.5 meters
BrickCellSize=4x4x4
// 4x4x4 bricks in one brickgroup, fast access for uint64
BrickGroupSize=4x4x4
// 16x16x16 meters, one brick is 4 meters, four brick is 16 meters
BrickGroupCellSize=16x16x16
// one group 16 meters, 32x32x8 brick groups multiply 16 meters is 512x512x128 meters
BrickGroupBitMask=32x32x8
```

Here are the default config variables.

- No Lightmap
- Hierarchical Voxelization Global Illumination
 - BrickGI
 - Low Memory Usage
 - GPU Friendly Algorithm
 - Voxelization for Clusters
- Gathering lighting data from Bricks
 - Tracing from BrickGroup => Bricks => Voxels
 - Get lighting data from BrickVisFacesAtlas

```

BrickGroupTexcoord = (SamplePoint - BrickGroupStartPos) / BrickGroupCellSize;
// check any brick exist or not in 4x4x4 bricks
BrickGroupMask = BrickGroupBitMask(BrickGroupTexcoord);
uint64 bAnyBrickExists = BrickGroupMask != 0;
if (!bAnyBrickExists)
{
    // ray marching on bricks and check visibility with BrickGroupMask
    BrickTexcoord = RayMarchingInBricks(SamplePoint, RayDirection, BrickGroupTexcoord, BrickGroupMask);
    // get brick storage in BrickMappingAtlas
    uint2 AtlasBias = BrickTexture(BrickTexcoord);
    // check any voxel exist or not in 4x4x4 voxel groups
    BrickMask = BrickBitMask(AtlasBias);
    bAnyVoxelExists = BrickMask != 0;
    if (!bAnyVoxelExists)
    {
        // raymarching with current brick voxels with BrickMask
        SampleVoxelTexcoord = RayMarchingInVoxels(SamplePoint, RayDirection, BrickGroupTexcoord, BrickTexcoord, BrickMask);
        SampleVoxelTexcoord += AtlasBias + BrickSize;
        Opacity = BrickOpacityAtlas(SampleVoxelTexcoord);
        if (Opacity != 0)
        {
            // 8x8x8 voxels / brick
            PackedVoxelVisFaceInfo = BrickMappingAtlas(SampleVoxelTexcoord);
            BiasX, BiasY, FaceIndex = UnpackVisFaceInfo(PackedVoxelVisFaceInfo);
            VisFaceAtlasTexcoord = BiasY * VisFaceAtlasSize.x + BiasX + FaceIndex;
            // we get lighting data finally.
            VoxelLighting = VoxelVisFaceLightingAtlas(VisFaceAtlasTexcoord);
        }
    }
}
    
```

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

Here is the Pseudocode how we get lighting data from brick and visfacelightingatlas. When performing ray tracing, the first step involves quickly querying brickgroupbitmask to identify which bricks within a group are valid. This involves a single uint64 sampling operation. If the mask is not zero, it indicates that the current brick group contains valid bricks, and we can determine which specific bricks are valid by checking which bits in the mask are set. After locating the brickindex through the combination of the bitmask and the sampleposition, it is straightforward to convert this index into bricktexcoord, allowing us to query the actual bias stored for the brick in the bricktexture. Subsequently, we utilize both the brickbitmask and the sampleposition to continue with a rapid ray marching process, which identifies the valid voxelcoord within the brick. Then, in the brickmappingatlas, we locate the offset of the voxel within the visfaceatlas, and proceed to read the relevant lightingdata.

- No Lightmap
- Hierarchical Voxelization Global Illumination
 - BrickGI
 - Low Memory Usage
 - GPU Friendly Algorithm
 - Voxelization for Clusters
- Gathering lighting data from Bricks
 - Tracing from BrickGroup => Bricks => Voxels
 - Get lighting data from BrickVisFacesAtlas

```
int3 RayMarchingInBricks(SamplePoint, SampleDir, BrickGroupTexcoord, BrickGroupMask)
{
    step = 0;
    while(step < 64)
    {
        SampleBrickTexcoord = (SamplePoint - BrickGroupTexcoord * BrickGroupCellSize) / BrickCellSize;
        BrickIndex = BrickTexcoord.x + (BrickTexcoord.y + BrickTexcoord.z * 4) * 4;
        if( (1 << BrickIndex) & BrickGroupMask)
        {
            return SampleBrickTexcoord;
        }
        else
        {
            SamplePoint += SampleDir * BrickCellSize;
        }
    }
    return int3(-1,-1,-1);
}
```

Here is the Pseudocode how we get valid brick from brickmask.

ADDITIONAL: GI FOR SEAMLESS RENDERING



- No Lightmap
- Hierarchical Voxelization Global Illumination
 - BrickGI
 - Low Memory Usage
 - GPU Friendly Algorithm
 - Voxelization for Clusters
- Compare with lightmap
 - Area 512x512x64m, 100+ Instances
 - Lightmap resolution is 1024x1024
 - High-end device (Immortalis-G715)
 - Shading pass is used for gathering lighting data from Bricks

GI Method	Memory	GPU Time
Lightmap	400MB	< 1ms
BrickGI + Shading	30MB	< 3 ms

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

Let's compare the performance between brickgi and lightmap on mobile devices. In this test, we added over 100 objects, each with a lightmap resolution of 1024x1024, within a scene area of 512x512 meters. Our findings show that lightmap requires at least 160MB of memory overhead, whereas brickgi only needs 30MB. Additionally, while the GPU time for brickgi is slightly higher than that of lightmap, 3ms is still considered an acceptable latency for most mobile games.

ADDITIONAL: GI FOR SEAMLESS RENDERING



- No Lightmap
- Hierarchical Voxelization Global Illumination
 - BrickGI
 - Low Memory Usage
 - GPU Friendly Algorithm
 - Voxelization for Clusters

- Area 512x512x128m
- Brick memory usage: ~30MB

Texture	Size	Format	Memory
BrickGroupBitMask	32x64x8	R32_UINT	64KB
BrickTexture	128x128x32	R16G16_UINT	2MB
BrickBitMask	64x256	R32_UINT	16KB
BrickOpacityAtlas	512x1024x8	R8	2MB
BrickMappingAtlas	512x1024x8	R32_UINT	8MB
BrickVisFacesAlbedoAtlas	1024x1024	R32_UINT	4MB
BrickVisFacesNormalAtlas	1024x1024	R32_UINT	4MB
BrickVisFacesLightingAtlas	1024x1024	R11G11B10_FLOAT	4MB
BrickVisFacesEmissiveAtlas	1024x1024	R11G11B10_FLOAT	4MB
Total	/	/	28MB

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

This table shows the details of 30MB memory allocated to cover an area within 512 meters. Moreover, high-precision clusters aren't necessary for voxelization calculations, keeping the overall voxelization cost low.

- No Lightmap
- Hierarchical Voxelization Global Illumination
 - BrickGI
 - Low Memory Usage
 - GPU Friendly Algorithm
 - Voxelization for Clusters
 - Coarser Clusters
 - Precision suffice
 - Easy to voxelize

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

Moreover, high-precision clusters aren't necessary for voxelization calculations, keeping the overall voxelization cost low. Voxelization inherently has limited precision, and clusters with lower precision are more easily voxelized through rasterization

ADDITIONAL: SHADOW FOR SEAMLESS RENDERING



- Clusters LODs
 - Coarser Clusters will suffice
 - Acceptable effects
 - Drawing fewer triangles
- Shadow Cache
 - An effective method

finest cluster shadow rendering (BiasLODError = 0)

```
MobileSceneRender.ShadowDepths.Atlas0 6144x2048
vkCmdBeginRenderPass2(C=Don't Care, D=Load)
vkCmdDrawIndirect(1) => <384, 205>
vkCmdEndRenderPass(C=Store, D=Store)
```

coarser cluster shadow rendering(BiasLODError = 2)

```
MobileSceneRender.ShadowDepths.Atlas0 6144x2048
vkCmdBeginRenderPass2(C=Don't Care, D=Load)
vkCmdDrawIndirect(1) => <384, 117>
vkCmdEndRenderPass(C=Store, D=Store)
```

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

Regarding direct lighting calculations, our tests indicate that high-precision clusters for generating shadow depths aren't essential. Using coarser approximations has minimal impact on visual quality but significantly enhances performance. For example, the upper image uses a more detailed rendering requiring over 200 drawcalls, while the lower image achieves similar results with around 100 drawcalls.

ADDITIONAL: SHADOWS FOR SEAMLESS RENDERING

- Clusters LODs
 - Coarser Clusters will suffice
 - Acceptable effects
 - Drawing fewer triangles
- Shadow Cache
 - An effective method

finest cluster shadow rendering (BiasLODError = 0)



coarser cluster shadow rendering(BiasLODError = 2)



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

The upper image here shows the effect with 200 drawcalls, and the lower image displays the shadowmap with 100 drawcalls. The difference is barely noticeable. Additionally, one of the most effective methods is shadow caching.



Let's look at seamless rendering's performance, comparing effects, framerates, power consumption, package size, and production efficiency.

- Benefits

- Quality
 - More than 80 million triangles
- FPS
 - 60 FPS On Mobile
- Power
 - ~ 128 mW
- Package Size
 - Only 1 / 4 vs Manual Mesh Level LOD
- Productivity
 - More 4x efficiency

Scene

Type	Mesh count	Triangles of Mesh	Total
Rock	40	2 million	80 million
Landscape	4096	1024	4 million
Grass	10000+	20	2 million



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

Quality-wise, we can render 80 million triangles simultaneously, unimaginable on mobile, especially retail devices. With this many triangles, we maintain stable framerates and acceptable power consumption. We also save on package size. Plus, we don't need to manually create LODs; just LOD0 assets suffice.

FPS ON MOBILE (Immortalis-G715)



Here's the framerate on a mobile device. Blue and green lines show the rendering effect after using seamless mobile rendering. Notice the blue curve, maintaining around 60fps with minimal fluctuations.

PROFILING ON MOBILE

On high-end mobile device (Immortalis-G715)

Pass Stage	FrameTime(ms)
Instance Culling	0.11
Cluster Culling	0.07
Binning	0.2
Rasterization	0.72
Material Classify	0.92
Material Shading	0.87
Total Cost	~3ms

On low-end mobile device (5 years ago)

Pass Stage	FrameTime(ms)
Instance Culling	1.08
Cluster Culling	1.79
Binning	3.25
Rasterization	2.67
Material Classify	3.42
Material Shading	7.31
Total Cost	~20ms

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

Take a look at the detailed performance data on different level mobile device. On high-end device ,it only take about 3ms gpu time,while on low-end mobile device which is 5 years ago,it will be about 20ms.

PROFILING ON MOBILE (Immortalis-G715)



Frame Time (~3ms ,which is sum of 0.11 + 0.07 + 0.2 + 0.72+0.92+0.87)

	Start Frame Time(ms)	Cost Time(ms)	Primitives	Delta Primitives	Draw Call	Delta Draw Call
NanoMesh Off	5.32	-	116.7k	-	89	-
Instance Culling	5.43	0.11	116.7k	0	89	0
Cluster Culling	5.5	0.07	116.7k	0	89	0
Binning	5.7	0.2	116.7k	0	89	0
Rasterization	6.42	0.72	116.7k	0	90	1
Material Classify	7.34	0.92	116.7k	0	102	2
Material Shading	8.21	0.87	116.7k	0	190	7

IO/Streaming Time (~1.5ms)

	Start Frame Time(ms)	Cost Time(ms)	Memory(MB)	Delta Memory(MB)
Streaming Off	6.7	-	623.66	-
Streaming On	8.25	1.55	660.22	36.56

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

We can also take a look at the detailed performance data on high-end mobile device. In the table above, I've broken down each rendering phase and monitored the changes in GPU time consumption step by step. If enable instance culling, it will take 0.11 ms, and then turn on cluster culling, it take 0.07ms, and then turn on binning pass, it take 0.2 ms, then rasterization pass which take 0.72ms, then material classify pass take 0.92ms, the last pass is material shading which take 0.87ms. We can observe that the time consumption for culling is not significant, and in this scenario, rasterization and material shading is the more time-consuming part. The table below shows the time consumption of IO Streaming, and we can see that the GPU's IO performance is also acceptable.

PROFILING ON MOBILE (LOW-END/ 5 YEARS AGO)



Frame Time (~20 ms, which is sum of 1.08 + 1.79 + 3.25 + 2.67+3.42+7.31)

	Start Frame Time(ms)	Cost Time(ms)	Primitives	Delta Primitives	Draw Call	Delta Draw Call
NanoMesh Off	24.98	-	116.7k	-	89	-
Instance Culling	26.06	1.08	116.7k	0	89	0
Cluster Culling	27.85	1.79	116.7k	0	89	0
Binning	31.1	3.25	116.7k	0	89	0
Rasterization	33.77	2.67	116.7k	0	90	1
Material Classify	37.19	3.42	116.7k	0	102	2
Material Shading	44.5	7.31	116.7k	0	190	88

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

We can also take a look at the detailed performance data on low-end mobile device which is 5 year ago. We can observe that the time consumption for culling is not significant, and in this scenario, material shading is the most time-consuming part.

POWER ON MOBILE (Immortalis-G715)

Power (in 10 minutes)

Off On Adaptive



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

This is power consumption over 10 minutes on a mobile device. Seamless rendering on mobile (blue) performs well. We found that after enabling this, the power consumption is significantly lower than that of traditional mesh rendering mode.



These are our research and validation results, but there's much room for improvement. This is just the beginning.

- Optimization
 - Lower Memory Usage
 - Software VRS with Visibility Buffer
 - DirectStorage IO Optimization
 - 1.5ms is still a bit little heavy
- Features
 - RayTracing Geometry Support
 - MeshShader Usage for High-end Devices
 - Bindless Material Shading

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

We'll continue optimizing performance and iterating features based on advancements in rendering technology and hardware capabilities. such as memory usage, software vsr with buffer, io optimization, 1.5ms is still a bit little heavy. For features, We are not supporting raytracing now. And with mobile hardware capabilities development, we will add meshshader and bindless material shading support.



Thank you for your time. Feel free to reach out to me if you have any further questions.