

# Grappling With Performance: Rendering Optimization Strategies In Rumbleverse

Jon Moore  
Graphics Engineer  
Iron Galaxy Studios

[jonmanatee.com/s/GDC2023.pdf](https://jonmanatee.com/s/GDC2023.pdf)



Hello everyone, I'm Jon Moore, I work as a Graphics Engineer at Iron Galaxy Studios, and I'd like to welcome you to my talk "Grappling With Performance: Rendering Optimization Strategies in Rumbleverse"

# Rumbleverse

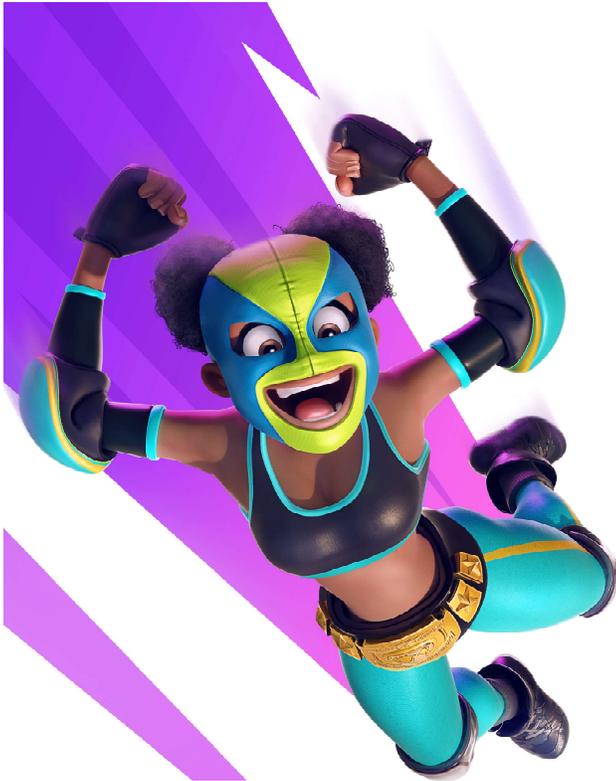
- 40-Player Battle Royale
- Melee Combat
- 1 km x 1 km Island Arena
- Performance target: 1080p  
60 FPS on PS4
- Shipped on Unreal 4.27.1



If you're not familiar with the game, Rumbleverse is a 40 player battle royale game featuring primarily melee combat on a 1km squared island arena.

Our gold standard performance target for the game is 1080p + 60 FPS on a base PS4 – with the game shipping across both generations of Sony and Microsoft consoles.

And an important detail to note – we are an Unreal licensee, and the game shipped using Unreal version 4.27.1, but using an engine that is not developed internally at IGS doesn't discourage me from finding plenty of optimization opportunities for the game, which is the motivation for this talk today



# Optimization Philosophy

- Choice of Engine doesn't affect optimization potential – your content creates opportunities and exposes limitations of the engine
- The fastest wavefront is the one you never launch
- Everything worth doing in life is at least an 0.1 ms speedup

Before getting into specifics, I think I should preface with my personal philosophy on optimization, especially as it relates to work done on the GPU. These are ideas that formed before I started working on Rumbleverse, but my experiences here have only further solidified these beliefs. (CLICK)

First and foremost, I think there is optimization potential anytime an engine is not built from scratch for a particular game. It is always useful to reflect on how your specific content might allow the engine to be modified/configured to run it most efficiently. This is very true with a widely licensed engine like Unreal, but I've seen this with shared engine tech used internally between teams, or even just when building off the engine used for a previous game on a new project that is not a direct sequel. (CLICK)

Secondly, I have over time come to appreciate that culling out redundant/unneeded work on the GPU is something that can be gone back to again and again when looking for gains. In the simplest sense, it's making sure that every wavefront running on the GPU is actually contributing to the final image. (CLICK)

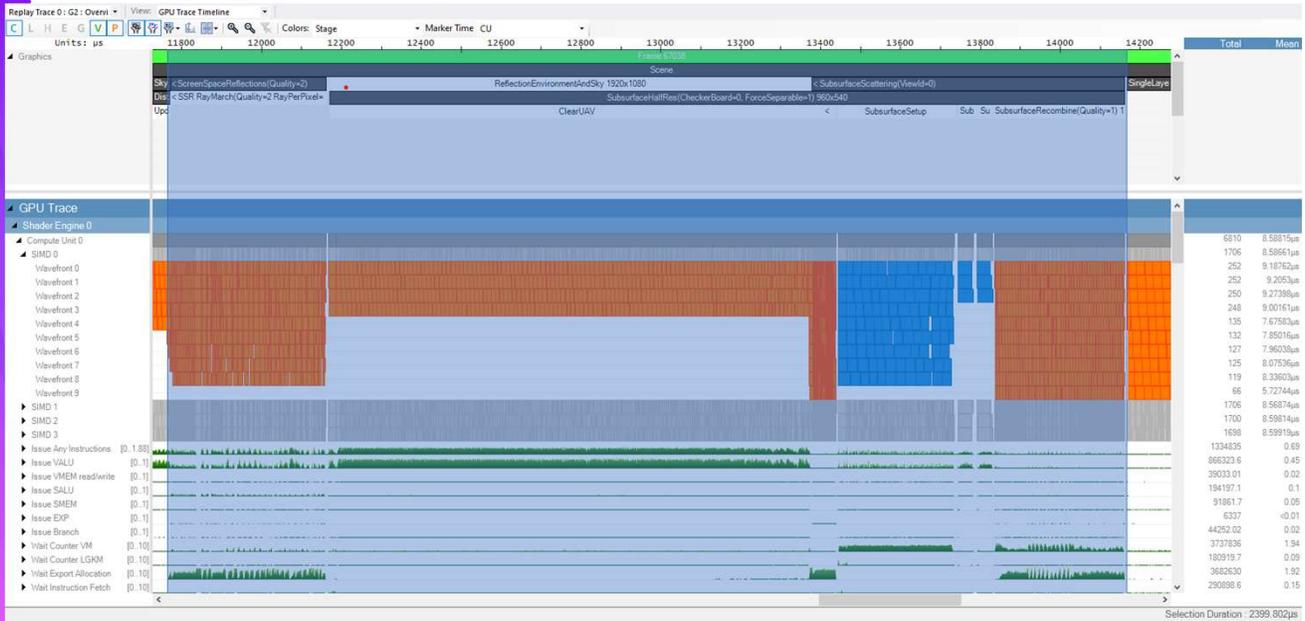
Finally, whenever I am working on a 60 hz video game, I generally use 0.1 ms as my measuring stick for if an optimization is worth doing. Less than that, and I will sit on a change as not being worth the risk of modifying the engine. 0.1 ms is often my sweetspot of feeling like a change is worthwhile and enough 0.1 ms changes added together will eventually make a big impact.



# Part 1: Proactive Optimizations

From day 1 of my time on Rumbleverse – rendering was my primary focus, and I wanted to make an impact on the game to ensure we could push our look to the best it could be. This was in service of allowing artists to build richer environments, have more detailed materials, and avoid dynamic resolution drops as much as possible. This first section is focused on optimizations made proactively in pursuit of this goal as we set-up our initial configuration of rendering features across our target platforms.

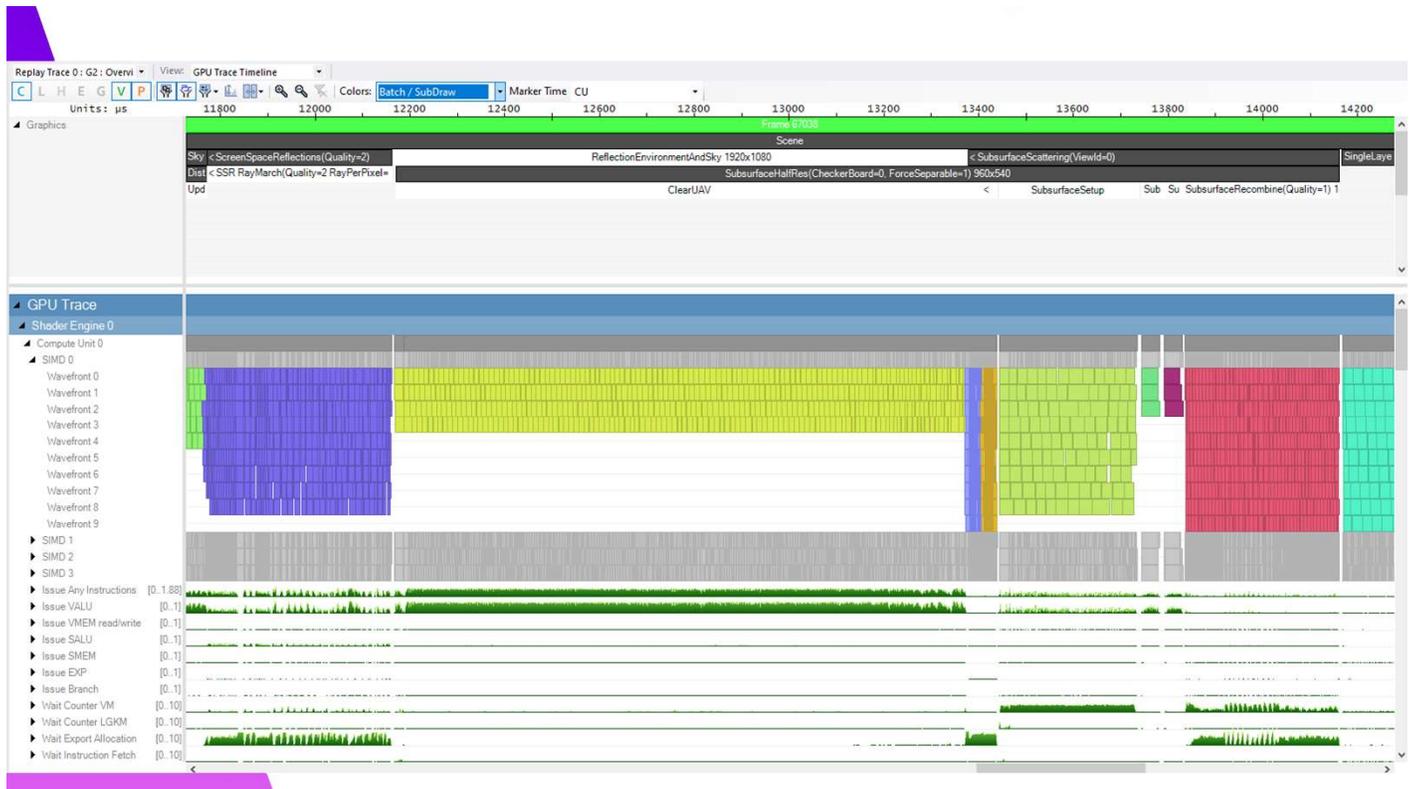
# Reflections + Subsurface



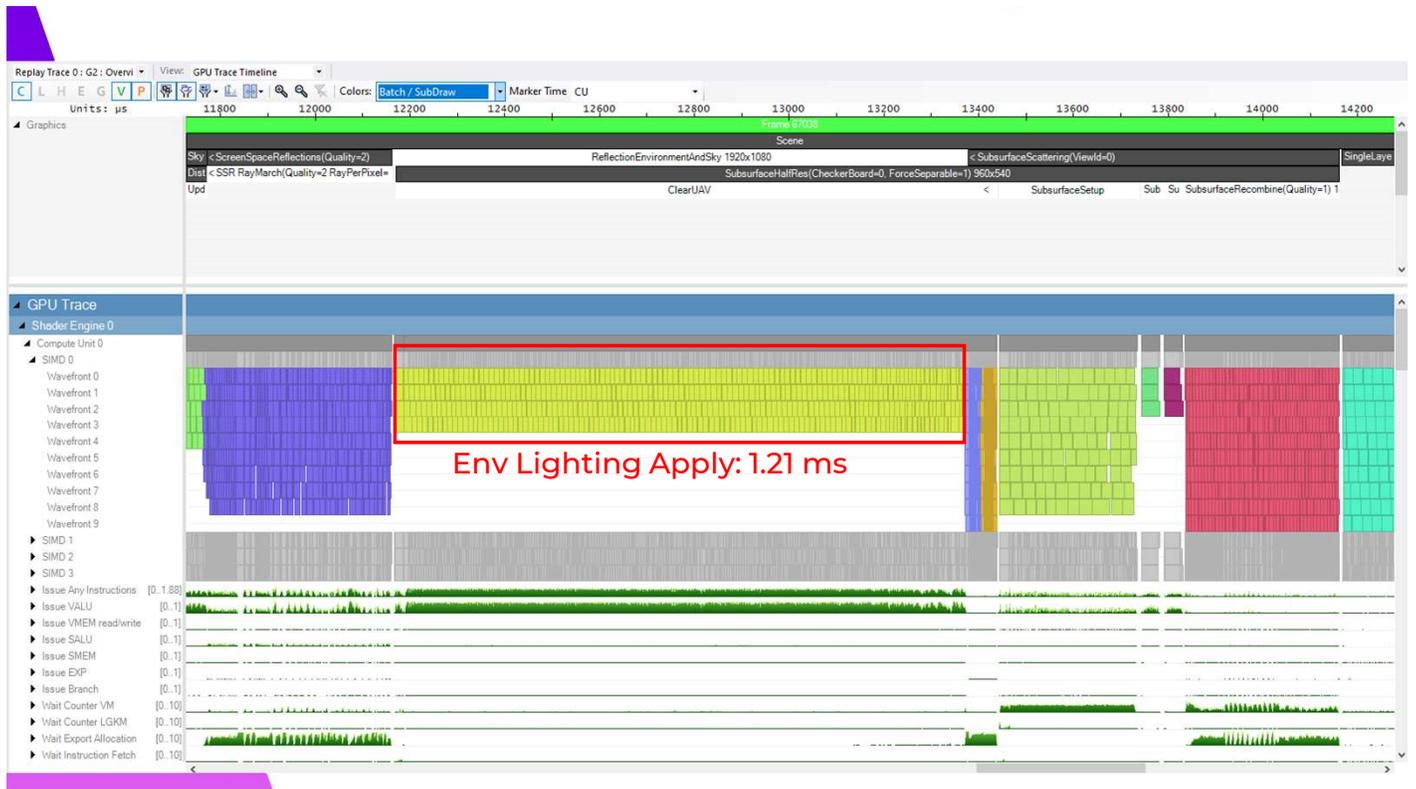
Let's start by looking at the 2.4 ms stretch of frame time in Razor on PS4 that I found myself looking at early on in development, where reflections and subsurface scattering are handled in stock UE4. All parts of this time were decided to be pretty important to the look of our game, and I wanted to try to reduce its cost as that is 14% of a 16.6 ms frame time.



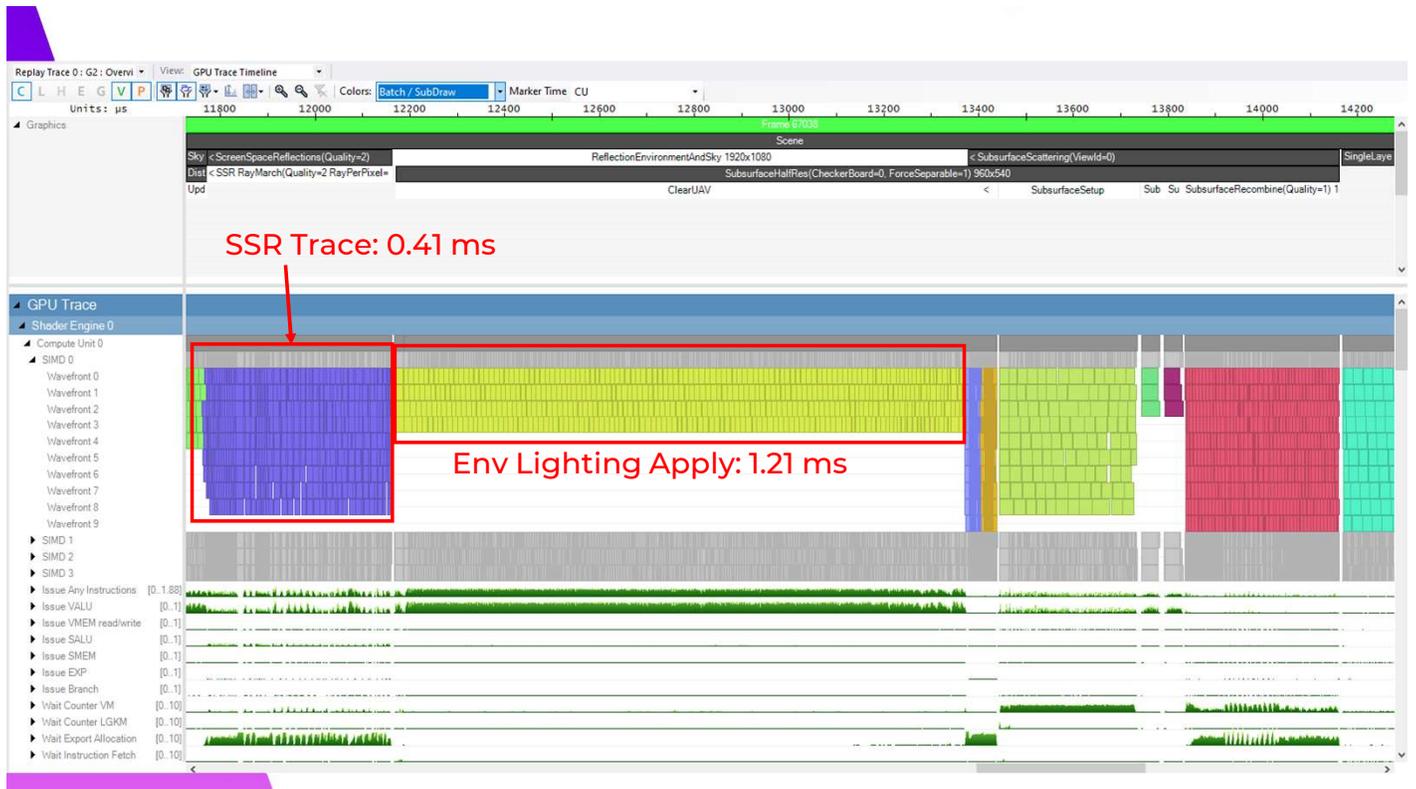
Here is the scene in question that the trace is from. I'll revisit this in a bit, but it's pretty standard for us: we have some foliage, a character, some sky, some shiny metallic objects, and some reflective windows.



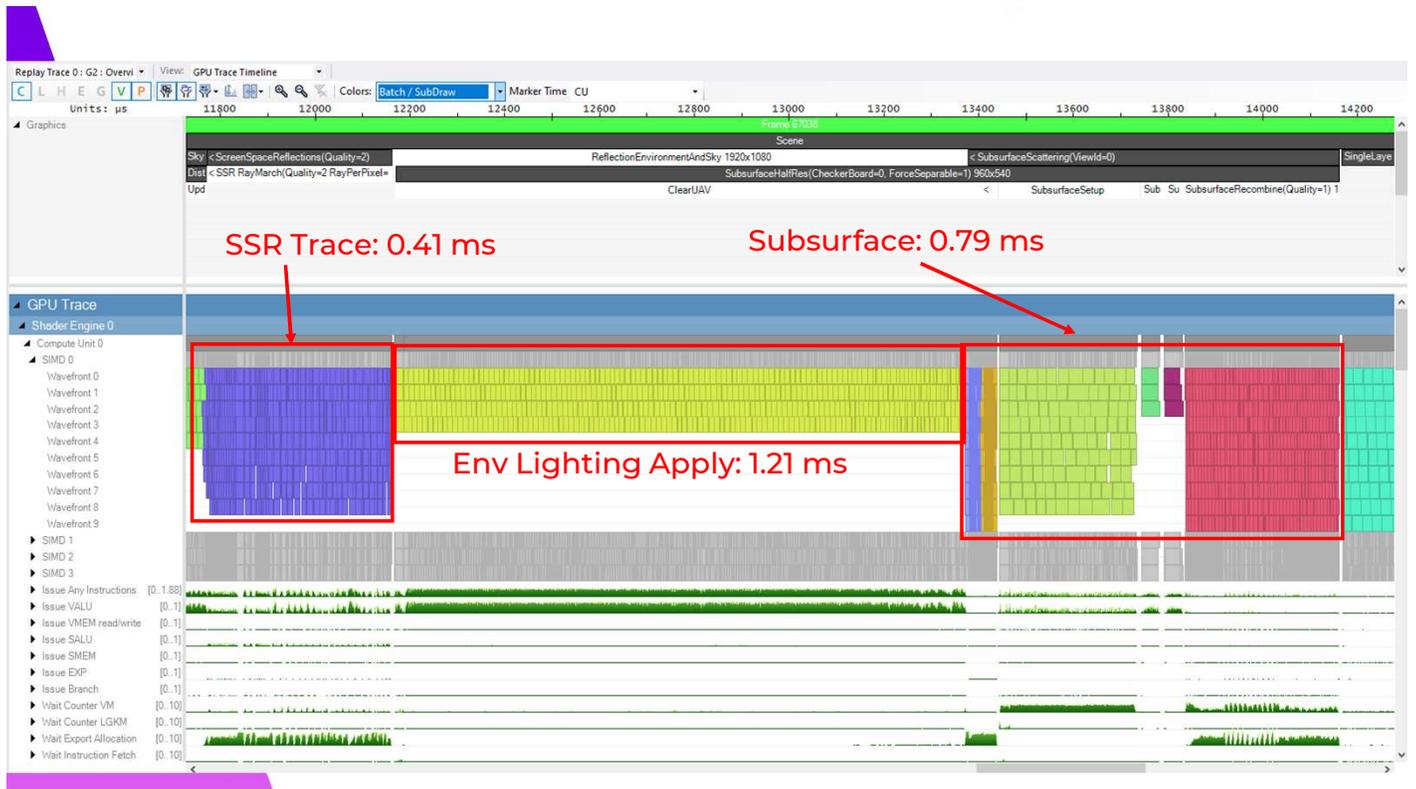
Let's go back to that stretch of time in Razor with the wavefronts colored by batch.



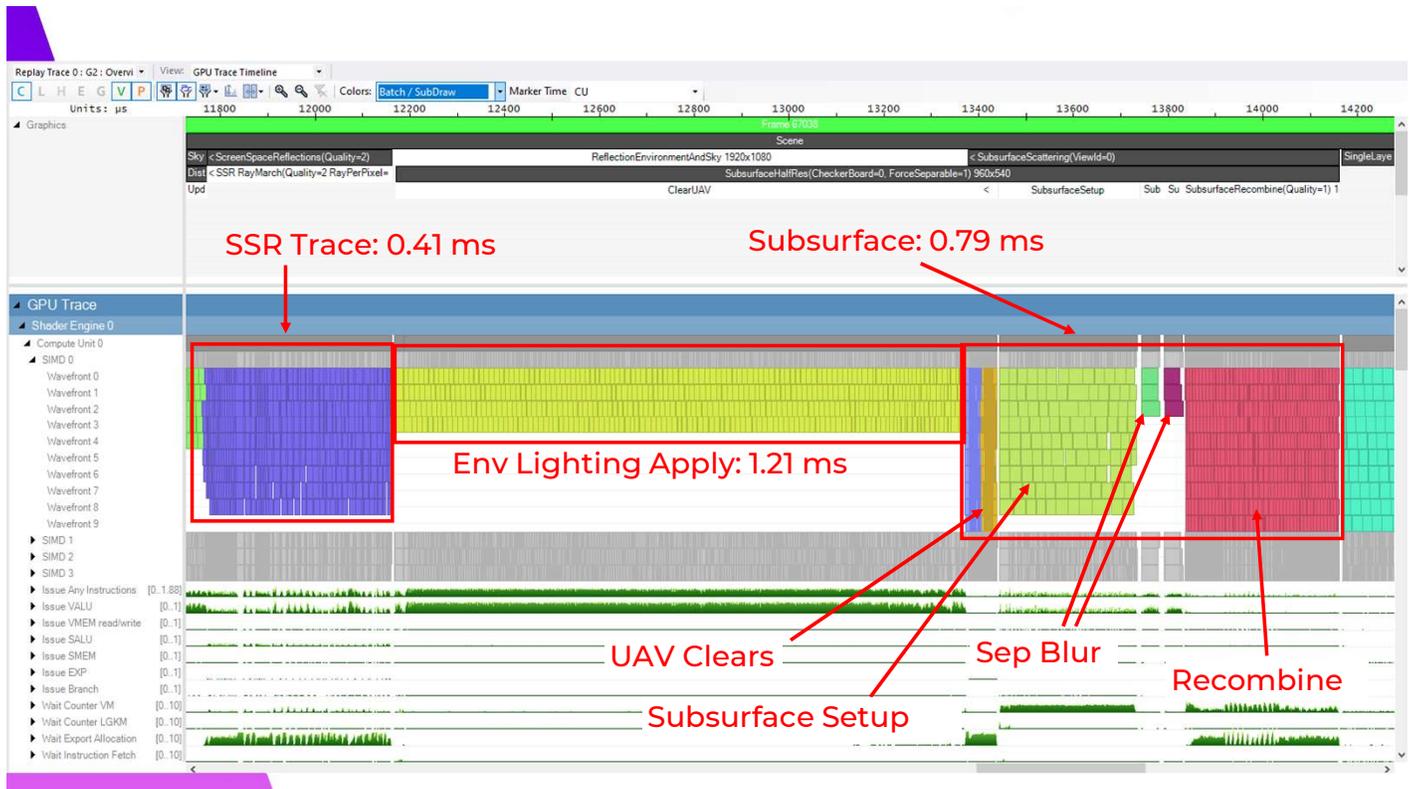
This long section in the middle is the Env Lighting Apply pixel shader, clocking in at 1.21 ms. This is where indirect diffuse, specular, and skylighting is combined together after the direct lighting is handled by the deferred lighting loop.



Before that, Screen Space Reflection tracing takes 4 tenths of a ms



And this time after Env Lighting Apply is 0.79 ms on subsurface scattering



If you're not familiar, SSS in this version of Unreal clears some UAVs, does SubsurfaceSetup which does tile classification and downsamples the subsurface to half res, and then dispatches compute shader batches for tiles based on the subsurface algorithm selected by tile classification.

And then finally there is a recombine with the Scene Color after the blur. Classifying by algorithm was added to Unreal after our character look was developed, so our tiles all fall into the path using Screen Space Separable Subsurface – based on the technique developed by Jorge Jimenez

*Separable Subsurface Scattering and Eye Rendering* by Jorge Jimenez:  
[http://advances.realtimerendering.com/s2012/activision/Jimenez-Separable\\_Subsurface\\_Scattering\\_and\\_Eye\\_Rendering\(Siggraph2012\).pptx](http://advances.realtimerendering.com/s2012/activision/Jimenez-Separable_Subsurface_Scattering_and_Eye_Rendering(Siggraph2012).pptx)



# Tile Classification

- Original Idea: improve occupancy in reflection apply by using tile classification
- Inspired by Ramy El Garawany's presentation: *Deferred Lighting in Uncharted 4*
- Algorithm:
  1. analyze GBuffer
  2. build lists of tiles based on material properties
  3. render each using different shader permutations + DispatchIndirect

Mentioning tile classification on SSS is an interesting launching point to how I approached improving these systems. Naughty Dog first presented a version of the tile classification technique for reducing per-tile cost of their lighting in the talk *Deferred Lighting in Uncharted 4* from Siggraph 2016. (CLICK)

I thought that the heavy time spent on reflection apply could be optimized in a similar way here. I would write a tile classification shader that looked at the Gbuffer properties of an 8x8 group of pixels and build lists based on the materials present. Then each list can be rendered using DispatchIndirect with the different shader permutations bound to each dispatch.

*Deferred Lighting in Uncharted 4* by Ramy El Garawany:

[http://advances.realtimerendering.com/s2016/s16\\_ramy\\_final.pptx](http://advances.realtimerendering.com/s2016/s16_ramy_final.pptx)



So for example in this frame we can pretty clearly see groupings of pixels that are all default lit or twosided foliage.



And here you can see a visualization of the tile classification in practice. The green tiles are default lit, the blue is all foliage, and the red contains a shading paths that are “complex” and run the full shader



However, the most important thing this work took me to was the unlit tiles here, tiles that are culled entirely had the biggest impact on the running time. This got me thinking about how the tile classification could be used for culling workload from SSR + SSS as well, and that cost of running the classification would be shared across multiple steps of our rendering.

Remember, the fastest wavefront is the one you never launch!



# Tile Classification 2.0

- Better Idea: skip empty waves in multiple passes based on a single classify step
- Cull unlit (skybox tiles) from everything
- Skip SSR trace on tiles above roughness threshold
- Skip SSS on tiles with no skin materials in them, run required clears with simplified shader permutation

So here's the updated plan: cull unlit tiles from all of the passes, add a classification for if all the pixels are too rough to trigger SSR traces, and skip SSS on tiles with no skin materials on them, and run a simplified clear on tiles that need to be cleared but don't need the full SSS set-up.

# Tile Classify Shader

```
44  uint bAnySSSProfile = 0;
45
46  // loop over each 8x8 tile within the 16x16 pixel area
47  uint2 PixelOffsets[4] = { uint2(0,0), uint2(1,0), uint2(0,1), uint2(1,1) };
48  UNROLL
49  for (int i = 0; i < 4; ++i)
50  {
51      uint2 PixelPos = (DispatchThreadId.xy * 2 + ViewDimensions.xy);
52      FScreenSpaceData ScreenSpaceData = GetScreenSpaceDataUint(PixelPos + (PixelOffsets[i] * 8)); ← Read GBuffer
53      GBufferData InGBufferData = ScreenSpaceData.GBuffer;
54
55      uint bIsDefaultLit = (InGBufferData.ShadingModelID == SHADINGMODELID_DEFAULT_LIT) ? 1 : 0;
56      uint bIsFoliageLit = (InGBufferData.ShadingModelID == SHADINGMODELID_TWOSIDED_FOLIAGE) ? 1 : 0;
57      uint bIsComplexLit = (InGBufferData.ShadingModelID > SHADINGMODELID_DEFAULT_LIT) ? 1 : 0;
58      uint bIsSSSProfile = UseSubsurfaceProfile(InGBufferData.ShadingModelID) ? 1 : 0;
59
60      float Roughness = InGBufferData.Roughness;
61      float RoughnessFade = GetRoughnessFade(Roughness);
62      uint bSkipSSR = (RoughnessFade <= 0.0 || InGBufferData.ShadingModelID == SHADINGMODELID_UNLIT) && InGBufferData.ShadingModelID != SHADINGMODELID_CLEAR_COAT;
63
64      // OR results
65      uint MergedResult = (bIsSSSProfile << 2) | (bIsComplexLit << 1) | bIsDefaultLit;
66      MergedResult = WaveAllBitOr(MergedResult);
67      uint bAnyDefaultLit = MergedResult & (1 << 0); ← WaveOp Bitwise Or+And
68      uint bAnyComplexLit = MergedResult & (1 << 1);
69      bAnySSSProfile = bAnySSSProfile | (MergedResult & (1 << 2));
70
71      // AND results
72      MergedResult = (bSkipSSR << 2) | (bIsFoliageLit << 1) | bIsDefaultLit;
73      MergedResult = WaveAllBitAnd(MergedResult);
74      uint bAllDefaultLit = MergedResult & (1 << 0); ← WaveOp Bitwise Or+And
75      uint bAllFoliageLit = MergedResult & (1 << 1);
76      uint bAllSkipSSR = MergedResult & (1 << 2);
77
78      // select which permutation
79      uint PermutationIndex = NUM_PERMUTATIONS;
80      if (bAllFoliageLit)
81      {
82          PermutationIndex = 4;
83      }
84      else if (bAllDefaultLit)
85
```

let's walk through a little bit of what is happening in the Tile Classify shader code. Classification happens on 8x8 tiles, but each group covers a 16x16 area because subsurface scattering is happening at half res. After sampling the gbuffer properties with UE4's GetScreenSpaceDataUint function, I use wave ops to merge the bitmasks for each 8x8 tile together. (CLICK)

You can see that in the code this happening with the UE4 shader API commands WaveAllBitOr and WaveAllBitAnd. After these wave ops, each thread in the wavefront is going to hold the same mask value in MergedResult.

One benefit to using wave ops is that the logic following the wave commands is all going to be scalar ALU, as the compiler knows MergedResult is going to be uniform across the wave

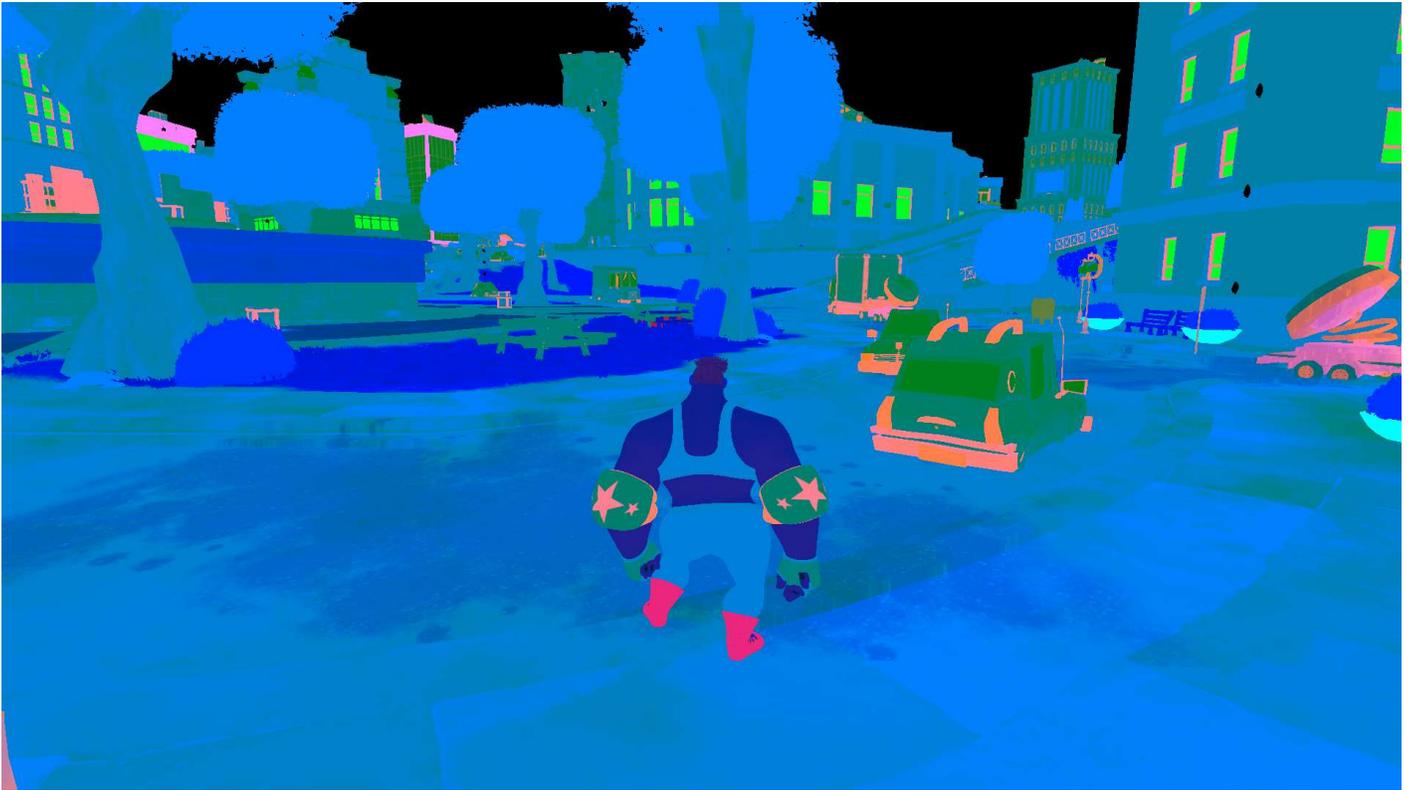
# Tile Classify Shader

```
77
78 // select which permutation
79 uint PermutationIndex = NUM_PERMUTATIONS;
80 if (bAllFoliageLit)
81 {
82     PermutationIndex = 4;
83 }
84 else if (bAllDefaultLit)
85 {
86     PermutationIndex = 6;
87 }
88 else if (bAnyComplexLit)
89 {
90     PermutationIndex = 0;
91 }
92 else if (bAnyDefaultLit)
93 {
94     PermutationIndex = 2;
95 }
96
97 // odd half of permutations lack SSR completely
98 if (bAllSkipSSR)
99 {
100     PermutationIndex += 1;
101 }
102
103 // write out the 8x8 data
104 // first thread does atomic increment and write, fully unlit tiles are skipped entirely
105 if (GroupIndex == 0 && PermutationIndex < NUM_PERMUTATIONS)
106 {
107     uint TileIndex;
108     InterlockedAdd(RWTileDispatchCounts[PermutationIndex * 3], 1, TileIndex);
109
110     uint TileLocationID = TileIndex + (PermutationIndex * NumTiles);
111     uint2 TileLocation = (GroupID * 2) + PixelOffsets[i];
112     RWTileLocationsBuffer[TileLocationID] = TileLocation.x | (TileLocation.y << 16);
113 }
114 }
115
116 // SSR permutations go beyond the end of the normal reflection tile permutations
117 uint SSSPermutationIndex = NUM_PERMUTATIONS + 1;
118 if (bAnySSRPerfFile)
```

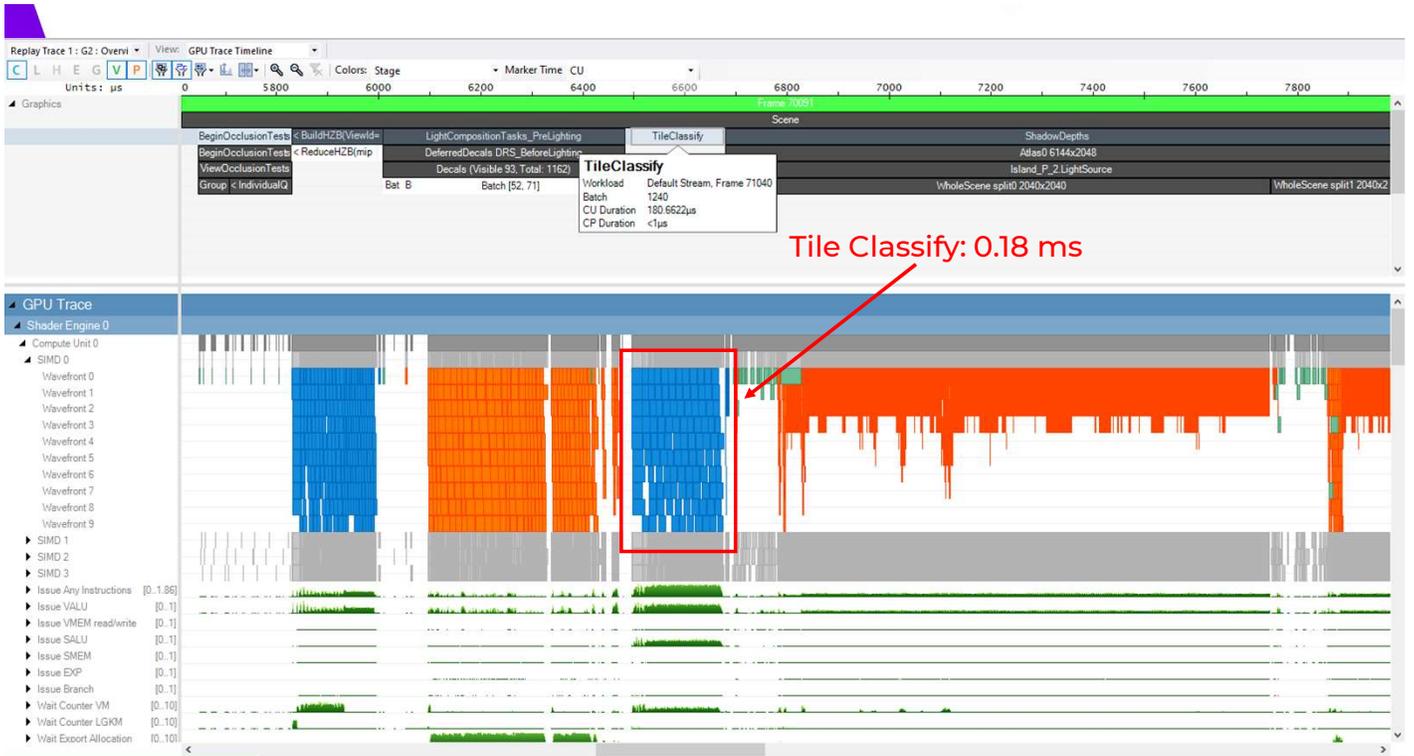
Then the shader permutation for the tile is selected based on the bits held in the MergedResult across the whole wave, and the result is written out by the first thread. (CLICK) An interlocked add occurs on the counts which gets a unique index for the tile that maps to the Tile Locations Buffer, which holds the pixel location for a given tile on the screen and will be used to reconstruct the pixel locations for each tile in the apply shaders.

Note that I selected 8x8 tiles because those are the size of 1 wavefront (64 threads) on GCN. Uncharted 4 used 16x16 tiles – which cuts the memory needed for the tile lists to 25%. Tile location lists require memory equal to max tile count \* permutation count. 8x8 tiles will allow tighter bounds on expensive material paths. I've opted for 8x8 partly because our permutation count is more limited. For example, it's in my backlog to try adding a path for tiles that contain foliage and default lit, as that is our most common boundary case falling into the complex path, but that would require more memory.

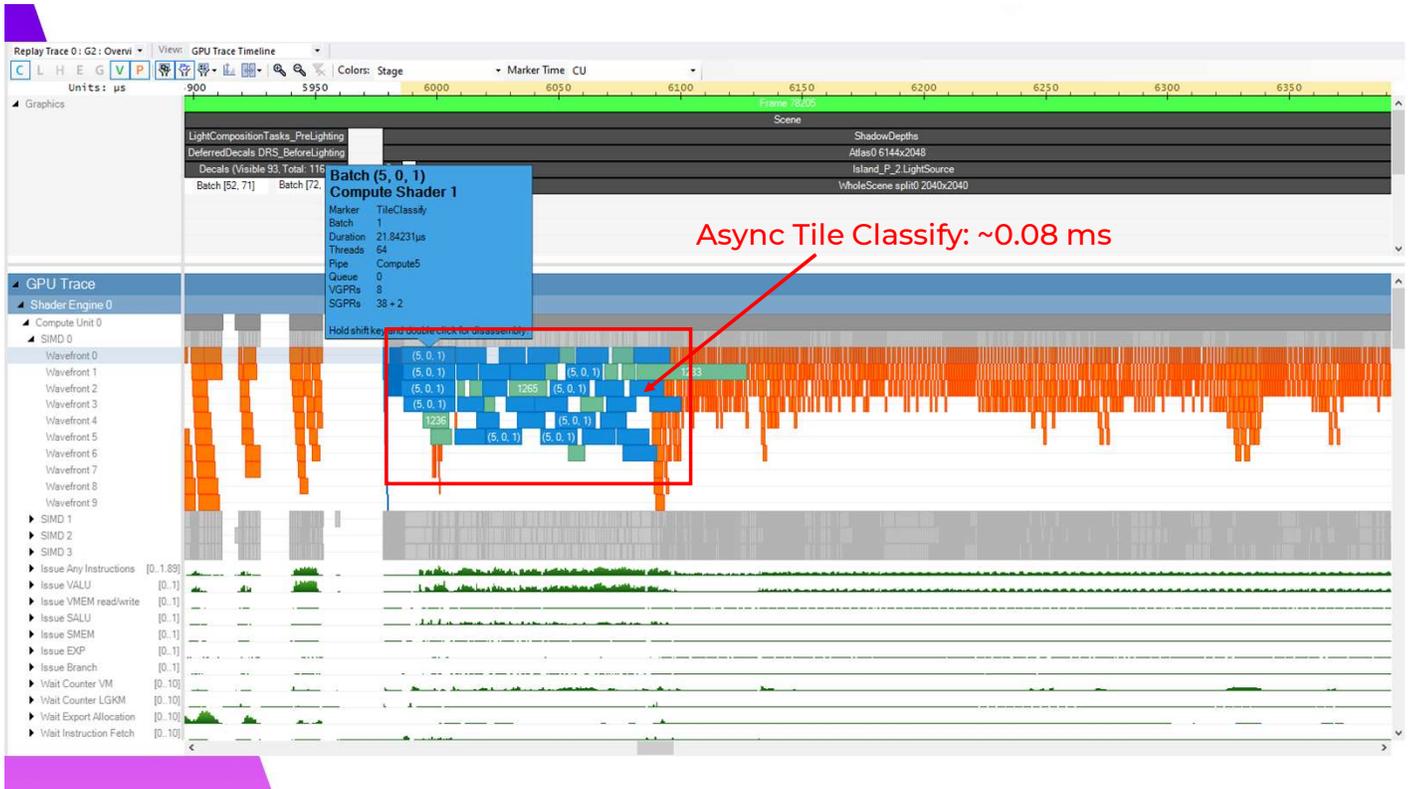
We have 10 shader permutations currently – which results in a 1.296 mb tile locations buffer for 8x8 tiles at 1080p. I should be able to reclaim 48 kb if I didn't over-allocate for the half-res tile lists, but most of the memory is coming from the 8 permutations used for SSR+Reflection Apply



It's worth noting – one really great thing about this shader is that those calls to sample the gbuffer properties all map back to just a single texture read. Epic has conveniently already packed all that information into just one Gbuffer target for us, which holds both Roughness and Material ID, which I'm showing for my shot here.



Now back in razor let's look at the cost of running this classification shader – it takes a modest 0.18 ms on a base PS4 at 1080p, and can execute as soon as the decals are done modifying the gbuffer



And we can actually do better, this classification job very nicely overlaps with shadow depth rendering using async compute, here you can see it overlapping with some vertex shading work before some pixel shader waves for masked materials run. This is frame dependent, but I generally see running it async save ~0.1 ms of frame time which makes the cost approximately 0.08 ms on a PS4.

# Env Apply Shader

```
487 // compute version of reflection and skylighting for dispatching tiles classified by shader features needed
488 [numthreads(8, 8, 1)]
489 void ReflectionEnvironmentSkyLightingCS(
490     uint3 GroupId : SV_GroupID,
491     uint3 DispatchThreadId : SV_DispatchThreadId, // DispatchThreadId = GroupId * int2(dimx,dimy) + GroupThreadId
492     uint3 GroupThreadId : SV_GroupThreadId, // 0..THREADGROUP_SIZEX 0..THREADGROUP_SIZEY
493     uint GroupIndex : SV_GroupIndex) // SV_GroupThreadId // SV_GroupThreadId.z*dimx*dimy + SV_GroupThreadId.y*dimx + SV_GroupThreadId.x
494 {
495     // lookup into tile data with group ID
496     uint TileLocationData = TileLocationsBuffer[GroupId.x + TILE_PERMUTATION * NumTiles];
497     // unpack tile location
498     uint2 PixelPos = 0;
499     PixelPos.x = (TileLocationData & 0xFFFF) * 8 + GroupThreadId.x;
500     PixelPos.y = (TileLocationData >> 16) * 8 + GroupThreadId.y;
501     PixelPos += ViewDimensions.xy;
502
503     float4 UVAndScreenPos;
504     UVAndScreenPos.xy = (float2(PixelPos.xy) + .5f) / (ViewDimensions.zw - ViewDimensions.xy);
505     UVAndScreenPos.zw = float2(2.0f, -2.0f) * UVAndScreenPos.xy + float2(-1.0f, 1.0f);
506
507     float4 SvPosition = float4(PixelPos.x, PixelPos.y, 0.f, 1.f);
508     float2 BufferUV = UVAndScreenPos.xy;
509     float2 ScreenPosition = UVAndScreenPos.zw;
510
511     // Sample scene textures.
512     GBufferData GBuffer = GetGBufferDataFromSceneTextures(BufferUV);
513
514     // Sample the ambient occlusion that is dynamically generated every frame.
515     float AmbientOcclusion = AmbientOcclusionTexture.SampleLevel(AmbientOcclusionSampler, BufferUV, 0).r;
516
517     // override GBuffer data if all pixels have same type
518     #if ALL_DEFAULT_LIGHTING
519         GBuffer.ShadingModelID = SHADINGMODELID_DEFAULT_LIT;
520     #elif ALL_FOLIAGE_LIGHTING
521         GBuffer.ShadingModelID = SHADINGMODELID_TWOSIDED_FOLIAGE;
522     #elif HAS_COMPLEX_LIGHTING
523         // if no complex lighting pixels we can do this clamp as a hint that everything is either unlit or default lit
524         GBuffer.ShadingModelID = clamp(SHADINGMODELID_UNLIT, SHADINGMODELID_DEFAULT_LIT, GBuffer.ShadingModelID);
525     #endif
526 }
```

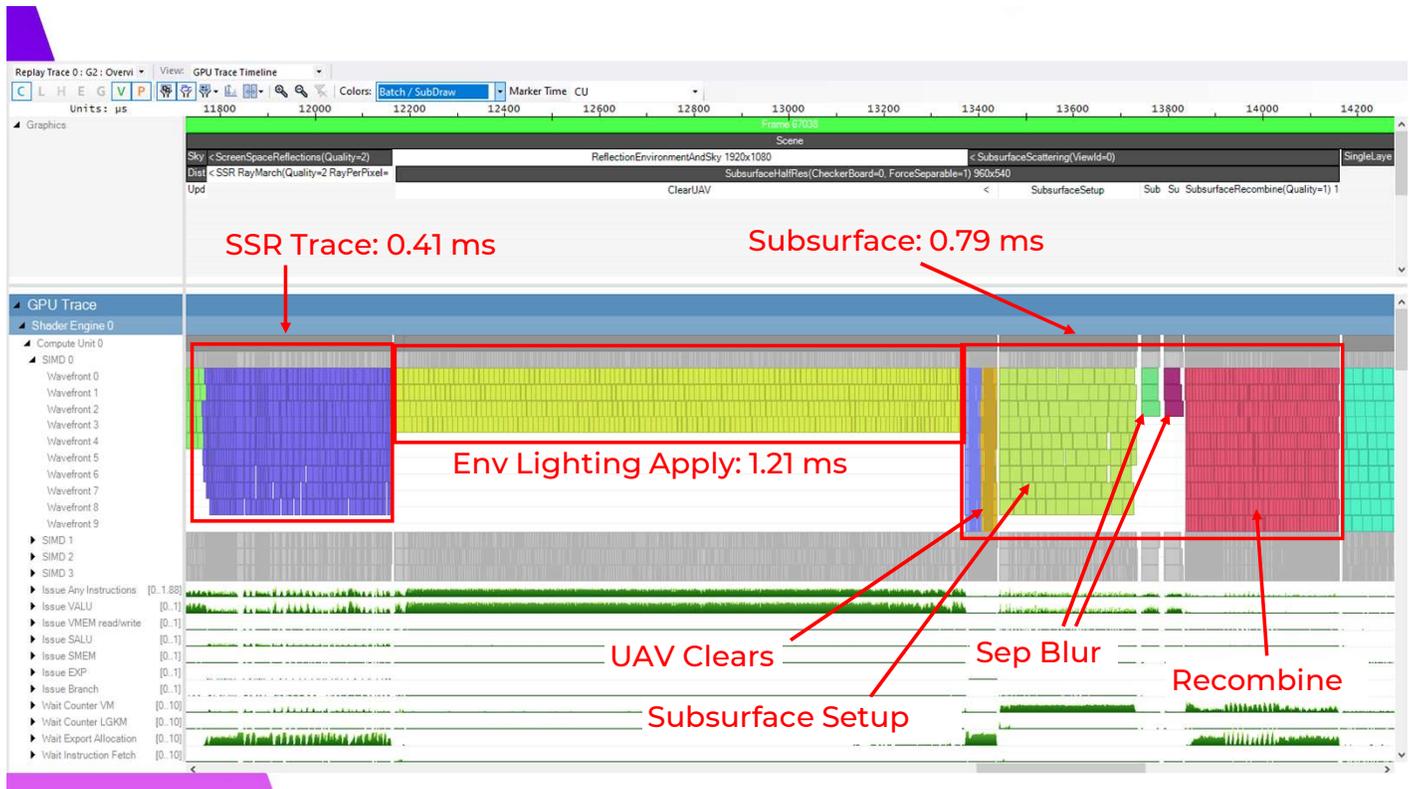
Reconstruct Tile Location

Read GBuffer

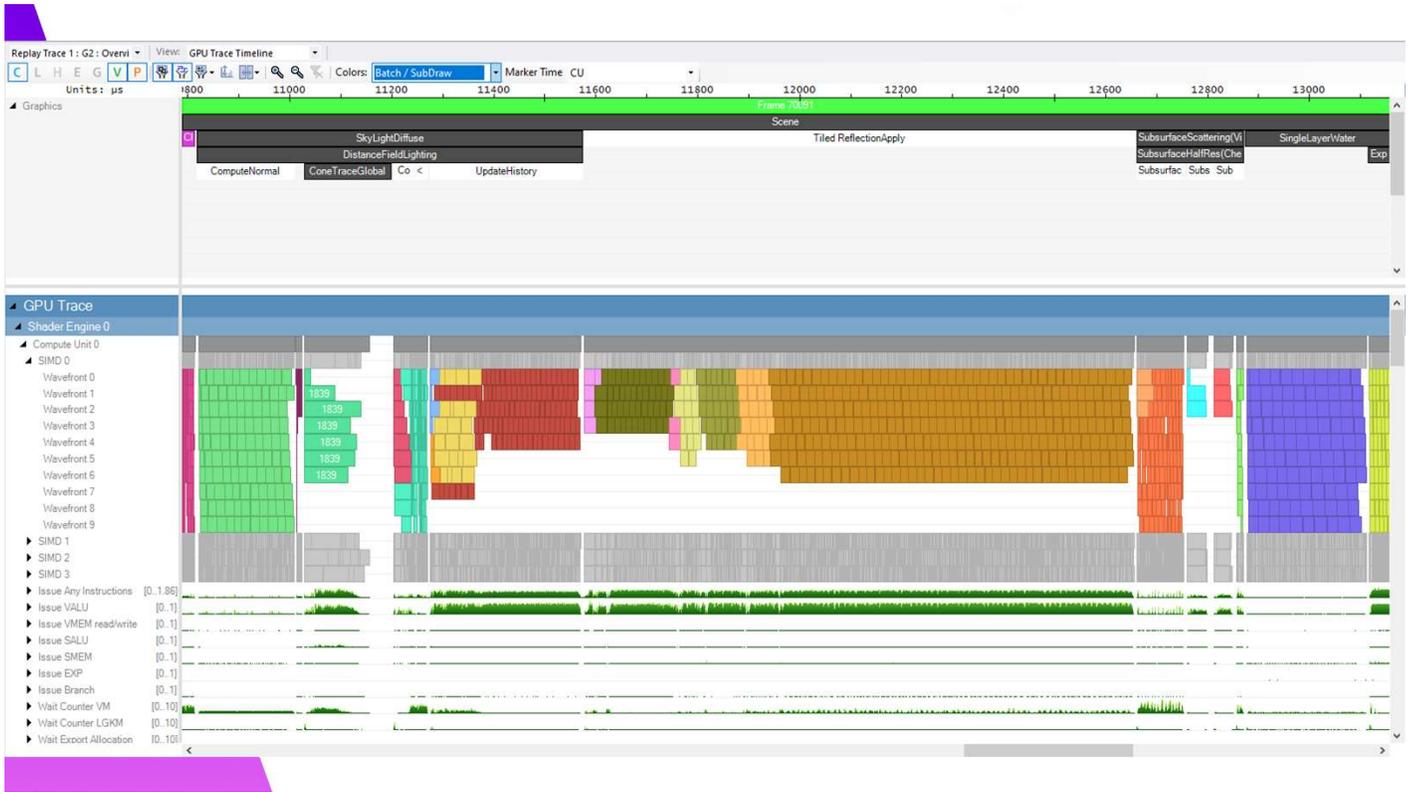
Override Shading Model ID

Before we look at the performance of the apply steps – let’s look at just the compute shader apply for Environment Lighting apply. This was just a full screen pixel shading pass in the original implementation, and this compute shader path runs using repeated calls to DispatchIndirect with different shader permutations.

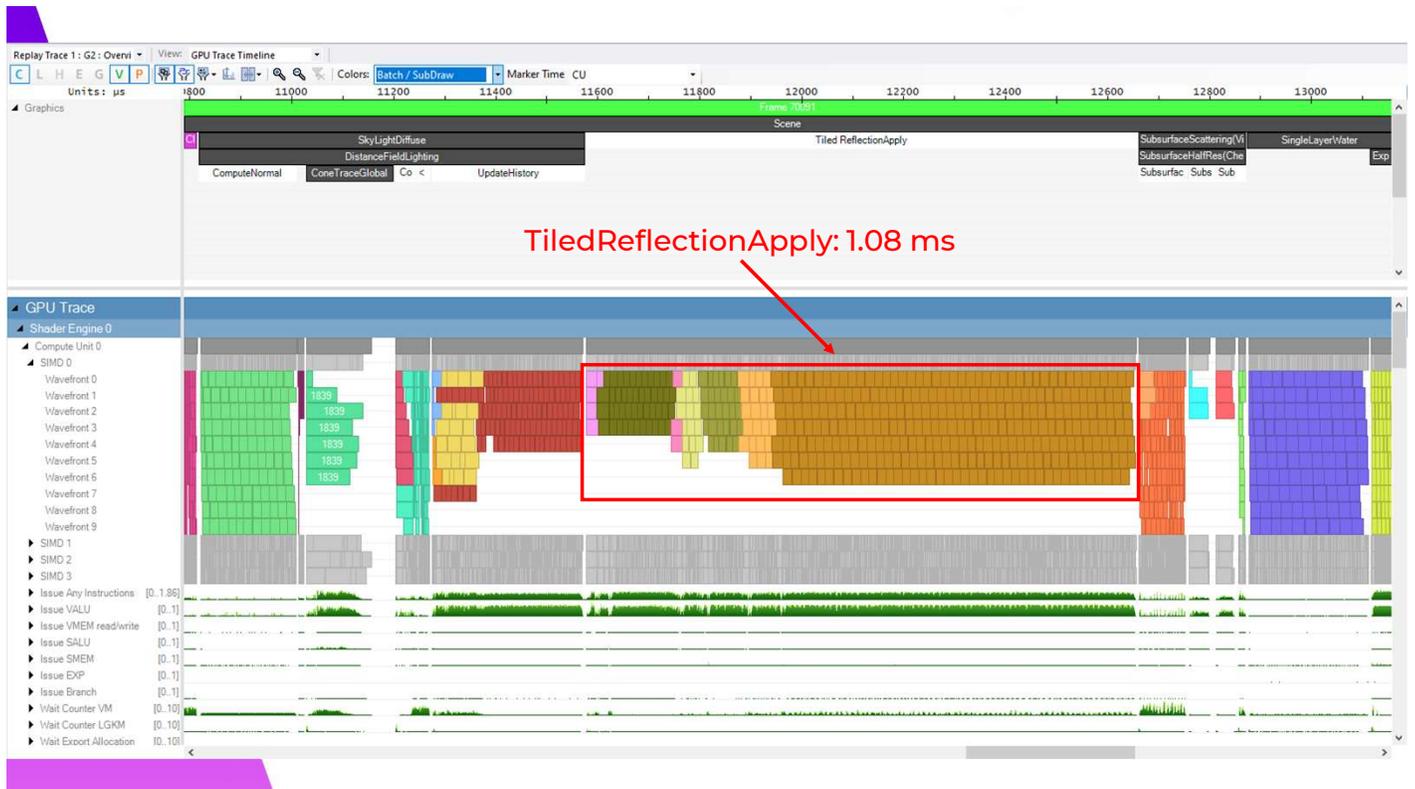
The shader begins by using the GroupId to look up into the tile locations buffer and then unpacking to an individual pixel location based on GroupThreadId. (CLICK ) You can see here that after the Gbuffer is read, an overwrite of the ShadingModelID will allow the optimizer to perform dead code elimination based on the preprocessor macros defined based on the shader permutation.



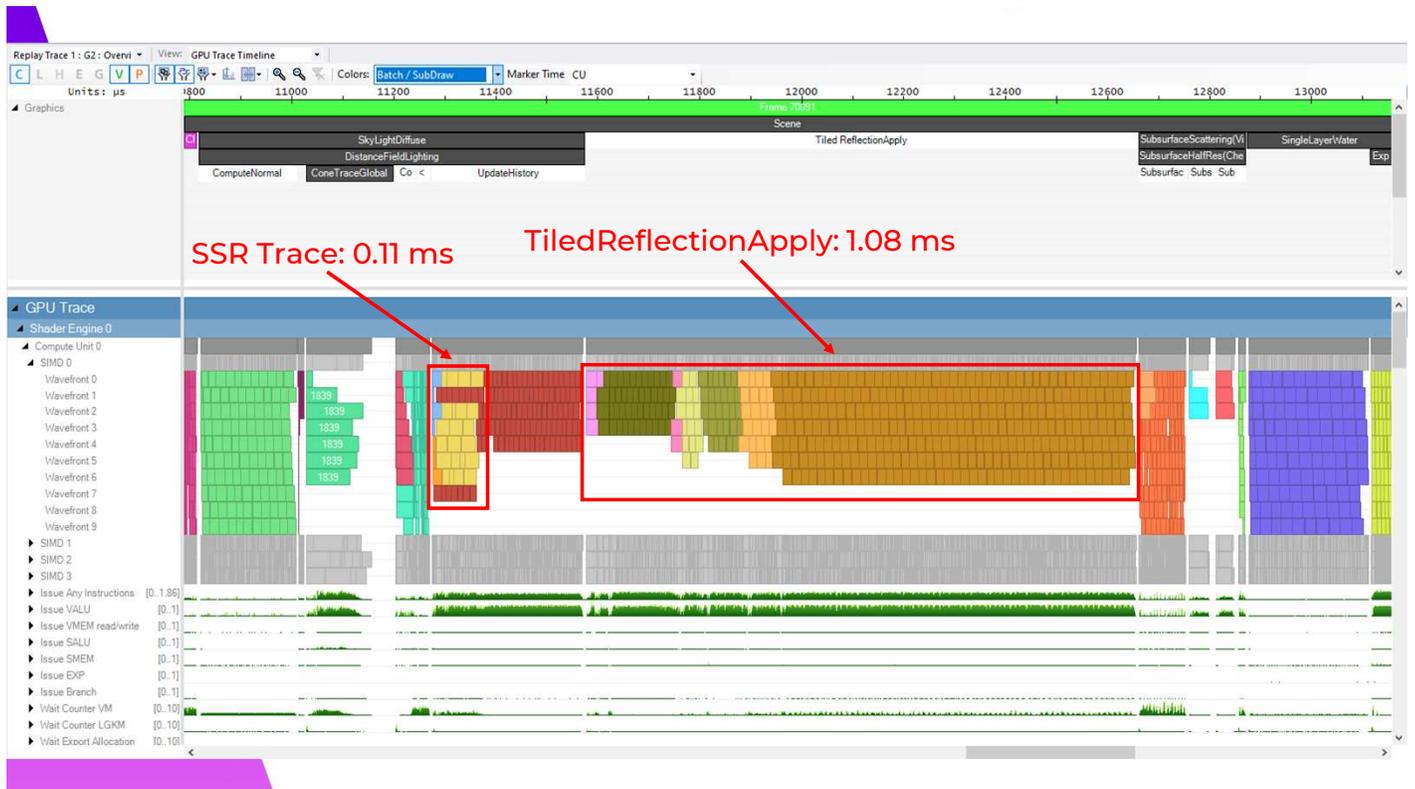
Now we need to look at what this 0.08 ms is buying is in the application of SSR, Reflection Environment, and SSS. Here is the original sequence of shading that I showed before.



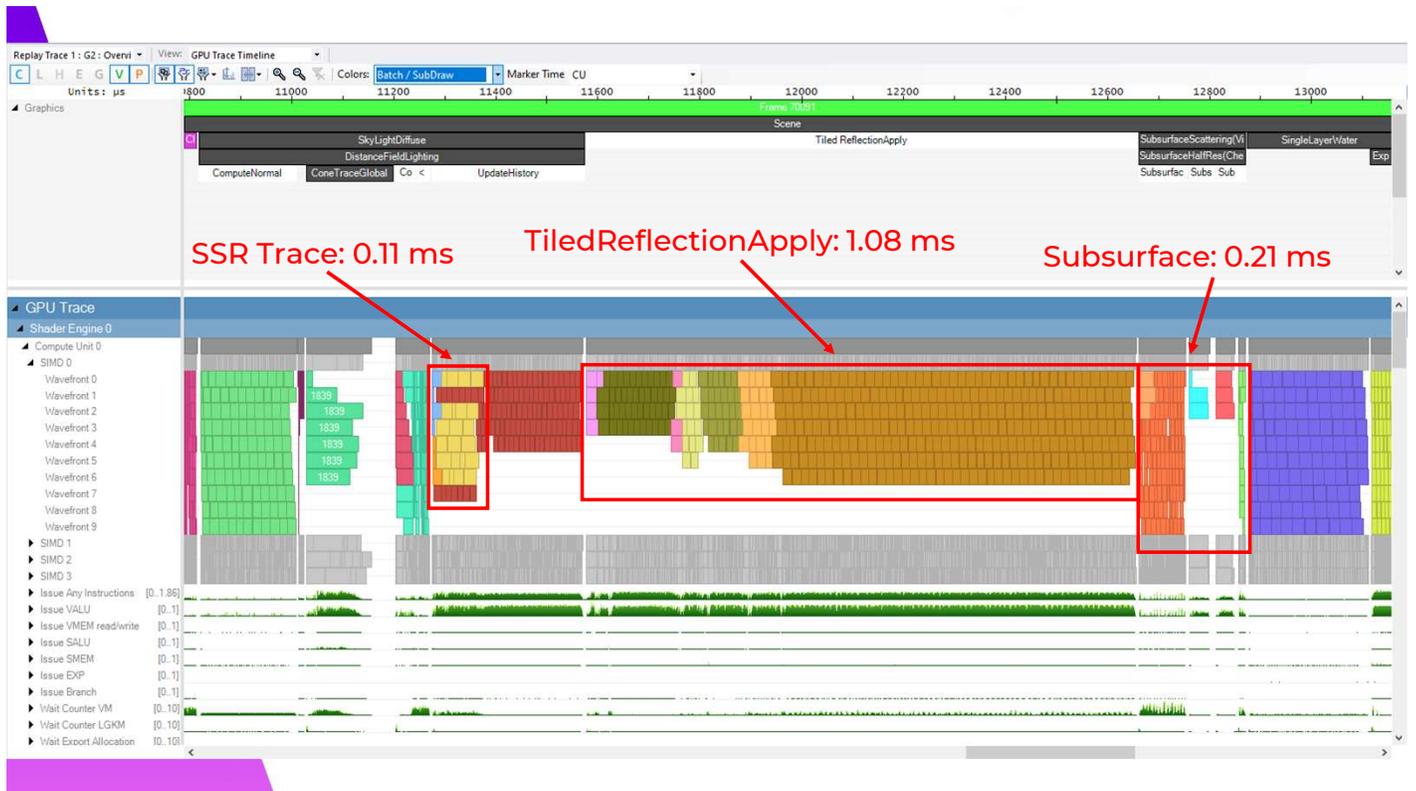
And here is our new frame.



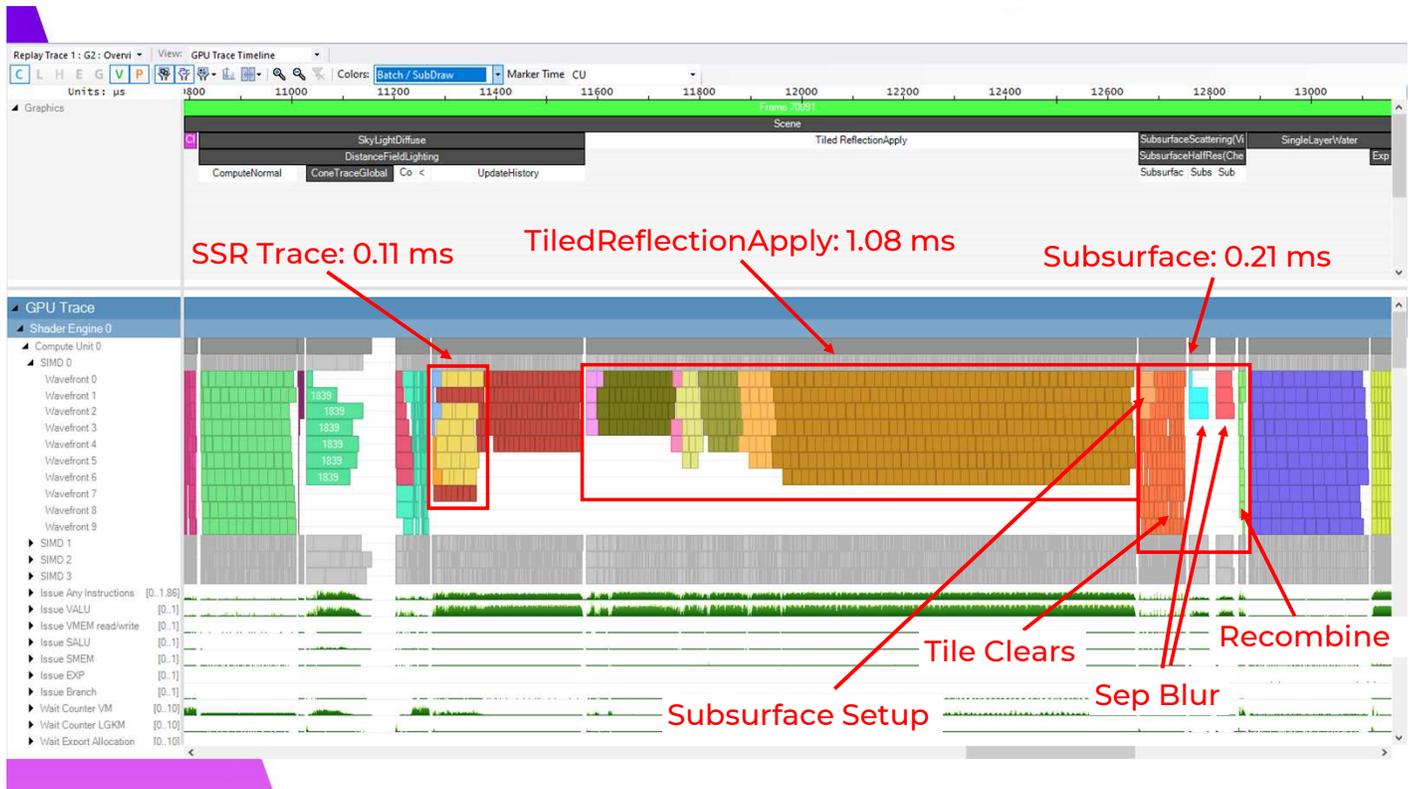
Here we have the Tiled Reflection apply shader at 1.08 ms, an 0.13 ms improvement. About half of this benefit is from culling sky pixels in this frame, so not really worthwhile in frames without sky pixels. One micro-optimization I would point out here is that I have slower waves with lower occupancy ordered first after the initial barrier and then have the fastest waves at the end. This helps reduce cracks between the batches – believe this is from lower occupancy waves waiting on more registers to become available. The lower occupancy waves also tend to be longer running, so putting the fastest batch last helps to ensure work drains quickly before the next barrier.



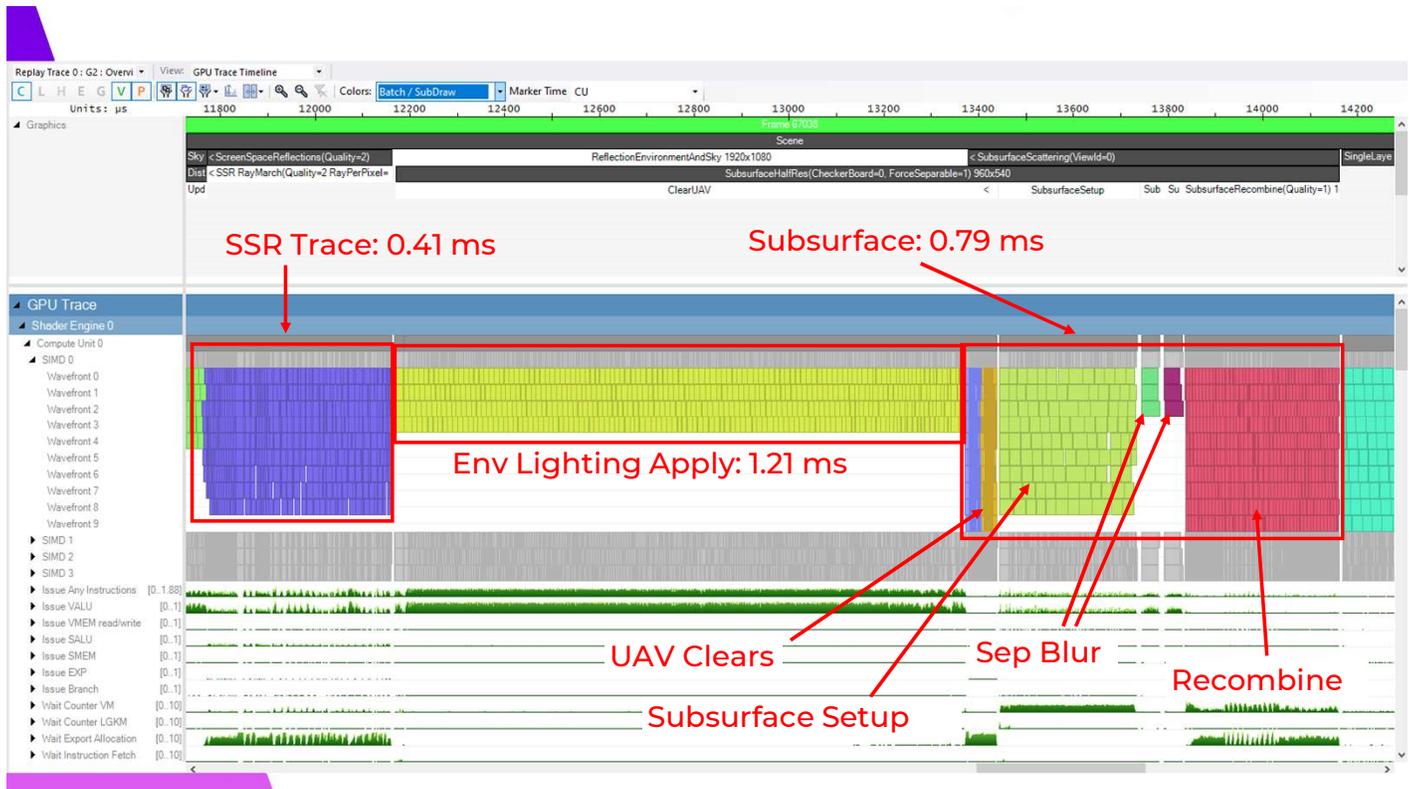
While TiledReflection apply is getting a small benefit, Screen Space Reflections is conversely seeing a really good benefit. 0.3 ms better and that is actually a conservative improvement because the waves are getting better overlap with DFAO history update. There is no need for a barrier as these are writing to separate buffers that both feed into the reflection apply. These results with SSR were what first gave me confidence that this was going to be a worthwhile optimization



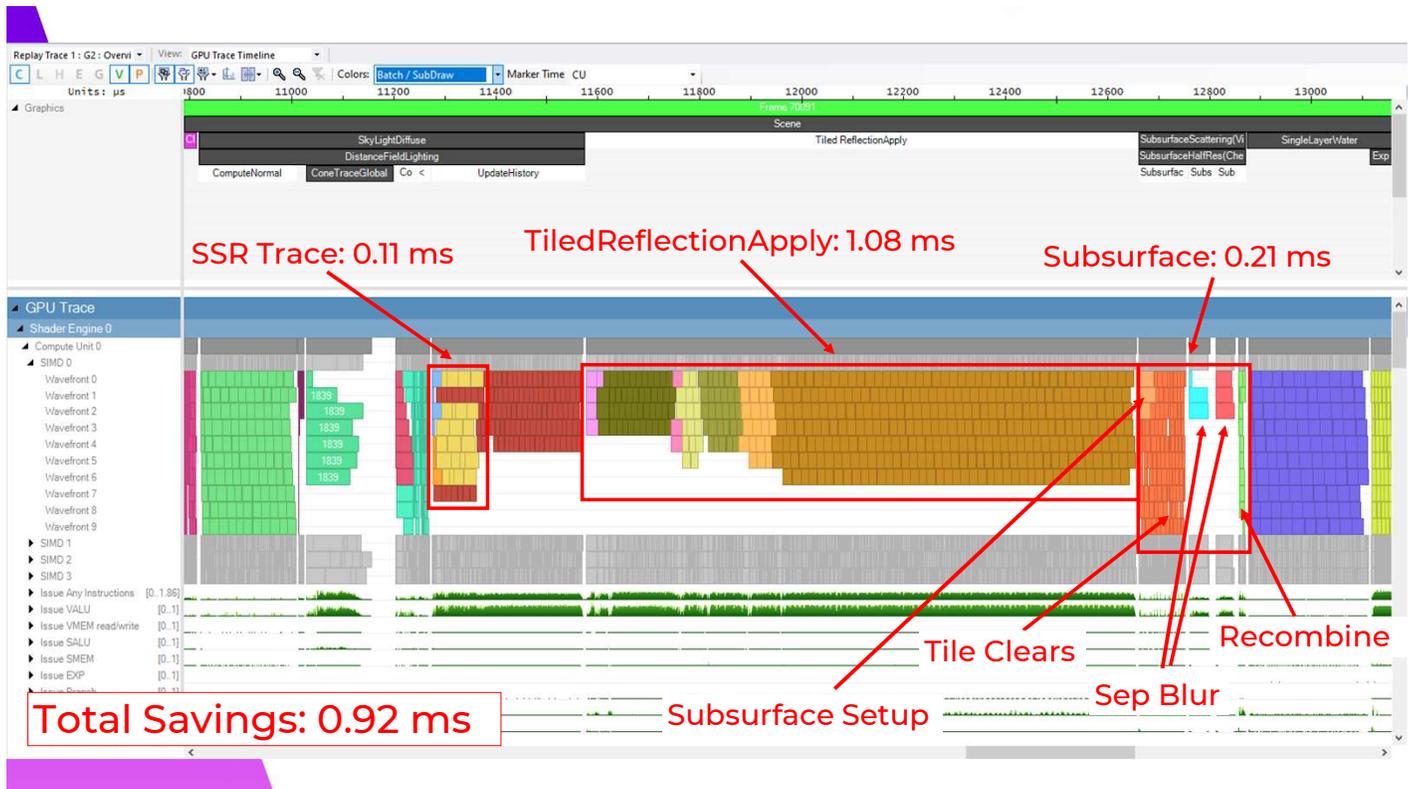
And then Subsurface after the reflection apply, which also sees a huge improvement, 0.58 ms better.



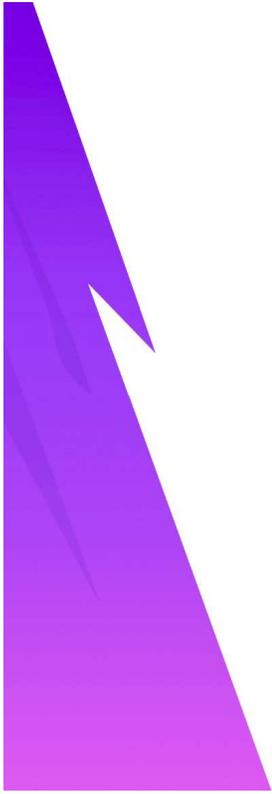
Here you can see the set-up and tile clears get really nice overlap with the tile clears having much shorter waves than the full setup shader. The blur steps are similar from the stock tile classify and the recombine is very fast as only tiles with skin in them are run.



Now that I've laid all this out – just a reminder again with what these passes were before.



And back to our results. This is a total savings of  $\sim 1$  ms from these passes, although obviously those benefits will vary based on the scene composition, 0.92 ms total for this shot once you subtract the classification cost

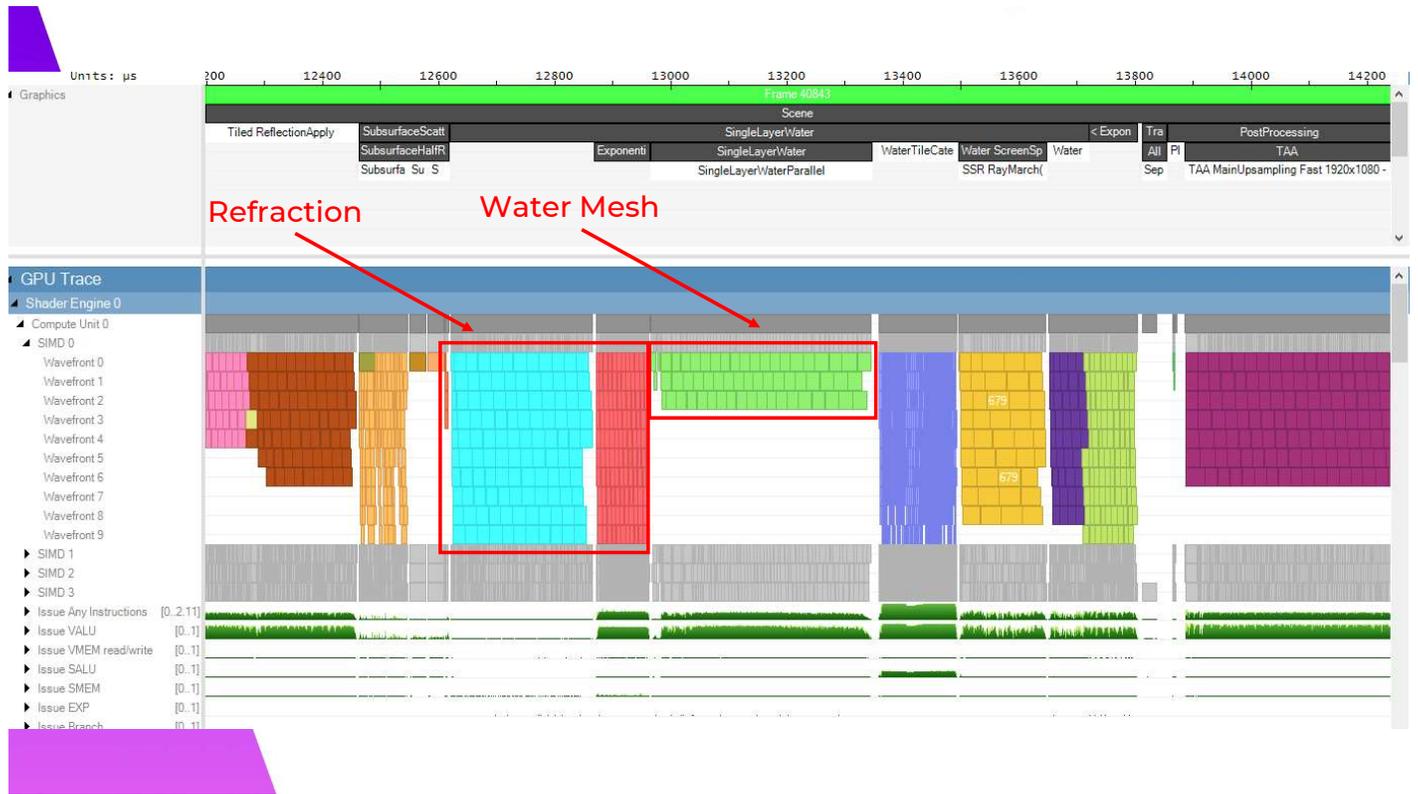


# Water Stenciling

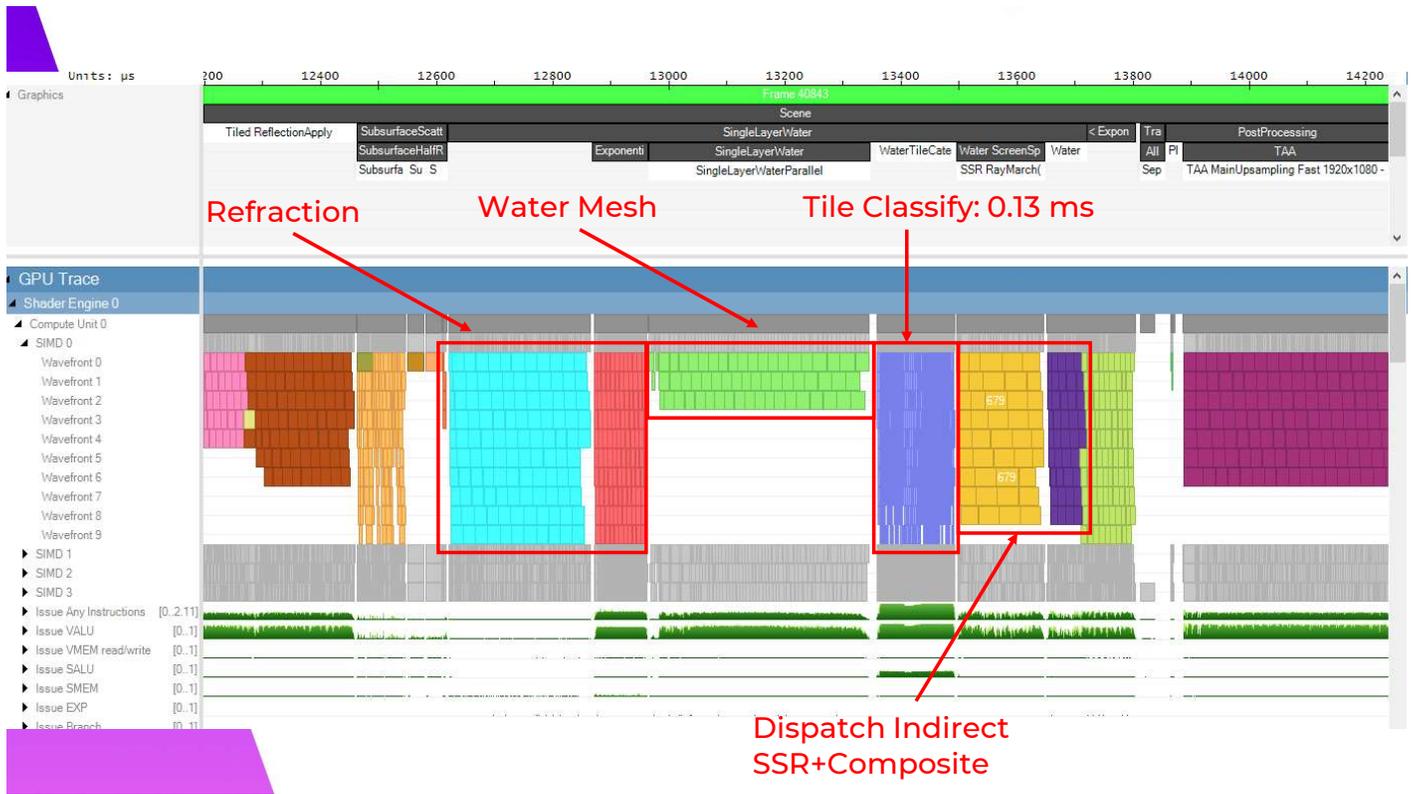
I just spent an amount of time talking up the virtues of tile classification, but sometimes it's just not the right tool for the job. I want to take a moment to show a small example where just using the stencil buffer is a better idea.



Consider this shot on PS4 at the edge of our Island. UE4's water system is really nice out of the box, and is taking up a reasonable number of pixels in this frame. The water is taking 1.1 ms to render on PS4.

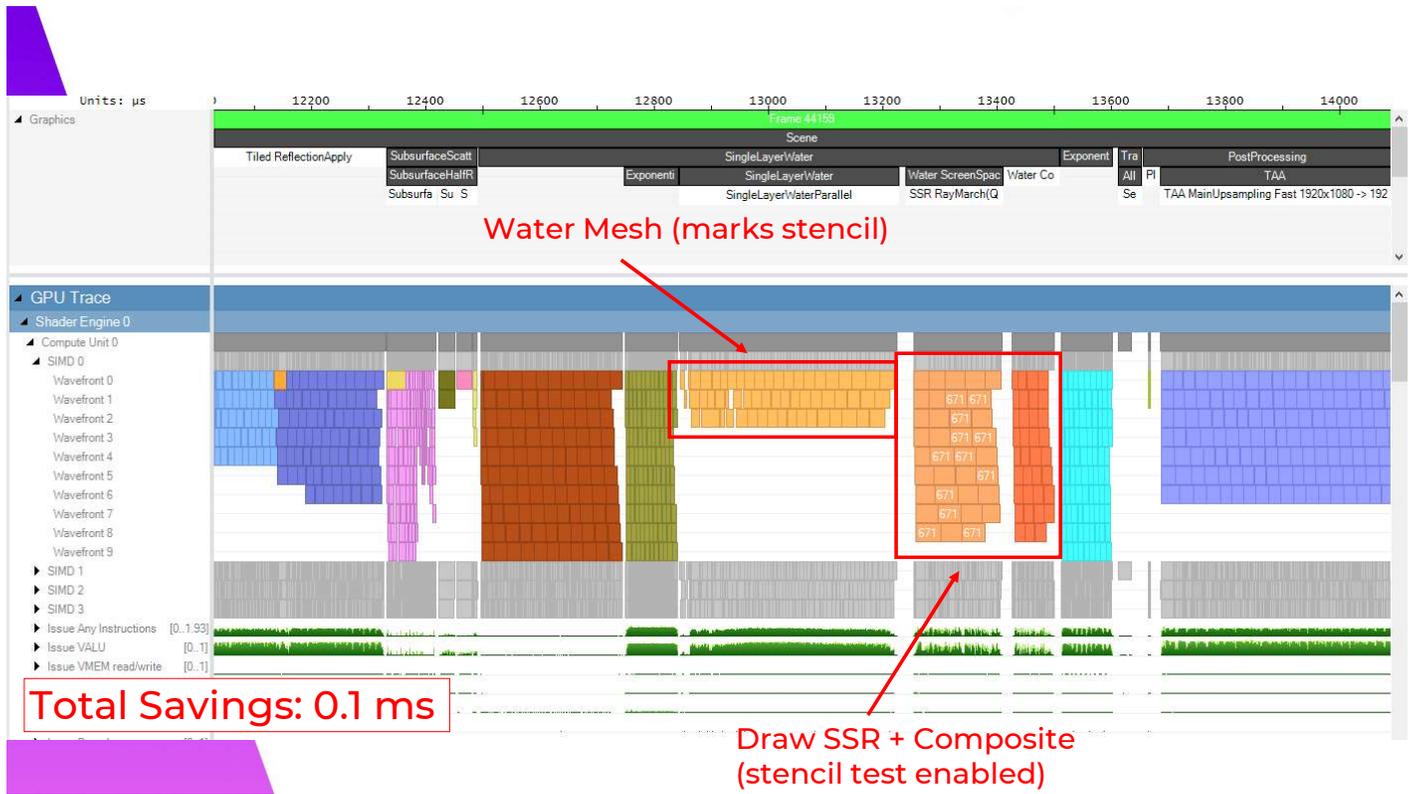


The water system first copies the scene color out to a secondary buffer and applies fog – this is for refraction. Then the water mesh renders to the full gbuffer like a regular opaque draw. Having an updated depth buffer has benefits for TAA in particular.



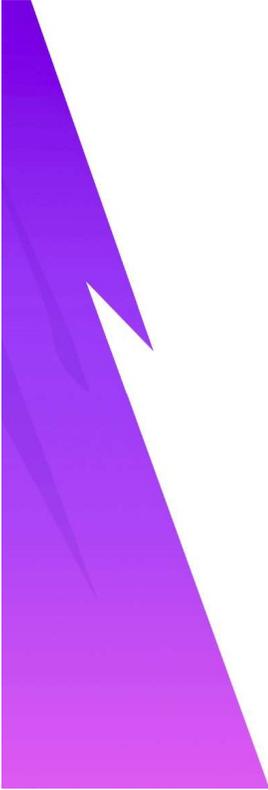
Then a classify is done on the gbuffer, very similar to what I did with reflections and subsurface, and just like before, Dispatch Indirect is used to run SSR and Composite only the tiles that were classified as having water in them





Here's what happens in razor - the tile classify is now missing between rendering the water mesh and handling SSR+Compositing. This saves 0.1 ms for very little effort, we don't get the full 0.13 ms of the tile classify back because we do end up with a barrier before the next pass now.

Credit to Oleksii from Dragon's Lake for doing this implementation



# Player Occlusion + Outline

My enthusiasm to use all parts of the stencil buffer and efficient EarlyZ utilization extends to our development of new effects for the game as well. I collaborated with my colleague Karinne Lorig to develop an outline and occlusion tint effect for our co-op modes.



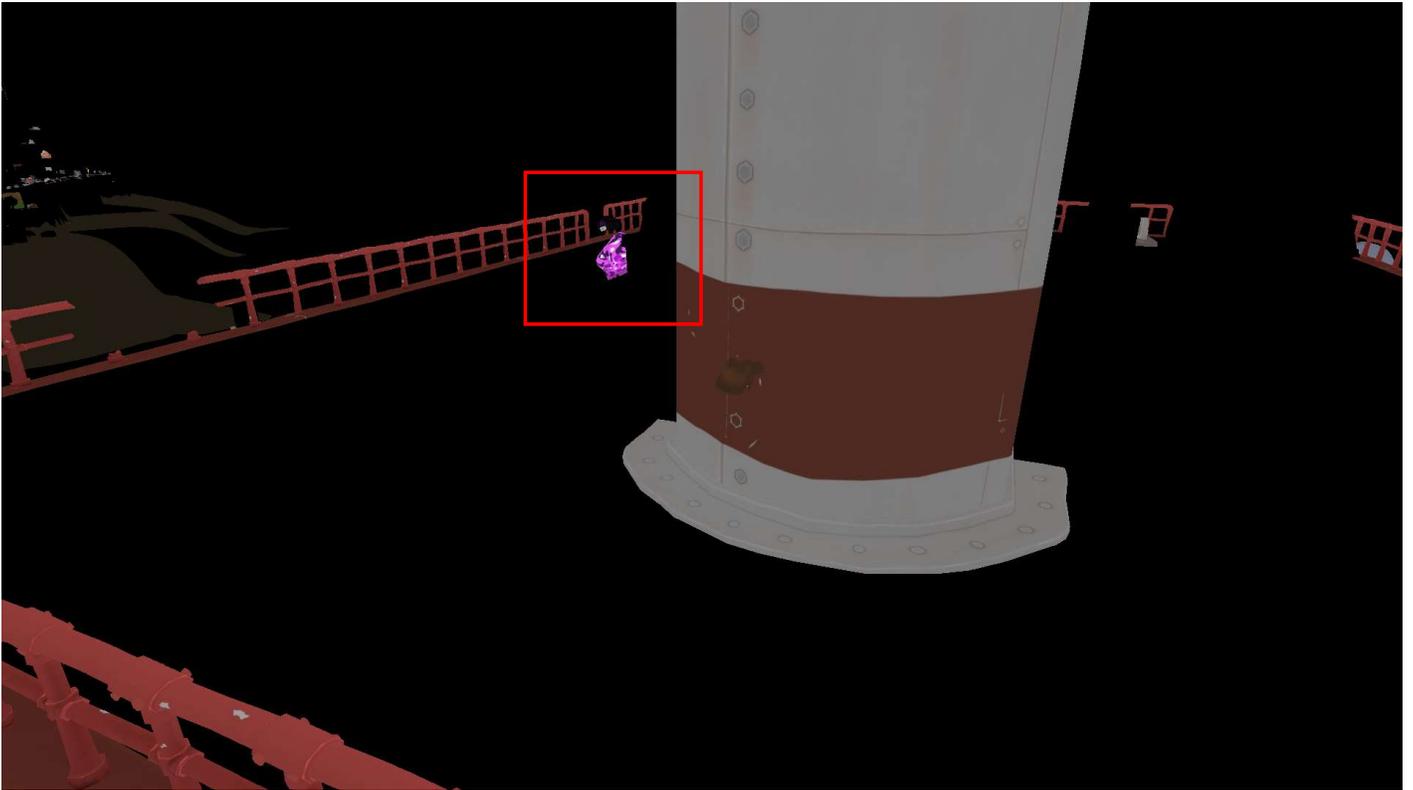
Here is a shot where you can see both parts of the effect on my teammate, which makes them easily identifiable in a brawl with another team, and a more subtle color of the occlusion effect is used on the active player to make sure their silhouette is always visible even if there is an object between the camera and the player.



# Player Occlusion + Outline

- Originally implemented with Custom Depth + Post Process Material – very slow
- Instead write stencil pass+fail for the characters during base pass and velocity pass
- Only run coloration+outline on pixels that are valid
- Stencil must be preserved from base pass through postprocessing, requires a few small engine modifications

This effect was originally prototyped by artists using UE4's Custom Depth feature and a post process material draw. This has a number of drawbacks: the extra depth buffer consumes memory and requires a full clear, and the characters must be rendered an additional time to it. Furthermore, the apply step runs as a full screen pass. We instead moved to setting stencil bits in the base pass and velocity pass, and then ran a custom shader on only the pixels needed. We did have to make few small engine mods needed to preserve the stencil through the frame, and we have reclaimed some stencil bits from other features that we were not using in the engine.



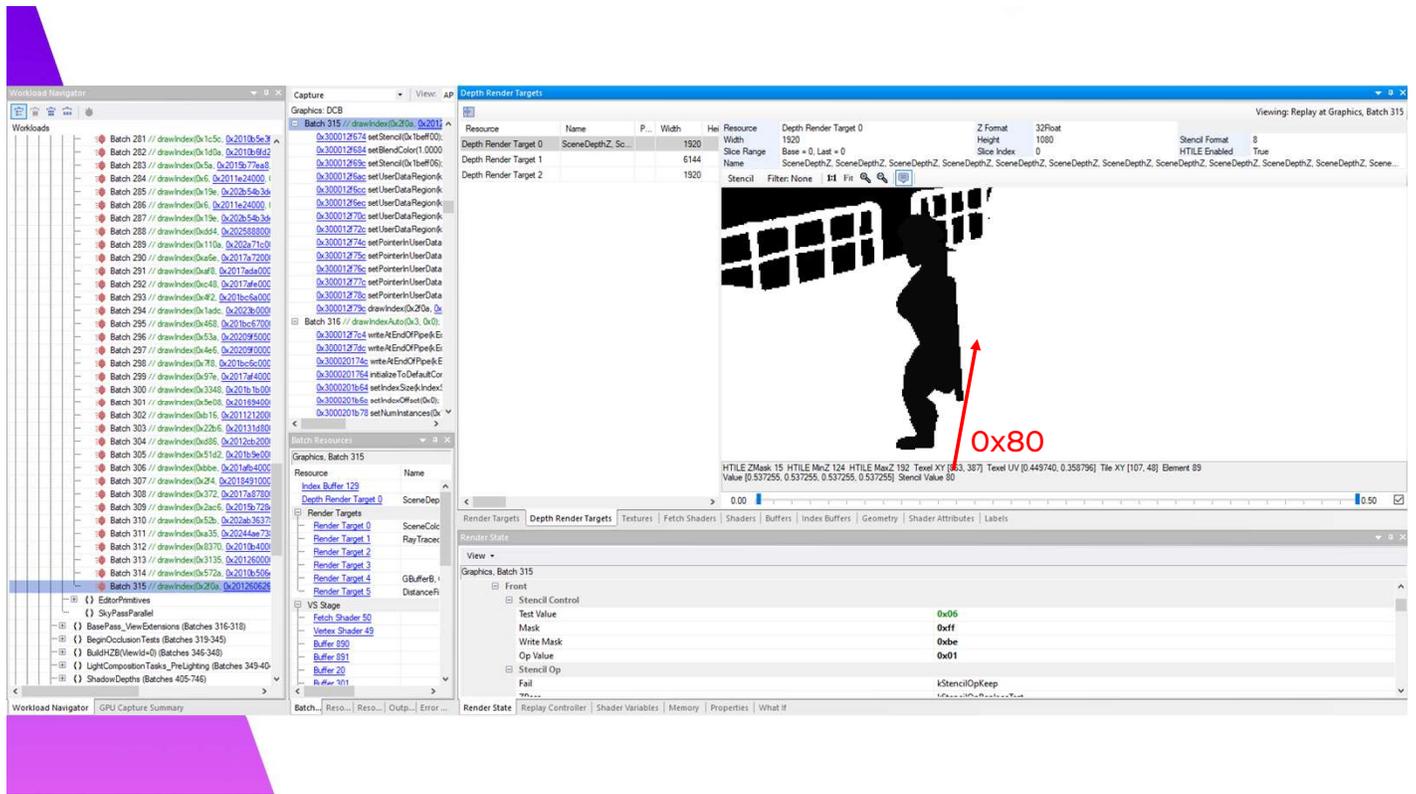
Let's walk through the stenciling passes - paying attention to the rendering of the pink shirt – I'm showing the albedo from the base pass here to make it easy to spot.

The screenshot displays the NVIDIA Nsight Graphics interface during a GPU capture. The main window shows the 'Depth Render Targets' configuration for 'Batch 238'. The viewport shows a white grid on a black background, with a red arrow pointing to the value '0x06' written on the grid. The 'Render State' window shows the 'Stencil Control' settings for the 'Front' pass, with 'Test Value' set to '0x06' and 'Write Mask' set to '0x06'.

Resource	Name	P...	Width	Height	Resource	Depth Render Target 0	Z Format	32Float
Depth Render Target 0	SceneDepth_Z_Sc...		1920		Width	1920	Height	1080
Depth Render Target 1			6144		Stencil Format	8	HTILE Enabled	True
Depth Render Target 2			1920		Base = 0, Last = 0			

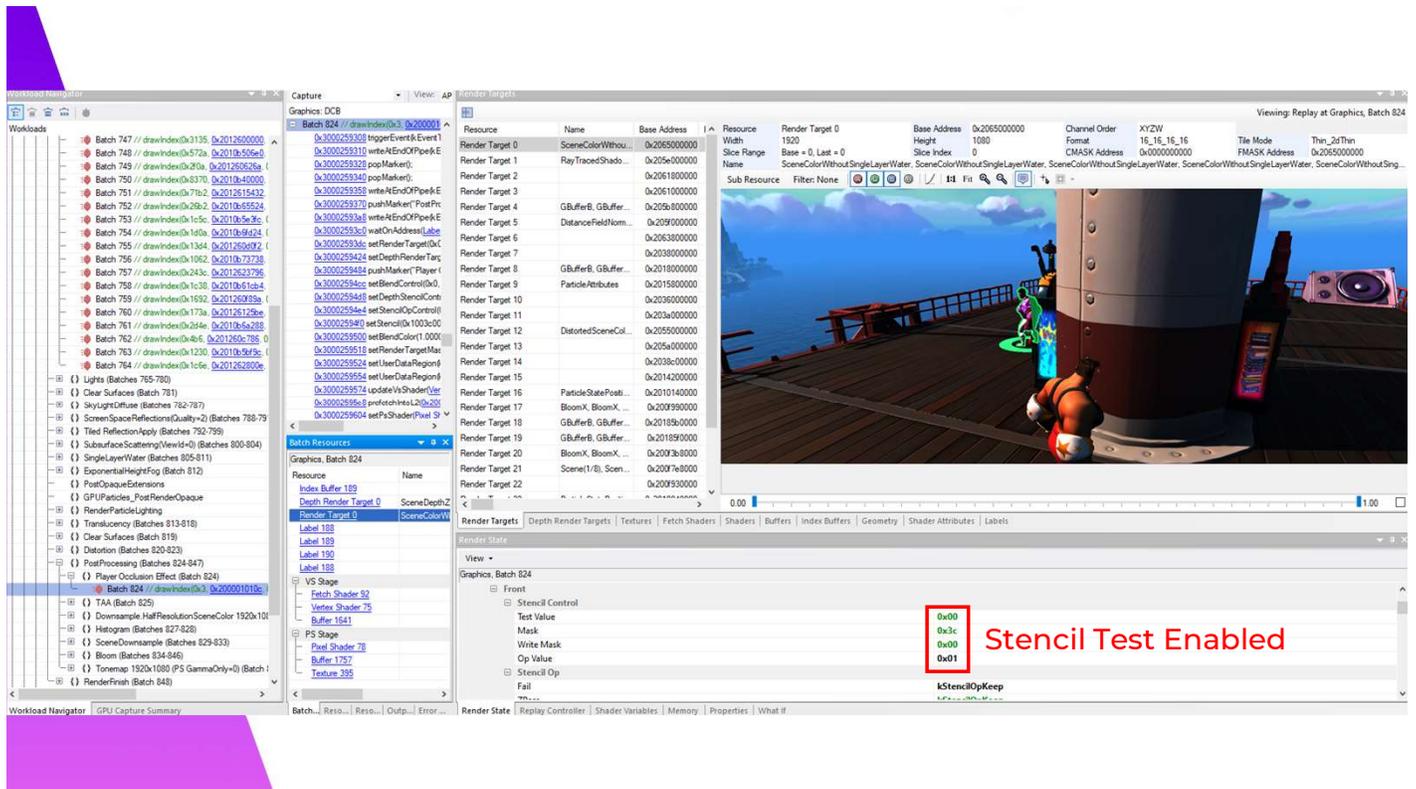
Resource	Name	Test Value	Write Mask	Op Value	Stencil Op
Stencil Control		0x06	0x06	0x01	kStencilOpKeep

And in Razor you can see in this particular shot that the stencil bits are set and the shirt writes out 0x06 for fragments that pass the depth test

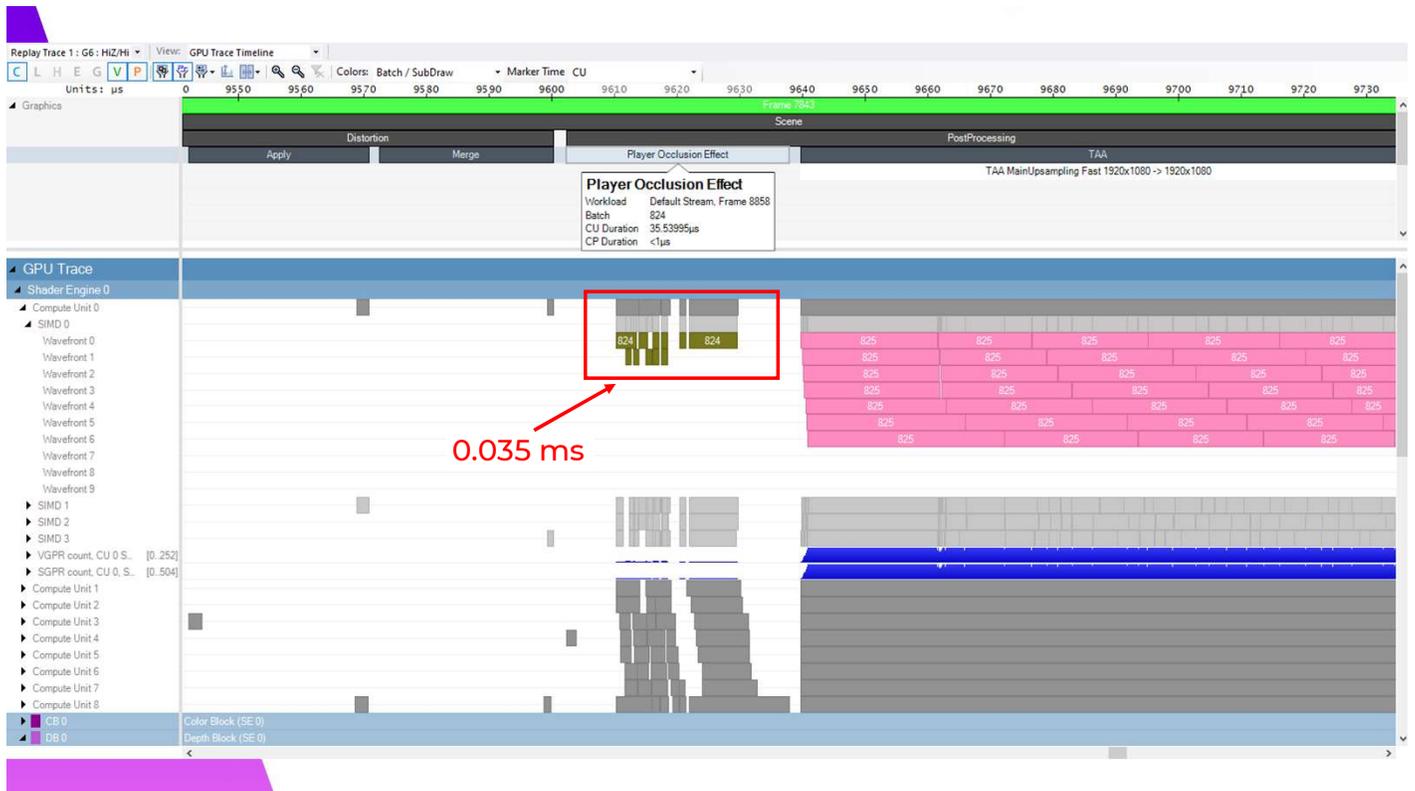


At the end of the base pass, you can see that the stencil is filled out with the remaining draws. The non-character stencil is set to 0x80, you can see a nice outline of the environment from the Receives Decals bit being set.

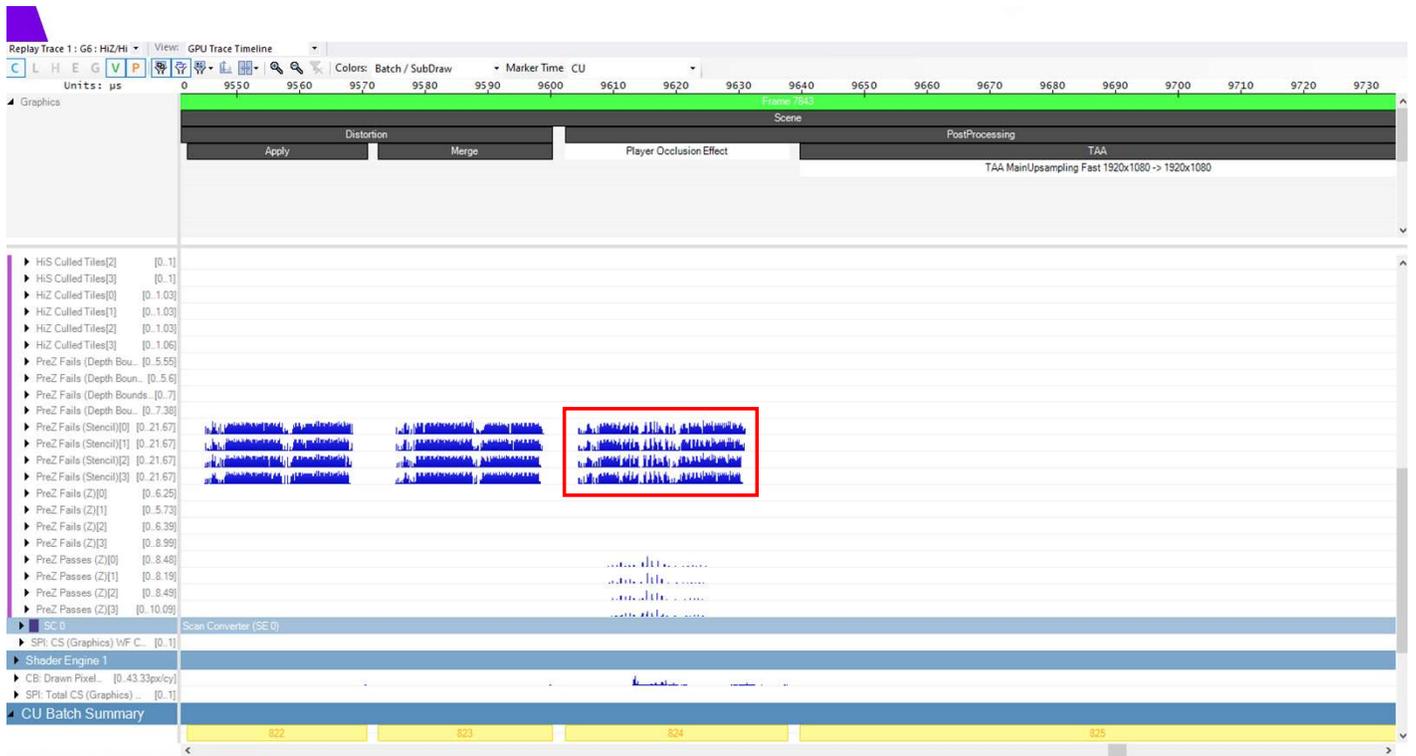




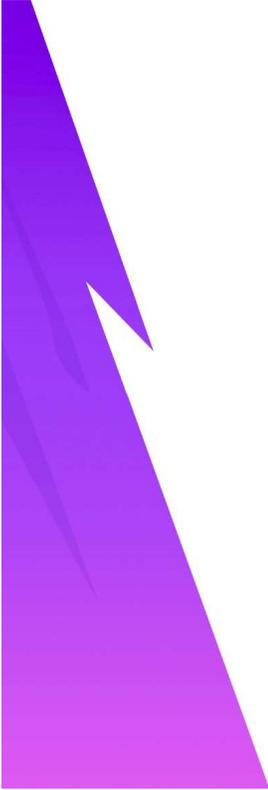
Finally, the post process effect is run before TAA on the fragments that the player rendered to. The shader samples neighboring pixels to render outlines at the edge of the character, and the center sample determines the player color and if the occlusion tint should be applied. The shader has to be run before TAA or else the stencil positions will not be valid for the un-jittered camera matrix. However, this is actually ideal, because TAA provides nice anti-aliasing on the effect. Because we only render on the inside of the character silhouette, we do not fight against the TAA algorithm as the outline effectively moves with the depth values of the character.



To illustrate how efficient this is, here is the razor trace for that frame. You can see the occlusion effect only take 0.035 ms, with very few pixel shader waves actually running. I have worked on a number of projects that only use 1-2 bits of the stencil buffer, so I am very happy that we frequently use all 8 bits of the buffer in a single frame to cull workloads with EarlyZ.

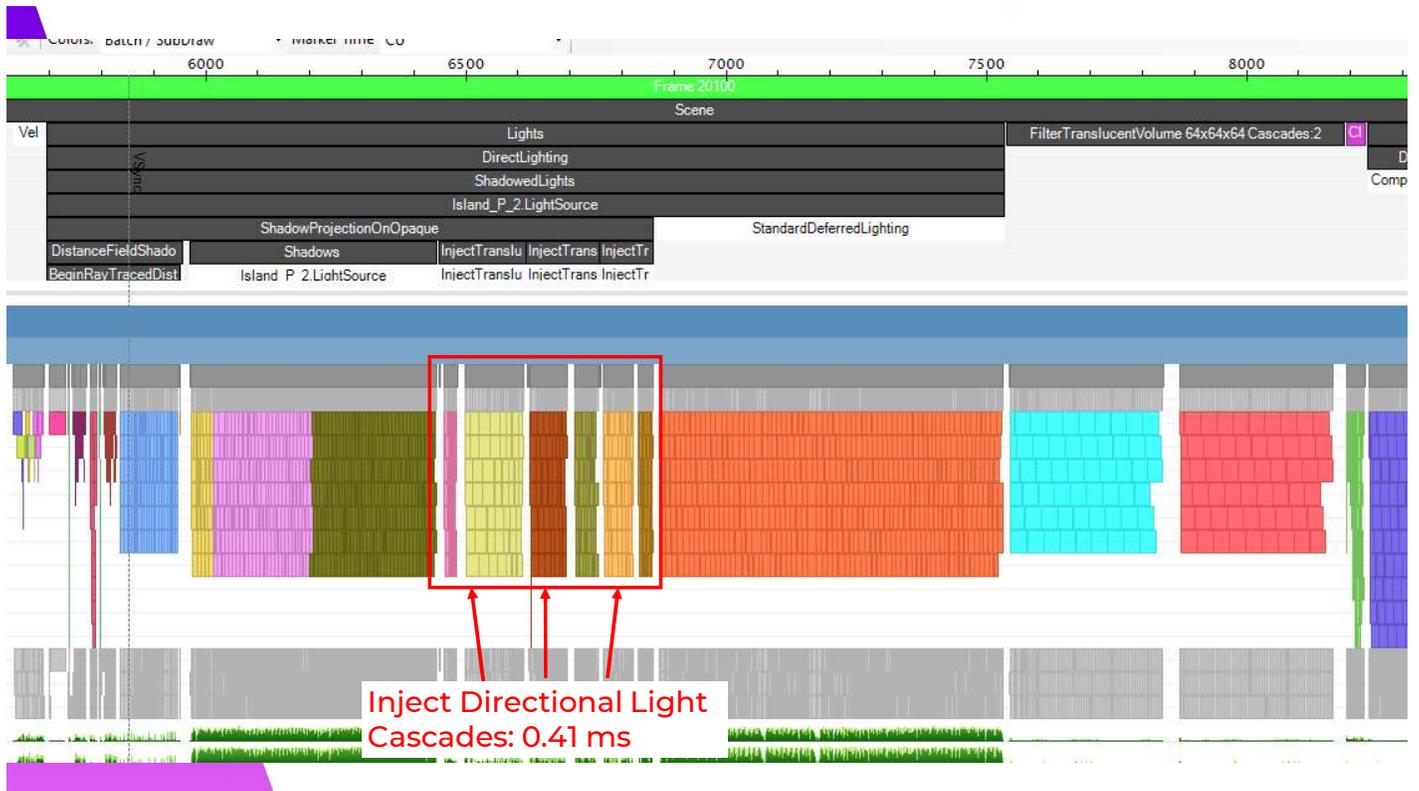


And here I have scrolled down to show that a large amount of this effect is just spent with EarlyZ rejecting the stencil bits. We have not configured HiStencil but theoretically this could be a use case that would benefit from it, but the time is so minimal that is not a priority and enabling HiStencil does have trade-offs with HiZ quality on GCN.

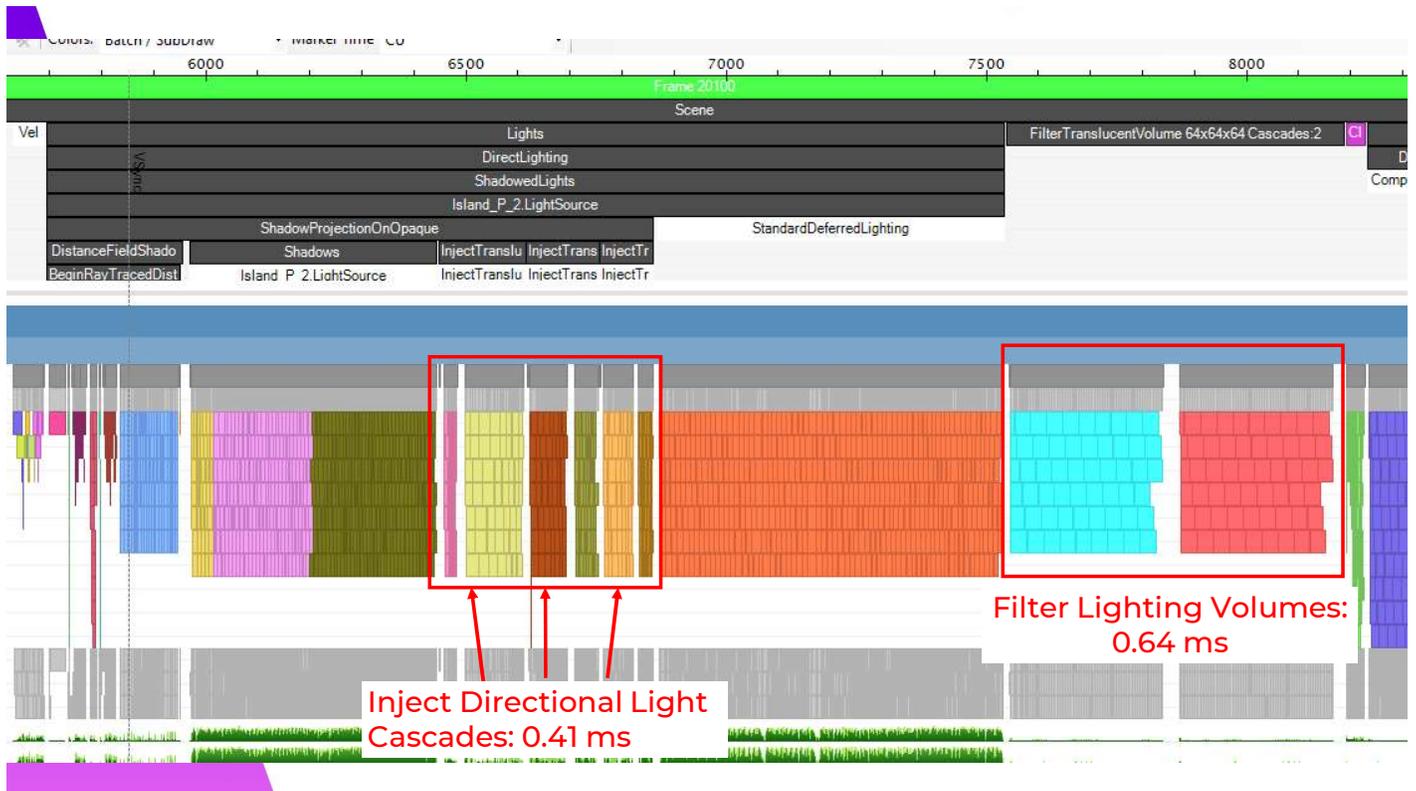


# Translucent Lighting

In the spirit of removing unnecessary work from our frame, I wanted to change how lit translucency was handled for our game compared to what UE4 does by default.



Stock UE4 accumulates lighting in a volume texture and then translucent particles and meshes can read from it to cheaply apply shading. Each light must be injected – you can see that happening once for each shadow cascade on our sunlight. There are two cascades on the translucent lighting volume, so there are a total of 6 injections happening here. More injections have to occur if there are additional dynamic lights rendering in the bounds of the lighting volume.

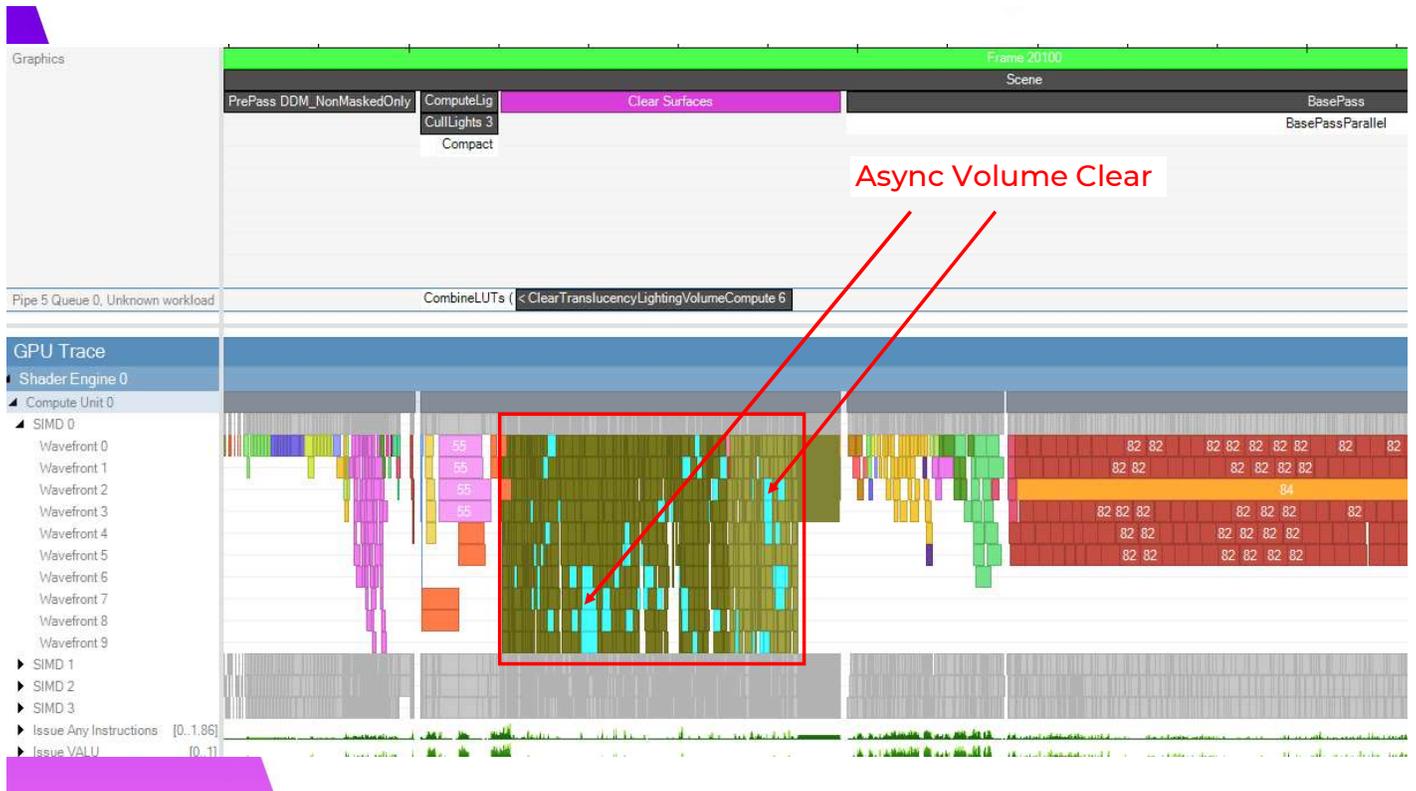


After lights are injected as part of the main lighting loop, the volume gets filtered to soften the shadows and make aliasing at the low resolution of the volume less noticeable. This takes 0.6 ms for two  $64^3$  lighting volume cascades.

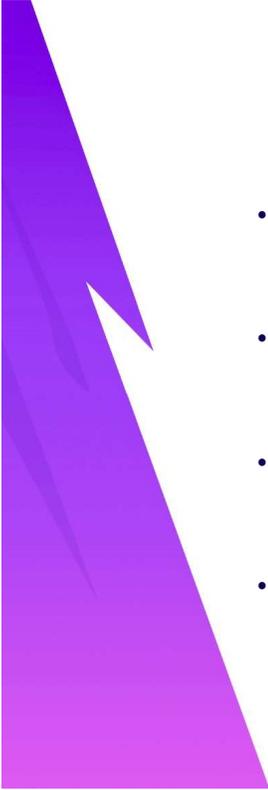
The design consequence here is that the system is designed to have a high cost on the number of lights present in a scene, but a very low cost per lit translucent pixel processed.

In total, the translucent volume is taking  $\sim 1$  ms per frame to support particle lighting. The only option to disable this in the engine is to resort to only using per-pixel forward shading on any lit translucency, which is quite expensive on something like smoke particles.

This is a lot to be paying for something that frequently doesn't contribute to the look of the final image if there aren't actually any lit particles or translucency firing, and our game is a little goofy, we don't have \*that\* much lit translucency. My goal was to eliminate all of these fixed costs from our frame.



The volume must also be cleared each frame, this happens in async compute, although you can see it's not behaving super well in this trace and it is overlapped with other DRAM heavy work in the surface clears. It can be effectively free with the right scheduling though.



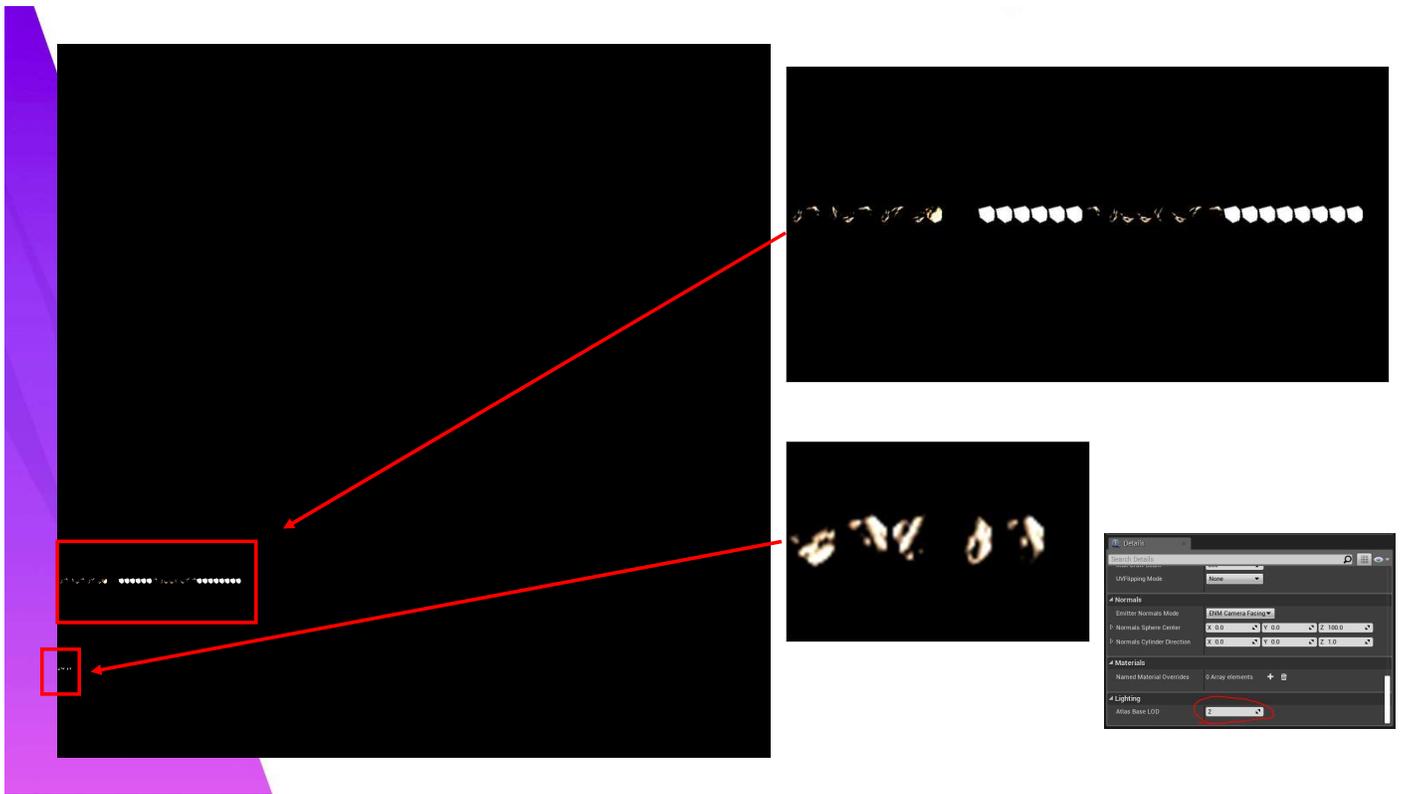
# Particle Lightmaps

- Originates from *The Devil is in the Details: idTech 666* by Tiago Sousa and Jean Geffroy
- Accumulate particle lighting into an atlas rendertarget
- Particles read from lightmap atlas during actual shading
- Can support LODs at various resolutions in the atlas

In pursuit of this goal, I worked with my colleague Rusty Swain to implement a system similar to what was used in DOOM 2016, a particle lightmap system that particles accumulate their lighting into at less than the final resolution of the particles on screen, but still higher quality than what you would get from just shading the particle vertices.

We select an LOD resolution for an effect based on distance from camera, and Doom could even reduce the shading frequency in the time dimension or have the atlas update in async compute. We only do LODs and update the atlas every frame for simplicity, but if art direction started pushing more lit particles sprites we would definitely look into additional optimization.

[http://advances.realtimerendering.com/s2016/Siggraph2016\\_idTech6.pdf](http://advances.realtimerendering.com/s2016/Siggraph2016_idTech6.pdf)

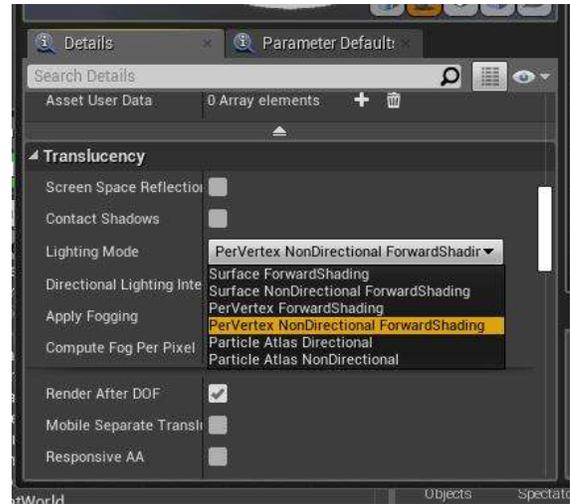


Here you can see a particle effect going into our atlas - this takes the sprite's position and normal map and calculates a lighting result at each texel. This is some smoke coming off of an impact crater of another rumbler that landed nearby me in the match. Updating this particular effect in the atlas took 20 microseconds, and gets better results than the translucent lighting volume. Having lit particles render to that atlas every frame does simplify managing allocations within it because the particles can be full repacked at each LOD every frame.

#### EXTRA NOTE:

Artists can tune a base LOD bias for if they want it to bet more texels per quad, so by default sprites will not take up more than 32 x 32 pixels for lighting when the emitter is close to the camera. LOD0 takes 128x128 pixels, but is almost never used by our artists. Emitters will drop an LOD level for ever 25% reduction in estimated screen size of the effect so there is dynamic LOD selection occurring for us within the atlas.

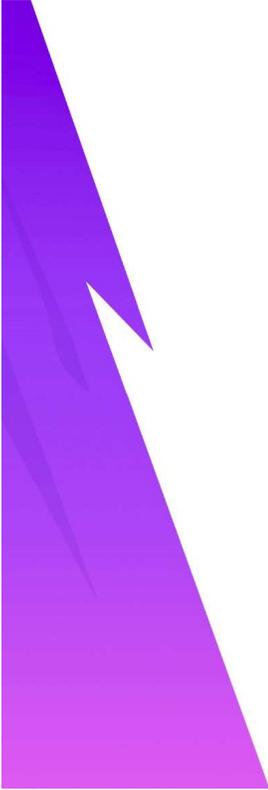
# Mesh Lighting



We only support sprite based particles in the lighting atlas. Meshes need to go down a forward shading path – it's okay for them to receive per-pixel forward shading when needed, but we added support for vertex shading which is frequently used on mesh particles in our game. We also added support for NonDirectional or “Wrap” shading in both the per-pixel and per-vertex paths.

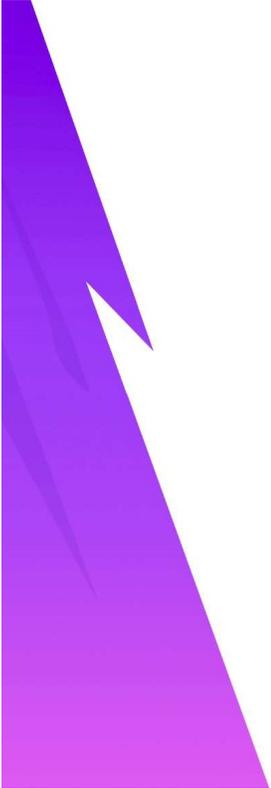
UE4 has a clustered lighting data structure for forward shading – the vertex shading path reads the light data at the cluster it intersects and accumulates it into a spherical harmonic. This is an old trick at this point, but this allows us to keep per-pixel normal map variation on a vertex lit object, which fits really nicely with our art style.

In most frames we see 0.5 ms improvement in frame time, and in many frames we see a full 1 ms better performance than using the translucent lighting volume path.



## Part 2: Reactive Optimizations

Now, while all that work was proactively done, we also did a lot to respond to performance problems we did not fully expect until content reached the limits of what we could do.



# Background: Lighting

I need to set the stage a little first with some decisions we made early on that would impact our performance later, specifically with how we built our environment



Early on we decided to lean in heavily on Unreal's distance field features. Our medium to far shadows are all traced against mesh distance fields that are calculated offline for static meshes, which allocate into a shared atlas texture accessed by compute shaders on the GPU. In the images I've included here, you can see the world without the distance field shadows in the top left, and a visualization of the per-object distance fields in the bottom left, together which make the final shot in the right with far shadows.

#### EXTRA NOTE:

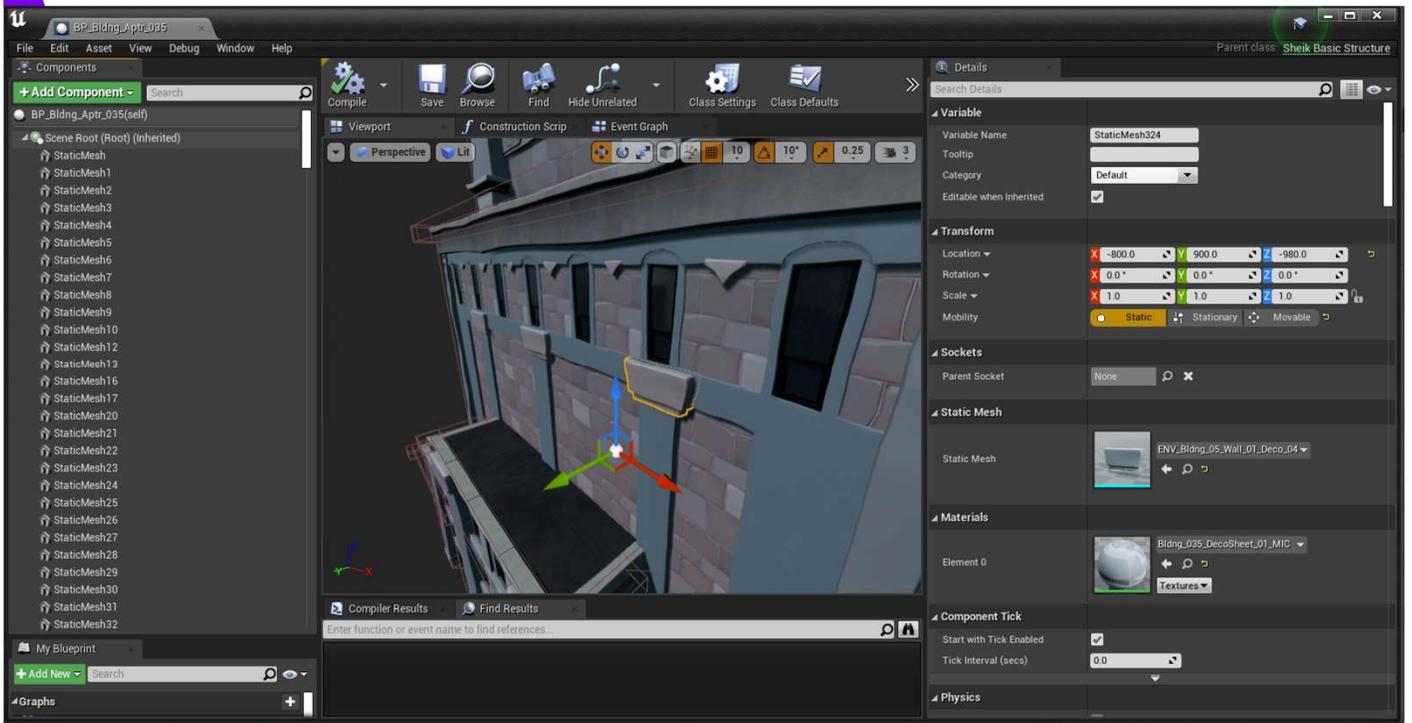
The terrain supports using a heightfield representation, but we are able to get away with disabling heightfield shadows on last-gen consoles because our artists and designers almost always choose to place decorator skirt meshes near changes in elevation like cliffs. The artists like it aesthetically and the designers like it because the game's traversal system works really well with wall climbing.



We also use distance field ambient occlusion, which traces against a lower quality clip map that composites the individual distance fields into a global distance field representation. This clipmap can be extremely low quality but DFAO still successfully adds a lot of definition to our shadowed areas. The only other ambient occlusion comes from a channel in the Gbuffer that artists can write out to from material graphs. This AO is pretty important to our look because we don't have any global illumination baked into lightmaps – just an ambient skylight that has a diffuse component (spherical harmonic at runtime) and a detailed specular environment map.



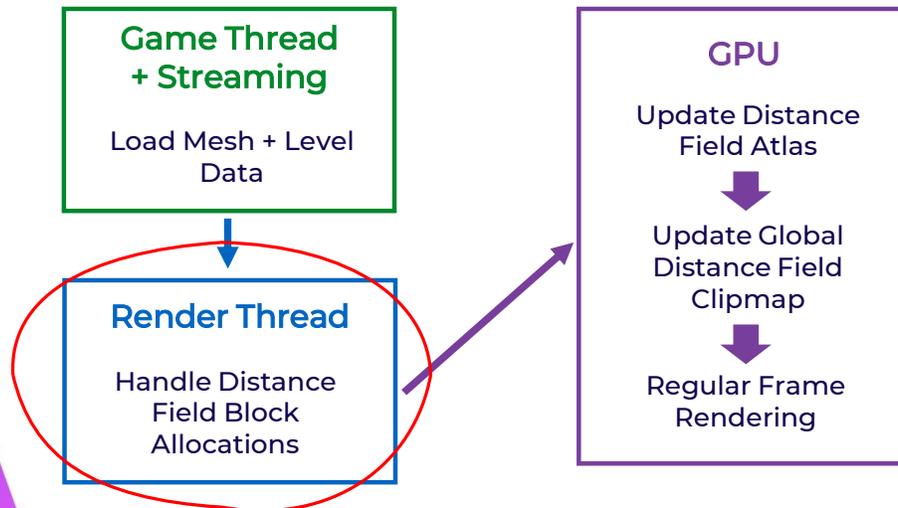
One minor detail I'll point out while we're on the subject. The ambient shadowing under the characters is just a simple decal that modifies only the ambient occlusion gbuffer. Super simple, but didn't work out of the box in Unreal and our designers really appreciate the characters always having some sort of shadow anchoring them to the ground.



A consequence of the distance field data is that it is all stored in volume textures – meaning that it really likes being associated with static meshes that have minimal empty space, for example we don't want the volumes to include the empty space on the interiors of buildings.

We really leaned into Unreal's dynamic instancing system and built our assets out of many smaller mesh components from very early on in development, this minimizes texels wasted on empty space. These smaller meshes are grouped together into entities we call "basic structures" - you can see that this piece of a structure is its own mesh which is then instanced multiple times on the building, and the building has many static mesh instances making it up in the component panel on the left.

# Distance Field Upload

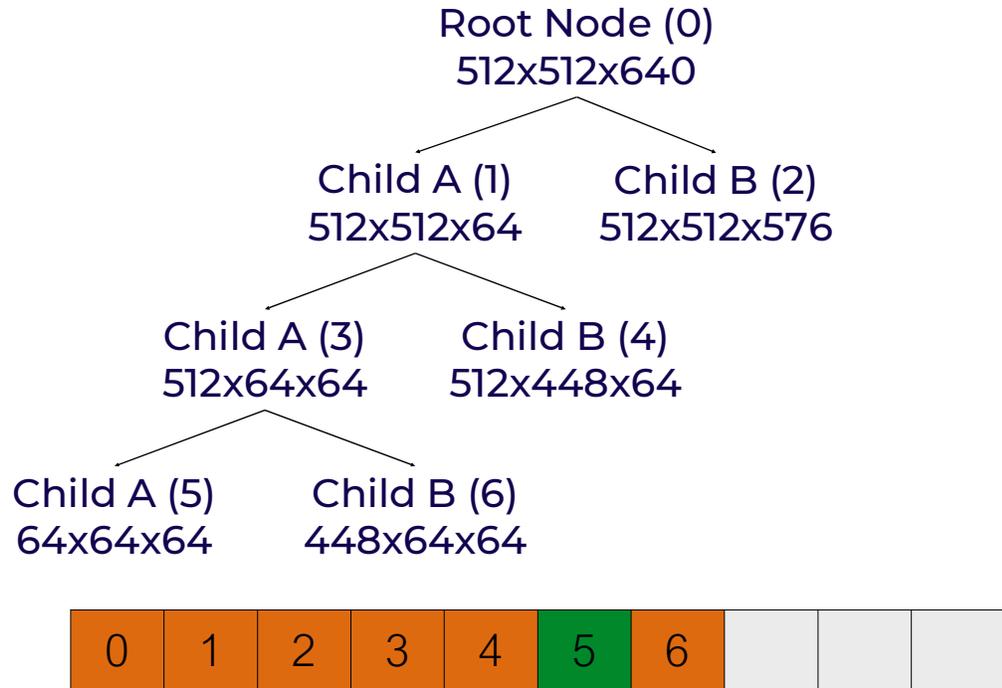
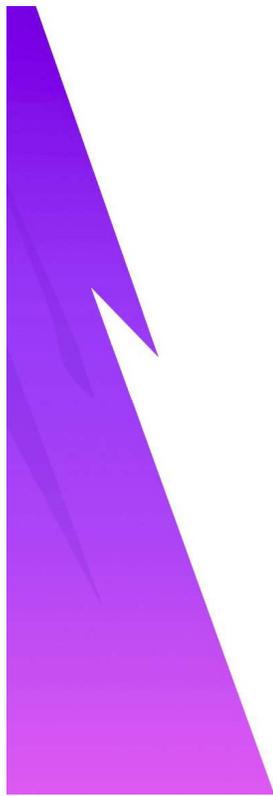


The distance field data has a number of systems it has to pass through while loading in before it is ready for rendering. This is after the main game thread runs it through the regular streaming process on the mesh and gameplay data.

The first is that the render thread manages the allocations within the distance field atlas with a block allocator that tracks where a particular mesh's data will be in the atlas and tries to minimize fragmentation within that resource.

Finally, there are multiple operations that must happen over on the GPU - it copies from the individual distance field memory into their assigned places in the atlas, and finally any regions of the global clipmap that need to be dirtied are reprocessed, including scrolling of the clipmap as the player moves through the world. CLICK

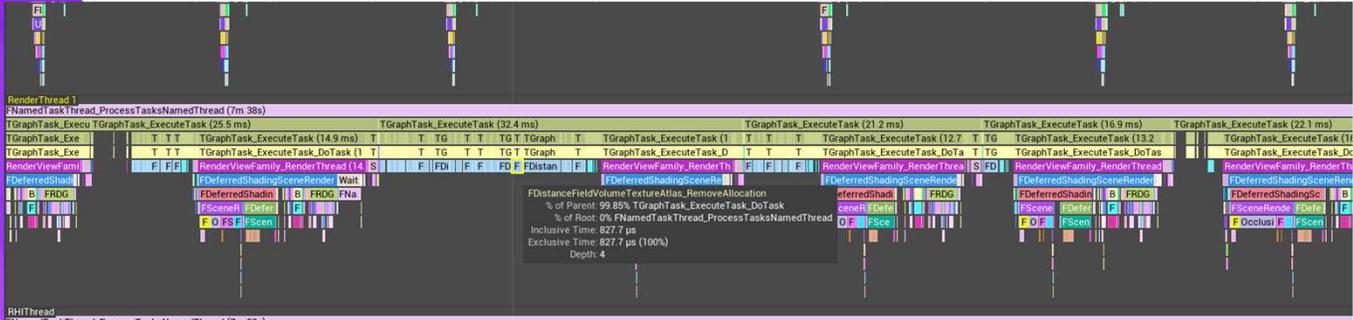
In the interest of time I'm going to cover my optimization effort on the render thread cost – but I've included details on GPU improvements in the appendix for these slides for those interested.



Now, recall that I mentioned that the quality and atlas memory usage encouraged the artists to use many small meshes. This pushed the render thread atlas manager to the limit. The block allocator constructs a binary tree of allocations with each node splitting the space in one of the 3 dimensions, which I've illustrated here for a single 64x64x64 allocation. The actual nodes are stored in a linear array, which I've marked the indices of. You can see that Node 5 is the only leaf node in this example.

Due to the increasingly large number of distance fields the artists were fitting into the atlas, removals and additions to the atlas became very expensive on the render thread. This time would probably be inconsequential in a more limited tree, but it became a problem for us as the distance field count in any given scene grew.

# Removal Cost



Name	Count	Incl	Excl
CPU (1 / 755)	801	320.2 ms	320.2 ms
FDistanceFieldVolumeTextureAtlas_Rer	801	320.2 ms	320.2 ms

Average Incl  
Time: 0.40 ms

Removals were the first problem I tackled – here you can see that after a garbage collect the render thread is processing a ton of expensive removals in the atlas for multiple frames as incremental gc is processed. The cost to do a removal is variable but in one section I profiled the average cost was 0.40 ms. Even with throttling that seems expensive.

Let's investigate the code and look for improvements. This is going to be an exercise of algorithm optimization – think Big O notation

```
bool RemoveElement(uint32 ElementBaseX, uint32 ElementBaseY, uint32 ElementBaseZ)
{
    int32 FoundNodeIndex = INDEX_NONE;
    // Search through nodes to find the element to remove
    // @todo - traverse the tree instead of iterating through all nodes
    for (int32 NodeIndex = 0; NodeIndex < Nodes.Num(); NodeIndex++)
    {
        FTextureLayoutNode3DB& Node = Nodes[NodeIndex];

        if (Node.MinX == ElementBaseX
            && Node.MinY == ElementBaseY
            && Node.MinZ == ElementBaseZ
            && Node.SizeX == ElementSizeX
            && Node.SizeY == ElementSizeY
            && Node.SizeZ == ElementSizeZ)
        {
            FoundNodeIndex = NodeIndex;
            break;
        }
    }
}
```

```
int32 FindNode(int32 CurrentNodeIndex, uint32 ElementBaseX, uint32 ElementBaseY, uint32 ElementBaseZ, uint32 ElementSizeX, uint32 ElementSizeY, uint32 ElementSizeZ)
{
    FTextureLayoutNode3DB Node = Nodes[CurrentNodeIndex];

    // if the node fits perfectly it is a match, this is our node
    if (Node.MinX == ElementBaseX
        && Node.MinY == ElementBaseY
        && Node.MinZ == ElementBaseZ
        && Node.SizeX == ElementSizeX
        && Node.SizeY == ElementSizeY
        && Node.SizeZ == ElementSizeZ)
    {
        return CurrentNodeIndex;
    }

    // check if it could fit children
    if (ElementBaseX >= Node.MinX
        && ElementBaseY >= Node.MinY
        && ElementBaseZ >= Node.MinZ)
    {
        uint32 SlackWithMinX = ElementBaseX - Node.MinX;
        uint32 SlackWithMinY = ElementBaseY - Node.MinY;
        uint32 SlackWithMinZ = ElementBaseZ - Node.MinZ;

        if ((SlackWithMinX + ElementSizeX) <= Node.SizeX
            && (SlackWithMinY + ElementSizeY) <= Node.SizeY
            && (SlackWithMinZ + ElementSizeZ) <= Node.SizeZ)
        {
            // Check each child depth-first
            if (Node.ChildA != INDEX_NONE)
            {
                int32 Result = FindNode(Node.ChildA, ElementBaseX, ElementBaseY, ElementBaseZ, ElementSizeX, ElementSizeY, ElementSizeZ);
                if (Result != INDEX_NONE)
                {
                    return Result;
                }
            }

            if (Node.ChildB != INDEX_NONE)
            {
                int32 Result = FindNode(Node.ChildB, ElementBaseX, ElementBaseY, ElementBaseZ, ElementSizeX, ElementSizeY, ElementSizeZ);
                if (Result != INDEX_NONE)
                {
                    return Result;
                }
            }
        }
    }

    // otherwise this is a dead-end
    return INDEX_NONE;
}
```

Recurse Tree

- First up – there are a number of linear searches with todo's in the code by the original authors (CLICK)
- it seems like a no-brainer to get those addressed. (CLICK)

This one is a linear search through the node array to try to find the allocation associated with the distance field we are removing, which we can change to a recursive traversal through the tree from the root node. Since we know the position and size of the node this is pretty close to going from linear to logarithmic in search time, depending on how balanced the tree is.

```
/**
 * Removes a previously allocated element from the layout and collapses the tree as much as possible,
 * in order to create the largest free block possible and return the tree to its state before the element was added.
 * @return True if the element specified by the input parameters existed in the layout.
 */
bool RemoveElement(uint32 ElementBaseX, uint32 ElementBaseY, uint32 ElementBaseZ, uint32 ElementSizeX, uint32 ElementSizeY, uint32 ElementSizeZ)
{
    int32 FoundNodeIndex = INDEX_NONE;
    // Search through nodes to find the element to remove
    // the first node is always the root
    FoundNodeIndex = FindNode(0, ElementBaseX, ElementBaseY, ElementBaseZ, ElementSizeX, ElementSizeY, ElementSizeZ);

    if (FoundNodeIndex != INDEX_NONE)
    {
        // Mark the found node as not being used anymore
        Nodes[FoundNodeIndex].Used = false;

        // Walk up the tree to find the node closest to the root that doesn't have any used children
        int32 ParentNodeIndex = FindParentNode(FoundNodeIndex);
        ParentNodeIndex = IsUsedNode(ParentNodeIndex) ? INDEX_NONE : ParentNodeIndex;
        int32 LastParentNodeIndex = ParentNodeIndex;
        while (ParentNodeIndex != INDEX_NONE)
        {
            // IsUsedNode(Nodes[ParentNodeIndex].ChildA)
            // IsUsedNode(Nodes[ParentNodeIndex].ChildB)
            {
                LastParentNodeIndex = ParentNodeIndex;
                ParentNodeIndex = FindParentNode(ParentNodeIndex);
            }
        }

        // Remove the children of the node closest to the root with only unused children,
        // which restores the tree to its state before this element was allocated,
        // And allows allocations as large as LastParentNode in the future.
        if (LastParentNodeIndex != INDEX_NONE)
        {
            RemoveChildren(LastParentNodeIndex);
        }
        return true;
    }
    return false;
}
```

```
/**
 * Removes a previously allocated element from the layout and collapses the tree as much as possible,
 * in order to create the largest free block possible and return the tree to its state before the element was added.
 * @return True if the element specified by the input parameters existed in the layout.
 */
bool RemoveElement(uint32 ElementBaseX, uint32 ElementBaseY, uint32 ElementBaseZ, uint32 ElementSizeX, uint32 ElementSizeY, uint32 ElementSizeZ)
{
    int32 FoundNodeIndex = INDEX_NONE;
    // Search through nodes to find the element to remove
    // the first node is always the root
    FoundNodeIndex = FindNode(0, ElementBaseX, ElementBaseY, ElementBaseZ, ElementSizeX, ElementSizeY, ElementSizeZ);

    if (FoundNodeIndex != INDEX_NONE)
    {
        // Mark the found node as not being used anymore
        Nodes[FoundNodeIndex].Used = false;

        // Walk up the tree to find the node closest to the root that doesn't have any used children
        int32 LastParentNodeIndex = INDEX_NONE;
        int32 ParentNodeIndex = Nodes[FoundNodeIndex].Parent;
        ParentNodeIndex = IsUsedNode(ParentNodeIndex) ? INDEX_NONE : ParentNodeIndex;
        LastParentNodeIndex = ParentNodeIndex;
        while (ParentNodeIndex != INDEX_NONE)
        {
            // IsUsedNode(Nodes[ParentNodeIndex].ChildA)
            // IsUsedNode(Nodes[ParentNodeIndex].ChildB)
            {
                LastParentNodeIndex = ParentNodeIndex;
                ParentNodeIndex = Nodes[ParentNodeIndex].Parent;
            }
        }

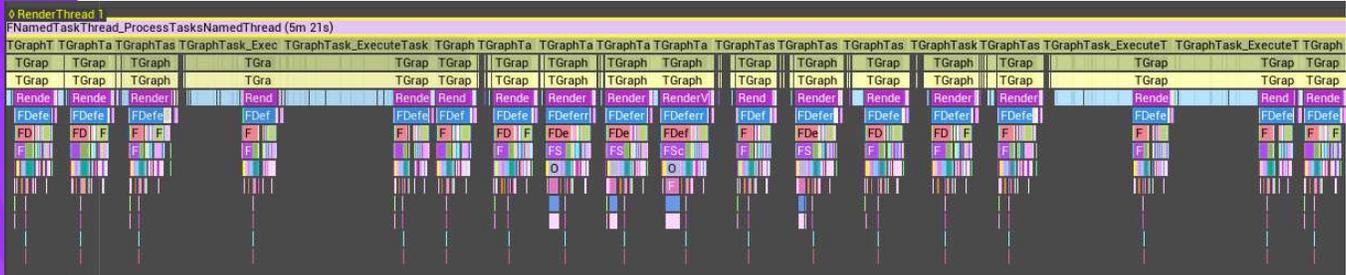
        // Remove the children of the node closest to the root with only unused children,
        // which restores the tree to its state before this element was allocated,
        // And allows allocations as large as LastParentNode in the future.
        if (LastParentNodeIndex != INDEX_NONE)
        {
            RemoveChildren(LastParentNodeIndex);
        }
        return true;
    }
    return false;
}
// @todo - could be a constant time search if the nodes stored a parent index
```

```
/** Returns the index into nodes of the parent node of SearchNode. */
int32 FindParentNode(int32 SearchNodeIndex)
{
    // @todo - could be a constant time search if the nodes stored a parent index
    for (int32 NodeIndex = 0; NodeIndex < Nodes.Num(); NodeIndex++)
    {
        FTextureLayoutNode3d& Node = Nodes[NodeIndex];
        if (Node.ChildA == SearchNodeIndex || Node.ChildB == SearchNodeIndex)
        {
            return NodeIndex;
        }
    }
    return INDEX_NONE;
}
```

After the block being unloaded is identified, the RemoveElement function needs to find which parent nodes need to be unloaded with it. It traverse up from the node repeatedly calling this FindParentNode function, which upon closer inspection also has a handy todo comment in it that it is once again doing an unnecessary linear search. (CLICK)

This is even simpler to resolve, we'll pay the extra 4 bytes per node to store the parent node index when blocks are inserted, to be able to turn these calls to FindParentNode into just an assignment

# Removal Cost



Name	Count	Incl	Excl
CPU (1 / 785)	692	227.3 ms	227.3 ms
FDistanceFieldVolumeTextureAtlas_Rer	692	227.3 ms	227.3 ms

Average Incl  
Time: 0.33 ms

Alright let's check in on how we're doing. Better, but still not great. We're still averaging a third of a ms per removal which does not scale to dozens of removals happening in a single frame. You can see we frequently are spending more than 16 ms a frame processing large numbers of removals, which takes me to looking at more code, guided by function sample profiling.

```
    if (!Nodes[ParentNodeIndex].ChildB)
    {
        LastParentNodeIndex = ParentNodeIndex;
        ParentNodeIndex = Nodes[ParentNodeIndex].Parent;
    }
}

// Remove the children of the node closest to the root with only unused children,
// which restores the tree to its state before this element was allocated,
// And allows allocations as large as LastParentNodeIndex in the future.
if (LastParentNodeIndex != INDEX_NONE)
{
    RemoveChildren(LastParentNodeIndex);
}
return true;
}
return false;
}
```

```
/** Recursively removes the children of a given node from the Nodes array and adjusts existing indices to compensate. */
void RemoveChildren(int32 NodeIndex)
{
    // Traverse the children depth first
    if (Nodes[NodeIndex].ChildA != INDEX_NONE)
    {
        RemoveChildren(Nodes[NodeIndex].ChildA);
    }
    if (Nodes[NodeIndex].ChildB != INDEX_NONE)
    {
        RemoveChildren(Nodes[NodeIndex].ChildB);
    }
    if (Nodes[NodeIndex].ChildA != INDEX_NONE)
    {
        // Store off the index of the child since it may be changed in the code below
        const int32 OldChildA = Nodes[NodeIndex].ChildA;
        // Remove the child from the Nodes array
        Nodes.RemoveAt(OldChildA);
        // Iterate over all the Nodes and fix up their indices now that an element has been removed
        for (int32 OtherNodeIndex = 0; OtherNodeIndex < Nodes.Num(); OtherNodeIndex++)
        {
            if (Nodes[OtherNodeIndex].ChildA >= OldChildA)
            {
                Nodes[OtherNodeIndex].ChildA--;
            }
            if (Nodes[OtherNodeIndex].ChildB >= OldChildA)
            {
                Nodes[OtherNodeIndex].ChildB--;
            }
            if (Nodes[OtherNodeIndex].Parent >= OldChildA)
            {
                Nodes[OtherNodeIndex].Parent--;
            }
        }
        // Mark the node as not having a ChildA
        Nodes[NodeIndex].ChildA = INDEX_NONE;
    }
    if (Nodes[NodeIndex].ChildB != INDEX_NONE)
    {
        const int32 OldChildB = Nodes[NodeIndex].ChildB;
        Nodes.RemoveAt(OldChildB);
        for (int32 OtherNodeIndex = 0; OtherNodeIndex < Nodes.Num(); OtherNodeIndex++)
        {
            if (Nodes[OtherNodeIndex].ChildA >= OldChildB)
            {
                Nodes[OtherNodeIndex].ChildA--;
            }
            if (Nodes[OtherNodeIndex].ChildB >= OldChildB)
            {
                Nodes[OtherNodeIndex].ChildB--;
            }
        }
    }
}
```

Following the search up the parent nodes to decide the sub-tree to remove, there's this RemoveChildren function that is called to actually remove those nodes. (CLICK)

Function sampling profiling told me that this was a hot function – so let's take a look at what's happening inside of it and see if we can identify some improvements.

```
    && !Nodes[nodeIndex].ChildB)
    {
        LastParentNodeIndex = ParentNodeIndex;
        ParentNodeIndex = Nodes[ParentNodeIndex].Parent;
    }
}

// Remove the children of the node closest to the root with only unused children.
// This restores the tree to its state before this element was allocated.
// No more allocations as large as lastParentNode in the future.
if (lastParentNodeIndex != INDEX_NONE)
{
    RemoveChildren(LastParentNodeIndex);
}
return true;
}

return false;
}

/** Recursively removes the children of a given node from the Nodes array and adjusts existing indices to compensate. */
void RemoveChildren(int32 NodeIndex)
{
    // Traverse the children depth first
    if (Nodes[NodeIndex].ChildA != INDEX_NONE)
    {
        RemoveChildren(Nodes[NodeIndex].ChildA);
    }
    if (Nodes[NodeIndex].ChildB != INDEX_NONE)
    {
        RemoveChildren(Nodes[NodeIndex].ChildB);
    }
}

if (Nodes[NodeIndex].ChildA != INDEX_NONE)
{
    // Store off the index of the child since it may be changed in the code below
    const int32 OldChildA = Nodes[NodeIndex].ChildA;

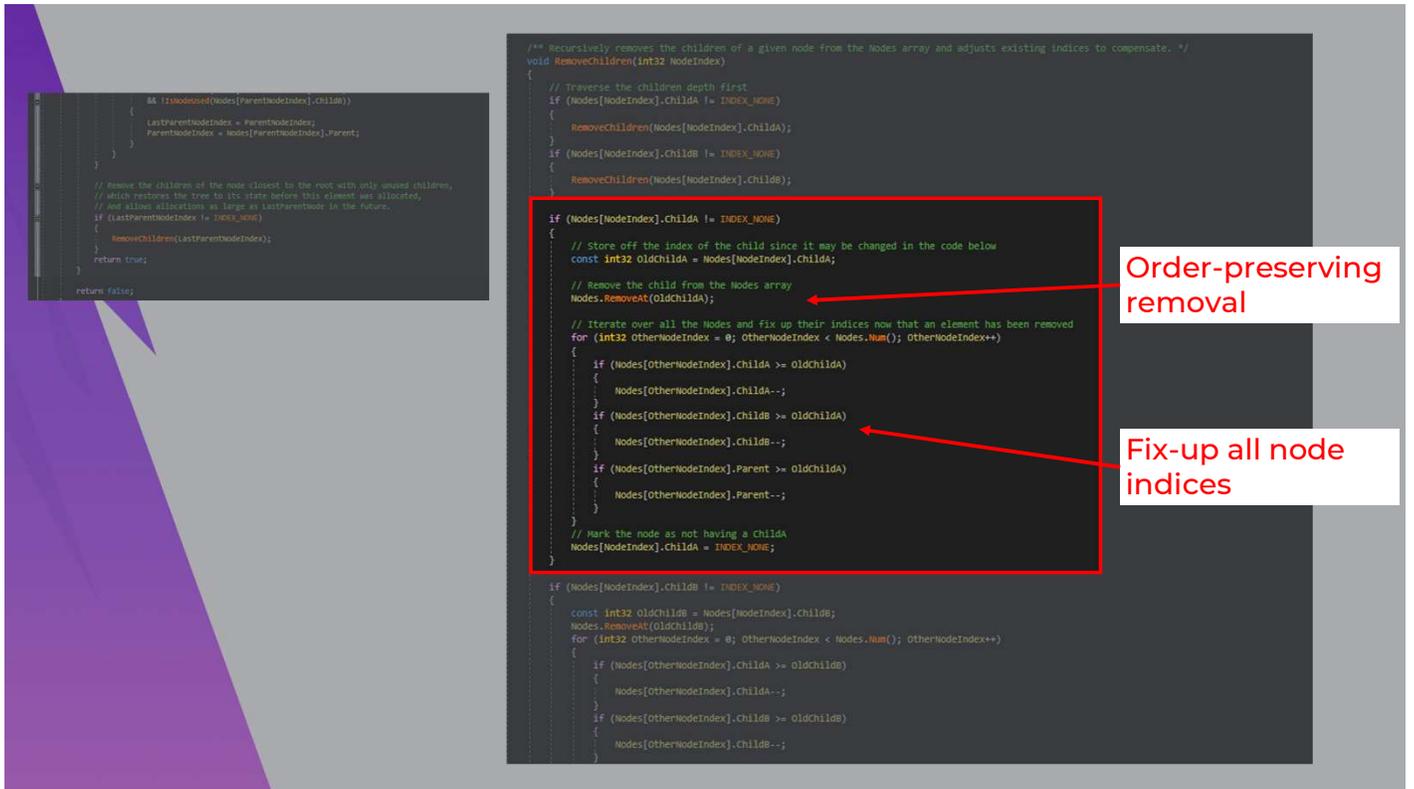
    // Remove the child from the Nodes array
    Nodes.RemoveAt(OldChildA);

    // Iterate over all the Nodes and fix up their indices now that an element has been removed
    for (int32 OtherNodeIndex = 0; OtherNodeIndex < Nodes.Num(); OtherNodeIndex++)
    {
        if (Nodes[OtherNodeIndex].ChildA >= OldChildA)
        {
            Nodes[OtherNodeIndex].ChildA--;
        }
        if (Nodes[OtherNodeIndex].ChildB >= OldChildA)
        {
            Nodes[OtherNodeIndex].ChildB--;
        }
        if (Nodes[OtherNodeIndex].Parent >= OldChildA)
        {
            Nodes[OtherNodeIndex].Parent--;
        }
    }

    // Mark the node as not having a ChildA
    Nodes[NodeIndex].ChildA = INDEX_NONE;
}

if (Nodes[NodeIndex].ChildB != INDEX_NONE)
{
    const int32 OldChildB = Nodes[NodeIndex].ChildB;
    Nodes.RemoveAt(OldChildB);
    for (int32 OtherNodeIndex = 0; OtherNodeIndex < Nodes.Num(); OtherNodeIndex++)
    {
        if (Nodes[OtherNodeIndex].ChildA >= OldChildB)
        {
            Nodes[OtherNodeIndex].ChildA--;
        }
        if (Nodes[OtherNodeIndex].ChildB >= OldChildB)
        {
            Nodes[OtherNodeIndex].ChildB--;
        }
    }
}
}
```

The first thing to note that this is once again a recursive function traversing down each child from the node we've selected for removal. So it makes sense that this function showed up in profiling.



Now let's look at the actual work is happening – each child is removed from the node array. This is handled by calling `RemoveAt` on the array holding all the nodes – which is using Unreal's `TArray` container class. This is an order-preserving removal, which internally means a large memmove of all the nodes following the removal index are copied down into the hole.

Then in the following loop, each node has the indices of their children iterated and adjusted to account for the shift in the array. Now, you may be familiar that there is a faster way to remove an item from an array if you *don't* care to preserve the order, which is to simply copy the final element into the hole that was just created. I'd like to do that here to avoid iterating every node inside a sub-tree traversal by simply fixing up the child indices of the parent node – which is a constant time access due to the prior change I showed where we store that on each node now. However, this is tricky to do in practice because we are trying to modify the tree as we are traversing it.

```

}
// Remove the children of the node closest to the root with only unused children,
// which restores the tree to its state before this element was allocated,
// and allows allocations as large as LastParentNode in the future.
if (LastParentNodeIndex != INDEX_NONE)
{
    // provide a stack for bookkeeping when walking the tree
    Teregs(int32 ParentNodeStack);
    ParentNodeStack.Reserve(5); // reserve 5 based on empirical evidence that most removals take 2-3 recursions
    ParentNodeStack.Push(LastParentNodeIndex);
    RemoveChildren(ParentNodeStack);
}
return true;

```

Allocate node stack for top level call

```

{
    ParentNodeStack.Push(Nodes[ParentNodeStack.Top()].ChildA);
    RemoveChildren(ParentNodeStack);
}
if (Nodes[ParentNodeStack.Top()].ChildB != INDEX_NONE)
{
    ParentNodeStack.Push(Nodes[ParentNodeStack.Top()].ChildB);
    RemoveChildren(ParentNodeStack);
}
if (Nodes[ParentNodeStack.Top()].ChildA != INDEX_NONE)
{
    // Store off the index of the child since it may be changed in the code below
    const int32 OldChildA = Nodes[ParentNodeStack.Top()].ChildA;
    // Mark the node as not having a ChildA
    Nodes[ParentNodeStack.Top()].ChildA = INDEX_NONE;
    const int32 EndIndex = Nodes.Num() > 0 ? (Nodes.Num() - 1) : 0;
    if (OldChildA < EndIndex)
    {
        // copy end Node into hole, no need to swap as we are about to trim the hole off the end
        FMemory::Memcpy(&Nodes[OldChildA], &Nodes[EndIndex], sizeof(FTextureLayoutNode3d));
        // update parent/children node tracking of node we just moved into the hole
        int32 ParentIndex = Nodes[OldChildA].Parent;
        if (ParentIndex != INDEX_NONE)
        {
            if (Nodes[ParentIndex].ChildA == EndIndex)
            {
                Nodes[ParentIndex].ChildA = OldChildA;
            }
            if (Nodes[ParentIndex].ChildB == EndIndex)
            {
                Nodes[ParentIndex].ChildB = OldChildA;
            }
        }
        if (Nodes[OldChildA].ChildA != INDEX_NONE)
        {
            Nodes[Nodes[OldChildA].ChildA].Parent = OldChildA;
        }
        if (Nodes[OldChildA].ChildB != INDEX_NONE)
        {
            Nodes[Nodes[OldChildA].ChildB].Parent = OldChildA;
        }
    }
    // we need to take care if we moved a node as we walk the stack, scan our list and make any corrections needed
    for (int32 StackIndex = 0; StackIndex < ParentNodeStack.Num(); StackIndex++)
    {
        if (ParentNodeStack[StackIndex] == EndIndex)
        {
            ParentNodeStack[StackIndex] = OldChildA;
        }
    }
}
// Remove the node moved into the hole from the end of the array
Nodes.RemoveAt(EndIndex, 1, false);
}

```

To get around this problem, I allocate and pass around a stack of parent nodes as we recurse back down the tree. (CLICK)

This allows us to perform fixup on the stack and avoid getting lost during our traversal as we remove nodes. This can be seen here in the updated logic inside of the RemoveChildren function.

```
    }
    // Remove the children of the node closest to the root with only unrooted children,
    // which restores the tree to its state before this element was allocated,
    // and allows allocations as large as lastParentNode in the future.
    if (LastParentNodeIndex != INDEX_NONE)
    {
        // avoid a stack for bookkeeping when walking the tree
        TArray<int32> ParentNodeStack;
        ParentNodeStack.Reserve(5); // Reserve 5 based on empirical evidence that most removals take 2-3 recursions
        ParentNodeStack.Push(LastParentNodeIndex);

        RemoveChildren(ParentNodeStack);
    }
    return true;
}

ParentNodeStack.Push(Nodes[ParentNodeStack.Top()].ChildA);
RemoveChildren(ParentNodeStack);
}
if (Nodes[ParentNodeStack.Top()].ChildB != INDEX_NONE)
{
    ParentNodeStack.Push(Nodes[ParentNodeStack.Top()].ChildB);
    RemoveChildren(ParentNodeStack);
}
if (Nodes[ParentNodeStack.Top()].ChildA != INDEX_NONE)
{
    // Store off the index of the child since it may be changed in the code below
    const int32 OldChildA = Nodes[ParentNodeStack.Top()].ChildA;
    // Mark the node as not having a ChildA
    Nodes[ParentNodeStack.Top()].ChildA = INDEX_NONE;

    const int32 EndIndex = Nodes.Num() > 0 ? (Nodes.Num() - 1) : 0;
    if (OldChildA < EndIndex)
    {
        // copy end Node into hole, no need to swap as we are about to trim the hole off the end
        FMemory::Memcpy(&Nodes[OldChildA], &Nodes[EndIndex], sizeof(FTextureLayoutNode3D));

        // update parent/children node tracking of node we just moved into the hole
        int32 ParentIndex = Nodes[OldChildA].Parent;
        if (ParentIndex != INDEX_NONE)
        {
            if (Nodes[ParentIndex].ChildA == EndIndex)
            {
                Nodes[ParentIndex].ChildA = OldChildA;
            }
            if (Nodes[ParentIndex].ChildB == EndIndex)
            {
                Nodes[ParentIndex].ChildB = OldChildA;
            }
        }
        if (Nodes[OldChildA].ChildA != INDEX_NONE)
        {
            Nodes[Nodes[OldChildA].ChildA].Parent = OldChildA;
        }
        if (Nodes[OldChildA].ChildB != INDEX_NONE)
        {
            Nodes[Nodes[OldChildA].ChildB].Parent = OldChildA;
        }
    }

    // we need to take care if we moved a node as we walk the stack, scan our list and make any corrections needed
    for (int32 StackIndex = 0; StackIndex < ParentNodeStack.Num(); StackIndex++)
    {
        if (ParentNodeStack[StackIndex] == EndIndex)
        {
            ParentNodeStack[StackIndex] = OldChildA;
        }
    }

    // Remove the node moved into the hole from the end of the array
    Nodes.RemoveAt(EndIndex, 1, false);
}
}
```

Allocate node stack for top level call

Copy into hole

Fix-up child and parent pointers

Here is the operation where we take the node at the end of the array and copy it into the location where we are doing a removal. Then we fix-up the child and parent node references to account for the node we just moved.

The image shows two snippets of C++ code with red arrows pointing to specific lines and white text boxes with red arrows pointing to those lines. The code is related to a tree traversal algorithm that uses a stack to manage nodes.

**Left Snippet:**

```

}
// Remove the children of the node closest to the root with only unvisited children,
// which restores the tree to its state before this element was allocated,
// and allows allocations at large as lastParentNode in the future.
if (LastParentNodeIndex != INDEX_NONE)
{
    // avoid a stack for bookkeeping when walking the tree
    Teryy(int32 ParentNodeStack);
    ParentNodeStack.Reserve(1); // Reserve 1 based on empirical evidence that most removals take 2-3 recursions
    ParentNodeStack.Push(LastParentNodeIndex);
}
RemoveChildren(ParentNodeStack);
return true;

```

**Annotations for Left Snippet:**

- Allocate node stack for top level call:** Points to the `ParentNodeStack` stack creation and `Reserve` call.
- Fix-up node stack:** Points to the `ParentNodeStack.Push(LastParentNodeIndex);` line.

**Right Snippet:**

```

ParentNodeStack.Push(Nodes[ParentNodeStack.Top()].ChildA);
RemoveChildren(ParentNodeStack);
}
if (Nodes[ParentNodeStack.Top()].ChildB != INDEX_NONE)
{
    ParentNodeStack.Push(Nodes[ParentNodeStack.Top()].ChildB);
    RemoveChildren(ParentNodeStack);
}
if (Nodes[ParentNodeStack.Top()].ChildA != INDEX_NONE)
{
    // Store off the index of the child since it may be changed in the code below
    const int32 OldChildA = Nodes[ParentNodeStack.Top()].ChildA;
    // Mark the node as not having a ChildA
    Nodes[ParentNodeStack.Top()].ChildA = INDEX_NONE;
    const int32 EndIndex = Nodes.Num() > 0 ? (Nodes.Num() - 1) : 0;
    if (OldChildA < EndIndex)
    {
        // copy end Node into hole, no need to swap as we are about to trim the hole off the end
        FMemory::Memcpy(&Nodes[OldChildA], &Nodes[EndIndex], sizeof(FTextureLayoutNode3D));
    }
    // update parent/children node tracking of node we just moved into the hole
    int32 ParentIndex = Nodes[OldChildA].Parent;
    if (ParentIndex != INDEX_NONE)
    {
        if (Nodes[ParentIndex].ChildA == EndIndex)
        {
            Nodes[ParentIndex].ChildA = OldChildA;
        }
        if (Nodes[ParentIndex].ChildB == EndIndex)
        {
            Nodes[ParentIndex].ChildB = OldChildA;
        }
    }
    if (Nodes[OldChildA].ChildA != INDEX_NONE)
    {
        Nodes[Nodes[OldChildA].ChildA].Parent = OldChildA;
    }
    if (Nodes[OldChildA].ChildB != INDEX_NONE)
    {
        Nodes[Nodes[OldChildA].ChildB].Parent = OldChildA;
    }
}
// we need to take care if we moved a node as we walk the stack, scan our list and make any corrections needed
for (int32 StackIndex = 0; StackIndex < ParentNodeStack.Num(); StackIndex++)
{
    if (ParentNodeStack[StackIndex] == EndIndex)
    {
        ParentNodeStack[StackIndex] = OldChildA;
    }
}
// Remove the node moved into the hole from the end of the array
Nodes.RemoveAt(EndIndex, 1, false);

```

**Annotations for Right Snippet:**

- Copy into hole:** Points to the `Memcpy` call.
- Fix-up child and parent pointers:** Points to the `if (Nodes[ParentIndex].ChildA == EndIndex)` and `if (Nodes[ParentIndex].ChildB == EndIndex)` blocks.
- Fix-up node stack:** Points to the `ParentNodeStack[StackIndex] = OldChildA;` line in the `for` loop.

Following this I want to call attention to this code – here we need to look at our stack of nodes we’re maintaining and make sure that if one of the nodes in our traversal path was the one we just moved, we need to fix-up the location in our stack so we don’t get lost as we continue our traversal



# Higher Level Code

```
void FDistanceFieldVolumeTextureAtlas::RemoveAllocation(FDistanceFieldVolumeTexture* Texture)
{
    SCOPED_NAMED_EVENT(FDistanceFieldVolumeTextureAtlas_RemoveAllocation, FColor::Emerald);
    InitializeIfNeeded();
    PendingAllocations.Remove(Texture);
    AllocatedCPUDataInBytes -= Texture->VolumeData.CompressedDistanceFieldVolume.GetAllocatedSize();

    if (FailedAllocations.Remove(Texture) > 0)
    {
        FailedAllocatedPixels -= Texture->VolumeData.Size.X * Texture->VolumeData.Size.Y * Texture->VolumeData.Size.Z;
    }

    if (!CurrentAllocations.Contains(Texture))
    {
        return;
    }

    FIntVector Size = Texture->SizeInAtlas;
    int PixelAreaSize = Size.X * Size.Y * Size.Z;

    const FIntVector Min = Texture->GetAllocationMin();
    verify(BlockAllocator.RemoveElement(Min.X, Min.Y, Min.Z, Size.X, Size.Y, Size.Z));
    CurrentAllocations.Remove(Texture);
    AllocatedPixels -= PixelAreaSize;

    FIntVector RemainingSize = Size;
}
```

```
void FDistanceFieldVolumeTextureAtlas::RemoveAllocation(FDistanceFieldVolumeTexture* Texture)
{
    SCOPED_NAMED_EVENT(FDistanceFieldVolumeTextureAtlas_RemoveAllocation, FColor::Emerald);
    InitializeIfNeeded();
    // @igs(jmoore) - START - [IGRender][IGOptimization]
    // prefer RemoveSwap to Remove for TArray where the order does not currently matter
    PendingAllocations.RemoveSwap(Texture, false);
    AllocatedCPUDataInBytes -= Texture->VolumeData.CompressedDistanceFieldVolume.GetAllocatedSize();
    if (FailedAllocations.RemoveSwap(Texture, false) > 0)
    {
        FailedAllocatedPixels -= Texture->VolumeData.Size.X * Texture->VolumeData.Size.Y * Texture->VolumeData.Size.Z;
    }

    FIntVector Size = Texture->SizeInAtlas;
    int PixelAreaSize = Size.X * Size.Y * Size.Z;
    const FIntVector Min = Texture->GetAllocationMin();
    // eliminate redundant check if Contains by just checking if RemoveSwap removed anything before continuing
    if (CurrentAllocations.RemoveSwap(Texture, false) == 0)
    {
        return;
    }
    AllocatedPixels -= PixelAreaSize;

    verify(BlockAllocator.RemoveElement(Min.X, Min.Y, Min.Z, Size.X, Size.Y, Size.Z));
    // @igs(jmoore) - END - [IGRender][IGOptimization]
}
```

Name	Count	Incl	Excl
CPU (1 / 774)	1,064	10.9 ms	10.9 ms
FDistanceFieldVolumeTextureAtlas_RemoveAllocation	1,064	10.9 ms	10.9 ms

Average Incl  
Time: 0.01 ms

There's one more thing improvement that I want to mention that is a bit of a pet peeve of mine. Don't do a search and remove separately! The top level management that uses the block allocator I was optimizing has several lists that can use RemoveSwap instead of Remove because order doesn't matter, and then the call to Contains (which will do a linear search) can be eliminated by just using the fact that RemoveSwap will return the number of elements removed.

This halves our removal time again – although this was a just a small portion of the cost before our other improvements were made. This is a 40x improvement overall. Again, I'm sure that this time is inconsequential for many games, the game we were building coupled open world level streaming with high numbers of small distance field allocations, it became very worthwhile for us to spend time optimizing this code.

# Insertion Cost

```
bool AddElement(uint32& OutBaseX, uint32& OutBaseY, uint32& OutBaseZ, uint32 ElementSizeX, uint32 ElementSizeY, uint32 ElementSizeZ)
{
    SCOPED_NAMED_EVENT(TextureLayout3D_AddElement, FColor::Emerald);

    if (ElementSizeX == 0 || ElementSizeY == 0 || ElementSizeZ == 0)
    {
        OutBaseX = 0;
        OutBaseY = 0;
        OutBaseZ = 0;
        return true;
    }

    if (bAlignByFour)
    {
        // Pad to 4 to ensure alignment
        ElementSizeX = (ElementSizeX + 3) & ~3;
        ElementSizeY = (ElementSizeY + 3) & ~3;
        ElementSizeZ = (ElementSizeZ + 3) & ~3;
    }

    // Try allocating space without enlarging the texture.
    int32 NodeIndex = AddSurfaceInner(0, ElementSizeX, ElementSizeY, ElementSizeZ, false);
    if (NodeIndex == INDEX_NONE)
    {
        // Try allocating space which might enlarge the texture.
        NodeIndex = AddSurfaceInner(0, ElementSizeX, ElementSizeY, ElementSizeZ, true);
    }
}
```

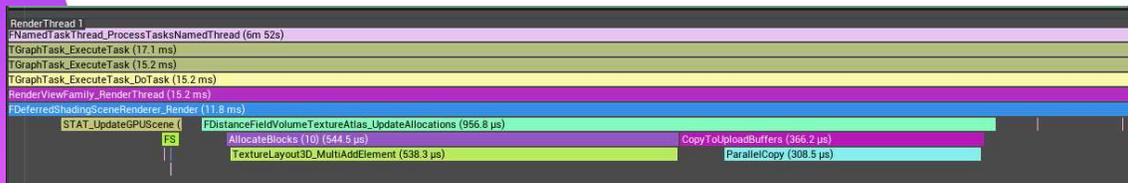


Besides removals from the block allocator, insertions also had lots of room for improvement – these were walking the tree looking for a suitable hole to allocate the DF into, but each insertion would start over at the root node. This seemed redundant – in the trace associated with this code you can see AddElement is called 10 times, and the tree is potentially walked twice by each call – this is taking 2.5 ms in a single frame

# Insertion Cost

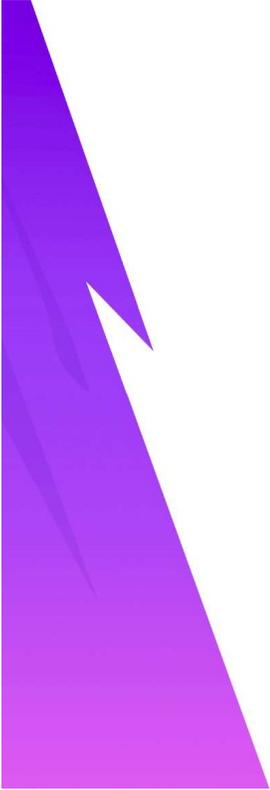
```
void MultiAddElement(const int32 ElementCount, TArrayuint16& OutResults, TArrayuint32& OutBaseX, TArrayuint32& OutBaseY, TArrayuint32& OutBaseZ, TArrayuint32& OutBaseW, TArrayuint32& ElementSizeX, TArrayuint32& ElementSizeY, TArrayuint32& ElementSizeZ, TArrayuint32& ElementSizeW, TArrayuint32& ElementSizeH, TArrayuint32& ElementSizeD)  
{  
    SCOPED_BMIED_EVENT(TextureLayout3D_MultiAddElement, FColor::Emerald);  
  
    for (int32 Index = 0; Index < ElementCount; ++Index)  
    {  
        if (ElementSizeX[Index] == 0 || ElementSizeY[Index] == 0 || ElementSizeZ[Index] == 0)  
        {  
            OutBaseX[Index] = 0;  
            OutBaseY[Index] = 0;  
            OutBaseZ[Index] = 0;  
            OutBaseW[Index] = 0;  
            OutBaseH[Index] = 0;  
            OutBaseD[Index] = 0;  
        }  
  
        if (!AlignByFour)  
        {  
            // Pad to 4 to ensure alignment  
            ElementSizeX[Index] = (ElementSizeX[Index] + 3) & ~3;  
            ElementSizeY[Index] = (ElementSizeY[Index] + 3) & ~3;  
            ElementSizeZ[Index] = (ElementSizeZ[Index] + 3) & ~3;  
            ElementSizeW[Index] = (ElementSizeW[Index] + 3) & ~3;  
            ElementSizeH[Index] = (ElementSizeH[Index] + 3) & ~3;  
            ElementSizeD[Index] = (ElementSizeD[Index] + 3) & ~3;  
        }  
  
        TArray<int32> NodeIndices;  
        NodeIndices.Reserve(ElementCount);  
        for (int32 Index = 0; Index < ElementCount; ++Index)  
        {  
            NodeIndices.Add(INDEX_NONE);  
        }  
  
        // Try allocating space without enlarging the texture.  
        MultiAddSurfFaceInner(0, NodeIndices, ElementSizeX, ElementSizeY, ElementSizeZ, false);  
        if (NodeIndices.Contains(INDEX_NONE))  
        {  
            // Try allocating space which might enlarge the texture.  
            MultiAddSurfFaceInner(0, NodeIndices, ElementSizeX, ElementSizeY, ElementSizeZ, true);  
        }  
    }  
}
```

Batched Insertion



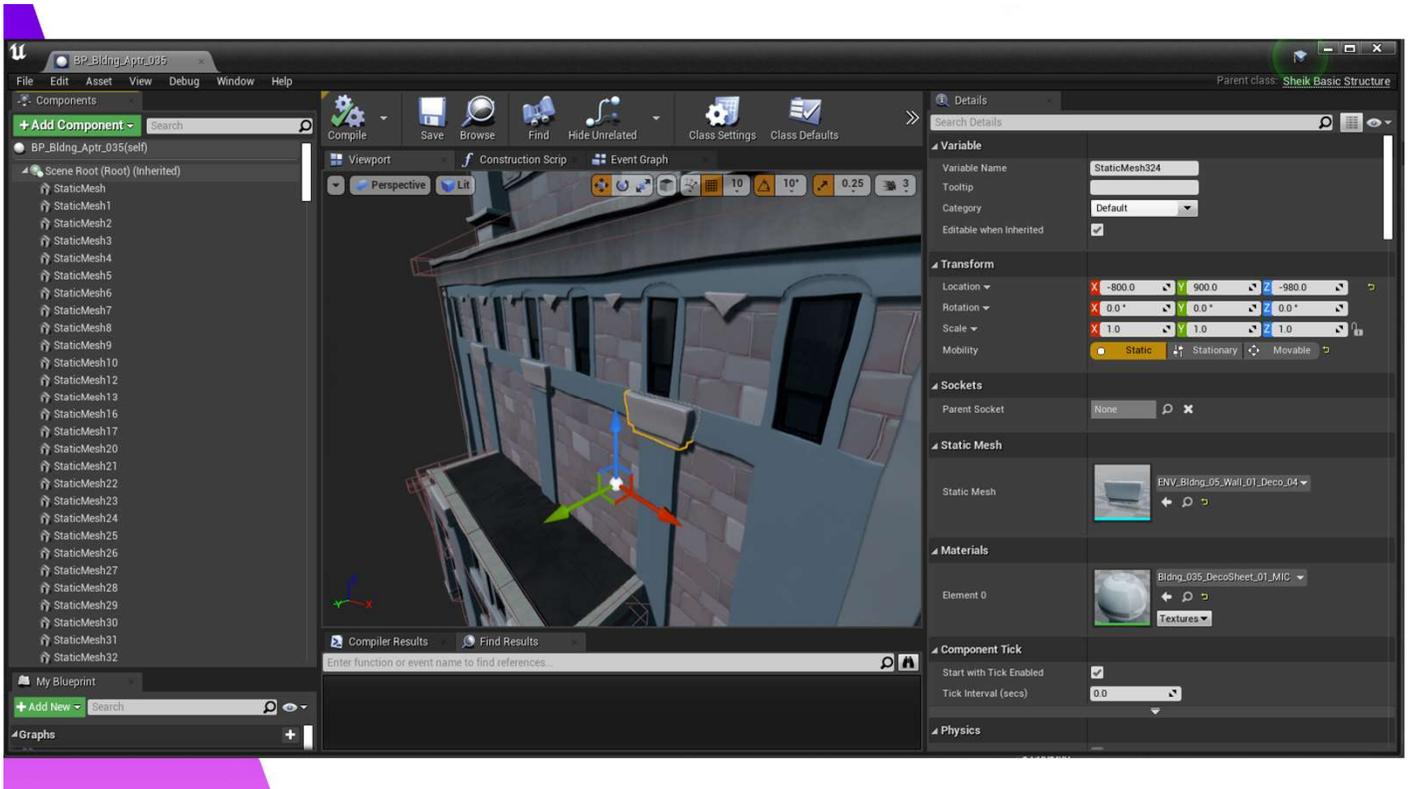
I got much better performance passing around a reference to the full list of desired allocations and checking each node against that list. This means more work being done when a node is visited, but less times traversing parts of the tree redundantly.

Here you can see that for a frame with 10 insertions, it is 2 ms faster than before.



# Reactive Optimization: Too Many Primitives

I've talked extensively now about the pain points caused by our large numbers of distance fields, but let's look at another render thread challenge we encountered – also related to our decisions with how we built the environment

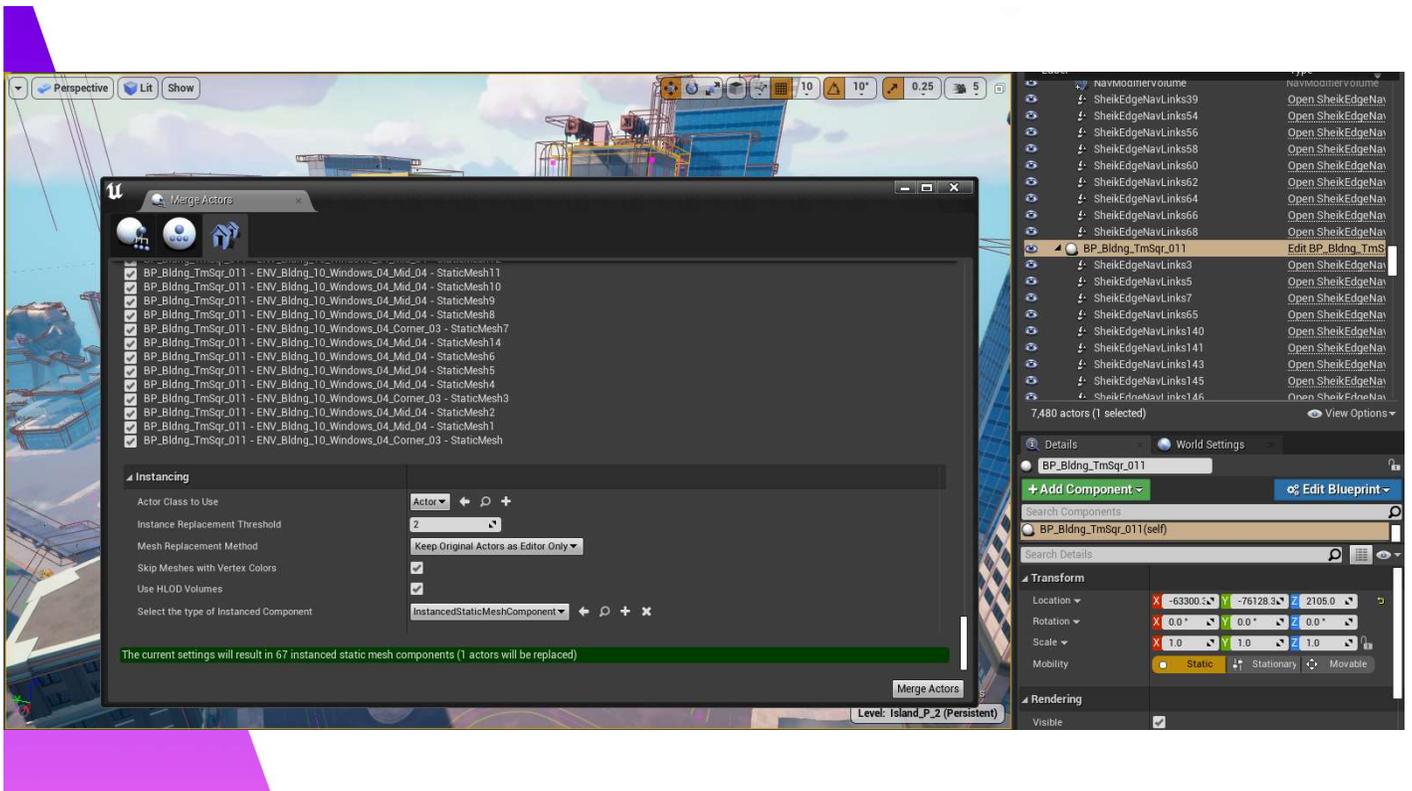


As a reminder – our assets are constructed out of many smaller mesh components kit bashed together.

The scope of the project had expanded after our vertical slice milestone and our art team pushed for even more polish and detail than was in our original plan. We started hitting the limits of dynamic instancing and became largely bottlenecked on what is known as “Init Views” in Unreal Engine – this is where frustum and occlusion culling takes place.



Here is an example of a challenging shot in front of some buildings that are made up of many primitives that are processed individually before instancing. You can see here that `InitViews` is taking up a hefty 7.7 ms, which is approaching half of our total frame time.



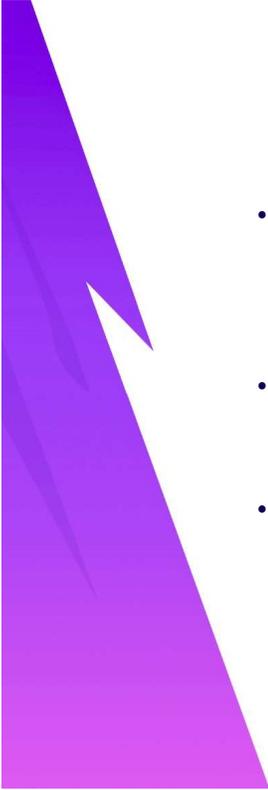
Now – Unreal 4 has an older system for managing costs of large numbers of primitives, which is to merge compatible components into a single “instanced static mesh component.” These manually instance their meshes, ignoring the dynamic instancing system, and go through the renderer as a single primitive.



# ISM Component Drawbacks

- LOD selection overly conservative due to inflated bounds, hurting GPU performance
- Large bounds also degrade frustum and occlusion culling quality
- Hard to work with in editor – burden placed on artists

However, stock Instanced Static Mesh components have a number of drawbacks compared to using individual mesh components that dynamically instance. The first two are that GPU performance degrades – the larger bounds on the merged primitive results in an overly conservative LOD selection and culling result. It also places a lot of burden on the artists to ask them to do this manually, which is orthogonal to the performance characteristics but still a concern.

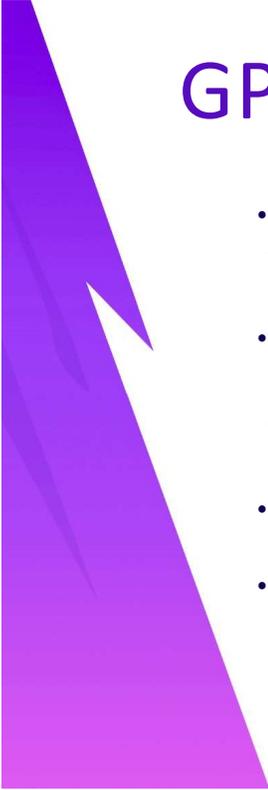


# GPU-driven rendering

- GPU driven rendering of meshes with MultiDrawIndirect can allow us to handle culling + LOD selection on the GPU
- Seb Aaltonen first presented many of these ideas in the talk *GPU-Driven Rendering Pipelines*
- But we don't have the resources to rewrite UE4!

The Cadillac solution is to instead go down the path of full GPU-driven rendering. Issue just a small number of commands on the CPU and handle LOD selection and culling on the GPU instead of dealing with expanded bounds on the primitives on the CPU. Seb Aaltonen pioneered a lot of these ideas in his part of the talk *GPU-Driven Rendering Pipelines* from Siggraph 2015. But our little didn't have the resources to do that ambitious of a re-write of the UE4 renderer.

[http://advances.realtimerendering.com/s2015/aaltonenhaar\\_siggraph2015\\_combined\\_final\\_footer\\_220dpi.pptx](http://advances.realtimerendering.com/s2015/aaltonenhaar_siggraph2015_combined_final_footer_220dpi.pptx)



# GPU-driven ISM Component

- Go halfway – modify ISM component to render with a special ISM GPU Scene
- Convert compatible components on Basic Structure actors into ISM components at cook time
- Render each ISM with DrawIndirect
- Keep ISM's grouped into the Basic Structures (do not merge buildings) so coarse CPU LOD selection and culling is still effective

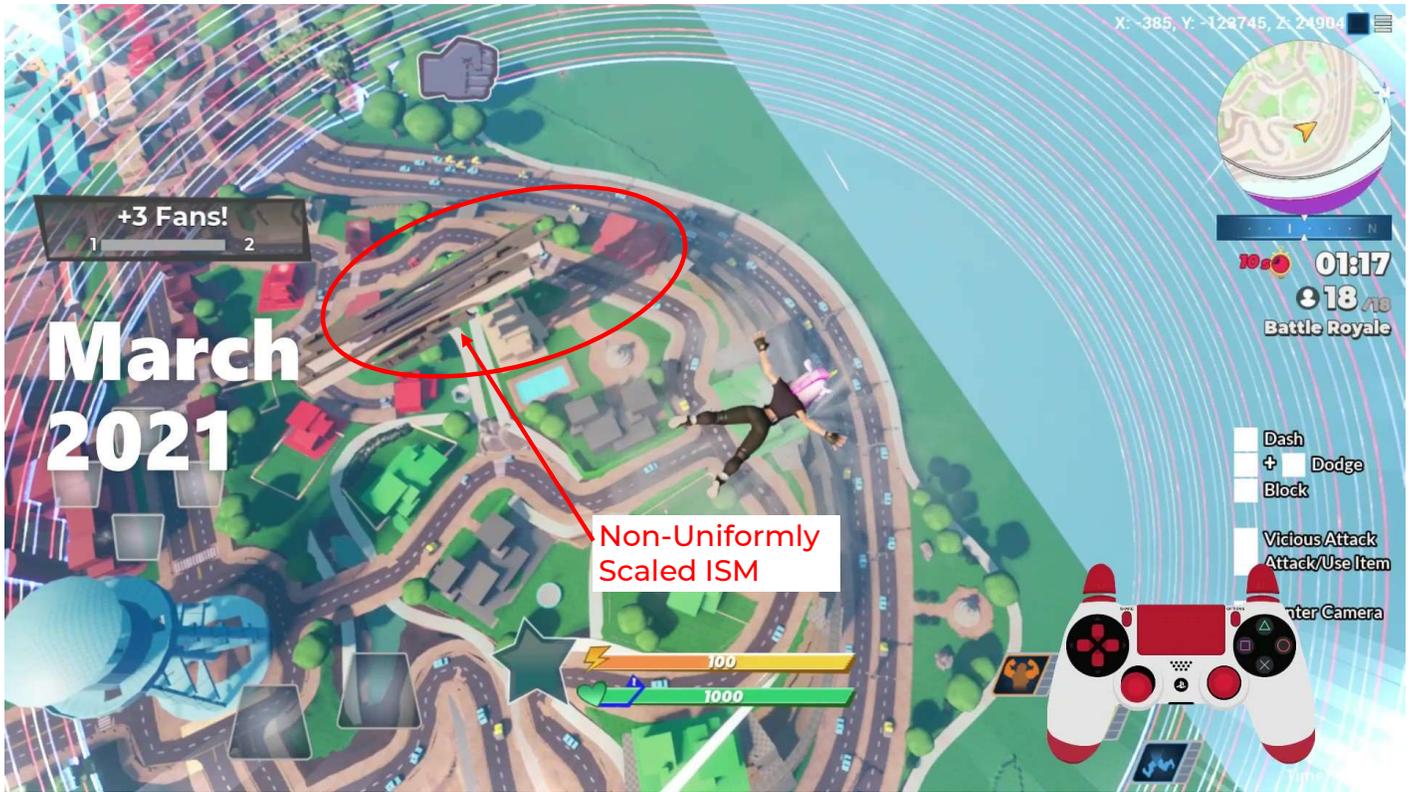
My colleague Nate Mefford figured out a really great compromise that fit our project, to go halfway and modify ISM component to render with a special ISM GPU Scene. We convert compatible components on Basic Structure actors into ISM components at cook time, which keeps the burden off of the artists. We keep ISM's grouped into their owning Basic Structures, so components that could have dynamically instanced between two neighboring buildings within view do not do so. This is because buildings offer a nice unit of coarse CPU-side LOD selection and culling which helps keep our GPU overhead low.

# CPU LOD Selection

```
1435 bool FInstancedStaticMeshSceneProxy::IsUsingCustomLODRules() const
1436 {
1437     return UseInstancedStaticMeshGPUScene();
1438 }
1439
1440 FLODMask FInstancedStaticMeshSceneProxy::GetCustomLOD(const FSceneView& InView, float InViewLODScale, int32 InForcedLODLevel, float& OutScreenSizeSquared) const
1441 {
1442     SCOPE_CYCLE_COUNTER(STAT_GetCustomLODTime);
1443
1444     if (!GInstancedStaticMeshGPUScene_EnableCoarseCPUculling)
1445     {
1446         FLODMask NOCPUcullingLODMask;
1447         NOCPUcullingLODMask.DitheredLODIndices[0] = 0;
1448         NOCPUcullingLODMask.DitheredLODIndices[1] = FInstancedStaticMeshShaderData::MaxLODS;
1449         return NOCPUcullingLODMask;
1450     }
1451
1452     check(PerInstanceShaderData.Num());
1453     check(PerInstanceShaderData.Num() == InstanceRenderData.PerInstanceRenderData->InstanceBuffer.GetNumInstances());
1454
1455     int8 MinLOD = MAX_int8;
1456     int8 MaxLOD = 0;
1457
1458     // NOTE: ComputeStaticMeshLOD *internally* takes into account View.LODDistanceFactor, so do *NOT* preemptively it in here like happens in other places in the code
1459     float LODDistanceScale = GetCachedScalabilityCVars().StaticMeshLODDistanceScale;
1460
1461     // @todo: This function can get kind of expensive. However this function only needs to return an FLODMask that contains a conservative range of LODIndices.
1462     // That implies an alternative faster way to do this function might be as follows:
1463     // -Store off the min / max radius of all instances
1464     // -set MinLOD to the LOD that would occur by placing the largest instance radius at a location of Proxy.WorldLocation + ProxyRadius * ViewDirProxy
1465     // -set MaxLOD to the LOD that would occur by placing the smallest instance radius at a location of Proxy.WorldLocation + ProxyRadius * ViewDirProxy
1466     // That would of course lead to a few more cases of unused LOD draws being generated which is a problem for GPU perf, so for now I think this function is a worthwhile tradeoff of CPU/GPU perf
1467
1468     for (int32 InstanceIndex = 0; InstanceIndex < PerInstanceShaderData.Num(); ++InstanceIndex)
1469     {
1470         // @todo: The way this code is written right now, we are relying on this code to match up bit-for-bit exactly with the LOD calculation code in the GPU job. That is a bad idea
1471         // What we should do is compute the LOD with the bound radius both slightly shrunk and slightly enlarged by epsilon and keep the Min/Max of both of those calculations.
1472         int8 InstanceLOD = ComputeStaticMeshLOD(RenderData, PerInstanceShaderData[InstanceIndex].WorldBoundCenter, PerInstanceShaderData[InstanceIndex].BoundSphereRadius, InView, ClampedInLOD, LODDistanceScale);
1473
1474         MinLOD = FMath::Min(MinLOD, InstanceLOD);
1475         MaxLOD = FMath::Max(MaxLOD, InstanceLOD);
1476     }
1477
1478     FLODMask Result;
1479     Result.DitheredLODIndices[0] = MinLOD;
1480     Result.DitheredLODIndices[1] = MaxLOD;
1481
1482     return Result;
1483 }
```

Find Min+Max LOD range of instances

When I say that coarse CPU LOD selection is important, I mean that we can't just send a draw per LOD to the command processor as that will result in too many empty batches. This code is running per-instance LOD selection on the CPU. This is a CPU vs GPU balancing act and is one clear area where the full GPU driven rendering pipeline would allow bigger gains. Right now we have to decide between being overly conservative and sending extra batches, or spending more CPU time on LOD selection eating into our savings on the render thread.



I should also mention there were lots of fiddly bits with getting negative and non-uniform scale to work correctly in the conversion to Instanced Static Meshes. I'm a believer in supporting both in a renderer, but this work really challenged my faith in that and we hit a lot of under-exercised code paths getting everything to render correctly when converted from individual static meshes to ISM's. You can see a shot from a playtest where a storefront has stretched out into the sky above a neighborhood due to the static meshes going into the ISM conversion having non-uniform scale.

I've included Nate's Unreal-specific workarounds in the appendix of these slides for those interested.



Now before I get into some more implementation details, here are the results for that challenging shot, starting with what we had before – InitViews taking 7.7 ms



And here it is with ISM components on the Basic Structures. This reduces render thread time by 2.37 ms. It varies from scene to scene but it is generally in the 2-3 ms range in savings on a base PS4. And this particular shot sees GPU time virtually flat, although I see it sometimes take 0.1-0.25 ms longer GPU time, which I think is very acceptable given the CPU benefits.

# Instance ID Indirection

```
140 // @big(mefford) - START - [SheikRender][ISOPimization]
141 // NOTE: This really should be owned and controlled by FInstanceStaticMeshGPUScene.h, but it has to live here in the Engine so it is available for when the InstanceVertexFactory is initialized
142 // This is technically the same problem the GPUScene has with the PrimitiveBuffer. They solve that by having a setting a dummyBuffer in the vertexfactory initialization and then override it
143 // dynamically in FInstanceCommandSubmitDraw, but that's a lot of plumbing for something that is still reasonably elegantly handled this way...
144 class FInstanceIndirectionBuffer : public FVertexBufferWithSRV
145 {
146 public:
147
148 virtual void InitSRV() override
149 {
150 // create a static vertex buffer
151 FRHIResourceCreateInfo CreateInfo;
152 void* BufferData = nullptr;
153
154 uint32 TotalNumberOfIdsToAllocate = INSTANCE_ID_INDIRECTION_BUFFER_MAX_DIRECT_INSTANCES + INSTANCE_ID_INDIRECTION_BUFFER_NUM_INDIRECT_IDS;
155
156 VertexBufferRHI = RHICreateVertexBuffer(TotalNumberOfIdsToAllocate * sizeof(uint32), BUF_Static | BUF_UnorderedAccess | BUF_ShaderResource, CreateInfo, BufferData);
157 uint32* VertexBufferData = (uint32*)BufferData;
158
159 // For the first INSTANCE_ID_INDIRECTION_BUFFER_MAX_DIRECT_INSTANCES instances, write the sequential values 0,1,2,3... That way Instance Mesh Rendering will still work without
160 // going through InstanceStaticMeshGPUScene indirection (so long as no more than INSTANCE_ID_INDIRECTION_BUFFER_MAX_DIRECT_INSTANCES are ever rendered in a single call to DrawIndexedPrimitive)
161 for (uint32 InstanceId = 0; InstanceId < INSTANCE_ID_INDIRECTION_BUFFER_MAX_DIRECT_INSTANCES; ++InstanceId)
162 {
163 VertexBufferData[InstanceId] = InstanceId;
164
165 // From INSTANCE_ID_INDIRECTION_BUFFER_MAX_DIRECT_INSTANCES to the end of the buffer, write an invalid index
166 // this is the dynamic portion of the buffer which will be allocated and filled out from the GPU as part of the InstanceStaticMeshGPUScene's indirect rendering path
167 for (uint32 InstanceId = INSTANCE_ID_INDIRECTION_BUFFER_MAX_DIRECT_INSTANCES; InstanceId < TotalNumberOfIdsToAllocate; ++InstanceId)
168 {
169 VertexBufferData[InstanceId] = 0xdeadbeef;
170
171
172
173 }
174
175 RHIUnlockVertexBuffer(VertexBufferData);
176
177 // Create a view of the buffer
178 ShaderResourceViewRHI = RHICreateShaderResourceView(VertexBufferData, sizeof(uint32), PF_R32_UINT);
179 UnorderedAccessViewRHI = RHICreateUnorderedAccessView(VertexBufferData, PF_R32_UINT);
180 }
181
182 TGlobalResource<FInstanceIndirectionBuffer> GInstanceIndirectionBuffer;
183
184 // NOTE: Really it would make more sense if this buffer were managed by FInstanceStaticMeshGPUScene, but that lives in the Renderer and this buffer needs to be available
185 // here in the Engine when the InstanceStaticMeshVertexFactory is initialized. So this semi-awkward accessor function gives the scene access to this buffer when it needs it
186 FRHIUnorderedAccessView* GetInstanceStaticMeshIndirectionUAV()
187 {
188 return GInstanceIndirectionBuffer.UnorderedAccessViewRHI;
189 }
190 // @big(mefford) - END - [SheikRender][ISOPimization]
```

Sequential IDs for Legacy Path

Indirection Values Updated Each Frame

The ISM component proxy code has a key modification in that we bind an instance indirection buffer. Regular old ISM's in UE4 will just access 0 to N instance ID's sequentially. To continue supporting this, our indirection buffer sets up a number of ID's at the front, this code path is still important to us for our grass rendering in particular. But any ISM's using the GPU scene will instead read IDs from beyond that point into the portion of the buffer that is dynamically updated by our scene update compute shader that runs each frame.

# Update Args Shader

```
207 [mthreads(NUM_THREADS, 1, 1)]
208 void updateIndirectArgs(uint GroupId, uint GroupThreadId, uint DispatchThreadId : SV_DispatchThreadID)
209 {
210     uint InstanceStaticMeshIndex = GroupId;
211
212     InstanceStaticMeshData InstanceStaticMesh = LoadInstanceStaticMesh(InstanceStaticMeshIndex, GroupThreadId);
213
214     uint NumInstances = InstanceStaticMesh.NumInstances;
215     uint NumElements = InstanceStaticMesh.NumElements;
216     uint NumLODs = InstanceStaticMesh.NumLODs;
217     uint MinLOD = InstanceStaticMesh.MinLOD;
218
219     // Set up some group shared data and sync
220     if (GroupThreadId == 0)
221     {
222         // Allocate our "space" in the output id buffer
223         // NOTE that currently we're over-allocating here because we don't know how many of each LOD we'll need.
224         // Fixing this would be fairly tricky since we don't know how many of each LOD we need until we run through the actual LOD selection/culling code below, and these indices need to be contiguous per-LOD
225         // In practice though, this buffer should be + 500 even with this sloppy allocation scheme.
226
227         uint NumInstancesToAllocate = NumInstances * NumInstances * MAX_VIEWS;
228         InterlockedAdd(InstanceIndirectionBuffer_GlobalOffset[0], NumInstancesToAllocate, InstanceIndirectionBuffer_GroupOffset);
229
230     }
231
232     if (GroupThreadId < MAX_LODS * MAX_VIEWS)
233     {
234         NumInstancesToDraw_Shared[GroupThreadId] = 0;
235     }
236
237     // Note that this group sync is for "both" the InstanceIndirectionBuffer_GroupOffset "and" the NumInstancesToDraw_Shared setup above.
238     GroupMemoryBarrierWithGroupSync();
239
240     uint GroupInstanceOffset = 0;
241     while (GroupInstanceOffset < NumInstances)
242     {
243         uint InstanceId = GroupInstanceOffset + GroupThreadId;
244
245         bool InstanceIsHidden = true;
246         if (InstanceId < NumInstances)
247             InstanceIsHidden = GetInstanceIsHidden(InstanceStaticMeshIndex, InstanceId);
248
249         if (!InstanceIsHidden)
250         {
251             //PerInstanceData PerInstanceData = GetPerInstanceData(InstanceStaticMeshIndex, InstanceId);
252
253             //ToDo: Right now, always pass in MAX_LODS and letting compiler unroll loop, but might be worth passing in NumLODs?
254             uint LODIndex = ComputeLODForInstance(PerInstanceData.WorldBoundsCenter, PerInstanceData.BoundsSphereRadius, MAX_LODS, InstanceStaticMesh.LODScreenSizes, 0, MinLOD);
255
256             // Add 1 to the number of instances to draw for the computed LOD and remember our unique draw offset so we can write our InstanceId to it
257             {
258                 uint ViewLODDrawOffset = 0;
259                 uint ViewIndex = 0;
260                 InterlockedAdd(NumInstancesToDraw_Shared[LODIndex + (ViewIndex * MAX_LODS)], 1, ViewLODDrawOffset);
261
262                 // Write to the InstanceIndirectionBuffer "which" Instance should be rendered at this LOD and InstanceId
263                 uint IndirectionBufferWriteOffset = GetInstanceIndirectionBufferOffset(InstanceIndirectionBuffer_GroupOffset, ViewIndex, LODIndex, NumLODs, NumInstances, ViewLODDrawOffset);
264                 InstanceIndirectionBuffer[IndirectionBufferWriteOffset] = InstanceId;
265             }
266         }
267     }
268 }
```

Read ISM Component Data

Reserve ID Buffer Space for ISM

Read Per-instance Data

Per-Instance LOD Selection

Write Instance ID

Let's take a look at that shader. Each ISM component gets a single group (the size of 1 wavefront) that allocates space into the ID buffer. This is done atomically as multiple groups are in flight at once. The some LDS is cleared to track the IndirectArgs and a barrier is inserted. (CLICK)

After this point we loop over the instances and perform LOD selection. The appropriate indirect arg counter in LDS is atomically incremented, and the instance ID is written to the indirection buffer. You might notice that this has similarities to the tile locations buffer from the tile classification code.

# Update Args Shader

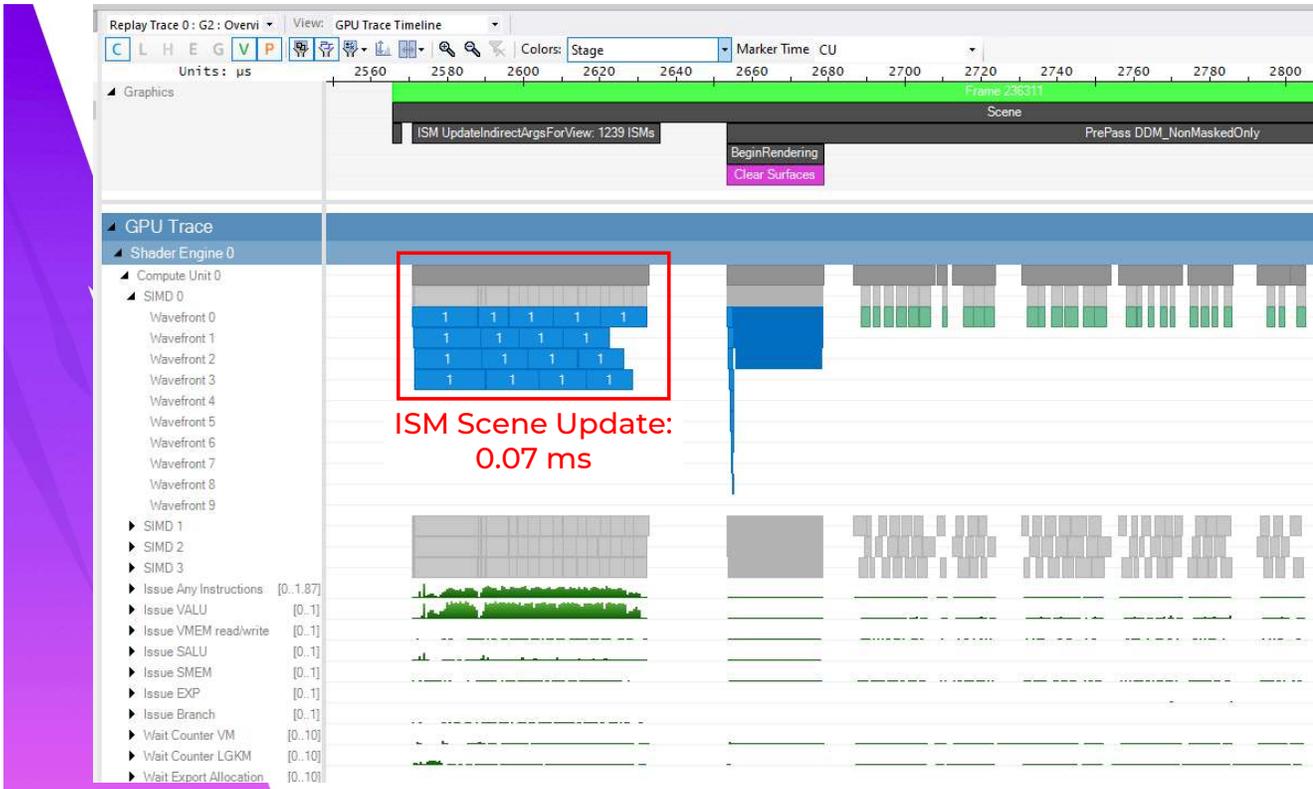
```
277     if (passedCulling)
278     {
279         uint LODRowOffset = 0;
280         uint ViewIndex = 1;
281         InterlockedAdd(numInstancesToDraw_Shared[LODIndex + (ViewIndex * MAX_LOD)], 1, LODRowOffset);
282
283         // write to the InstanceIndirectionBuffer which Instance should be rendered at this LOD and InstanceId
284         uint InstanceIndirectionBufferOffset = getInstanceIndirectionBufferOffset(InstanceIndirectionBuffer_GroupOffset, ViewIndex, LODIndex, NumLODs, NumInstances, LODRowOffset);
285         InstanceIndirectionBuffer[InstanceIndirectionBufferOffset] = InstanceId;
286     }
287
288     // call each shadow view
289     for (int ShadowIndex = 0; ShadowIndex < NUM_SHADOW_VIEWS - IGNORED_SHADOW_VIEWS; ShadowIndex++)
290     {
291         bool PassedCulling = true;
292         for (int PlaneIndex = 0; PlaneIndex < NUM_SHADOW_CULLING_PLANES; PlaneIndex++)
293         {
294             int ArrayIndex = PlaneIndex + ShadowIndex * NUM_SHADOW_CULLING_PLANES;
295             float DistanceToPlane = PerInstanceData.WorlBoundsCenter.x * ShadowViewCullingPlanes[ArrayIndex].x
296                 + PerInstanceData.WorlBoundsCenter.y * ShadowViewCullingPlanes[ArrayIndex].y
297                 + PerInstanceData.WorlBoundsCenter.z * ShadowViewCullingPlanes[ArrayIndex].z
298                 - ShadowViewCullingPlanes[ArrayIndex].w;
299             PassedCulling = PassedCulling && ((DistanceToPlane > PerInstanceData.BoundsSphereRadius));
300         }
301         if (PassedCulling)
302         {
303             uint LODRowOffset = 0;
304             uint ViewIndex = 2 + ShadowIndex;
305             InterlockedAdd(numInstancesToDraw_Shared[LODIndex + (ViewIndex * MAX_LOD)], 1, LODRowOffset);
306
307             // write to the InstanceIndirectionBuffer which Instance should be rendered at this LOD and InstanceId
308             uint InstanceIndirectionBufferOffset = getInstanceIndirectionBufferOffset(InstanceIndirectionBuffer_GroupOffset, ViewIndex, LODIndex, NumLODs, NumInstances, LODRowOffset);
309             InstanceIndirectionBuffer[InstanceIndirectionBufferOffset] = InstanceId;
310         }
311     }
312
313     groupDistanceOffset += THREAD_COUNT;
314
315 }
316
317 groupMemoryBarrierMemGroupSync();
318
319 // For each element, have a thread update their indirect args with the number of instances to draw and where in the InstanceIndirectionBuffer to read from
320 if (GroupThreadId < NumElements)
321 {
322     for (uint viewIndex = 0; viewIndex < MAX_VIEWS; viewIndex++)
323     {
324         PerElementData PerElementData = InstanceData.StaticMesh.PerElementData;
325
326         uint IndirectArgsIndex = PerElementData.IndirectArgsIndex;
327         uint LODIndex = PerElementData.LODIndex;
328
329         uint InstanceIndirectionBufferOffset = getInstanceIndirectionBufferOffset(InstanceIndirectionBuffer_GroupOffset, viewIndex, LODIndex, NumLODs, NumInstances, 0);
330
331         // each view is contiguously stored off of the base IndirectArgsIndex
332         writeIndirectArgsResults(IndirectArgsIndex + viewIndex, NumInstancesToDraw_Shared[LODIndex + (viewIndex * MAX_LOD)], InstanceIndirectionBufferOffset);
333     }
334 }
335 }
```

Main View  
Frustum Cull

Shadow Cascade  
Frustum Cull

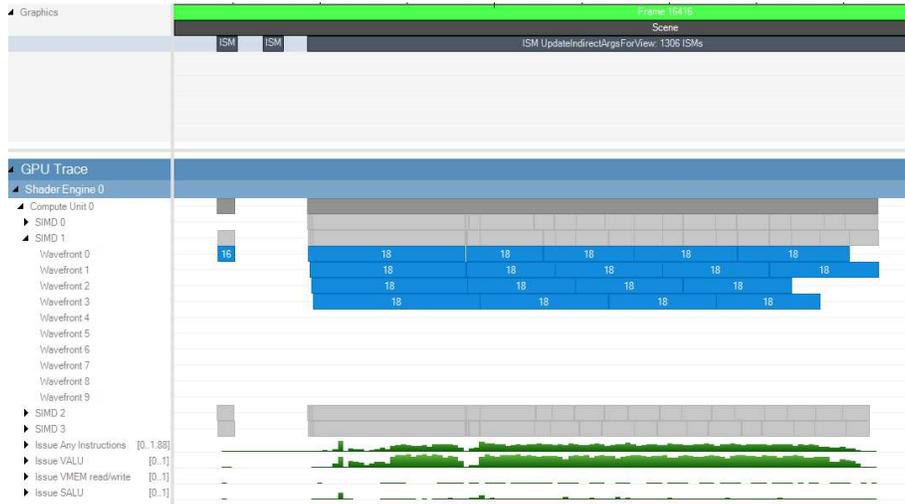
Write Indirect  
Args to DRAM

While the first arg holds all the instances for a given ISM, we also perform per-view culling for our main + shadow cascade views to further reduce wasted verts when possible. Finally we have one more barrier and write the indirect args from LDS into main memory.

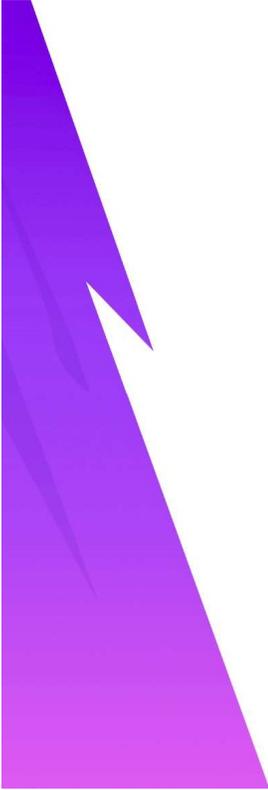


The ISM scene update shader is generally pretty fast to run, in this frame it only took 0.07 ms to run. We also have to handle scattered updates of the mesh description data when a new ISM component is added, but that is not a per-frame operation, and is also very fast, just 0.01 ms in traces I took recently.

While I'm very with the results for this system, we have ideas for further improvements – dynamic merging of loose StaticMeshActors could help improve prop rendering, and the system as it is currently only does occlusion culling on the bounds of the full structure, so some form of per-instance occlusion culling could be beneficial.



Here you can see one such trace where two small updates to the mesh descriptions are patched before the scene update shader



# Reactive Optimization: Decals

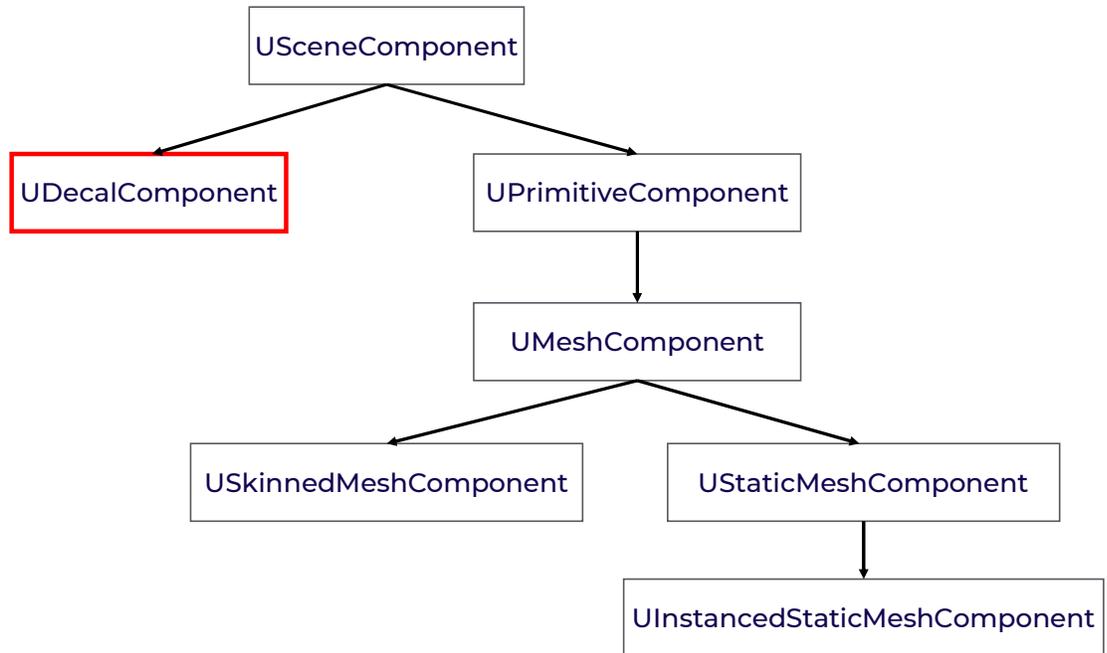
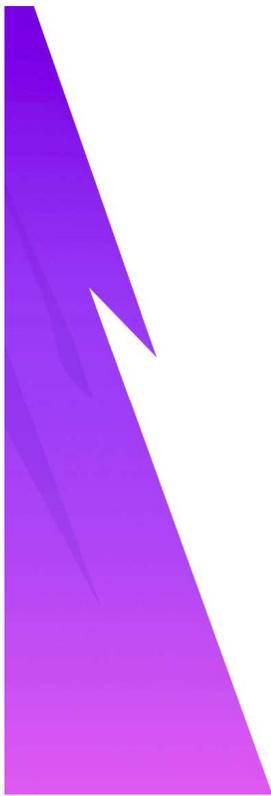
Finally, one last problem I wanted to touch on was unanticipated challenges with decals late in our development cycle. My colleague Karinne and I made a number of improvements in service of getting systems to behave as we expected them to.



The artists and I had a bit of contention in that they (rightly) felt the environment was too stiff in our original greenlight and vertical slice builds. They wanted more variety and character, which leaning heavily on instancing was not great at supporting.

I pushed for them to try to use decals more and more, and they finally did. You can see that the graffiti on the left side of this slide is overlapping multiple instanced wall meshes that are underneath it. And the cracks on the street in the lower left are breaking up the repetition in a road mesh that is generated from a spline.





I'm illustrating a partial chain of inheritance from scene component here. A scene component in unreal is just something that has a transform matrix that is placed in a level, while a primitive is a scene component that goes through the renderer.

The problem is that at some point in the history of the engine a decision was made that Decals should not be Primitives and handled separately. And likely separately from that, the Texture streaming system was engineered to operate on Primitives.

While my gut reaction was frustration that Decals were not inheriting from `UPrimitiveComponent` and that that was an oversight on the engine team's part – there *are* good reasons to not make Decals into Primitives. I already covered the efforts we went through to keep our primitive count lower with Instanced Static Mesh components. For example – decals never need to cast shadows so we don't need to waste time considering them for shadow casting like static meshes and skeletal meshes.

```

139     for (const AActor* Actor : Level->Actors)
140     {
141         if (!Actor) continue;
142
143         const bool bIsStaticActor = Actor->IsRootComponentStatic();
144
145         TInlineComponentArray<UPrimitiveComponent*> Primitives;
146         Actor->GetComponents<UPrimitiveComponent*>(Primitives);
147         for (const UPrimitiveComponent* Primitive : Primitives)
148         {
149             check(Primitive);
150             if (bIsStaticActor && Primitive->Mobility == EComponentMobility::Static)
151             {
152                 SetAsStatic(DynamicComponentManager, Primitive);
153                 UnprocessedComponents.Push(Primitive);
154             }
155             else
156             {
157                 SetAsDynamic(DynamicComponentManager, LevelContext, Primitive);
158             }
159         }
160
161         //@igs(klorig) - [IGRender] - START
162         TInlineComponentArray<UDecalComponent*> Decals;
163         Actor->GetComponents<UDecalComponent*>(Decals);
164         for (const UDecalComponent* Decal : Decals)
165         {
166             check(Decal);
167             if (bIsStaticActor && Decal->Mobility == EComponentMobility::Static)
168             {
169                 Decal->bAttachedToStreamingManagerAsStatic = true;
170                 UnprocessedDecals.Push(Decal);
171             }
172             //We don't currently handle dynamic decals at the moment.
173         }
174
175         NumStepsLeft -= (int64)FMath::Max<int32>((Primitives.Num() + Decals.Num()), 1);
176         //@igs(klorig) - [IGRender] - END
177     }
178
179     NumStepsLeft -= (int64)FMath::Max<int32>(Level->Actors.Num(), 1);
180

```

```

195
196 //@igs(klorig) - [IGRender] - START
197 //Render time logic moved from primitive component to support texture streaming for decals.
198 //
199 //
200 // The value of WorldSettings->TimeSeconds for the frame when this component was last rendered. This is written
201 // from the render thread, which is up to a frame behind the game thread, so you should allow this time to
202 // be at least a frame behind the game thread's world time before you consider the actor non-visible.
203 //
204 mutable float LastRenderTime;
205
206 /** Same as LastRenderTimeOnScreen but only updated if the component is on screen, used by the texture streamer. */
207 mutable float LastRenderTimeOnScreen;
208

```

The fix we went with is to have the streaming manager maintain a list of UDecalComponents separately from the primitives, and promoted a small amount of data up to scene component from primitive component. This allowed all of those 2kx2k graffiti pieces to drop down as low as 64x64 when they are offscreen or far away.

I wanted to call out this particular code change because I do not think this is what I would do if I was the maintainer of the engine. For example, we could avoid making Decal Components inherit from Primitive Component by creating an interface class that each implements. The reason *not* to do this is because doing so would increase our merge burden – as Epic will be fully unaware of our change unless we discuss a pull request.

I have always found myself on game teams looking for the simplest to maintain code without sacrificing performance, and often waiting for a request up to an engine team is not feasible in a lot of contexts when working on a deadline.



Now, If I had been thinking through all the consequences, I would have realized that everything I just said about Decals not being Primitives in the renderer meant that they also means that they likely don't get nearly as much attention regarding frustum culling performance on the render thread – and a little while later I started noticing that decal rendering time on the render thread was suddenly quite high as artists added more and more decals to our game. 2.2 ms in the location I showed with many decals in the surrounding area.

Each call to AddDeferredDecalPass is handling a different decal rendering stage is processing all of the decals in the scene separately and redoing frustum culling and distance fading – and it's happening serially on the main thread as well. Clearly, our artists were leaning more heavily into deferred decals than other games using the engine, which to be fair, is exactly what I had asked them to do.

```

624 // @igs(jmoore) - START - [IGRender][IGOptimization]
625 TArray<EDecalRenderStage> DecalStages;
626 TArray<FTransientDecalRenderDataList*> SortedDecals;
627 if (GUseParallelDecalCulling)
628 {
629     SCOPED_NAMED_EVENT(BuildVisibleDecalListParallel, FColor::Turquoise);
630     DecalStages.Reserve(3);
631     SortedDecals.Reserve(3);
632     DecalStages.Add(DRS_AfterBasePass);
633     SortedDecals.Add(GraphBuilder.AllocObject<FTransientDecalRenderDataList*>());
634     DecalStages.Add(DRS_BeforeLighting);
635     SortedDecals.Add(GraphBuilder.AllocObject<FTransientDecalRenderDataList*>());
636     DecalStages.Add(DRS_Emissive);
637     SortedDecals.Add(GraphBuilder.AllocObject<FTransientDecalRenderDataList*>());
638 }
639
640 const FScene& Scene = *(FScene*)ViewFamily.Scene;
641 const EShaderPlatform ShaderPlatform = View.GetShaderPlatform();
642
643 FDecalRendering::BuildVisibleDecalListParallel(Scene, View, DecalStages, SortedDecals);
644 }
645 // @igs(jmoore) - END - [IGRender][IGOptimization]
646
647 if (ViewFamily.EngineShowFlags.Decals && !ViewFamily.EngineShowFlags.ShaderComplexity)
648 {
649     // @igs(jmoore) - START - [IGRender][IGOptimization]
650     AddDeferredDecalPass(GraphBuilder, View, DecalPassTextures, DRS_AfterBasePass, GUseParallelDecalCulling ? SortedDecals[0] : nullptr);
651     // @igs(jmoore) - END - [IGRender][IGOptimization]
652 }
653
654 if (bDoDecal && !IsUsingGBuffers(View.GetShaderPlatform()))
655 {
656     // decals are before AmbientOcclusion so the decal can output a normal that AO is affected by
657     // @igs(jmoore) - START - [IGRender][IGOptimization]
658     AddDeferredDecalPass(GraphBuilder, View, DecalPassTextures, DRS_BeforeLighting, GUseParallelDecalCulling ? SortedDecals[1] : nullptr);
659     // @igs(jmoore) - END - [IGRender][IGOptimization]
660 }
661
662 if (bDoDecal && !IsSimpleForwardShadingEnabled(View.GetShaderPlatform()))
663 {
664     // DBuffer decals with emissive component
665     // @igs(jmoore) - START - [IGRender][IGOptimization]
666     AddDeferredDecalPass(GraphBuilder, View, DecalPassTextures, DRS_Emissive, GUseParallelDecalCulling ? SortedDecals[2] : nullptr);
667     // @igs(jmoore) - END - [IGRender][IGOptimization]
668 }
669

```

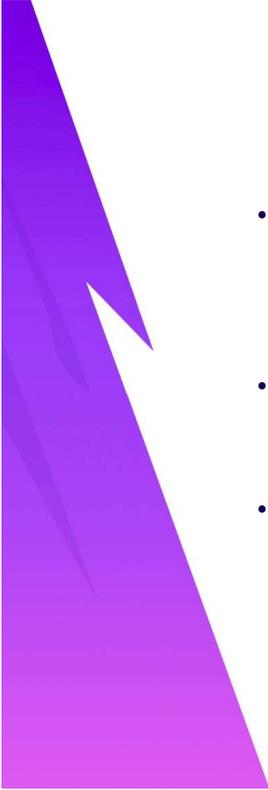
Parallel frustum cull

So the solution here is pretty simple – I build the three decal lists with a single frustum cull and distance fade on the task graph with a parallel For loop (very similar to primitive frustum culling in InitViews) and then I handle copying each decal that passes culling into the relevant list afterwards.

The resulting list is passed down into the different decal passes and the original per-pass processing is skipped.



In the area I showed this reduced processing down to 0.67 ms – a 1.53 ms improvement in render thread time.



# Final Thoughts

- Every game has optimization potential tied to the game you are making that any shared purpose engine may not anticipate
- Any technology you lean into will lead to finding its limits when you push it hard enough
- Search for solutions that fit the scope of your team

Alright so just a few thoughts to end on. I just want to reiterate that I believe you should always be thinking about how the engineers that you are downstream from are not necessarily going to anticipate the content you are making.

And related to that, any time you push into new territory and scope as a team – you are likely going to hit the limits of systems. A lot of what I showed improvements for today worked really well for a long time in development, but eventually became untenable when the content reached a critical mass.

And related to that – always search for the solutions that fit the scope of your team. Part of why I spend so much time figuring out engineering solution is because we don't have an army of artists to try to rework content causing performance problems. And I very much appreciate when they are able to help us solve problems on the content side – we're all working together to try to make a great game in the end.



# References

- [Deferred Lighting in \*Uncharted 4\*](#), Ramy El Garawany, Siggraph 2016
- [GPU-Driven Rendering Pipelines](#), Ulrich Haar, Sebastian Aaltonen, Siggraph 2015
- [Separable Subsurface Scattering and Photorealistic Eyes Rendering](#), Jorge Jimenez, Siggraph 2016
- [The Devil is in the Details: idTech666](#), Tiago Sousa, Jean Geffroy, Siggraph 2016
- [Dynamic Occlusion With Signed Distance Fields](#), Daniel Wright, Siggraph 2015



# Many Thanks

- Andreas Frederickson – the advisor for this talk
- Everyone that contributed to Rumbleverse rendering: especially Karinne Lorig, Nate Mefford, Rusty Swain, David Laskey, and our partners at Dragon's Lake
- The entire Rumbleverse development team – especially our wonderful artists and the team leadership
- My little family: Kelsey, Spaceman, and Nova



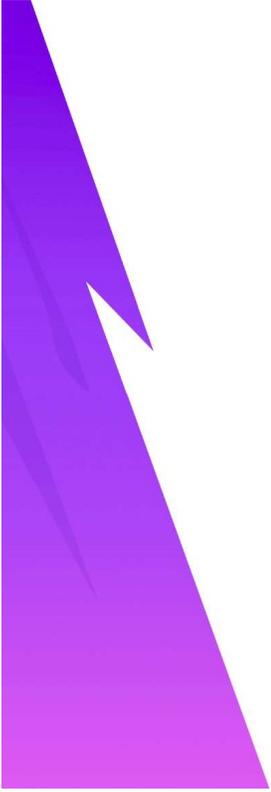
# Q&A

E-mail: [Jon@IronGalaxyStudios.com](mailto:Jon@IronGalaxyStudios.com)  
Twitter: [@JonManatee](https://twitter.com/JonManatee)

Slides:  
[jonmanatee.com/s/GDC2023.pdf](https://jonmanatee.com/s/GDC2023.pdf)

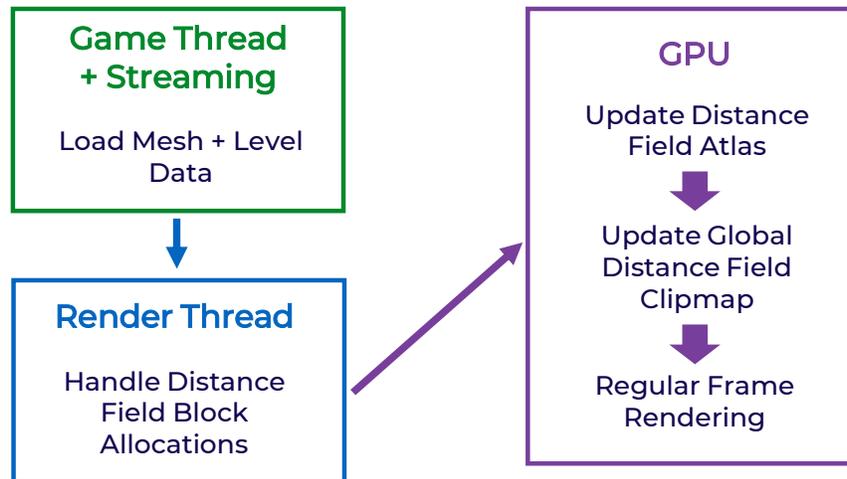
## RUMBLEVERSE





# Appendix 1: Distance Field GPU Optimizations

# Distance Field Upload



Let's go over and consider the work happening on the GPU. As a reminder we need to copy the distance field texels into the atlas texture at the locations decided by the allocator, and then we need to update dirty regions of the global distance field clipmaps

# Distance Field Atlas Barriers

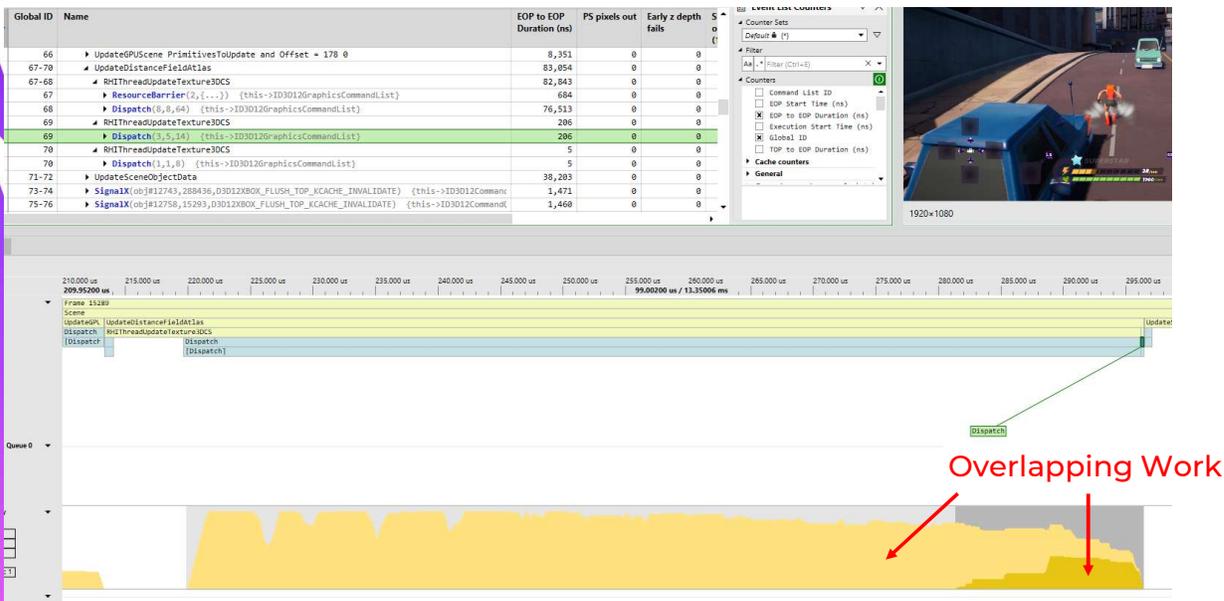


The first up is uploads into the distance field atlas CS. The consoles are using a compute shader path to update these – and it is inserting a barrier between each copy. Again, this is especially bad because the artists are motivated to keep the distance field memory on each object as small as possible so they can pack more into the atlas. This means that the ratio between barrier waits and actual work on the CUs gets worse as the artists do a better and better job.

These barriers would be correct to ensure deterministic results if the regions being updated within the volume overlapped, but given that there is no need to do that as this is a texture atlas where everything gets its own spot in the buffer, it is safe for these wavefronts to be executing at the same time. None of these are going to be writing to the same places in memory.

There are four updates happening in this frame in this capture I'm showing here, but there could be dozens in some frames.

# Distance Field Atlas Barriers



Artem from our co-dev partner Dragon's Lake got these barriers eliminated by implementing an interface in the RHI layer that hints that there will be multiple safe updates to the same resource and only places a barrier at the beginning + end, and then we put in place a per-frame limit of  $128^3$  texels worth of updates per-frame to keep things from getting unmanageable from just the total amount of bandwidth needed. This keeps our time spent updating the atlas to under 0.1 ms on any given frame.

In this updated capture you can see that three distance fields are being updated and two of the barriers have been eliminated, keeping the GPU fed during these updates.



# Global DF Update - Full



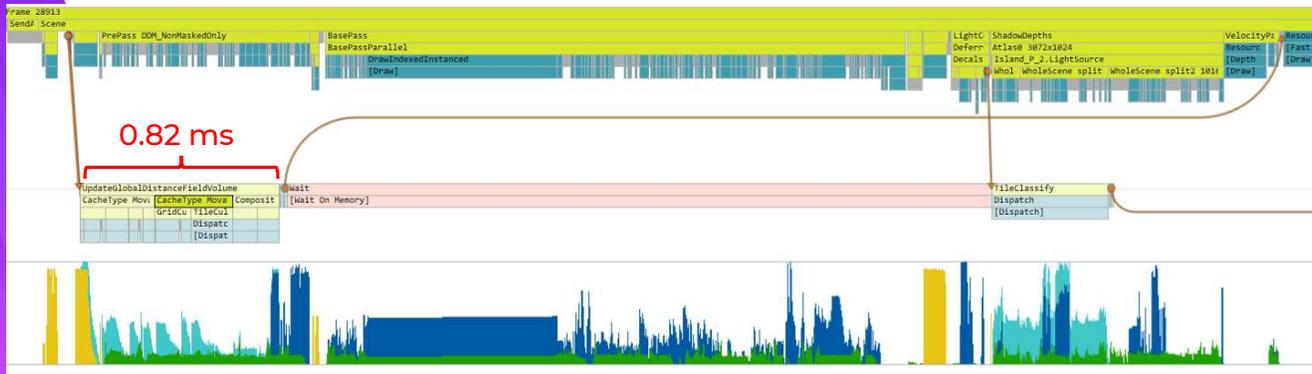
Surprisingly, we were hitting full clipmap updates, which can be quite slow – but not from my 1000 pending failsafe triggering.

It was happening from single large objects with distance fields that would invalidate the entire clip map at once. Throttling but update region count would never save us here and I was prepared to start implementing logic to slice up these regions into smaller sub sections over multiple frames, but I decided to look at what these objects were first.

It was coming from the fact that the artists liked to enable a distance field on LOD1 of many buildings levels so that the first LOD still casts a crude shadow, making the transition from the full LOD0 shadows less noticeable but still dropping down to 0 distance field memory at LOD2. These DFs that were the size of an entire buildings could easily trigger a full update of one of the clipmaps when running through the city. I'm showing that here – each of these clip maps takes >2 ms to update and we have 3 of them. We really want to avoid that case.

However- we really don't care about those LOD1 distance fields, and a really simple hack was put in place to eliminate that from happening. Which is that we simply don't process distance fields for the LOD1 meshes on the global clipmap. We know that the LOD0 meshes for that level are going to do a good job dirtying and updating the areas where buildings are placed anyways both on add and removal. The artists are generating the LOD1 distance field for the direct shadow tracing which is only traced against the distance field atlas and not the global clip map.

# Global DF Update - Async



And for good measure we moved this work to async compute – we don't overlap postprocessing with the next frame with async compute like some games do. So our early frame work and depth prepass could use some async compute and this is a perfect candidate. A normal frame is going to have to scroll the clip maps as the player moves and that takes ~0.5 ms on Xbox One, which stretches to 0.8 ms when run async and stays nicely overlapped with the vertex shading work in our prepass. When additional work is needed this can spill past the prepass, but overlapping with the base pass is generally not overly problematic for us and we still see an overall reduction in spikes.

The results of the distance field updates are not needed until after our velocity rendering completes, but we do have additional async work kicking off to overlap with shadow maps. Throttling of the global clipmap updates generally prevents that from ever running that deep into the frame.



## Appendix 2: UE4 Cook Conversion Static Mesh -> Instanced Static Mesh

# Cook-time Conversion

```
void AShedBasicStructure::ConvertToInstanceStaticMeshes()
{
    UE_LOG(LogShedkStructure, Log, TEXT("Converting Building Xs To ISM Components."), *GetFullName());
    // Get all the individual static mesh components in this building
    TArray<StaticMeshComponent> StaticMeshComponents;
    GetComponents<StaticMeshComponent>(StaticMeshComponents);
    // Map InstanceableMeshComponent to InstanceableMeshComponents;
    // For (UStaticMeshComponent* StaticMeshComponent : StaticMeshComponents)
    {
        // For now, ignore ISMCs. Eventually should possibly add capability to merge existing ISMCs with matching individual SMCs
        if (StaticMeshComponent->IsA(UInstanceableStaticMeshComponent::StaticClass()))
        {
            continue;
        }
        // Do not include hidden meshes
        if (!StaticMeshComponent->IsVisible())
        {
            continue;
        }
        TArray<StaticMeshComponent*> InstanceableMeshList = InstanceableMeshComponent->FindByAdd(StaticMeshComponent);
        InstanceableMeshList.Add(StaticMeshComponent);
    }
    UE_LOG(LogShedkStructure, Log, TEXT(" Contains Xsd Total Static Mesh Components, Converting to Xsd Instanceable Static Mesh Components."), StaticMeshComponents.Num(), InstanceableMeshComponents.Num());
    for (FMap<FInstanceableMeshComponent, TArray<UStaticMeshComponent*>>::Element UniqueMeshID : InstanceableMeshComponents)
    {
        FInstanceableMeshComponent InstanceableMeshData = UniqueMeshID.Key;
        TArray<UStaticMeshComponent*> InstanceableMeshComponentList = UniqueMeshID.Value;
        UE_LOG(LogShedkStructure, Verbose, TEXT(" Instance Count: Xsd : Static Mesh: Xs"), InstanceableMeshComponentList.Num(), *UniqueMeshID.Key.StaticMesh->GetName());
        // Guard against StaticMeshComponents that do not have a valid StaticMesh
        // @todo: Should we just remove these component though?
        if (InstanceableMeshData.StaticMesh == nullptr)
        {
            UE_LOG(LogShedkStructure, Warning, TEXT("Warning: Shedk Structure Xs Includes Static Mesh components that reference a null Static Mesh!"), *GetFullName());
            continue;
        }
        // Guard against StaticMeshes that have more LODs than the InstanceableStaticMeshScene supports
        // @note: An alternative would be to silently ignore LODs beyond FInstanceableStaticMeshSceneData::MaxLODs, and always just 'clamp' the last supported LOD...
        if (InstanceableMeshData.StaticMesh->GetNumLODs() > FInstanceableStaticMeshSceneData::MaxLODs)
        {
            UE_LOG(LogShedkStructure, Warning, TEXT("Warning: Shedk Structure Xs References Static Mesh Xs which contains Xd LODs. Meshes Must contain <= Xd LODs to use optimized GPU Scene rendering"), *GetFullName());
            *InstanceableMeshData.StaticMesh->GetName(),
            InstanceableMeshData.StaticMesh->GetNumLODs(),
            FInstanceableStaticMeshSceneData::MaxLODs);
            // This static mesh is not eligible for instancing because of the number of LODs, move to the next possible group of mesh components
            continue;
        }
        // InstanceableStaticMeshScene only supports InstanceableStaticMeshComponents with <= FInstanceableStaticMeshSceneData::MaxInstanceableInstances
        // In the situation where we have more than that many instance-able static meshes in a building, create multiple InstanceableStaticMeshComponents
        // Each with no more than FInstanceableStaticMeshSceneData::MaxInstanceableInstances
        uint32 ComponentToInstanceIndex = 0;
        uint32 NumRemainingComponentsToInstance = InstanceableMeshComponentList.Num();
    }
}
```

Skip Incompatible Components

Map Component to Static Mesh

Skip Incompatible Meshes

Here is the code in the basic structure class that converts the individual mesh components to ISM components. First we iterate all of the components and build a unique mapping to source meshes. We must filter out some incompatible meshes at this time as well. Then we iterate each unique mesh.

