

# Эффективное использование GPU на примере геймдева

**Ростислав  
Михеев**

Elverils

 @hacenator

 [rmikheev@elverils.com](mailto:rmikheev@elverils.com)

 **C++ Russia**  
2023



О нас



# Эффективное GPU

Использование ресурсов:

- памяти
- времени исполнения

# Архитектура GPU



CPU



GPU

RTX 4090

16384 cuda cores

- GPU это монстр с тысячами ядер
- Single Instruction Multiple Threads (SIMT)
- Различные типы памяти (registers, shared memory, local memory, L1 & L2 кеши, global memory)
- Вспомогательные блоки и ядра (Texture Processing, RT cores, Tensor cores, Rasterization)
- GPC <- TPC <- SM <- CUDA cores

# Инструменты анализ производительности

- Рантайм статистика: fps, ms на весь/часть кадра
- Снимок кадра: последовательность команд
- Анализ синхронизации GPU/CPU: выявление простоя (дыр) на GPU вследствие неправильного передачи команд (partial submit, threshold); анализ readback
- Профилирование вызовов отрисовки и шейдеров

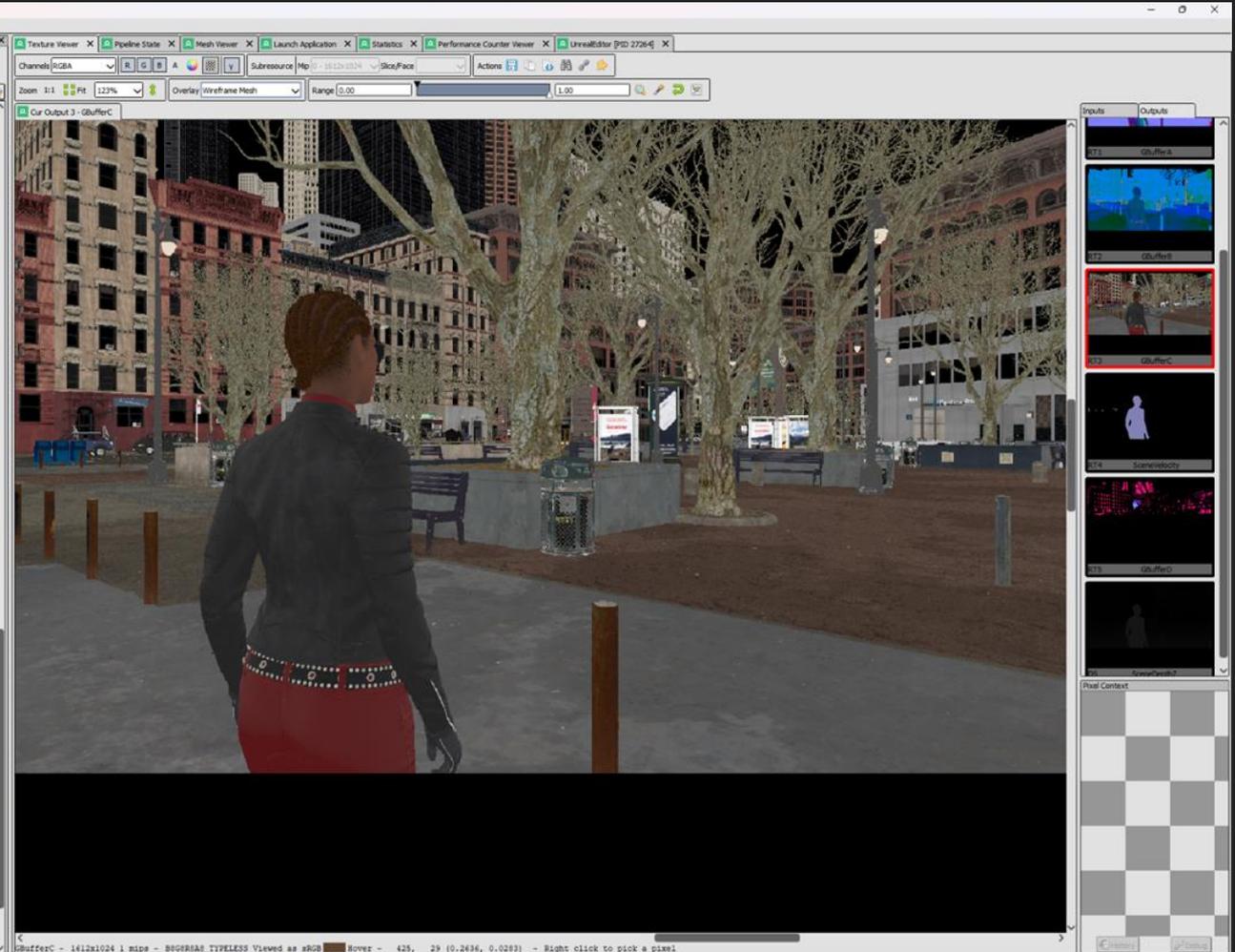
PC RUNNING WITH AI LOGGING ON!  
PC RUNNING WITH OC VERIFY ON!

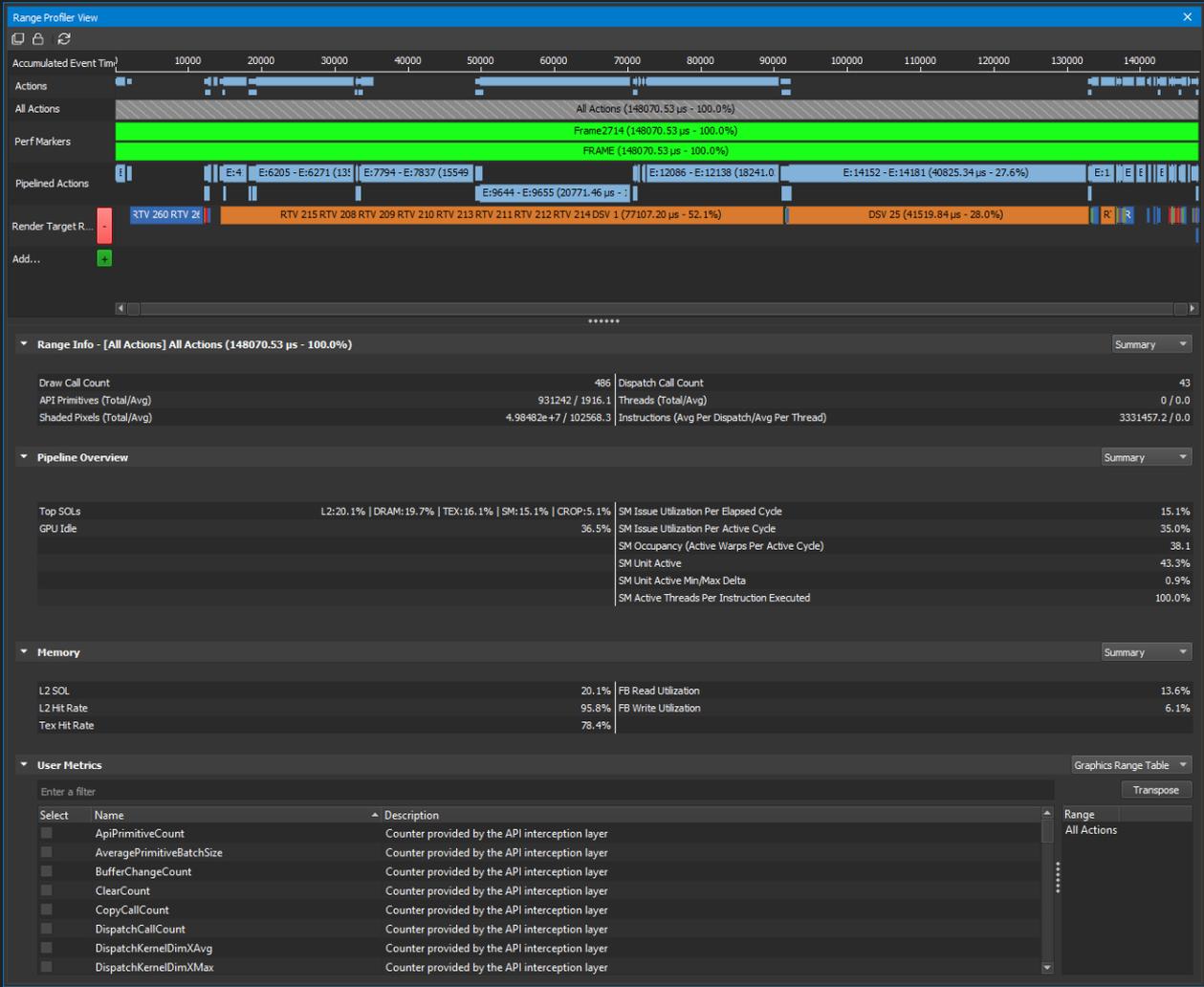
[TOTAL]	24.23	28.92	22.18
LumenScreenProbeGather	2.89	3.66	2.82
[unaccounted]	3.89	7.14	2.46
Nanite_BasePass	2.51	2.62	2.37
Shadow_Depth	2.63	4.19	2.10
Nanite_VisBuffer	1.70	1.78	1.65
LumenReflections	1.24	1.26	1.22
TemporalSuperResolution	1.15	1.91	1.13
Prepass	1.45	2.79	0.91
LumenSceneLighting	0.87	0.93	0.81
FrameRenderFinish	0.01	0.02	0.01
RayTracingScene	0.68	0.72	0.65
Basepass	0.75	1.00	0.70
Slate_UI	0.53	1.26	0.49
Postprocessing	0.50	0.52	0.49
Nanite_Readback	0.52	0.52	0.51
LumenSceneUpdate	0.29	0.89	0.18
Shadow_Projection	0.27	0.28	0.26
Translucency	0.32	0.39	0.27
RenderDeferredLighting	0.25	0.26	0.25
GPU_SkinCache	0.35	0.37	0.33
SubsurfaceScattering	0.16	0.18	0.15
MotionBlur	0.12	0.13	0.11
Lights	0.12	0.13	0.12
L-Shift Translucent_Lighting	0.11	0.12	0.11
BeginOcclusionTests	0.10	0.11	0.09
SkinnedGeometryUpdateBLAS	0.12	0.13	0.11
GPU_SceneUpdate	0.07	0.08	0.07
Render_Velocities	0.09	0.10	0.08
RayTracingGeometry	0.10	0.28	0.09
Niagara_GPU_Simulation	0.06	0.06	0.05
Niagara_GPU_Compute_All_Free_IDS	0.06	0.07	0.06
Composition_BeforeBasePass	0.05	0.06	0.04
HZB	0.05	0.05	0.05
PostRenderOpsFX	0.04	0.06	0.03
Hair_CardsInterpolation	0.04	0.05	0.04
SortLights	0.03	0.03	0.02
VisibilityCommands	0.02	0.05	0.01
Editor_Primitives			
UpdateLumenSceneBuffers	0.02	0.02	0.01
Morphi_Target_Compute	0.03	0.03	0.03
Fog	0.01	0.02	0.01
VirtualTextureUpdate	0.01	0.02	0.01
VirtualTexture	0.00	0.01	0.01
SkinnedGeometryBuildBLAS			
Nanite_Streaming			
DepthOfField			
Distortion	0.03	0.78	0.01
TAA			
Hair StrandsInterpolation	0.01	0.01	0.01
Hair StrandsGuideDeform			
Hair GuideInterpolation			
FXSystemPreRender	0.00	0.00	0.00

[10 more stats. Use the stats.MaxPerGroup CVar to increase the limit]

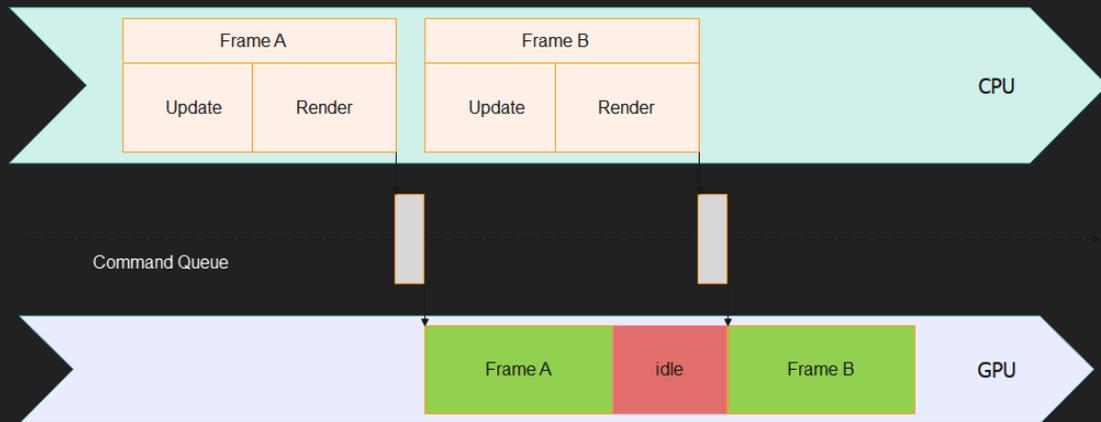
23.64 FPS  
42.30 ms  
Frame: 40.85 ms  
Game: 40.96 ms  
Draw: 35.70 ms  
GPU: 19.94 ms  
RHIT: 26.65 ms  
DynRes: Unsupported  
Draws: 5885  
Prims: 706.1K

ED	Name	Duration (ns)
728-33592	Scene	47.54368
750-738	VirtualTextureClear	0.00307
742-745	Nagara::GPUOffler_BeginFrame	0.00
740-740	Nagara::ExecuteTickFrame	0.00
752-739	FPUSystem::PreRender	0.00
743-754	FGPUStateManager::DrawBender	0.00
743-846	GPUScene::Update	0.02763
852-886	GPUScene::UploadDynamicPrimitiveShaderDataForView	0.01128
889-925	BulkRenderingCommandsDeferred(Culling-On)	0.02468
938-966	VirtualTexturePageTabletUpdates	0.00814
969-1027	HorzGridAndZorPosition	0.02927
1030-1055	UpdateDistanceFieldSettings	0.00838
1058-1112	UpdateDistanceFieldAtlas	0.11245
1115-1402	UpdateGlobalDistanceField	0.76662
1485-1481	> AccessModelAsGraphics (Textures: 7, Buffers: 15)	0.00
1494-1497	> ClearDepthStencil (SceneDepth)	0.00922
1500	> ExecuteCommands(2)(1): Close (Baked Command List 124888)	
1501	> ExecuteCommands(2)(2): Reset (Baked Command List 124888)	
1504	> ExecuteCommands(2)(2): Close (Baked Command List 124888)	
1509	> ExecuteCommands(2)(2): Reset (Baked Command List 124890)	
1513-2225	PrePass DDM: AllPasses (Forced by Nanite)	0.61728
2226-2244	> ResetOpenGLHistoryForCarCaptureAtlas	0.00922
2245-2263	UpdateRenderPipeline	0.00783
2271-2142	Nanite Visibility	1.49453
3145-3206	ComputeLightGrid	0.0256
3209-4125	> BeginOcclusionTests	0.37786
4126-4188	> BulkDraw(ViewId=0)	0.04492
4191-4193	Submit Commands	0.00
4196-5369	> ForceOcclusionTests	0.00
5172-5378	LumenSceneUpdate: 37 card captures 0.0294 levels	0.39787
5781-7003	Composition for Reflections	0.7639
7006-7014	> GPUBufferClear	0.01126
7017-7019	> GPUBufferClear	0.00
7015-20139	> GPUBufferClear	27.13107
20115-20139	> GPUBufferClear	0.00
20132-20138	CopyStencilLightingChannels	0.01743
20141-20722	FinalShadowMapArray: BulkPageAllocation	0.34243
20725-31095	Shadow Depths	2.09101
31098-31408	LumenSceneLighting	1.07347
31403-32140	> DispatchTable	0.00
31740-31741	> Nanite Readback	0.01338
31785-51813	> Readback	0.00205
31816-31826	LightCompositionTask: PreLighting	0.00
31829-31832	> ClearStencil (SceneDepth)	0.00512
31835-32496	DiffuseIndirectAO	6.83594
32499-33133	> CreateTranslucencyLightingVolumeCompute 64	0.01028
33111-33626	Light	0.49562
33631-33669	FilterTranslucentVolume (4x4x4x4 Cascades: 2)	0.10593
33672-32908	SubsurfaceScattering (ViewId=0)	0.20838
33911-33925	ExponentiateHeightMap	0.01126
33938-33941	SetSceneTexturesInVertexBuffer	0.00
33944-33950	FPUSystem::PostRenderOpaque	0.00
33953-33955	AsyncGPUTracerMember::PostRenderOpaque	0.00
33958-33959	Nagara::PostRenderFinish	0.00
33962-33963	UnsetSceneTexturesInVertexBuffer	0.00
33966-33967	FGPUStateManager::DrawPostRenderOpaque	0.00
33970	> ExecuteCommands(2)(1): Close (Baked Command List 124982)	
33971	> ExecuteCommands(2)(2): Reset (Baked Command List 124988)	
33974	> ExecuteCommands(2)(2): Close (Baked Command List 124988)	
33986	> ExecuteCommands(2)(1): Reset (Baked Command List 124995)	
33990	> ExecuteCommands(2)(1): Close (Baked Command List 124995)	
33993	> ExecuteCommands(2)(1): Reset (Baked Command List 124994)	
33996-34004	Final Passes	1.16624
34005-35172	> VirtualTextureFinalBackCopy	0.81024
35175-35685	PostProcessing	2.32723
35688-35690	ReleaseRayTracingResources	0.00
35693-35700	Render Finish	0.00
35703-35704	FinalPassInVertexBuffer	0.00



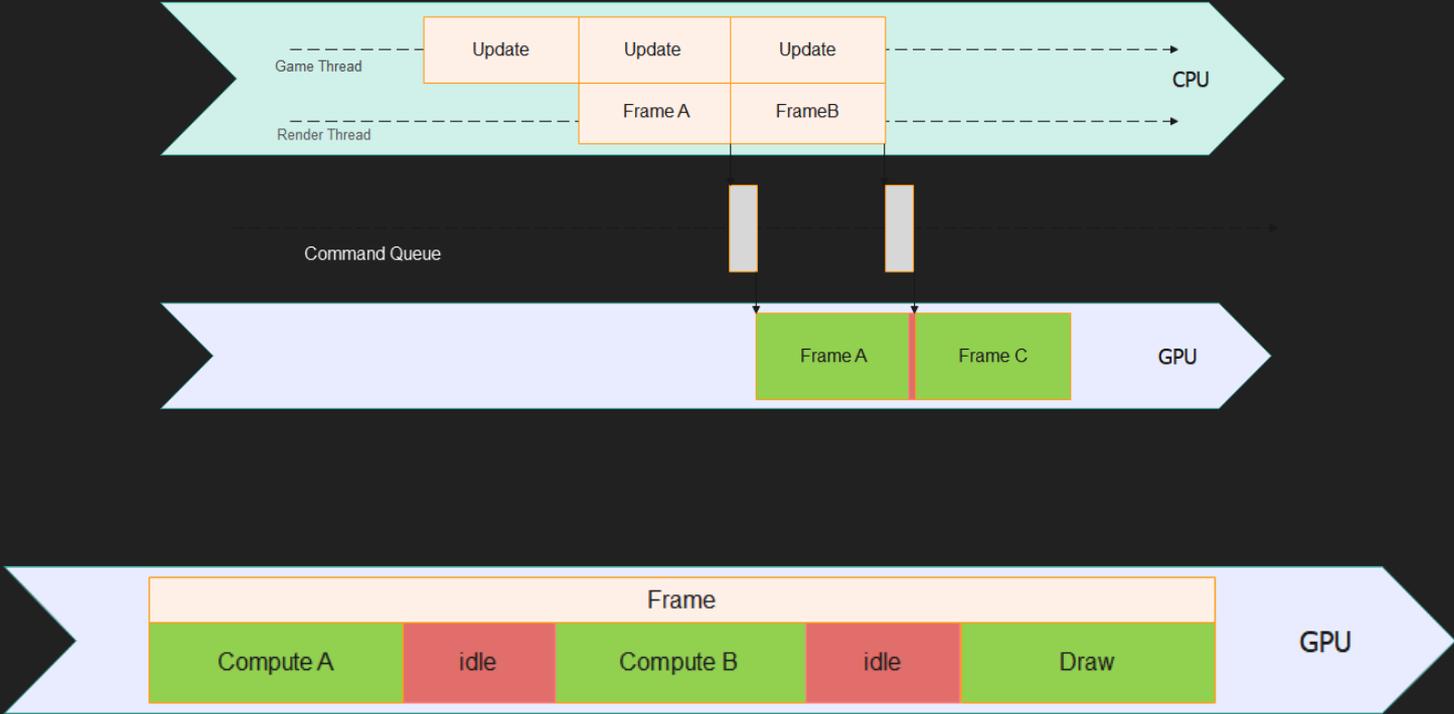


# GPU idles



- Render Thread
- Избегать readback в том же кадре

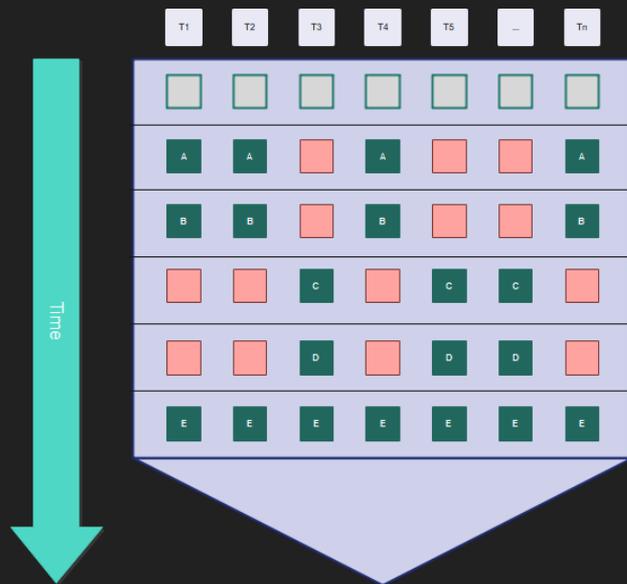
# GPU idles



# Ветвление на GPU

Неправильное или чрезмерное использование условий внутри шейдеров приводит к снижению производительности

```
if (rand() % 2) {  
    A();  
    B();  
} else {  
    C();  
    D();  
}  
E();
```



# Оптимизация ветвления

- Выносить общие куски за пределы блоков условий
- Заменять ветвление на арифметические операции (умножение, saturate, step)

```
float x = rand() % 2;  
A()*x+B()*x+C()*!x+D()*!x;  
E();
```

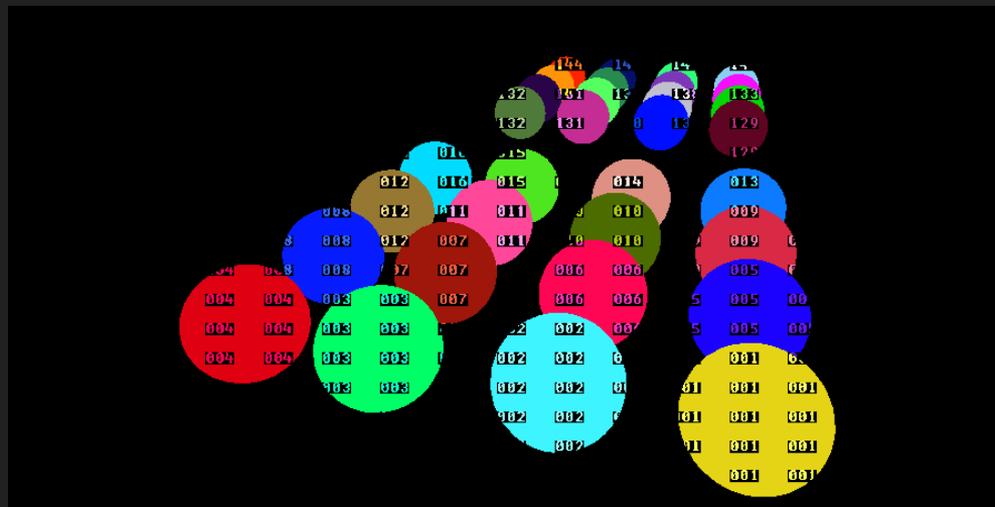
# Оптимизация ветвления

- Статическое бранчевание
- Использовать permutation через дефайны

```
#if BRANCH
    A();
    B();
#else
    C();
    D();
#endif
E();
```

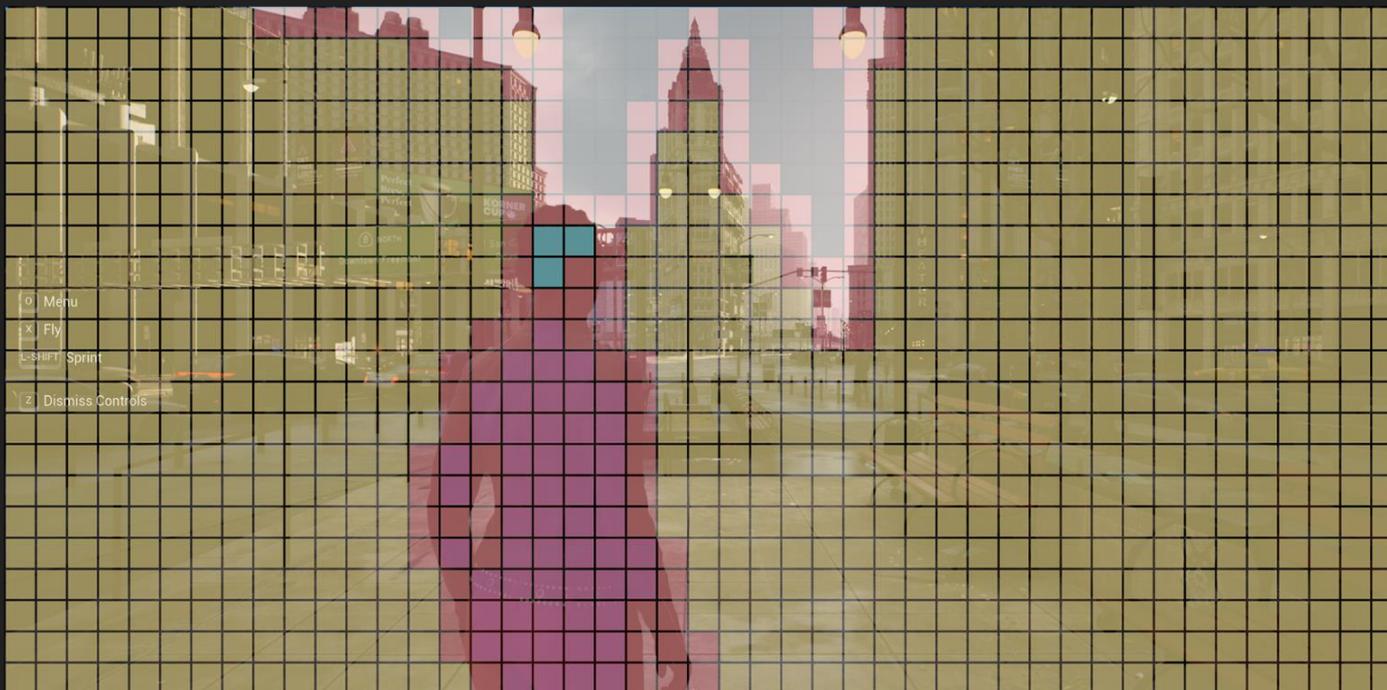
# Оптимизация ветвления

Использовать маску (через depth, stencil) для фрагментных шейдеров



# Потайловый компьютер

Потайловая проходка через ExecuteIndirect() для compute шейдеров



# Циклы

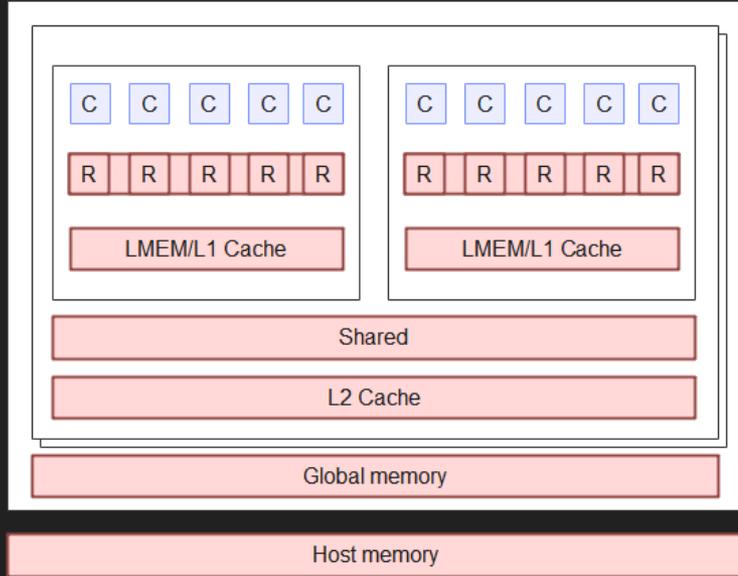
[unroll] / [loop]

```
uint count = 4;
float w = 0.0f;

[unroll]
for (uint i = 0; i < count && w <= 1.0f; ++i)
{
    A();
    w += B();
}
```



# Память GPU



Название	Объем	Область видимости	Скорость чтения и записи
Registers	Маленький	тред	высокая или очень высокая
Shared	Маленький	группа тредов	высокая
Local Memory	Средний	тред	средняя
Global memory	Большой	группа тредов	низкая или очень низкая

# Функции для Wave

- Wave Query - Общие операции для контроля (IsFirstLane / GetLaneCount / GetLaneIndex)
- Wave Vote - Операции сравнения значений между потоками (ActiveBallot / AllTrue / AnyTrue)
- Wave Broadcast - Получение значения из других тредов (ReadLaneFirst / ReadLaneAt )
- Wave Reduction - Арифметические и логические операции для активных тредов (ActiveMax / ActiveSum / BitOr и т.д.)

# Memory barriers

В условиях SIMT архитектуры GPU для правильной организации вычислений нам нужны синхронизации.

Барьеры в основном различаются по области памяти которую они синхронизируют - shared, device, all, none.

В некоторых случаях простая перестановка барьеров внутри кода шейдера может влиять на производительность.

# Пример использования

```
Shared mem:  groupshared float SharedMaxRadius[SIZE];

SharedMaxRadius[GroupThreadIndex] = MaxRadius;
GroupMemoryBarrierWithGroupSync();

[unroll]
for (uint i = 0; i < 5; i++)
{
    const uint ReduceSize = 32 >> i;
    if (GroupThreadIndex < ReduceSize)
    {
        MaxRadius = max(MaxRadius, SharedMaxRadius[GroupThreadIndex + ReduceSize]);
        SharedMaxRadius[GroupThreadIndex] = MaxRadius;
    }
}

TileMaxRadius = SharedMaxRadius[0];
```

# Пример использования

Atomic:

```
groupshared uint SharedMaxRadius;  
  
uint Unused;  
  
InterlockedMax(SharedMaxRadius, asuint(MaxRadius), Unused);  
GroupMemoryBarrierWithGroupSync();  
TileMaxRadius = asfloat(SharedMaxRadius);
```

Wave op:

```
TileMaxRadius = WaveActiveMax(MaxRadius);
```

# Упаковка

В ряде случаев точности float16 достаточно для проведения расчетов. Поэтому важно помнить о возможности упаковки данных.

- Bit-packing (f16tof32 <-> f32tof16) Можно запихивать в 2 раза больше параметров. Есть векторизация для float16 - т.е. некоторое железо может объединять две float16 операции в одной инструкции
- Текстуры - компоненты, битность. YCoCg color space
- Сферические гармоники - карты окружения в частности для лайт проб
- hardware компрессии - есть при растеризация. Поэтому иногда для полноэкранный отрисовки пиксельный шейдер будет предпочтительнее чем компьютер (ROP)

## Дополнительно

- fast-math оптимизации
- Избегать overdraw и контролировать blending
- Upscale/Downscale, Dynamic resolution, Checkerboard
- EarlyZ
- Threadgroup swizzling

# Автоматизация

Как и всегда частая проблема, что в активной фазе разработке изменения могут иметь деструктивный эффект, поэтому очень важно собирать и хранить статистику по билдам.

Нет возможности проверять билд целиком с каждой оптимизацией.

- скриншоты выбранных мест;
- сбор статистики в таблицы (DC, primitives count, FPS, ms);
- автоматизация проходов - пролеты камеры (можно использовать navmesh).
- анализ использования памяти

# ЧТИВО

- [nvidia-ada-gpu-architecture.pdf](#)
- [GDC2017 Wave Operations in DX SM6 and Vulkan \(gpuopen.com\)](#)
- [Optimizing Compute Shaders for L2 Locality using Thread-Group ID Swizzling | NVIDIA Technical Blog](#)