

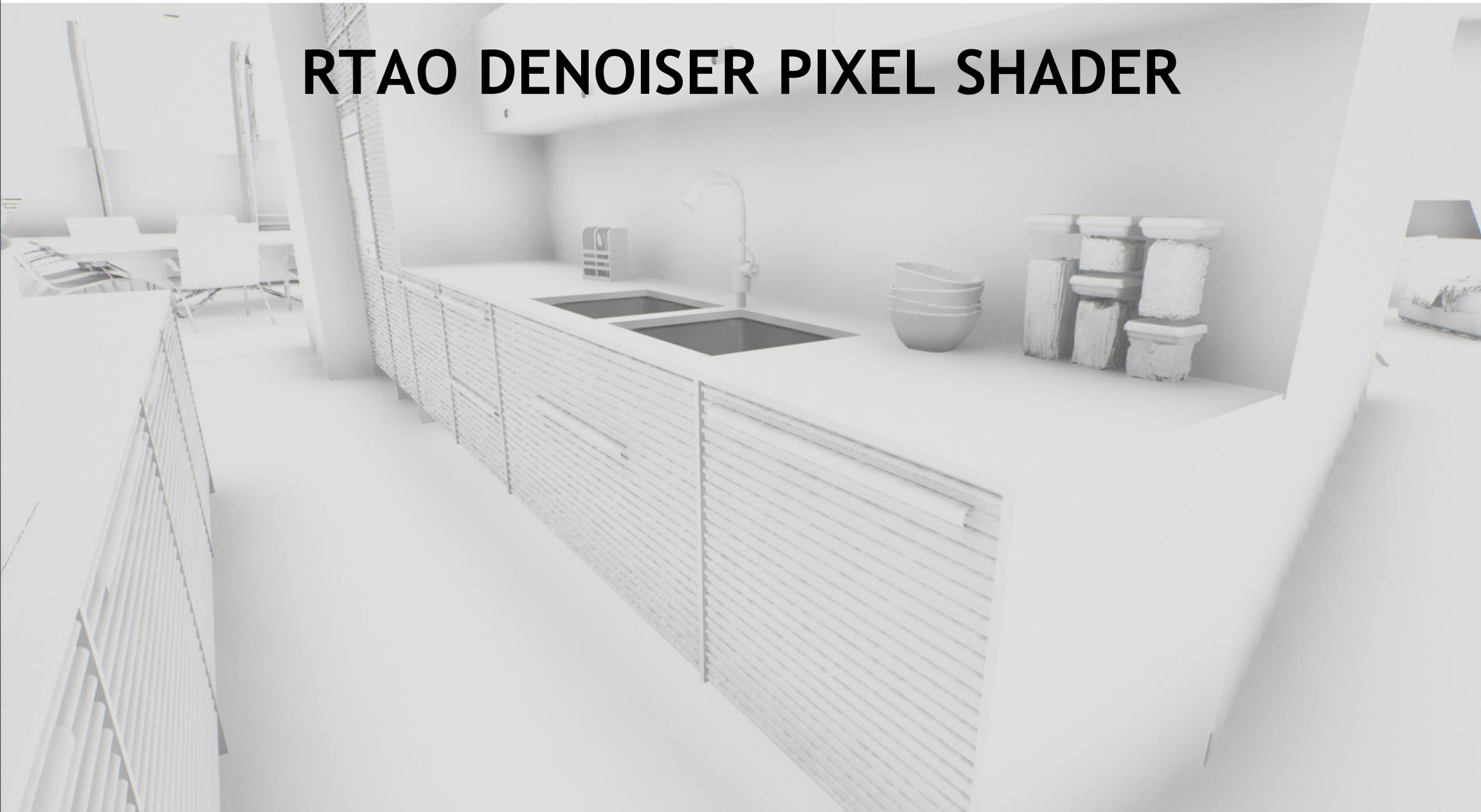


OPTIMIZING DX12/DXR GPU WORKLOADS

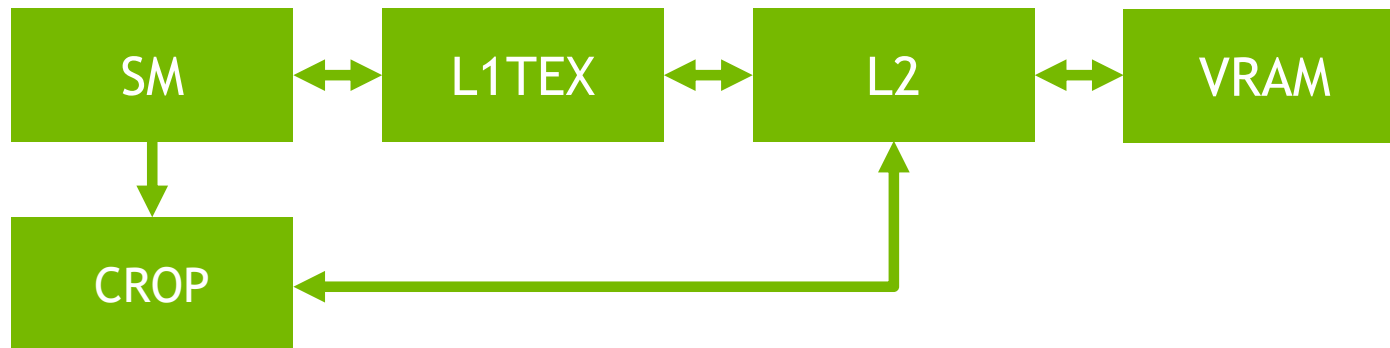
USING **NSIGHT GRAPHICS: GPU TRACE** AND
THE PEAK-PERFORMANCE-PERCENTAGE (P3) METHOD

Louis Bavoil | Principal Engineer, Developer Technology

RTAO DENOISER PIXEL SHADER



Full-Screen Pixel Shader



SM = Streaming Multiprocessor

L1TEX = Level 1 cache + Texture unit

L2 = Level 2 cache

VRAM = GDDR video-memory controller

CROP = Color ROP

Unit Throughput% Metrics

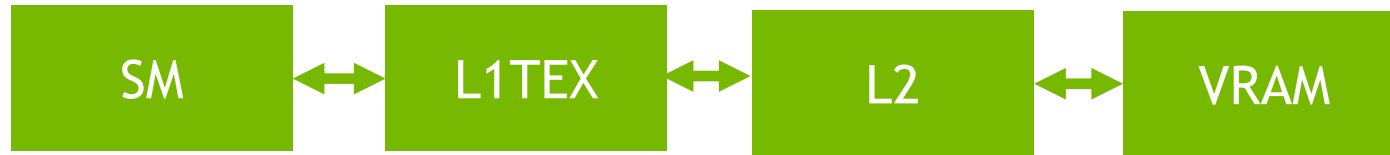


“Throughput%” = % of max theoretical throughput

Also known as:

- Speed Of Light% (SOL%)
- Peak-Perf%

Top Throughput% Units



Example:

SM: 70%

L1TEX: 57%

L2: 5%

VRAM: 2%

CROP: 2%



SM-throughput limited

SM Sub-Throughput Metrics

SM: 70%

L1TEX: 57%

L2: 5%

VRAM: 2%

CROP: 2%



SM FMA Pipe: 70%

SM SFU Pipe: 52%

SM ALU Pipe: 25%

SM FP16 Pipe: 0.0%



FMA-Pipe-Throughput
Limited

SM pipes on Turing GPUs

FMA: fp32 {FADD, FMUL, FMAD, ...} ops + int {IMUL, IMAD} ops

ALU: integer & logic ops

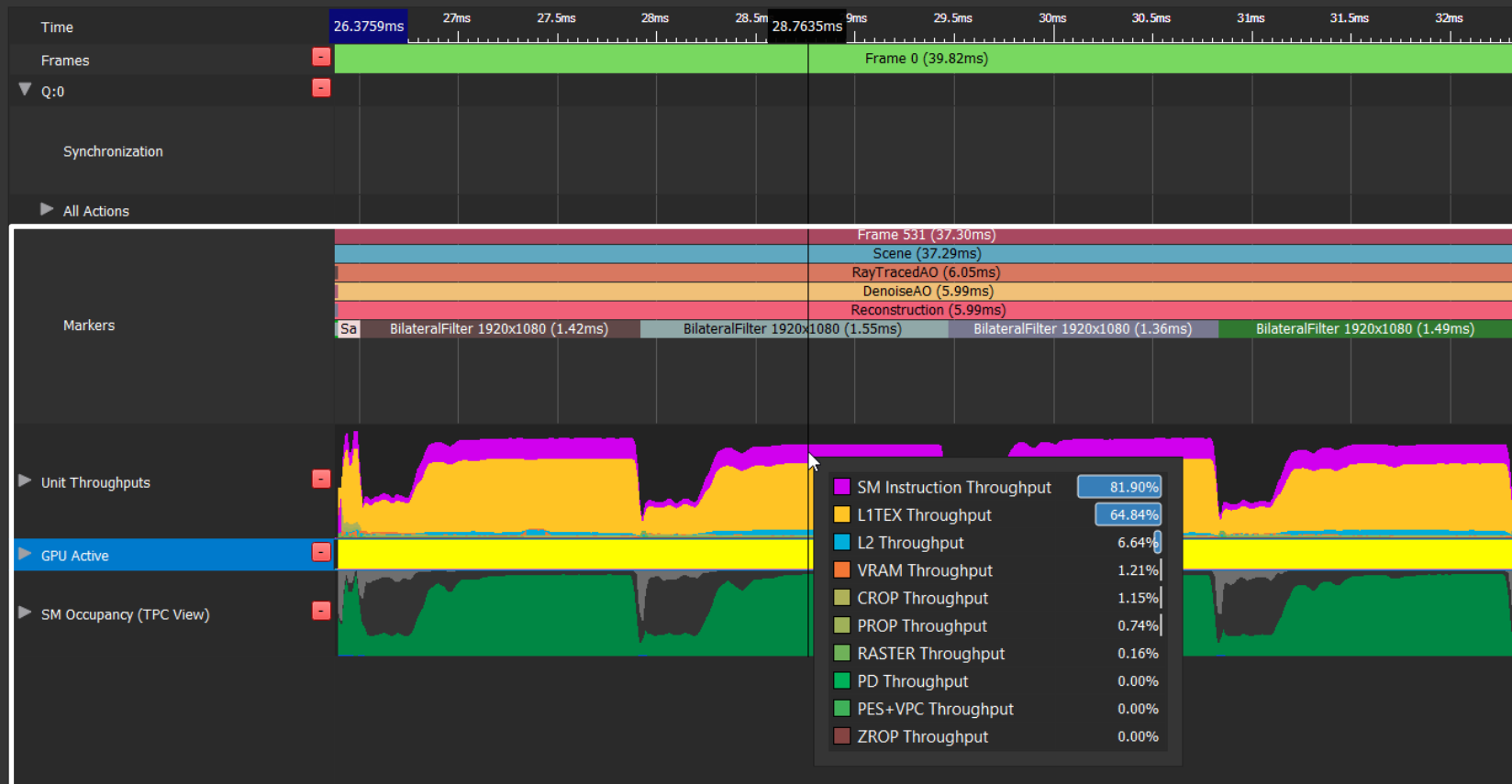
FP16: FP16 ops executed in pairs

SFU: transcendental ops (rsqrt, cos/sin, etc.)

Nsight Graphics: GPU Trace

DenoiseAO_Before.nsgt-gfxgpt X

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



Summary Metrics Detailed Metrics Capture Information

Start: 26.38 ms End: 32.36 ms Duration: 5.99 ms Range: All Visible

Unit Throughput	Value
SM Instruction Throughput	70.95%
SM FMA Pipe Throughput	69.98%
SM Issue Active	60.59%
L1TEX Throughput	56.57%
SM SFU Pipe Throughput	52.02%
SM ALU Pipe Throughput	25.18%
L2 Throughput	4.56%
CROP Throughput	2.25%
VRAM Throughput	2.20%
PROP Throughput	1.56%
RASTER Throughput	0.29%
PES+VPC Throughput	0.00%
PD Throughput	0.00%
ZROP Throughput	0.00%
SM FP16+Tensor/Pine Throughput	0.00%

Warp Occupancy	Value
Idle SM Unused Warp Slots	3.89%
Active SM Unused Warp Slots	20.25%
Compute Warps	0.00%
Pixel Warps	75.86%
Vertex+Tess+Geom Warps	0.00%

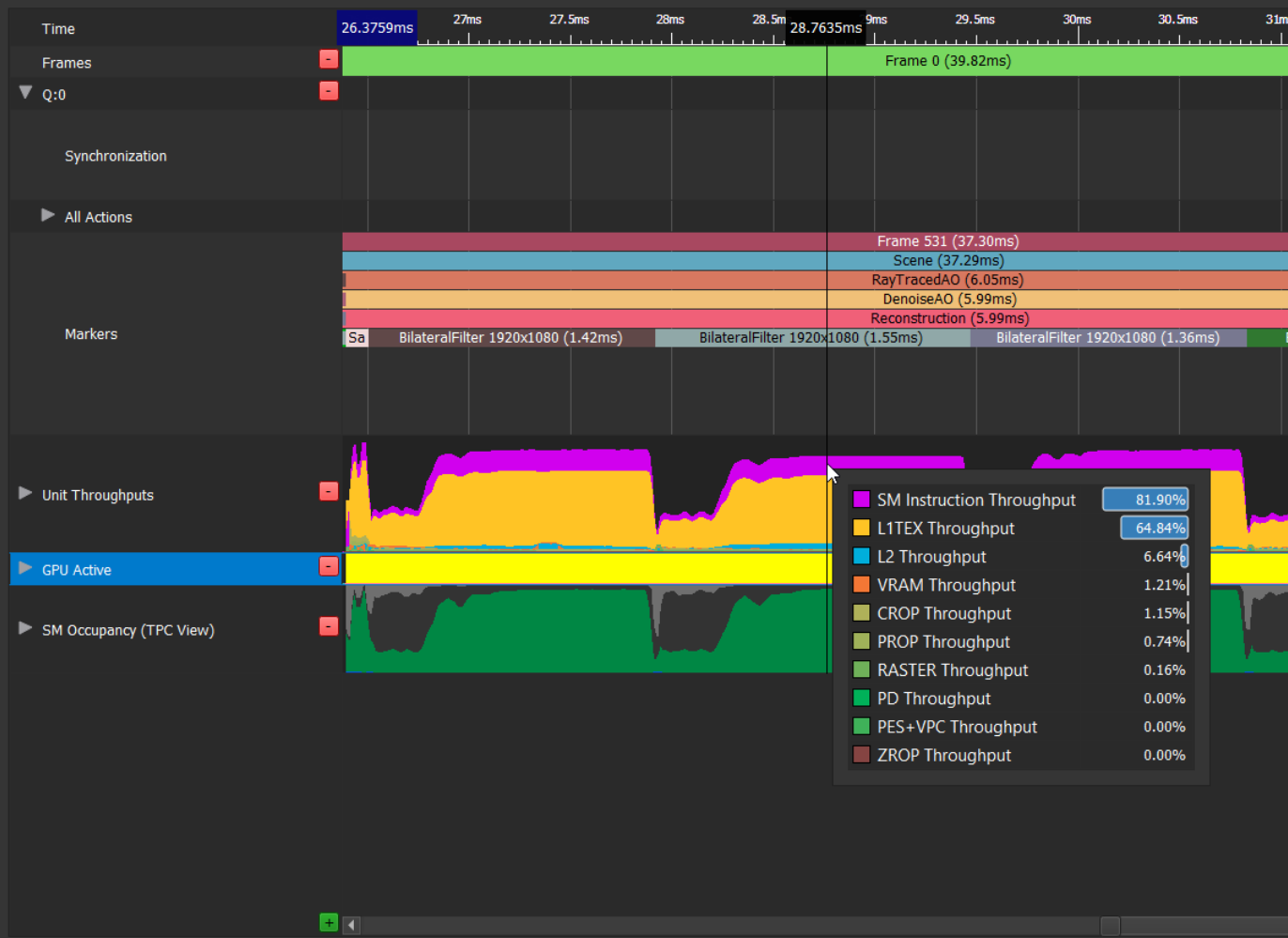


Metric Graphs (for >= Turing GPUs)

Nsight Graphics: GPU Trace

DenoiseAO_Before.nsgt-gfxgpt X

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



Summary Metrics Detailed Metrics Capture Information

Start: 26.38 ms **Duration:** 5.99 ms
End: 32.36 ms **Range:** All Visible

Unit Throughput	Value
SM Instruction Throughput	70.95%
SM FMA Pipe Throughput	69.98%
SM Issue Active	60.59%
L1TEX Throughput	56.57%
SM SFU Pipe Throughput	52.02%
SM ALU Pipe Throughput	25.13%
L2 Throughput	4.56%
CROP Throughput	2.25%
VRAM Throughput	2.20%
PROP Throughput	1.56%
RASTER Throughput	0.29%
PES+VPC Throughput	0.00%
PD Throughput	0.00%
ZROP Throughput	0.00%
SM FP16+Tensor Pine Throughput	0.00%

Warp Occupancy	Value
Idle SM Unused Warp Slots	3.89%
Active SM Unused Warp Slots	20.25%
Compute Warps	0.00%
Pixel Warps	75.86%
Vertex+Tess+Geom Warps	0.00%

Average Values for Current Range

FMA-THROUGHPUT LIMITED

Why?

For each sample, shaders reconstruct a 3D world-space position with:

```
float2 SampleScreenPosition = (SampleScreenUV.xy - View.ScreenPositionScaleBias.wz) /  
View.ScreenPositionScaleBias.xy;
```

```
float4 SampleHomogeneousWorldPosition = mul(float4(SampleScreenPosition *  
SampleDepth, SampleDepth, 1), View.ScreenToWorld); ← 4x4 matrix mul
```

```
float3 SampleWorldPosition = SampleHomogeneousWorldPosition.xyz /  
SampleHomogeneousWorldPosition.w;
```

FMA-REMOVAL EXPERIMENT:

```
#if 0
```

```
float4 SampleHomogeneousWorldPosition =  
    mul(float4(SampleScreenPosition * SampleZ, SampleZ, 1), View.ScreenToWorld);
```

```
#else
```

```
float4 SampleHomogeneousWorldPosition =  
    float4(SampleScreenPosition * SampleZ, SampleZ, 1);
```

```
#endif
```

FMA-REMOVAL EXPERIMENT:

4X4 MATRIX MUL -> NOP

	BEFORE	AFTER	RATIO
GPU Elapsed Time	5.99 ms	4.88 ms	1.23x Gain
Throughput: SM	71.0%	63.7%	0.90x
Throughput: L1TEX	56.6%	67.8%	1.20x
Throughput: L2	4.6%	5.5%	1.20x

On RTX 2080 with SetStablePowerState(TRUE)

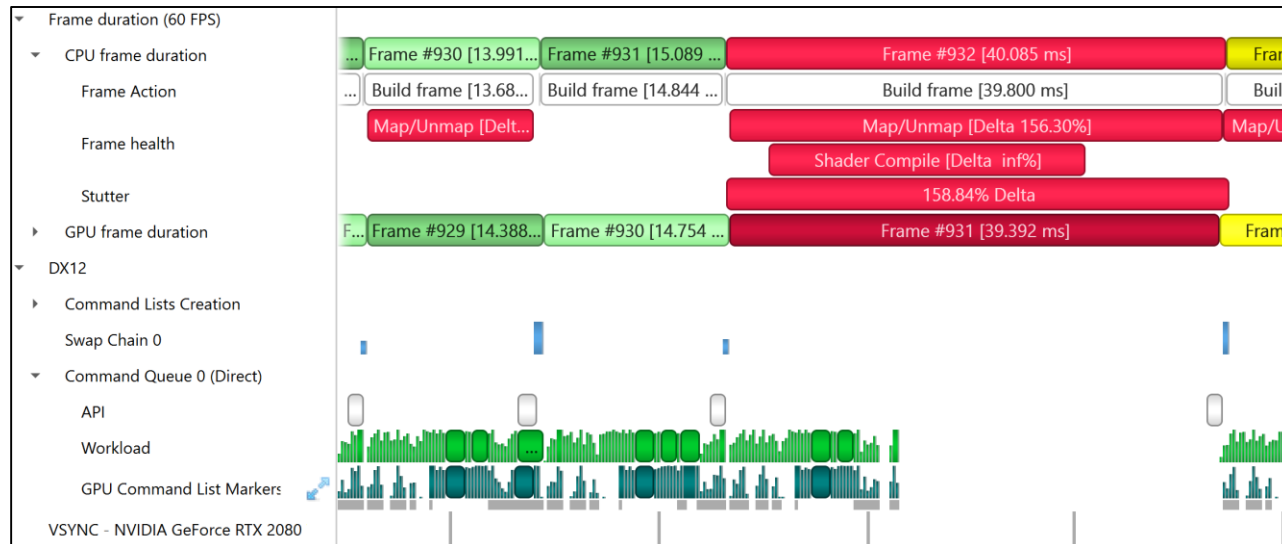
THE P3 (PEAK-PERF%) METHOD

START

GPU Active%

< 95%

Use
Nsight Systems

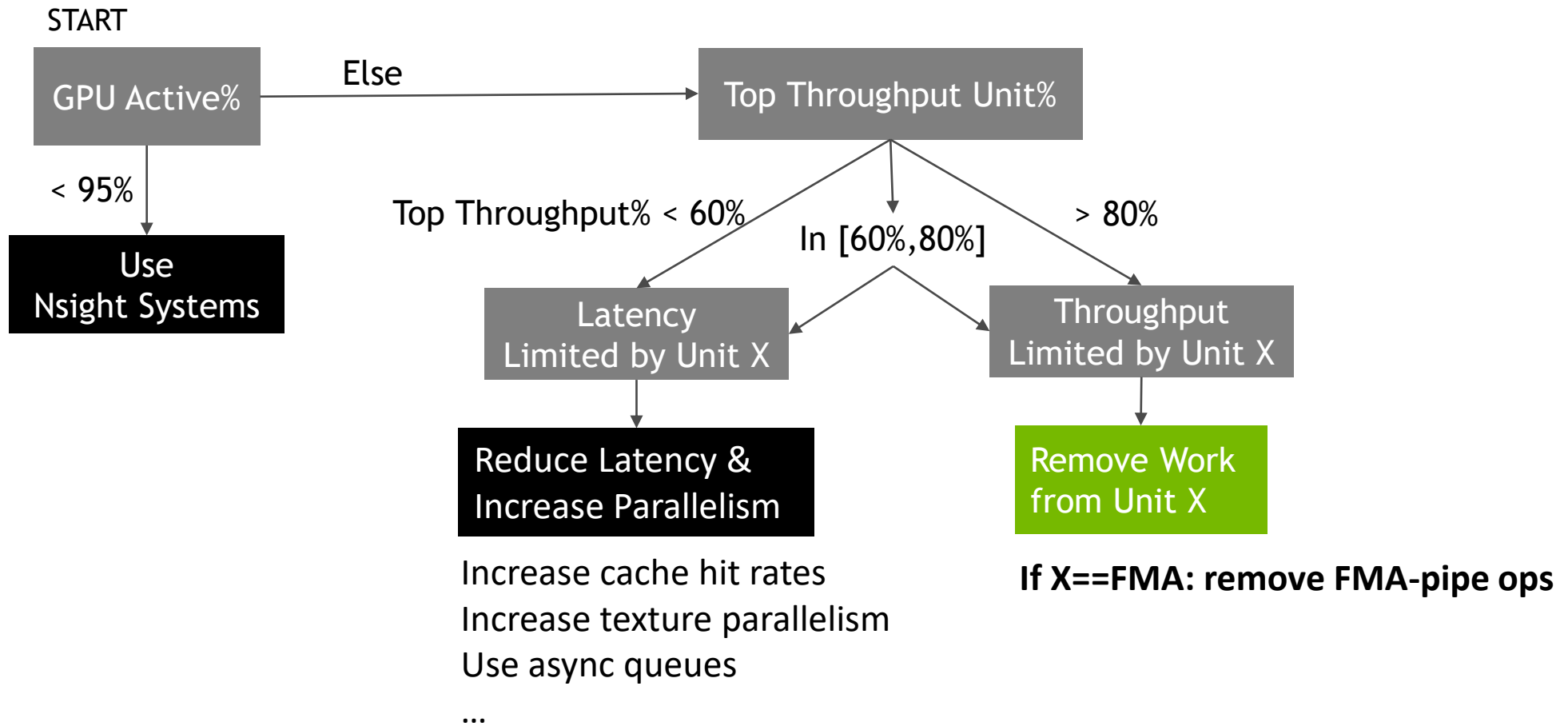


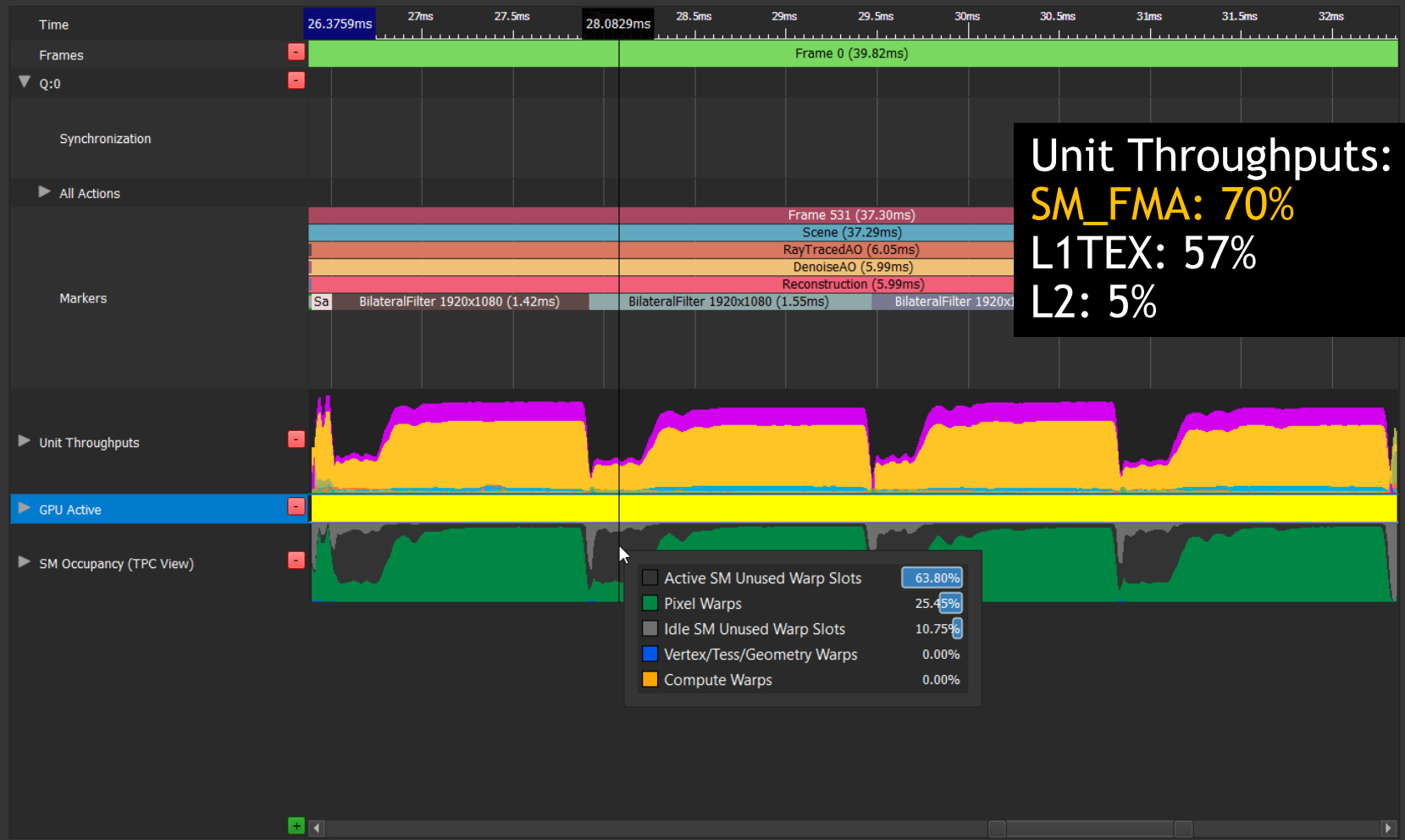
PSO Creation Stalls
on Critical Path



Video-Memory
Overcommitment

THE P3 (PEAK-PERF%) METHOD





Unit Throughputs:
SM_FMA: 70%
L1TEX: 57%
L2: 5%

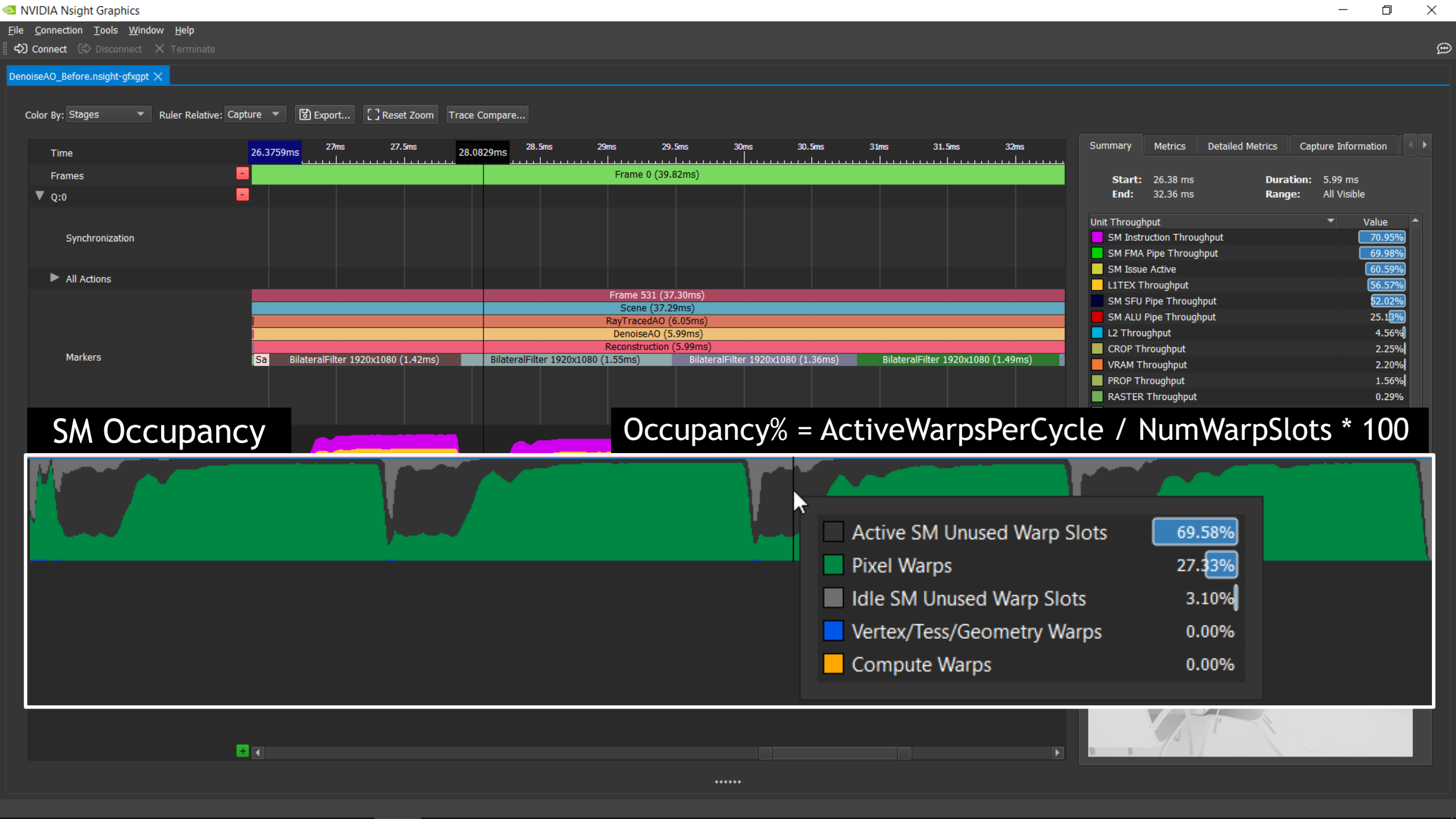
Summary	Metrics	Detailed Metrics	Capture Information
Start: 26.38 ms	Duration: 5.99 ms		
End: 32.36 ms	Range: All Visible		

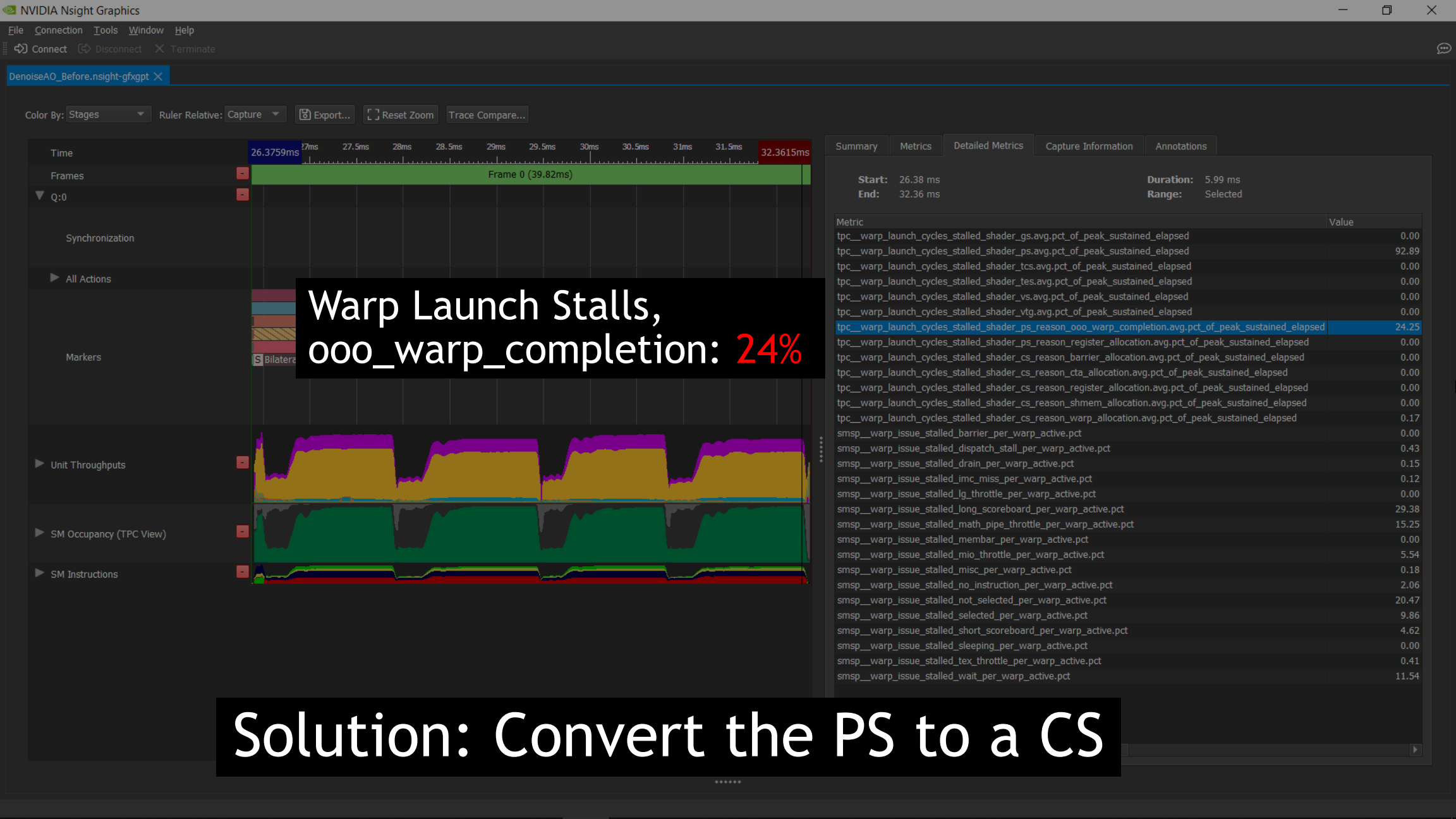
Unit Throughput	Value
SM Instruction Throughput	70.95%
SM FMA Pipe Throughput	69.98%
SM Issue Active	60.59%
L1TEX Throughput	56.57%
SM SFU Pipe Throughput	52.02%
SM ALU Pipe Throughput	25.18%
L2 Throughput	4.56%
CROP Throughput	2.25%
VRAM Throughput	2.20%
PROP Throughput	1.56%
RASTER Throughput	0.29%
PES+VPC Throughput	0.00%
PD Throughput	0.00%
ZROP Throughput	0.00%
SM FP16+Tensor Pine Throughput	0.00%

Warp Occupancy	Value
Idle SM Unused Warp Slots	3.89%
Active SM Unused Warp Slots	20.25%
Compute Warps	0.00%
Pixel Warps	75.86%
Vertex+Tess+Geom Warps	0.00%

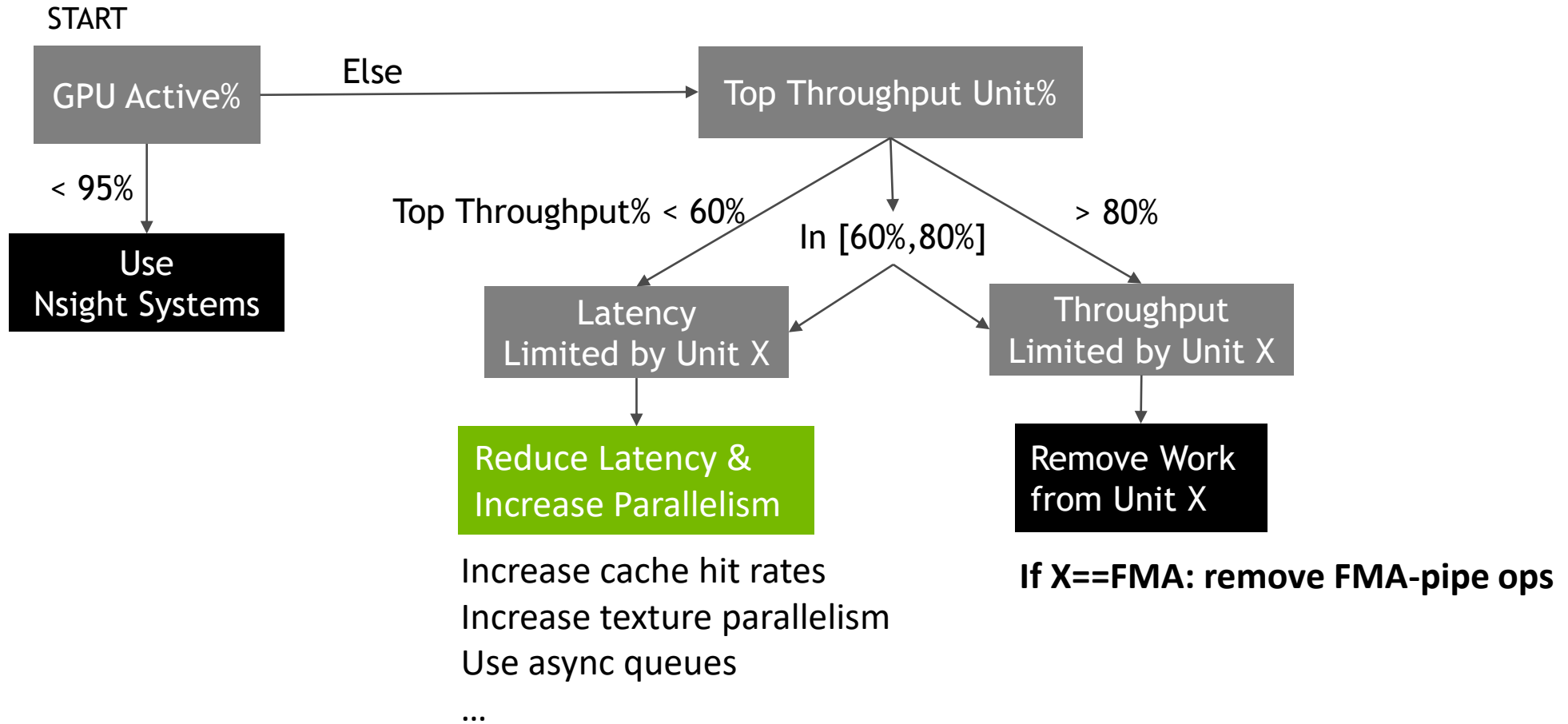


Active SM Unused Warp Slots	63.80%
Pixel Warps	25.45%
Idle SM Unused Warp Slots	10.75%
Vertex/Tess/Geometry Warps	0.00%
Compute Warps	0.00%





THE P3 (PEAK-PERF%) METHOD





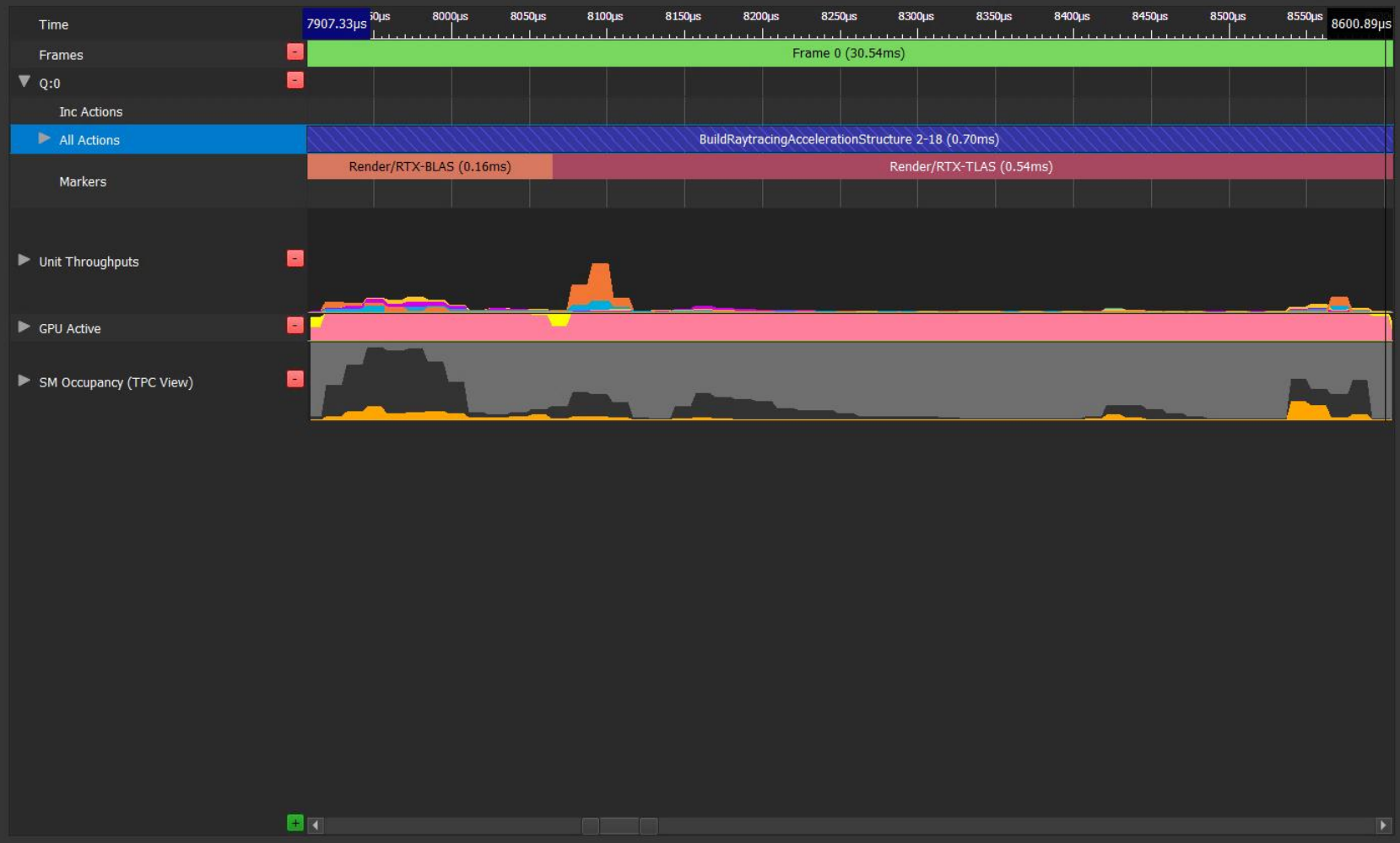
Case Study #1:
RTAS Updates // Async Compute
in Metro Exodus

RTAS Updates

Connect Disconnect Terminate

2-MetroBHV-Before.nsihgt-gfxgpt

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



Summary Metrics Detailed Metrics Capture Information

Start: 7.91 ms **Duration:** 0.70 ms
End: 8.61 ms **Range:** All Visible

Unit Throughput	Value
VRAM Throughput	3.03%
L1TEX Throughput	2.34%
SM Instruction Throughput	2.10%
SM Issue Active	2.05%
SM ALU Pipe Throughput	1.78%
L2 Throughput	1.59%
SM FMA Pipe Throughput	0.89%
SM SFU Pipe Throughput	0.72%
PROP Throughput	0.00%
PD Throughput	0.00%
PES+VPC Throughput	0.00%
RASTER Throughput	0.00%
ZROP Throughput	0.00%
CROP Throughput	0.00%
SM FP16+Tensor Pine Throughput	0.00%

Warp Occupancy	Value
Idle SM Unused Warp Slots	78.22%
Active SM Unused Warp Slots	18.35%
Compute Warps	3.43%
Pixel Warps	0.00%
Vertex+Tess+Geom Warps	0.00%

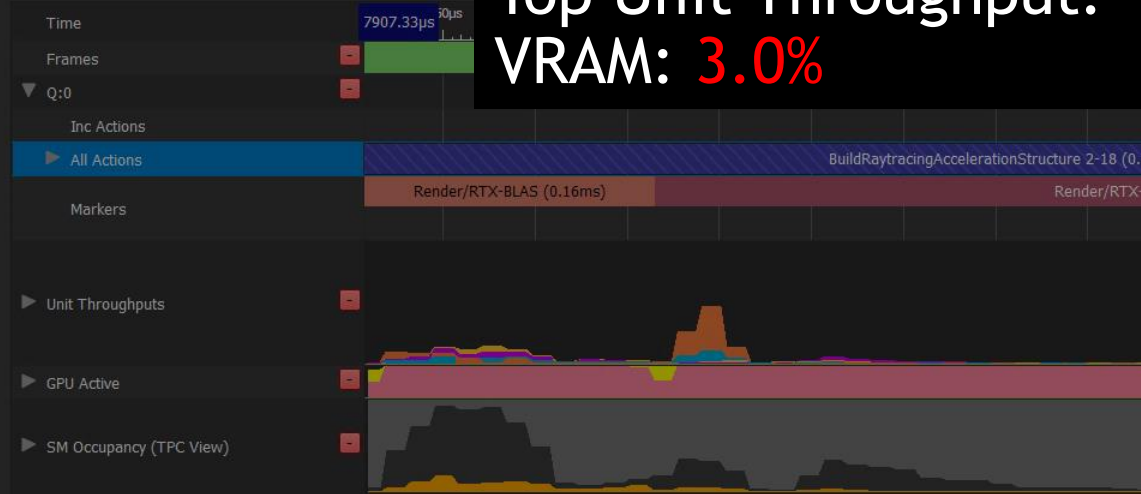
RTAS Updates

Connect Disconnect Terminate

2-MetroBHV-Before.nsisight-gfxgpt

Color By: Stages Ruler Relative: Capture

Top Unit Throughput:
VRAM: 3.0%



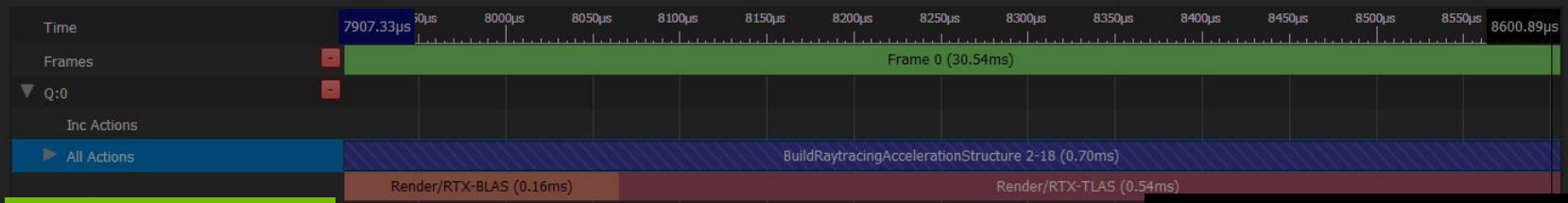
Unit Throughput	Value
VRAM Throughput	3.03%
L1TEX Throughput	2.34%
SM Instruction Throughput	2.10%
SM Issue Active	2.05%
SM ALU Pipe Throughput	1.78%
L2 Throughput	1.59%
SM FMA Pipe Throughput	0.89%
SM SFU Pipe Throughput	0.72%
PROP Throughput	0.00%
PD Throughput	0.00%
PES+VPC Throughput	0.00%
RASTER Throughput	0.00%
ZROP Throughput	0.00%
CROP Throughput	0.00%
SM FP16+Tensor Pipe Throughput	0.00%

RTAS Updates

Connect Disconnect Terminate

2-MetroBHV-Before.nsiht-gfxgpt

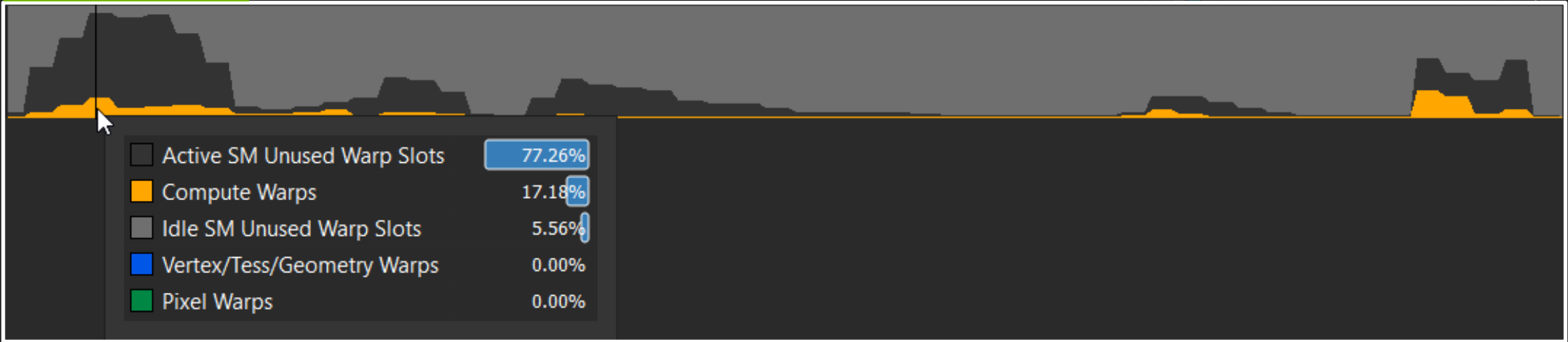
Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



Summary	Metrics	Detailed Metrics	Capture Information
Start: 7.91 ms	Duration: 0.70 ms		
End: 8.61 ms	Range: All Visible		
Unit Throughput		Value	
VRAM Throughput			3.03%
LITEX Throughput			2.34%

SM Occupancy

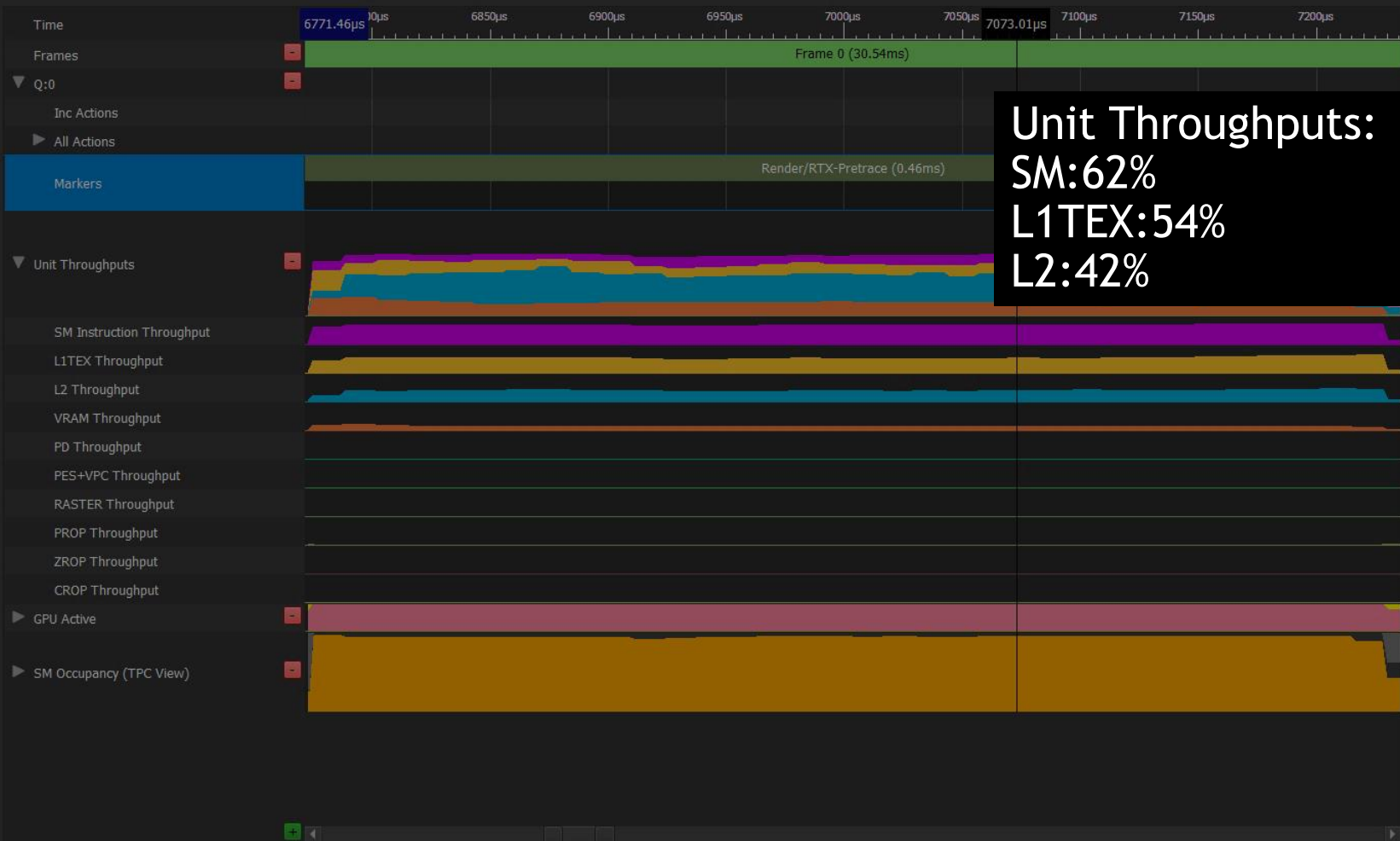
Average Warp Occupancy For Workload: 3.4%



Top Unit Throughput % << 60% + SM Occupancy% << 100% → Use Async Compute?

Independent Workload #1: Screen-Space PreTracing

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



Unit Throughputs:
SM:62%
L1TEX:54%
L2:42%

Summary Metrics Detailed Metrics Capture Information

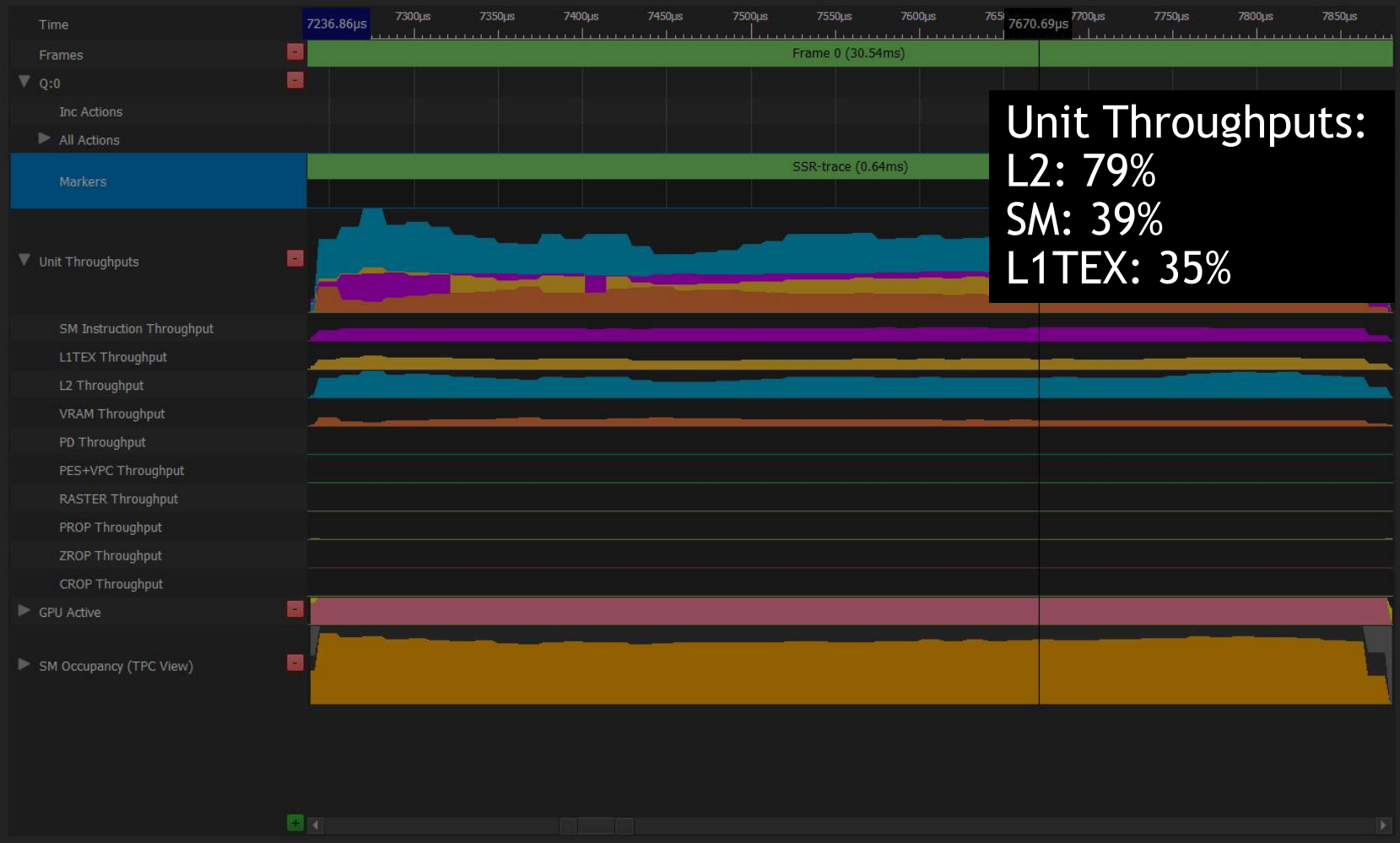
Start: 6.77 ms End: 7.24 ms Duration: 0.46 ms Range: All Visible

Unit Throughput	Value
SM Instruction Throughput	61.51%
SM FMA Pipe Throughput	61.46%
L1TEX Throughput	53.95%
SM Issue Active	49.80%
L2 Throughput	42.65%
SM SFU Pipe Throughput	33.86%
SM ALU Pipe Throughput	15.30%
VRAM Throughput	12.70%
PROP Throughput	0.00%
PD Throughput	0.00%
PES+VPC Throughput	0.00%
RASTER Throughput	0.00%
ZROP Throughput	0.00%
CROP Throughput	0.00%
SM FP16+Tensor Pipe Throughput	0.00%

Warp Occupancy	Value
Idle SM Unused Warp Slots	1.03%
Active SM Unused Warp Slots	4.46%
Compute Warps	94.51%
Pixel Warps	0.00%
Vertex+Tess+Geom Warps	0.00%

Independent Workload #2: SSR

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



Unit Throughputs:
L2: 79%
SM: 39%
L1TEX: 35%

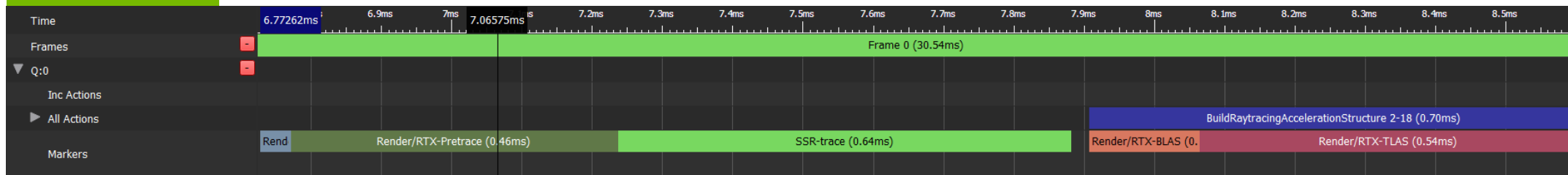
Unit Throughput	Value
L2 Throughput	78.76%
SM Instruction Throughput	38.82%
SM FMA Pipe Throughput	38.78%
SM Issue Active	35.99%
L1TEX Throughput	34.81%
SM ALU Pipe Throughput	20.80%
VRAM Throughput	18.17%
SM SFU Pipe Throughput	17.71%
PROP Throughput	0.00%
PD Throughput	0.00%
PES+VPC Throughput	0.00%
RASTER Throughput	0.00%
ZROP Throughput	0.00%
CROP Throughput	0.00%
SM FP16+Tensor Pipe Throughput	0.00%

Warp Occupancy	Value
Idle SM Unused Warp Slots	1.49%
Active SM Unused Warp Slots	17.88%
Compute Warps	80.62%
Pixel Warps	0.00%
Vertex+Tess+Geom Warps	0.00%

ASYNC COMPUTE DIFF

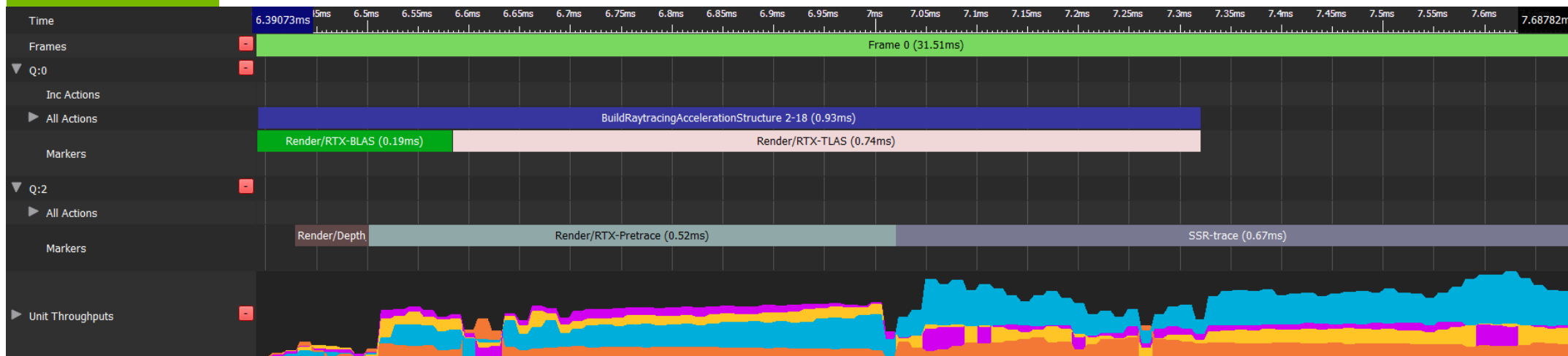
Serialized

1.83 ms



Overlapped

1.30 ms



ASYNC-COMPUTE OVERLAP

(RTAS Updates) // (Async Compute)

	BEFORE	AFTER	RATIO
GPU Elapsed Time	1.83 ms	1.30 ms	1.41x Gain
Throughput: L2	39.2%	54.8%	1.40x
Throughput: SM	30.1%	42.0%	1.40x
Throughput: L1TEX	26.9%	37.4%	1.39x
SM Occupancy	53.8%	78.2%	1.45x

On RTX 2080 with SetStablePowerState(TRUE)



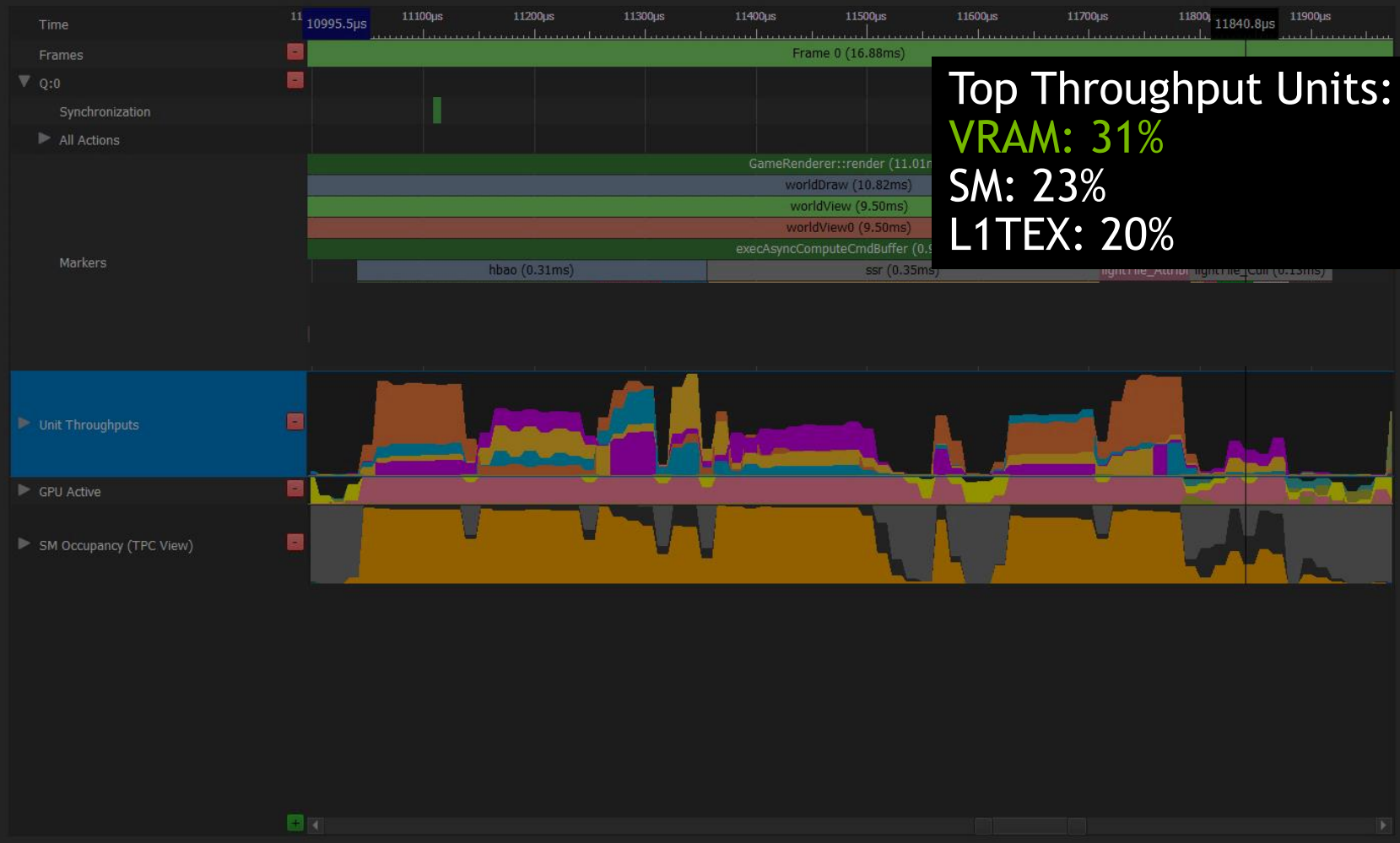
Case Study #2:
Shadow Maps // Async Compute

HBAO + SSR + Light Culling

Connect Disconnect Terminate

1-BF5_AsyncCompute-Before.nsisight-gfxgpt X

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



Top Throughput Units:
VRAM: 31%
SM: 23%
L1TEX: 20%

Summary Metrics Detailed Metrics Capture Information

Start: 11.00 ms End: 11.97 ms Duration: 0.98 ms Range: All Visible

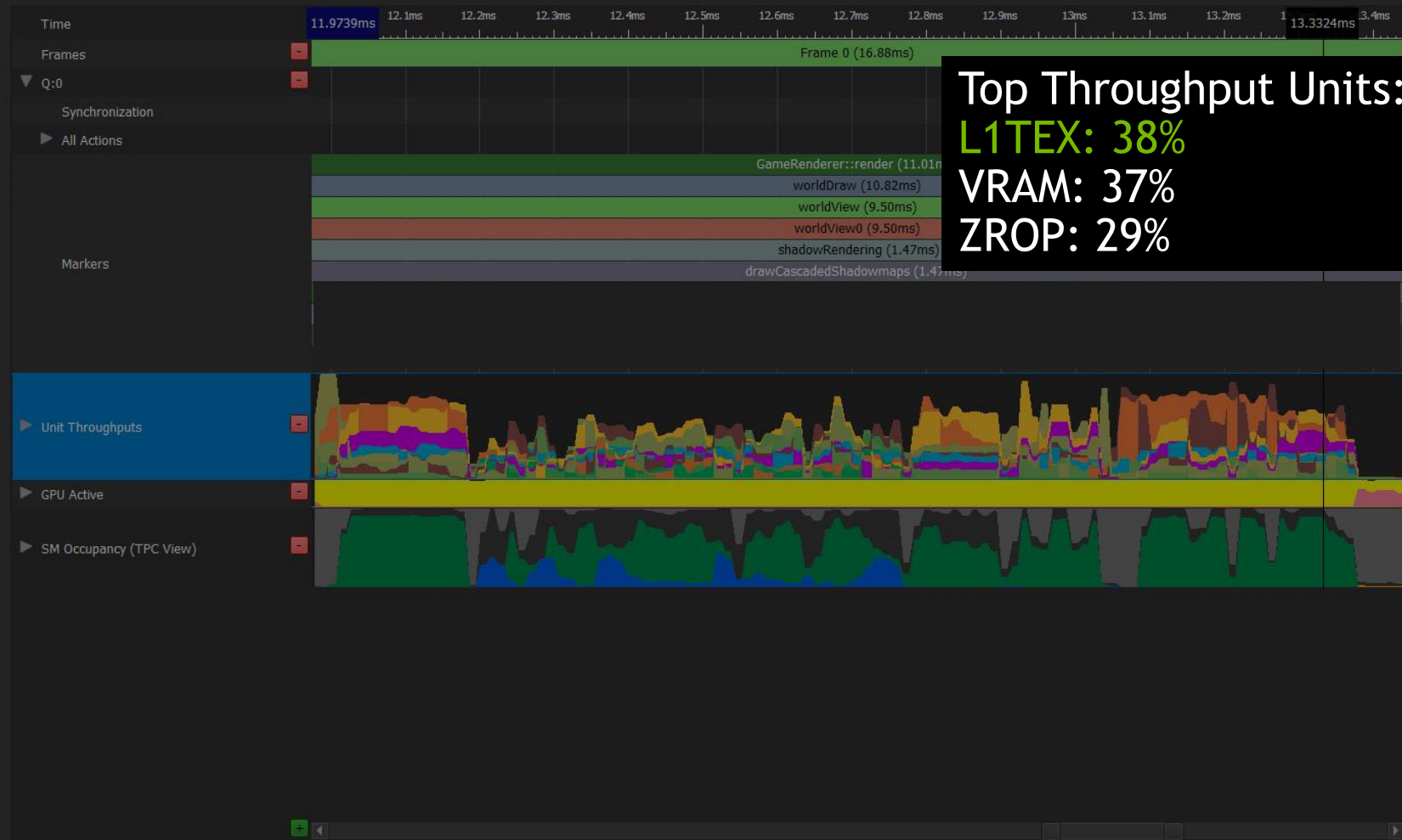
Unit Throughput	Value
VRAM Throughput	30.63%
SM Instruction Throughput	22.61%
L1TEX Throughput	19.96%
SM Issue Active	19.90%
L2 Throughput	18.30%
SM FMA Pipe Throughput	15.42%
SM SFU Pipe Throughput	12.62%
SM ALU Pipe Throughput	11.95%
CROP Throughput	0.30%
RASTER Throughput	0.15%
PROP Throughput	0.12%
ZROP Throughput	0.00%
PES+VPC Throughput	0.00%
PD Throughput	0.00%
SM FP16+Tensor Pine Throughput	0.00%

Warp Occupancy	Value
Idle SM Unused Warp Slots	28.20%
Active SM Unused Warp Slots	10.56%
Compute Warps	61.24%
Pixel Warps	0.00%
Vertex+Tess+Geom Warps	0.00%



Independent Workload: Shadow Maps

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



Top Throughput Units:
L1TEX: 38%
VRAM: 37%
ZROP: 29%

Summary Metrics Detailed Metrics Capture Information

Start: 11.97 ms End: 13.44 ms Duration: 1.47 ms Range: All Visible

Unit Throughput	Value
L1TEX Throughput	38.34%
VRAM Throughput	36.68%
ZROP Throughput	29.10%
L2 Throughput	22.04%
SM Instruction Throughput	21.77%
SM Issue Active	20.54%
PROP Throughput	20.48%
SM FMA Pipe Throughput	17.52%
RASTER Throughput	17.42%
CROP Throughput	12.11%
SM ALU Pipe Throughput	8.69%
SM SFU Pipe Throughput	7.62%
PES+VPC Throughput	5.61%
PD Throughput	4.31%
SM FP16+Tensor Pipe Throughput	0.00%

Warp Occupancy	Value
Idle SM Unused Warp Slots	20.09%
Active SM Unused Warp Slots	21.51%
Compute Warps	0.05%
Pixel Warps	51.97%
Vertex+Tess+Geom Warps	6.37%



ASYNC-COMPUTE OVERLAP

(SHADOW MAPS) // (HBAO+SSR+LIGHT-CULL)

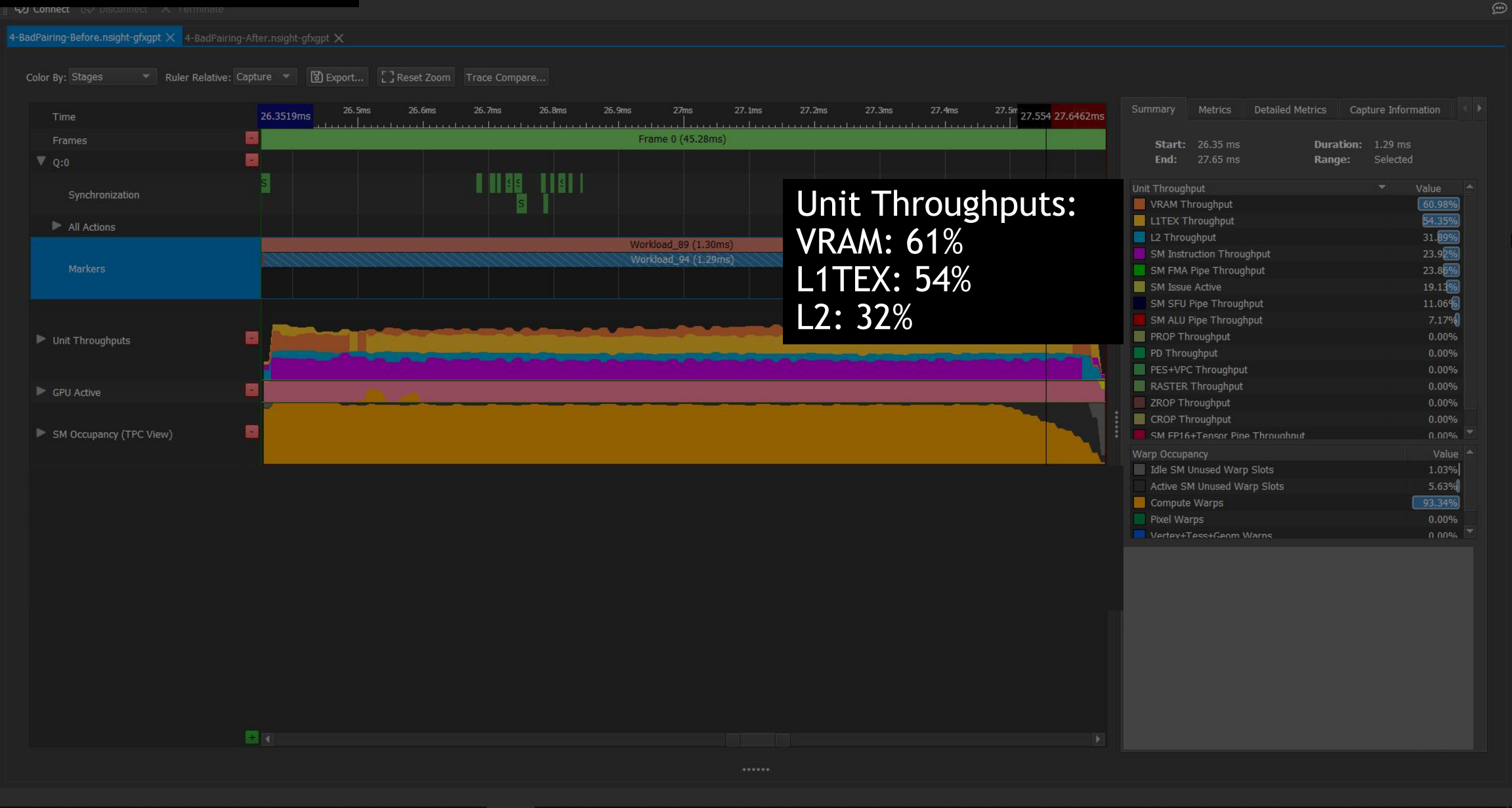
	BEFORE	AFTER	RATIO
GPU Elapsed Time	2.45 ms	2.15 ms	1.14x Gain
Throughput: VRAM	34.2%	40.8%	1.19x
Throughput: L1TEX	31.0%	34.0%	1.10x
Throughput: SM	22.1%	24.3%	1.10x
SM Occupancy	59.5%	70.6%	1.19x

On RTX 2080 + SetStablePowerState(TRUE)



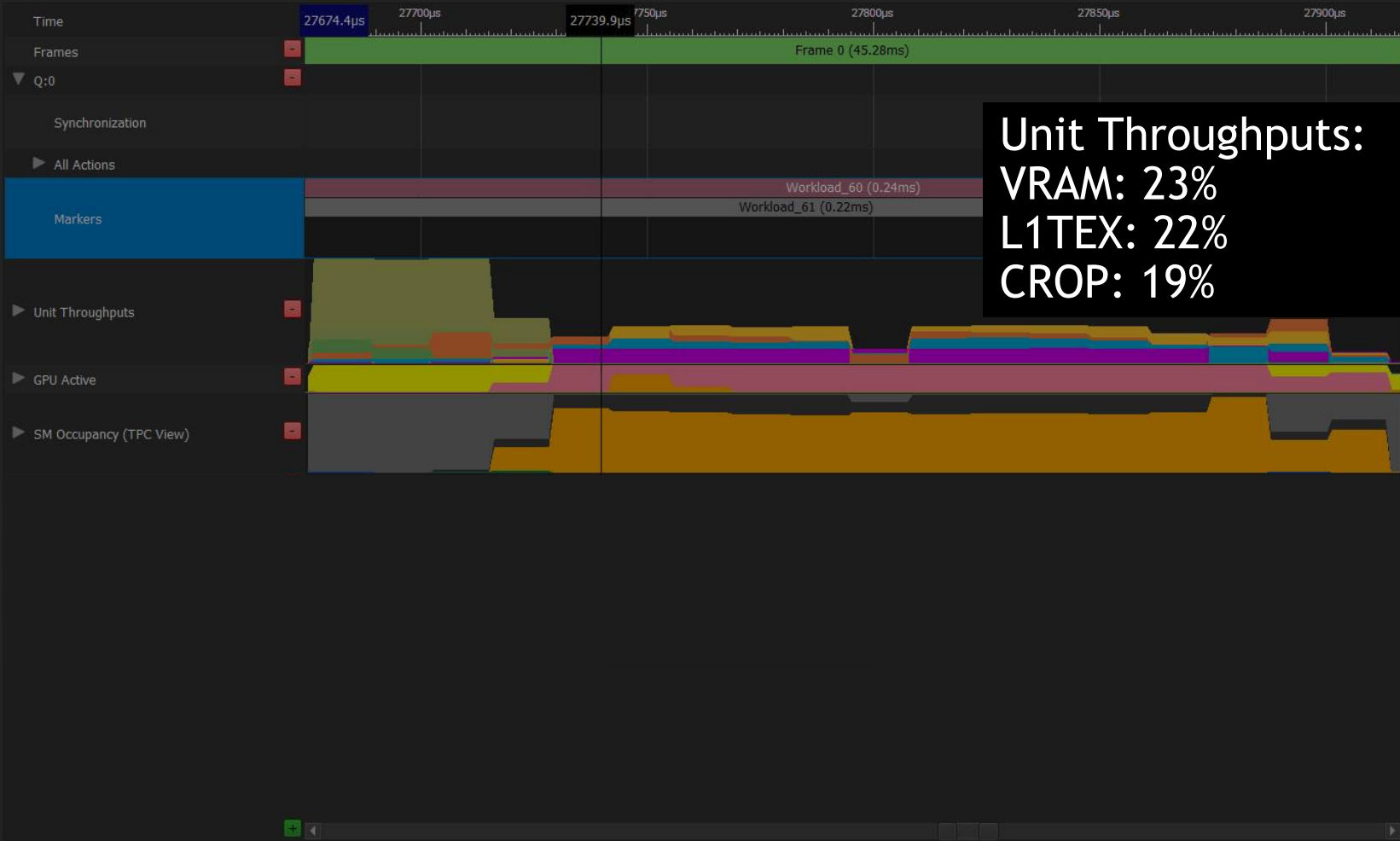
Case Study #3:
Bad Async-Compute Pairing

Blur Compute Shader



Independent Workload #1: Water Simulation

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



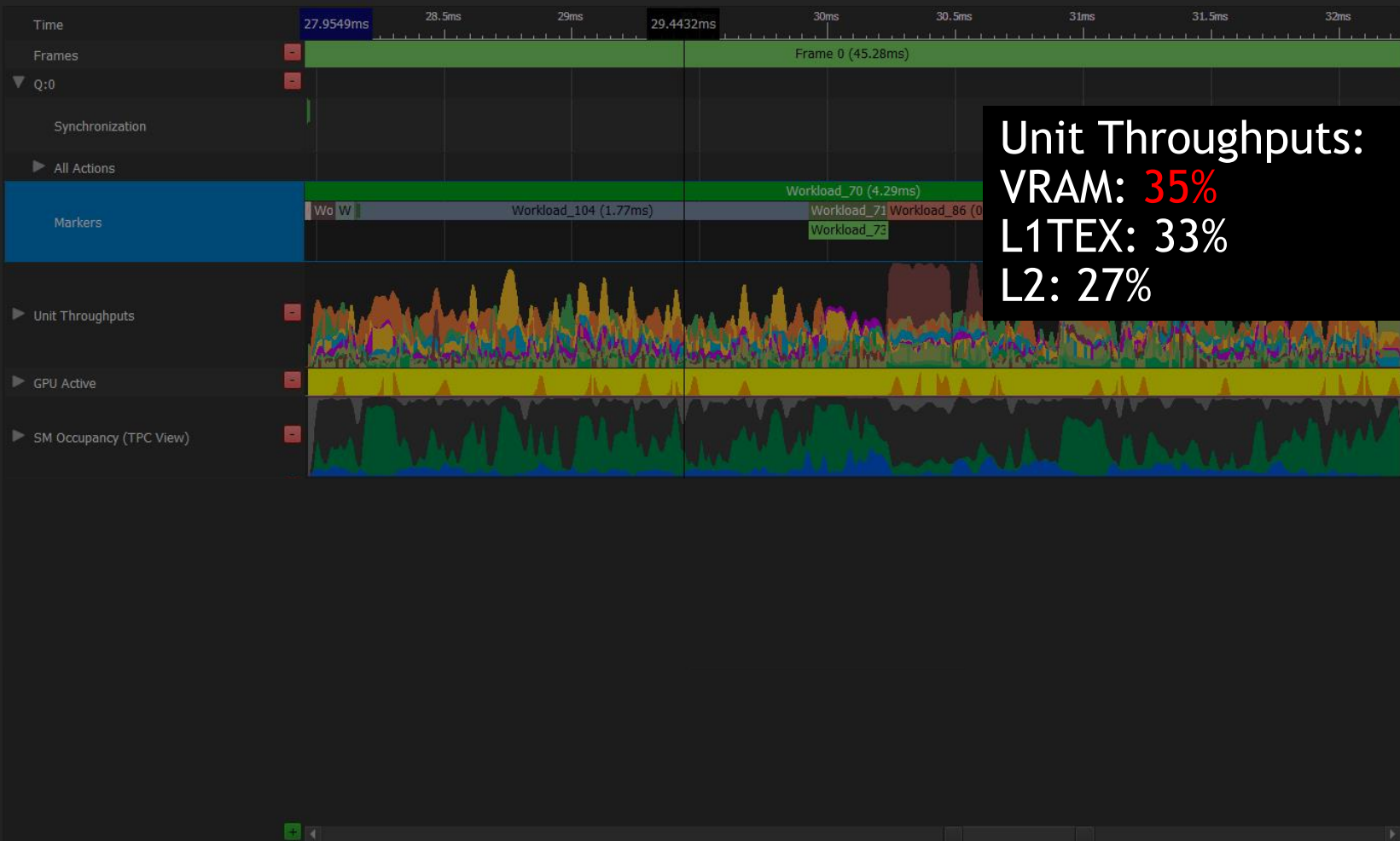
Unit Throughputs:
VRAM: 23%
L1TEX: 22%
CROP: 19%

Unit Throughput	Value
VRAM Throughput	22.96%
L1TEX Throughput	22.51%
CROP Throughput	18.58%
L2 Throughput	14.79%
SM Instruction Throughput	10.27%
SM Issue Active	10.27%
PROP Throughput	6.16%
SM ALU Pipe Throughput	5.01%
RASTER Throughput	2.29%
SM FMA Pipe Throughput	1.70%
SM SFU Pipe Throughput	1.08%
PES+VPC Throughput	0.00%
PD Throughput	0.00%
ZROP Throughput	0.00%
SM FP16+Tensor Pipe Throughput	0.00%

Warp Occupancy	Value
Idle SM Unused Warp Slots	26.14%
Active SM Unused Warp Slots	15.79%
Compute Warps	57.90%
Pixel Warps	0.17%
Vertex+Tess+Geom Warps	0.00%

Independent Workload #2: GBuffer Fill

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



Unit Throughputs:
VRAM: 35%
L1TEX: 33%
L2: 27%

Unit Throughput	Value
VRAM Throughput	34.57%
L1TEX Throughput	33.05%
L2 Throughput	26.56%
PES+VPC Throughput	22.22%
SM Instruction Throughput	22.07%
SM Issue Active	20.99%
ZROP Throughput	19.89%
SM FMA Pipe Throughput	19.50%
CROP Throughput	18.26%
RASTER Throughput	10.37%
PROP Throughput	9.39%
SM SFU Pipe Throughput	8.77%
SM ALU Pipe Throughput	8.12%
PD Throughput	7.13%
SM FP16+Tensor Pipe Throughput	0.00%

Warp Occupancy	Value
Idle SM Unused Warp Slots	5.74%
Active SM Unused Warp Slots	50.30%
Compute Warps	0.00%
Pixel Warps	36.02%
Vertex+Tess+Geom Warps	7.94%

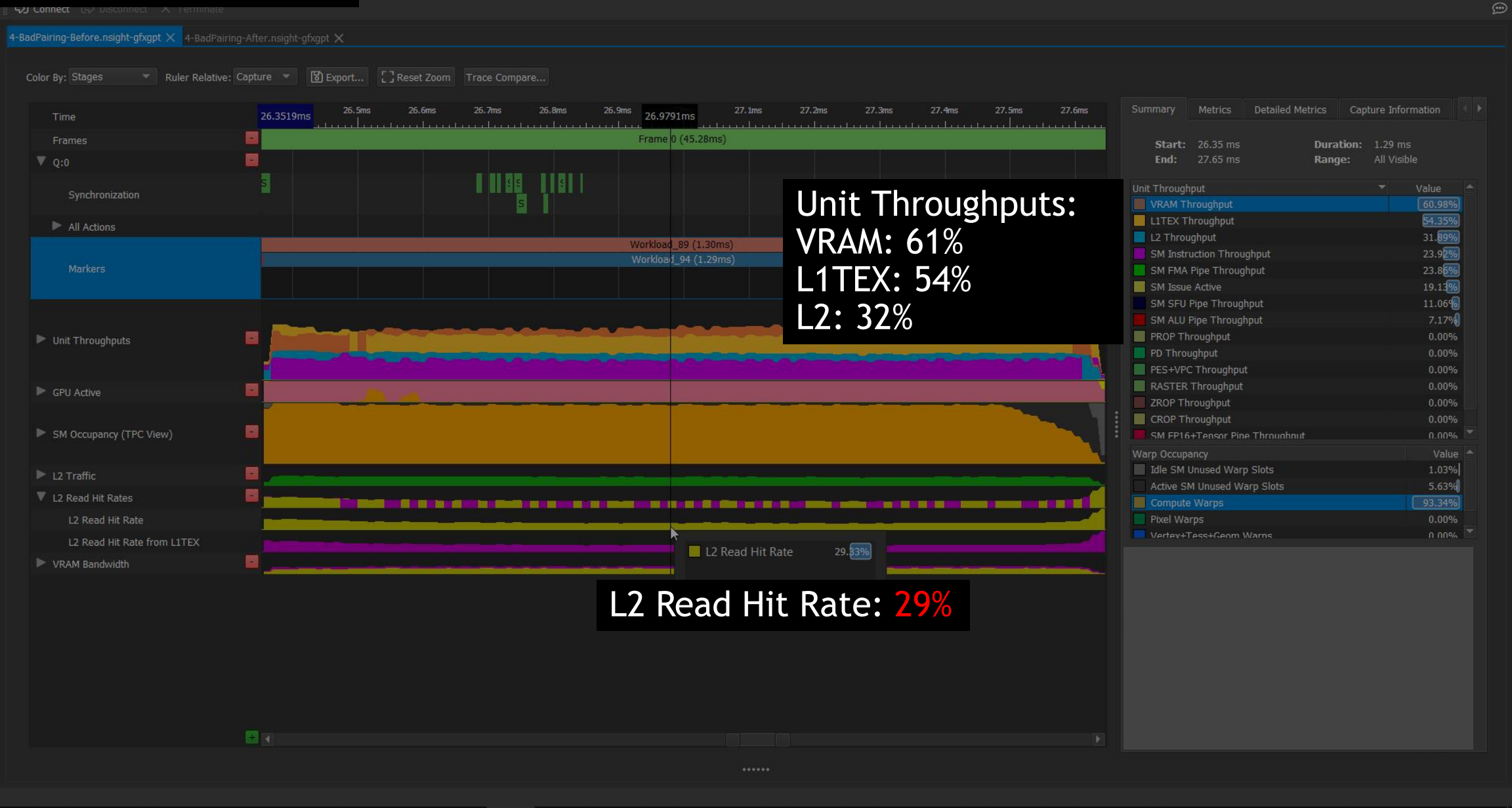
BAD ASYNC-COMPUTE PAIRING

(BLUR CS) // (GBUFFER + WATER SIM)

	BEFORE	AFTER	RATIO
GPU Elapsed Time	5.89 ms	6.12 ms	0.96x Loss
Throughput: VRAM	43.1%	47.1%	1.09x
Throughput: L1TEX	36.9%	35.4%	0.96x
Throughput: L2	27.0%	26.0%	0.96x
SM Occupancy	54.9%	57.5%	1.05x
L2 Read Hit Rate	52.3%	44.5%	0.85x

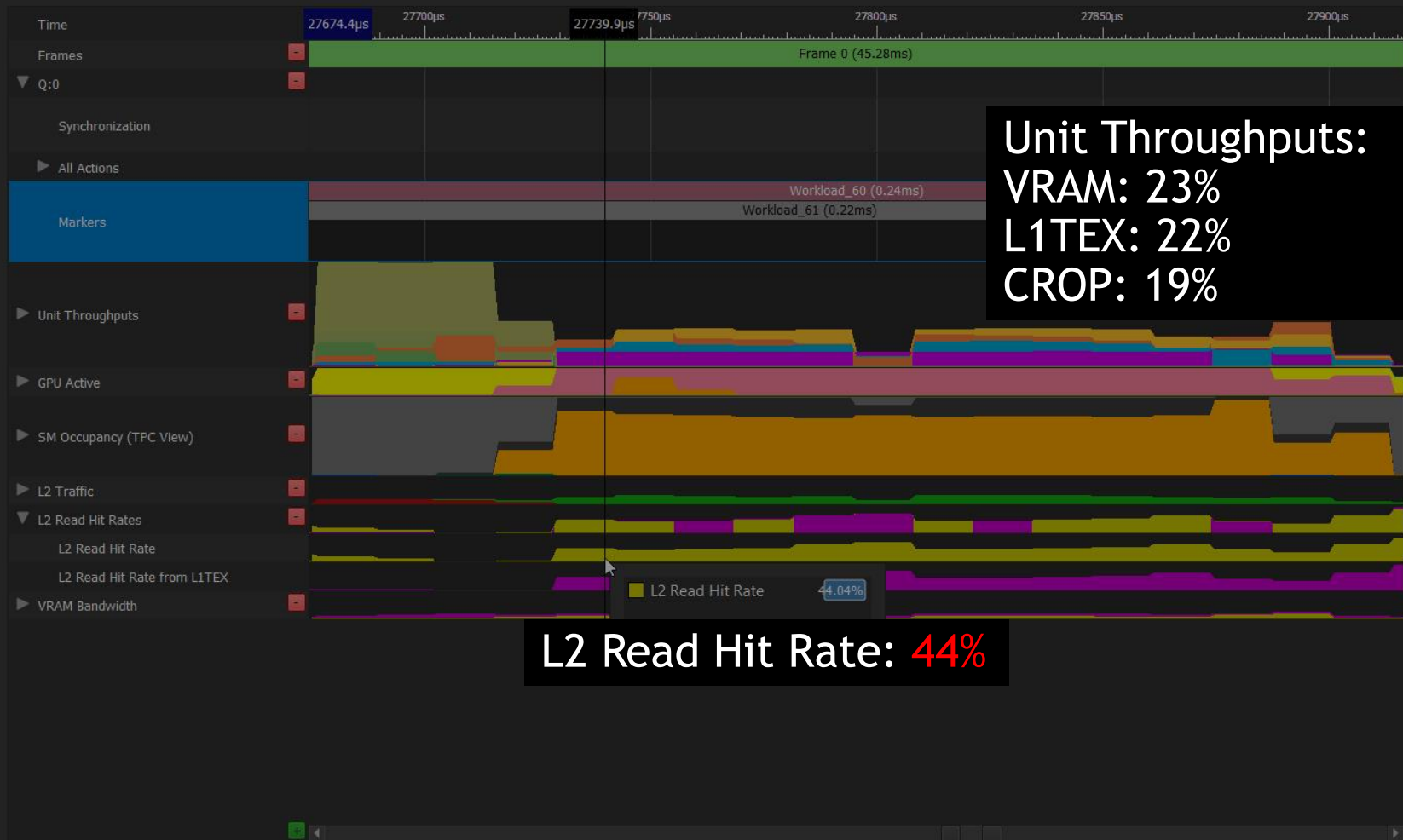
On RTX 2080 with SetStablePowerState(TRUE)

Blur Compute Shader



Independent Workload #1: Water Simulation

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



Unit Throughputs:
VRAM: 23%
L1TEX: 22%
CROP: 19%

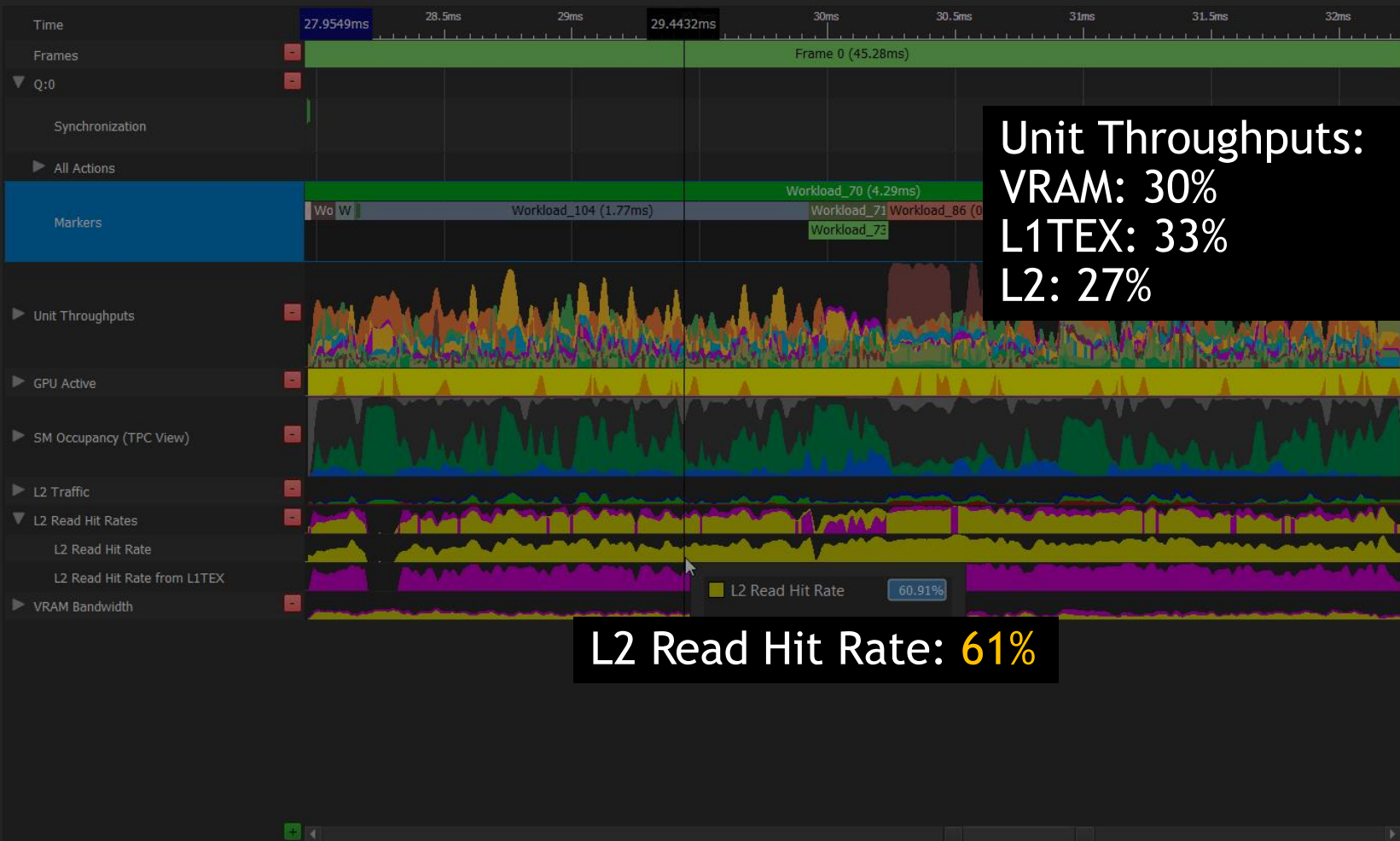
Unit Throughput	Value
VRAM Throughput	22.96%
L1TEX Throughput	22.51%
CROP Throughput	18.58%
L2 Throughput	14.79%
SM Instruction Throughput	10.27%
SM Issue Active	10.27%
PROP Throughput	6.16%
SM ALU Pipe Throughput	5.01%
RASTER Throughput	2.29%
SM FMA Pipe Throughput	1.70%
SM SFU Pipe Throughput	1.08%
PES+VPC Throughput	0.00%
PD Throughput	0.00%
ZROP Throughput	0.00%
SM FP16+Tensor Pipe Throughput	0.00%

Warp Occupancy	Value
Idle SM Unused Warp Slots	26.14%
Active SM Unused Warp Slots	15.79%
Compute Warps	57.90%
Pixel Warps	0.17%
Vertex+Tess+Geom Warps	0.00%

L2 Read Hit Rate: 44%

Independent Workload #2: GBuffer Fill

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



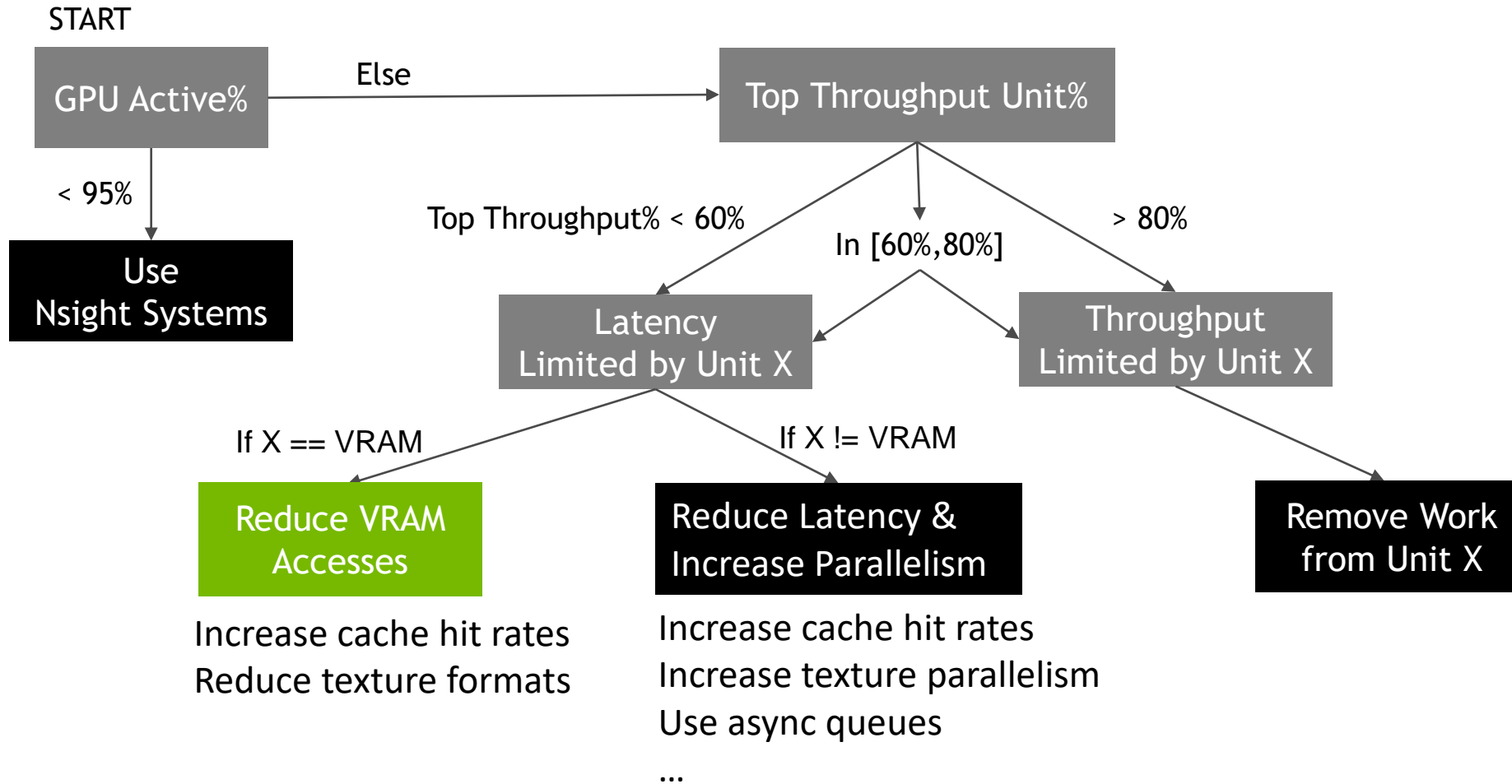
Unit Throughputs:
VRAM: 30%
L1TEX: 33%
L2: 27%

L2 Read Hit Rate: 61%

Unit Throughput	Value
VRAM Throughput	39.57%
L1TEX Throughput	33.05%
L2 Throughput	26.56%
PES+VPC Throughput	22.22%
SM Instruction Throughput	22.07%
SM Issue Active	20.99%
ZROP Throughput	19.89%
SM FMA Pipe Throughput	19.50%
CROP Throughput	18.26%
RASTER Throughput	10.37%
PROP Throughput	9.39%
SM SFU Pipe Throughput	8.77%
SM ALU Pipe Throughput	8.12%
PD Throughput	7.13%
SM FP16+Tensor Pipe Throughput	0.00%

Warp Occupancy	Value
Idle SM Unused Warp Slots	5.74%
Active SM Unused Warp Slots	50.30%
Compute Warps	0.00%
Pixel Warps	36.02%
Vertex+Tess+Geom Warps	7.94%

THE P3 (PEAK-PERF%) METHOD





Case Study #4:
VRAM-Limited Denoiser CS

Battlefield V, DXR, 1440p

W

149 m

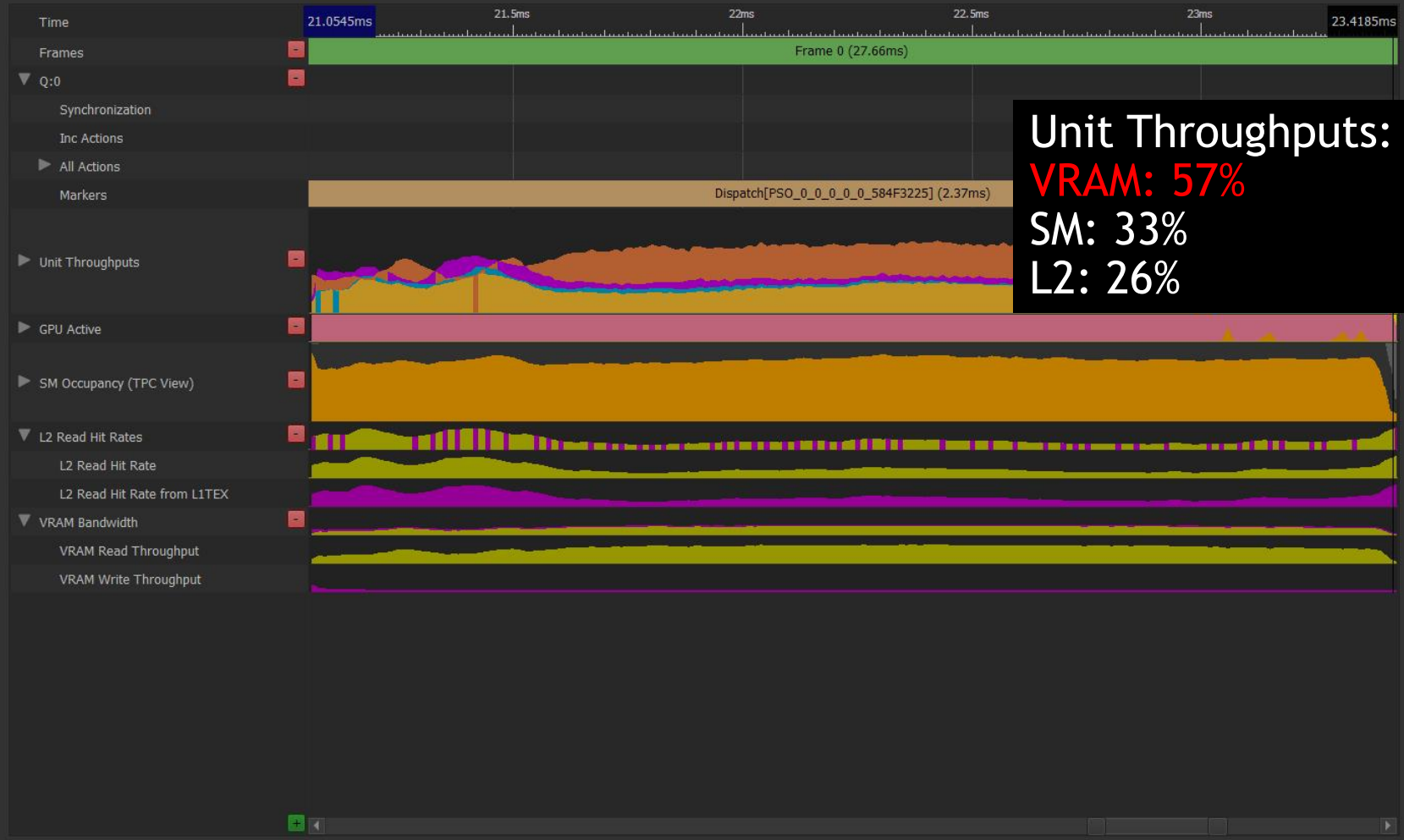
SECURE CHOKEPOINT KILO



+ 91

20 40 3

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



Unit Throughputs:
VRAM: 57%
SM: 33%
L2: 26%

Summary	Metrics	Detailed Metrics	Capture Information
Start:	21.05 ms	Duration:	2.37 ms
End:	23.43 ms	Range:	All Visible

Unit Throughput	Value
VRAM Throughput	56.09%
SM Instruction Throughput	33.76%
SM SFU Pipe Throughput	33.58%
SM Issue Active	29.74%
L2 Throughput	25.87%
SM FMA Pipe Throughput	25.31%
L1TEX Throughput	24.39%
SM ALU Pipe Throughput	18.03%
PROP Throughput	0.00%
PD Throughput	0.00%
PES+VPC Throughput	0.00%
RASTER Throughput	0.00%
ZROP Throughput	0.00%
CROP Throughput	0.00%
SM FP16+Tensor Pipe Throughput	0.00%

Warp Occupancy	Value
Idle SM Unused Warp Slots	0.50%
Active SM Unused Warp Slots	21.58%
Compute Warps	77.92%
Pixel Warps	0.00%
Vertex+Tess+Geom Warps	0.00%



Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



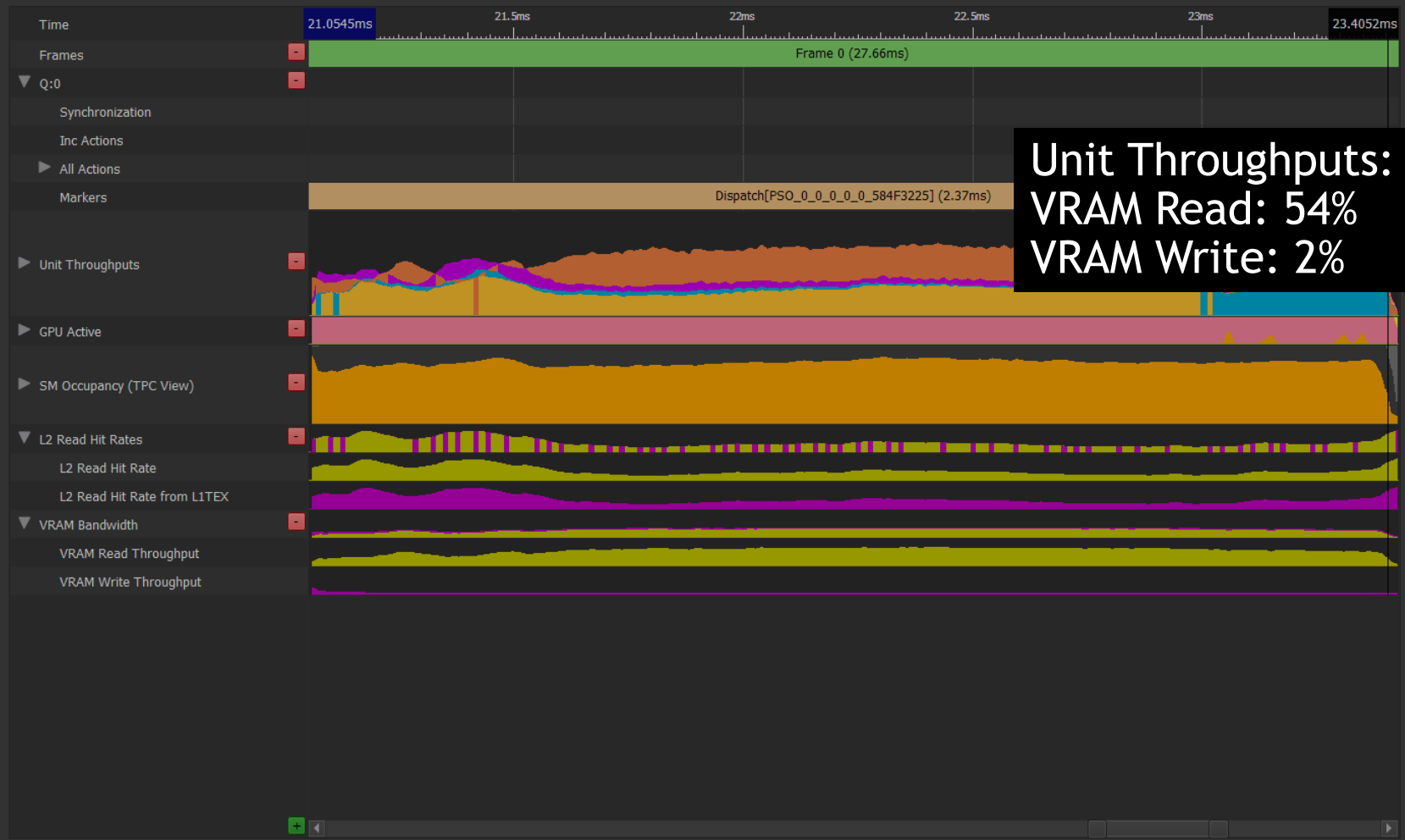
Summary Metrics Detailed Metrics Capture Information

Start: 21.05 ms End: 23.43 ms Duration: 2.37 ms Range: All Visible

Search...

Metrics	Value
Active SM Unused Warp Slots	21.56%
Async Compute In Flight	0.00%
Async Copy Engine Active	0.99%
Compute Warps	77.92%
CROP Throughput	0.00%
Dispatch Started	0.00%
Draw Started	0.00%
DXR Build	0.00%
DXR Dispatch	0.00%
GPU Active	50.13%
GPU Active	100.00%
GR Active	100.00%
Idle SM Unused Warp Slots	0.50%
L1TEX Local+Global Data Throughput	0.65%
L1TEX LSU Data Throughput	0.80%
L1TEX LSU Writeback Throughput	0.15%
L1TEX Shared+Attribute Data Throughput	0.15%
L1TEX Tag-Stage Throughput	23.88%
L1TEX Texture Data Throughput	24.25%
L1TEX Texture Filter Stage Throughput	3.53%
L1TEX Throughput	24.38%
L2 Bandwidth from CROP	0.00%
L2 Bandwidth from L1TEX	25.87%
L2 Bandwidth from ZROP	0.00%
L2 Read Hit Rate	32.42%
L2 Read Hit Rate from L1TEX	32.48%
L2 Read Hit Rates	32.45%
L2 Throughput	25.87%
PD Throughput	0.00%
PES+VPC Throughput	0.00%
Pixel Warps	0.00%
ZROP Throughput	0.00%

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



Unit Throughputs:
VRAM Read: 54%
VRAM Write: 2%

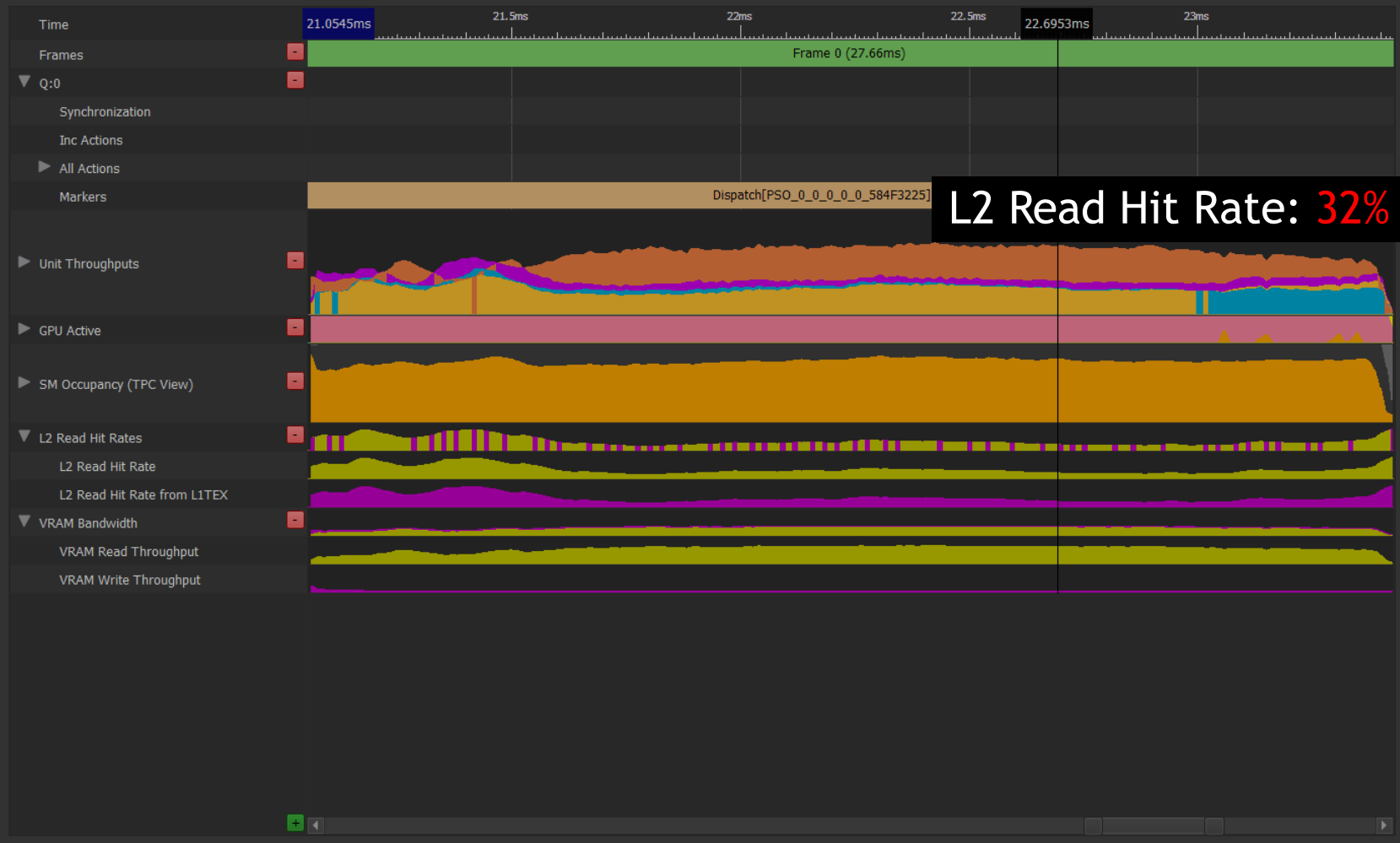
Summary Metrics Detailed Metrics Capture Information

Start: 21.05 ms End: 23.43 ms Duration: 2.37 ms Range: All Visible

VRAM

Metrics	Value
VRAM Bandwidth	28.05%
VRAM Read Throughput	54.19%
VRAM Throughput	56.09%
VRAM Write Throughput	1.90%

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...

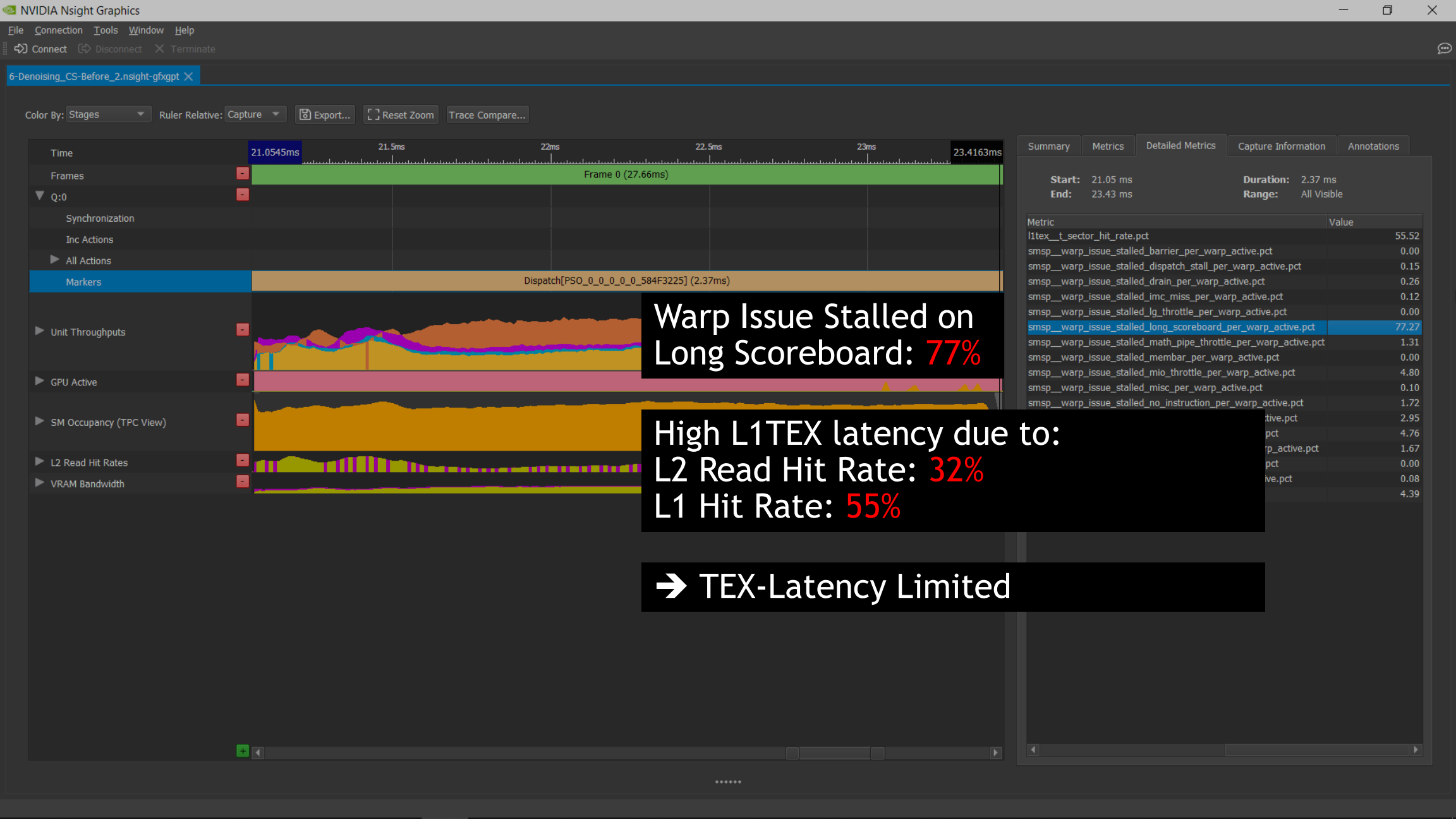


Summary Metrics Detailed Metrics Capture Information

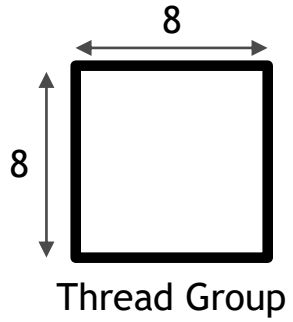
Start: 21.05 ms End: 23.43 ms Duration: 2.37 ms Range: All Visible

L2

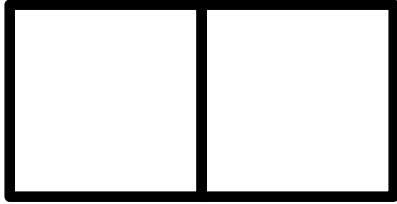
Metrics	Value
L2 Bandwidth from CROP	0.00%
L2 Bandwidth from LITEX	25.87%
L2 Bandwidth from ZROP	0.00%
L2 Read Hit Rate	32.42%
L2 Read Hit Rate from LITEX	32.48%
L2 Read Hit Rates	32.45%
L2 Throughput	25.87%



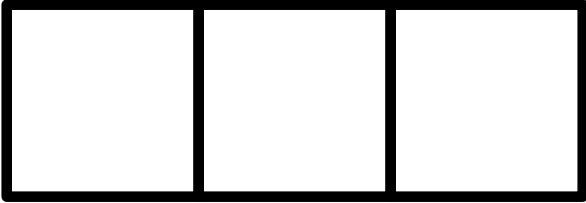
THREAD-GROUP LAUNCH ORDER



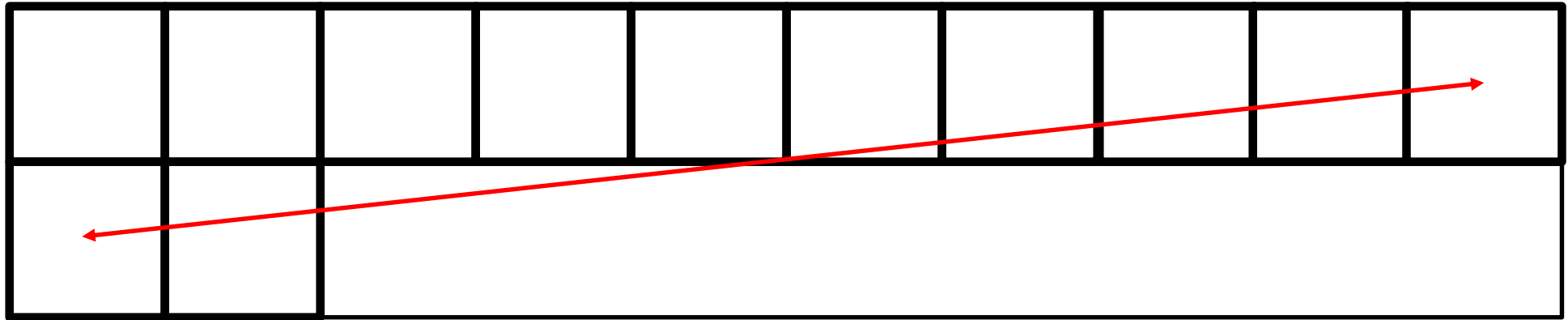
THREAD-GROUP LAUNCH ORDER



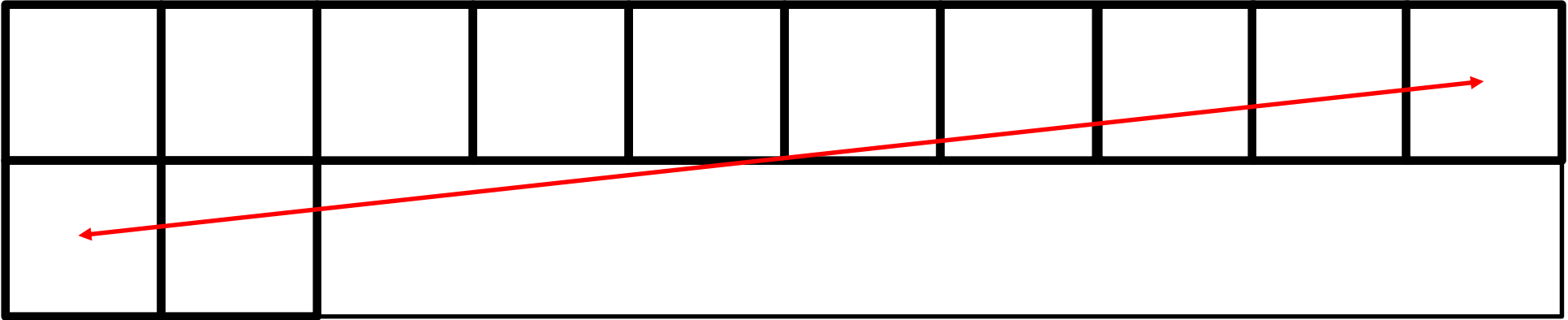
THREAD-GROUP LAUNCH ORDER



THREAD-GROUP LAUNCH ORDER



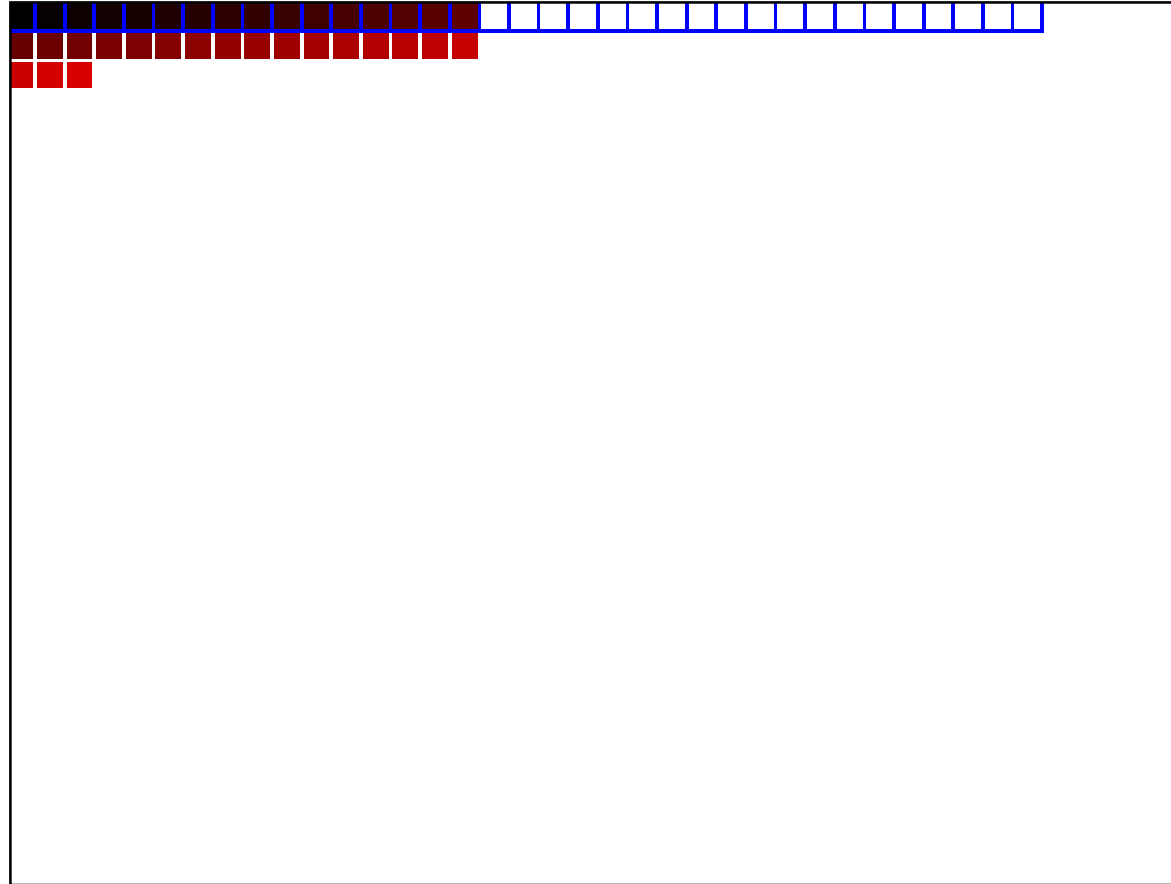
THREAD-GROUP LAUNCH ORDER



Thread groups launched sequentially have far-apart pixel coordinates

- ➔ Increases the size of the working set in L2
- ➔ Can cause poor L2 hit rate and long TEX miss latencies

THREAD-GROUP TILING



Divide the Dispatch grid into tiles of width= N

```

uint N = 16;
uint vThreadGroupIDFlattened = (Dispatch_Grid_Dim.x)*groupId.y + groupId.x;
uint Total_number_of_ThreadGroups_in_one_tile = N*(Dispatch_Grid_Dim.y);

uint Tile_ID_of_current_ThreadGroup = (vThreadGroupIDFlattened)/Total_number_of_ThreadGroups_in_one_tile;
uint Local_ThreadGroup_ID_flattened_within_current_tile =
(vThreadGroupIDFlattened)%Total_number_of_ThreadGroups_in_one_tile;
uint Local_ThreadGroup_ID_y_within_current_tile = (Local_ThreadGroup_ID_flattened_within_current_tile)/N;
uint Local_ThreadGroup_ID_x_within_current_tile = (Local_ThreadGroup_ID_flattened_within_current_tile)%N;

uint Tiled_vThreadGroupIDFlattened = Tile_ID_of_current_ThreadGroup*N +
Local_ThreadGroup_ID_y_within_current_tile*(Dispatch_Grid_Dim.x) + Local_ThreadGroup_ID_x_within_current_tile;

uint2 TiledvThreadGroupID;
TiledvThreadGroupID.y = Tiled_vThreadGroupIDFlattened/Dispatch_Grid_Dim.x;
TiledvThreadGroupID.x = Tiled_vThreadGroupIDFlattened%Dispatch_Grid_Dim.x;

uint2 TiledvThreadID;
TiledvThreadID.x = (ThreadGroup_Dim.x)*TiledvThreadGroupID.x + groupThreadIndex.x;
TiledvThreadID.y = (ThreadGroup_Dim.y)*TiledvThreadGroupID.y + groupThreadIndex.y;

```

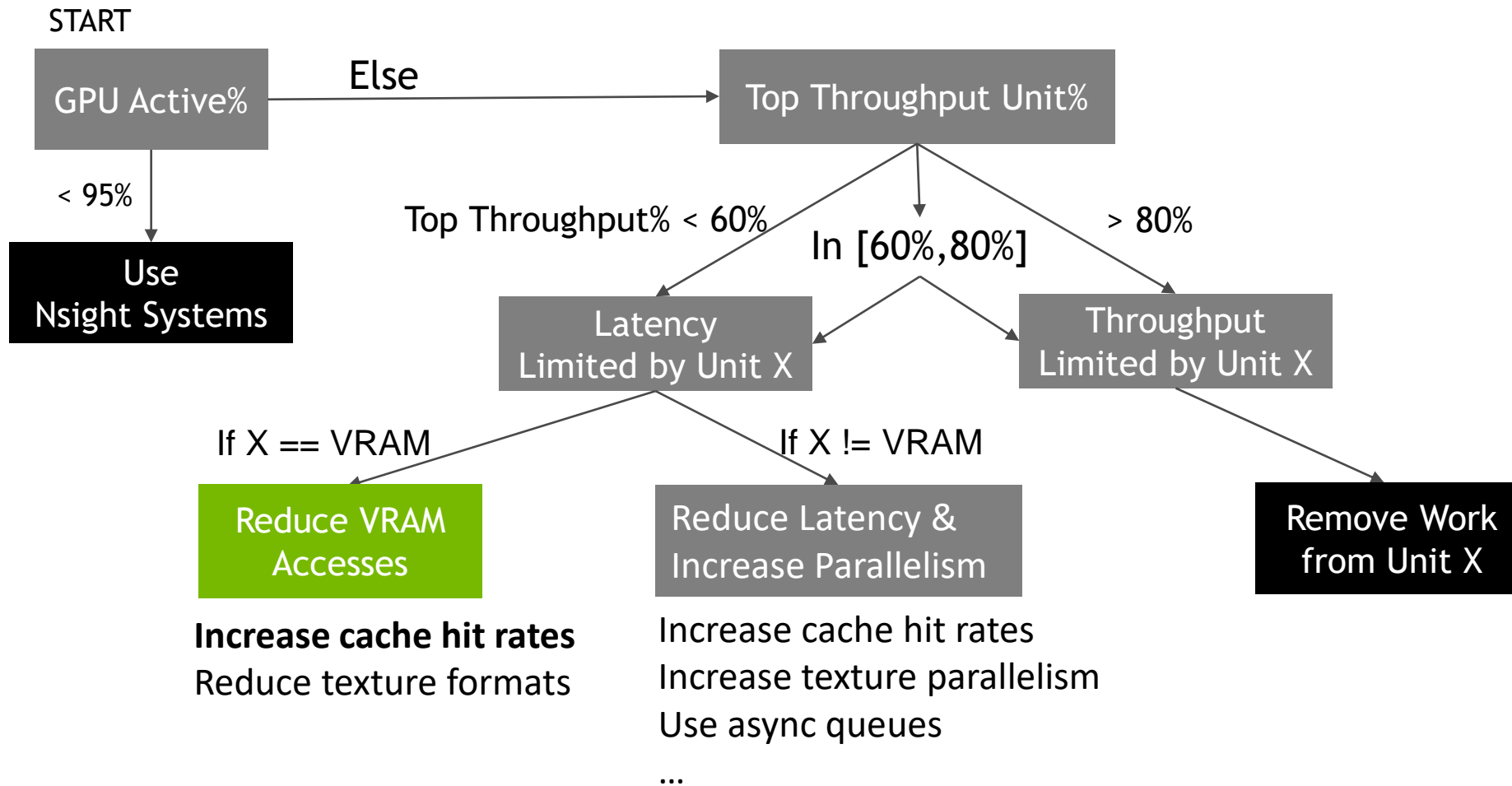
THREAD-GROUP TILING

Tile Size: [16, Dispatch_Grid_Dim.y]

	BEFORE	AFTER	RATIO
GPU Elapsed Time	2.36 ms	1.61 ms	1.47x Gain
Throughput: VRAM	56.3%	30.4%	0.54x
Throughput: SM	33.8%	51.3%	1.52x
L2 Read Hit Rate	62.5%	85.8%	1.37x
SM Warp Stalls on long_scoreboard	77.8%	58.6%	0.75x

On RTX 2080 with SetStablePowerState(TRUE)

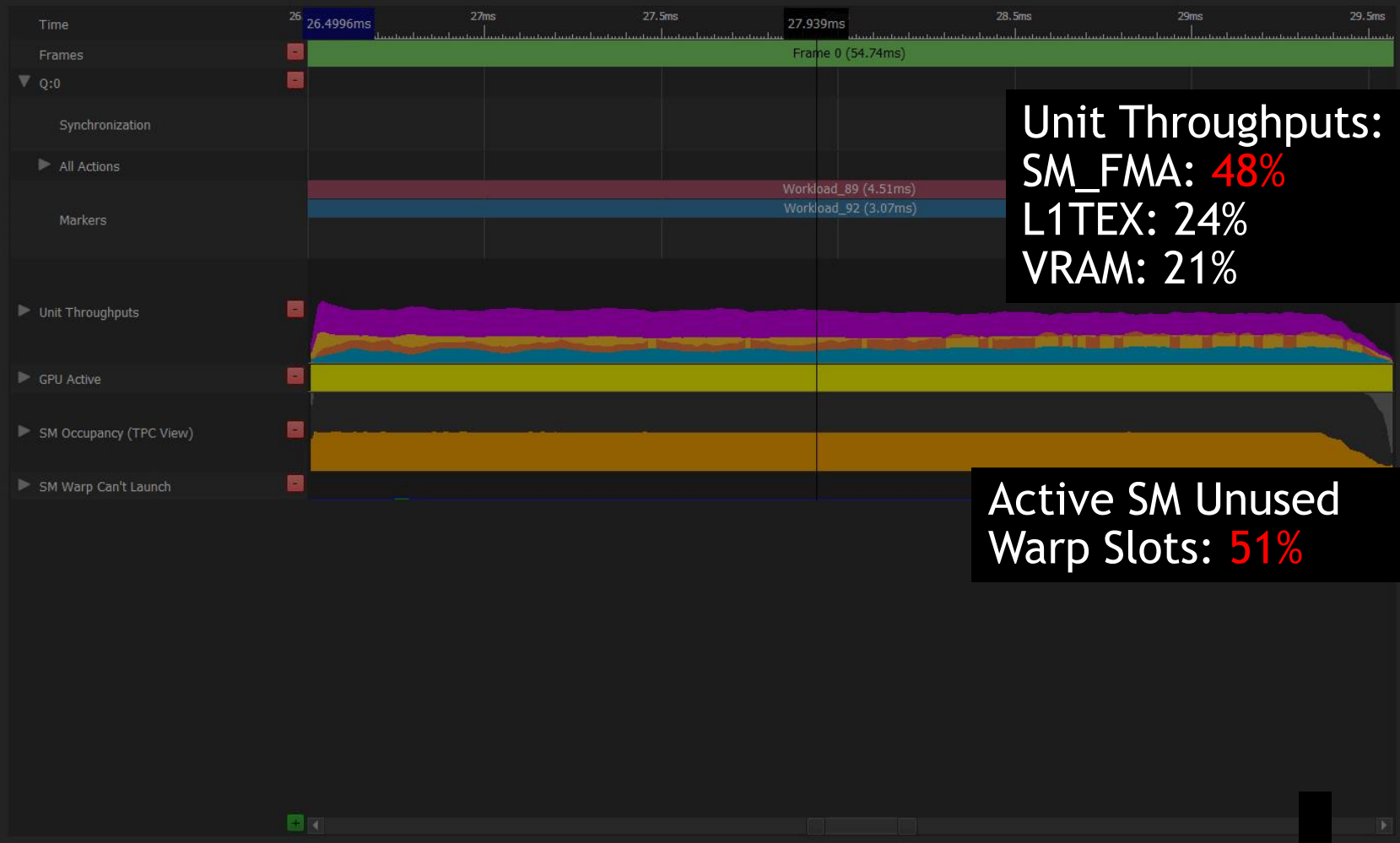
THE P3 (PEAK-PERF%) METHOD





Case Study #5:
Lighting CS using Shared Memory

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



Unit Throughputs:
SM_FMA: 48%
L1TEX: 24%
VRAM: 21%

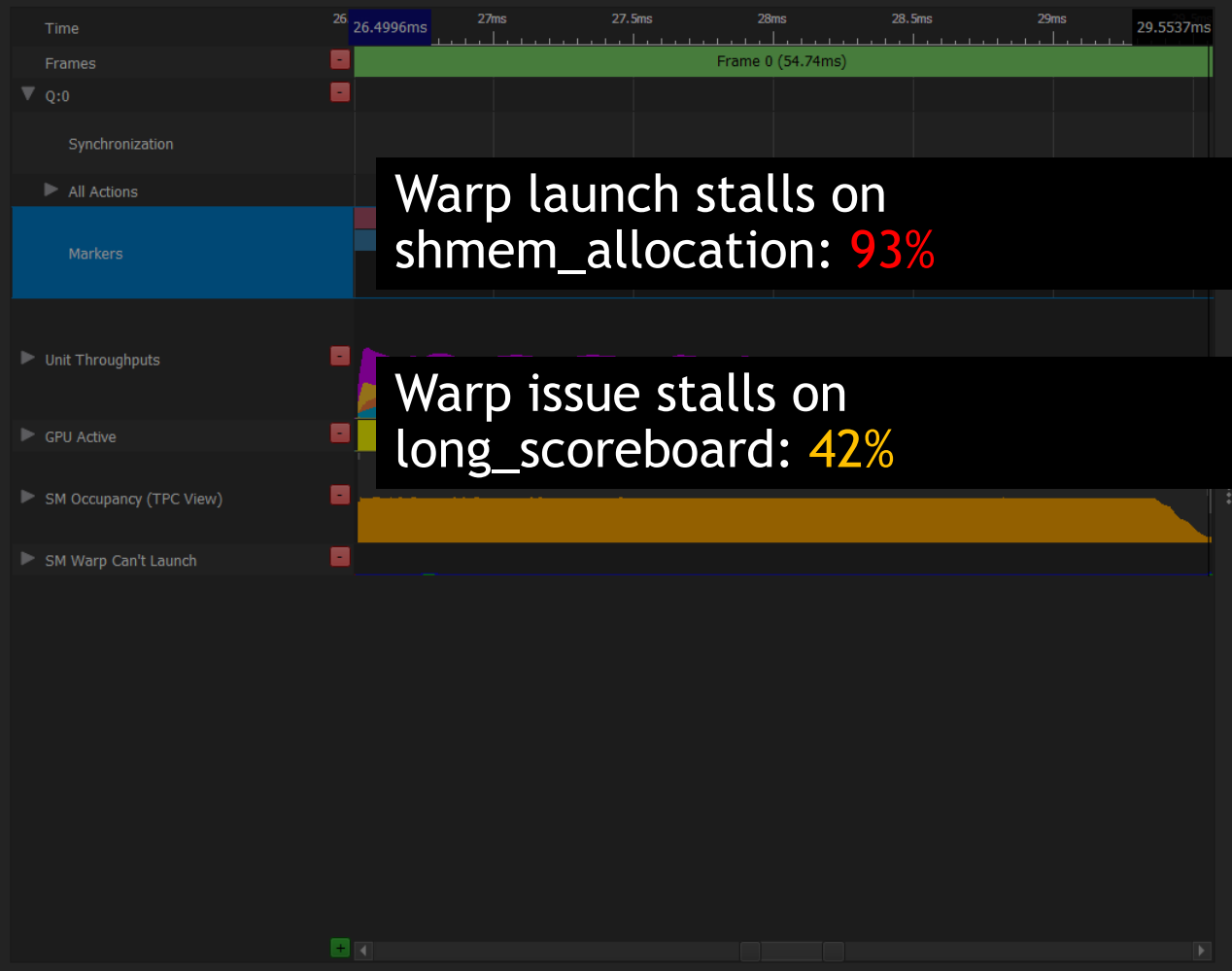
Active SM Unused Warp Slots: 51%

Start:	End:	Duration:	Range:
26.50 ms	29.57 ms	3.07 ms	All Visible

Unit Throughput	Value
SM Instruction Throughput	47.55%
SM FMA Pipe Throughput	47.53%
SM Issue Active	44.78%
L1TEX Throughput	23.92%
SM ALU Pipe Throughput	22.02%
VRAM Throughput	20.84%
SM SFU Pipe Throughput	15.82%
L2 Throughput	11.92%
PROP Throughput	0.00%
PD Throughput	0.00%
PES+VPC Throughput	0.00%
RASTER Throughput	0.00%
ZROP Throughput	0.00%
CROP Throughput	0.00%
SM FP16+Tensor Pipe Throughput	0.00%

Warp Occupancy	Value
Idle SM Unused Warp Slots	0.86%
Active SM Unused Warp Slots	51.05%
Compute Warps	48.09%
Pixel Warps	0.00%
Vertex+Tess+Geom Warps	0.00%

Color By: Stages Ruler Relative: Capture Export... Reset Zoom Trace Compare...



Summary Metrics Detailed Metrics Capture Information Annotations

Start: 26.50 ms Duration: 3.07 ms
 End: 29.57 ms Range: All Visible

Metric	Value
tpc_warp_launch_cycles_stalled_shader_cs_reason_barrier_allocation.avg.pct_of_peak_sustained_elapsed	0.00
tpc_warp_launch_cycles_stalled_shader_cs_reason_cta_allocation.avg.pct_of_peak_sustained_elapsed	0.00
tpc_warp_launch_cycles_stalled_shader_cs_reason_register_allocation.avg.pct_of_peak_sustained_elapsed	1.19
tpc_warp_launch_cycles_stalled_shader_cs_reason_shmem_allocation.avg.pct_of_peak_sustained_elapsed	93.00
tpc_warp_launch_cycles_stalled_shader_cs_reason_warp_allocation.avg.pct_of_peak_sustained_elapsed	0.00
l1tex_t_sector_hit_rate.pct	86.46
smsp_warp_issue_stalled_barrier_per_warp_active.pct	0.13
smsp_warp_issue_stalled_dispatch_stall_per_warp_active.pct	0.12
smsp_warp_issue_stalled_drain_per_warp_active.pct	0.01
smsp_warp_issue_stalled_imc_miss_per_warp_active.pct	3.80
smsp_warp_issue_stalled_lg_throttle_per_warp_active.pct	0.00
smsp_warp_issue_stalled_long_scoreboard_per_warp_active.pct	41.75
smsp_warp_issue_stalled_math_pipe_throttle_per_warp_active.pct	4.99
smsp_warp_issue_stalled_membar_per_warp_active.pct	0.00
smsp_warp_issue_stalled_mio_throttle_per_warp_active.pct	0.64
smsp_warp_issue_stalled_misc_per_warp_active.pct	0.16
smsp_warp_issue_stalled_no_instruction_per_warp_active.pct	7.22
smsp_warp_issue_stalled_not_selected_per_warp_active.pct	6.60
smsp_warp_issue_stalled_selected_per_warp_active.pct	11.79
smsp_warp_issue_stalled_short_scoreboard_per_warp_active.pct	4.43
smsp_warp_issue_stalled_sleeping_per_warp_active.pct	0.00
smsp_warp_issue_stalled_tex_throttle_per_warp_active.pct	0.00
smsp_warp_issue_stalled_wait_per_warp_active.pct	20.22

SHARED-MEM-SIZE REDUCTION

Before: store light data into shared mem

```
static groupshared LightData SharedMem[MaxLightCount];
```

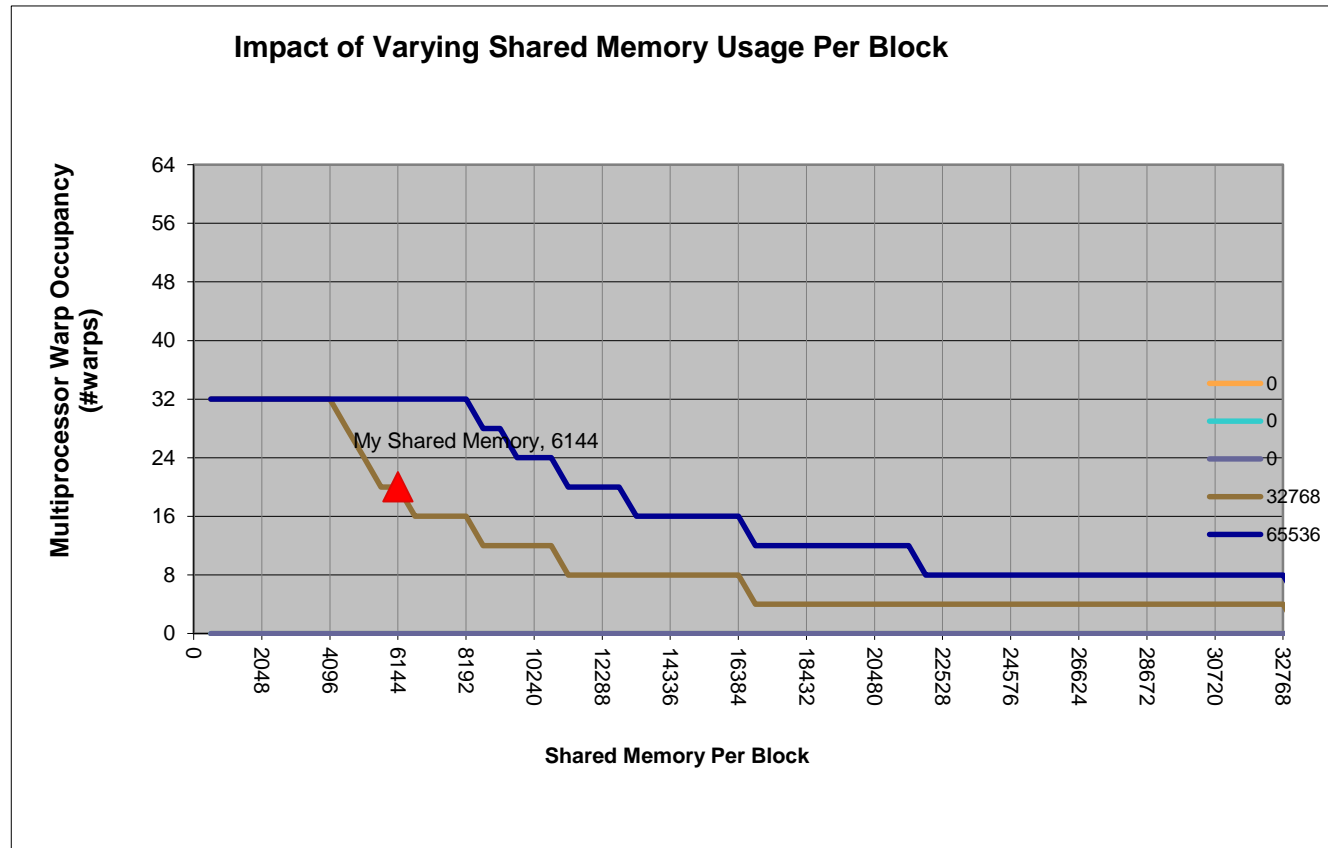
...

After: store light indices into shared mem

```
static groupshared uint SharedMem[MaxLightCount];
```

And load light data via non-divergent indexed constant-buffer loads

CUDA OCCUPANCY CALCULATOR



For Turing RTX GPUs (Compute Capability 7.5)

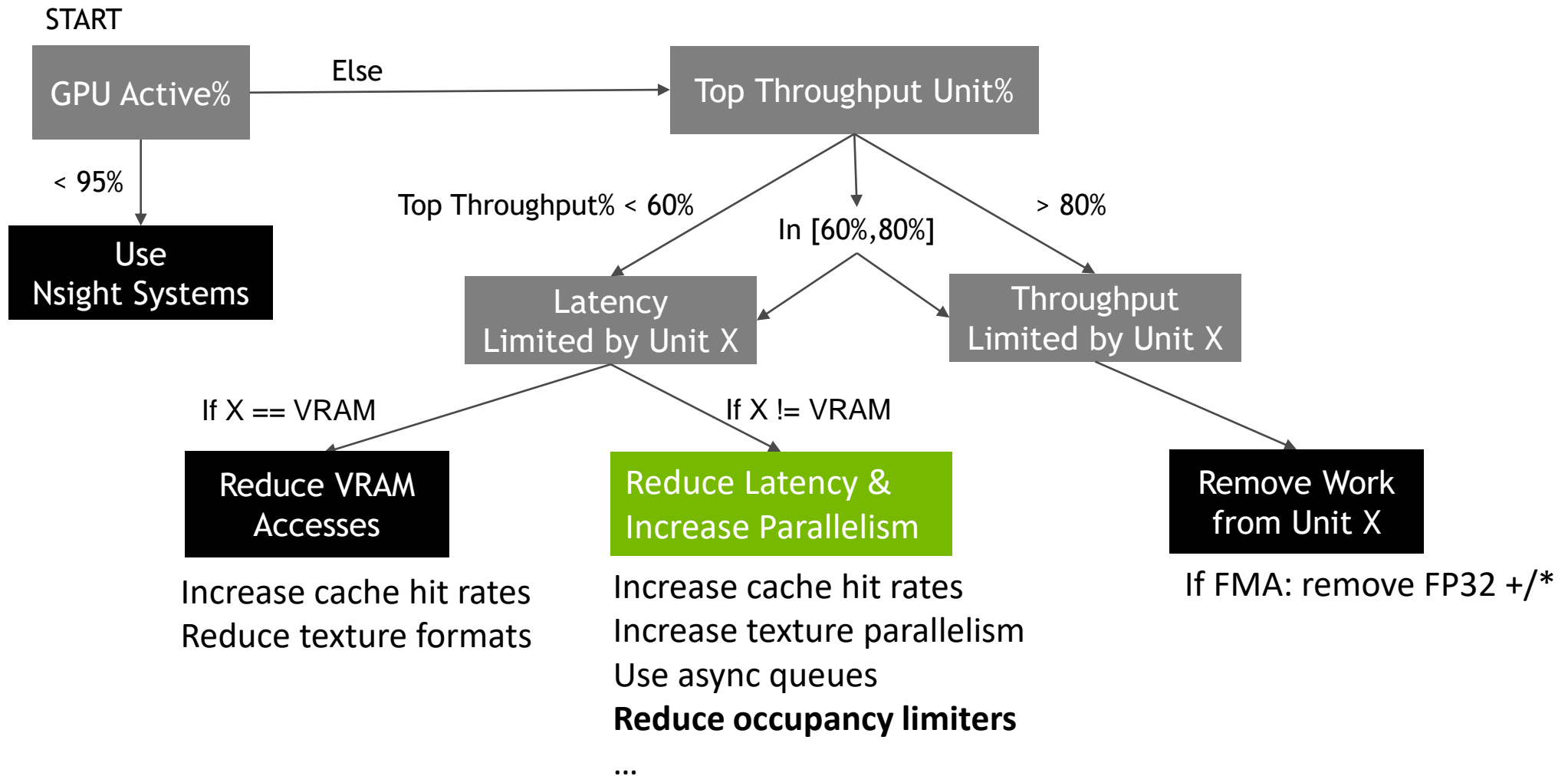
SHARED-MEM-SIZE REDUCTION:

6,144 -> 256 bytes per thread group

	BEFORE	AFTER	RATIO
GPU Elapsed Time	3.06 ms	2.21 ms	1.38x Gain
SM Throughput	47.5%	66.1%	1.39x
Warp launch stalls on shmem_allocation	94%	0%	+INF
Warp issue stalls on long_scoreboard	41.7%	36.7%	1.14x

On RTX 2080 with SetStablePowerState(TRUE)

THE P3 (PEAK-PERF%) METHOD





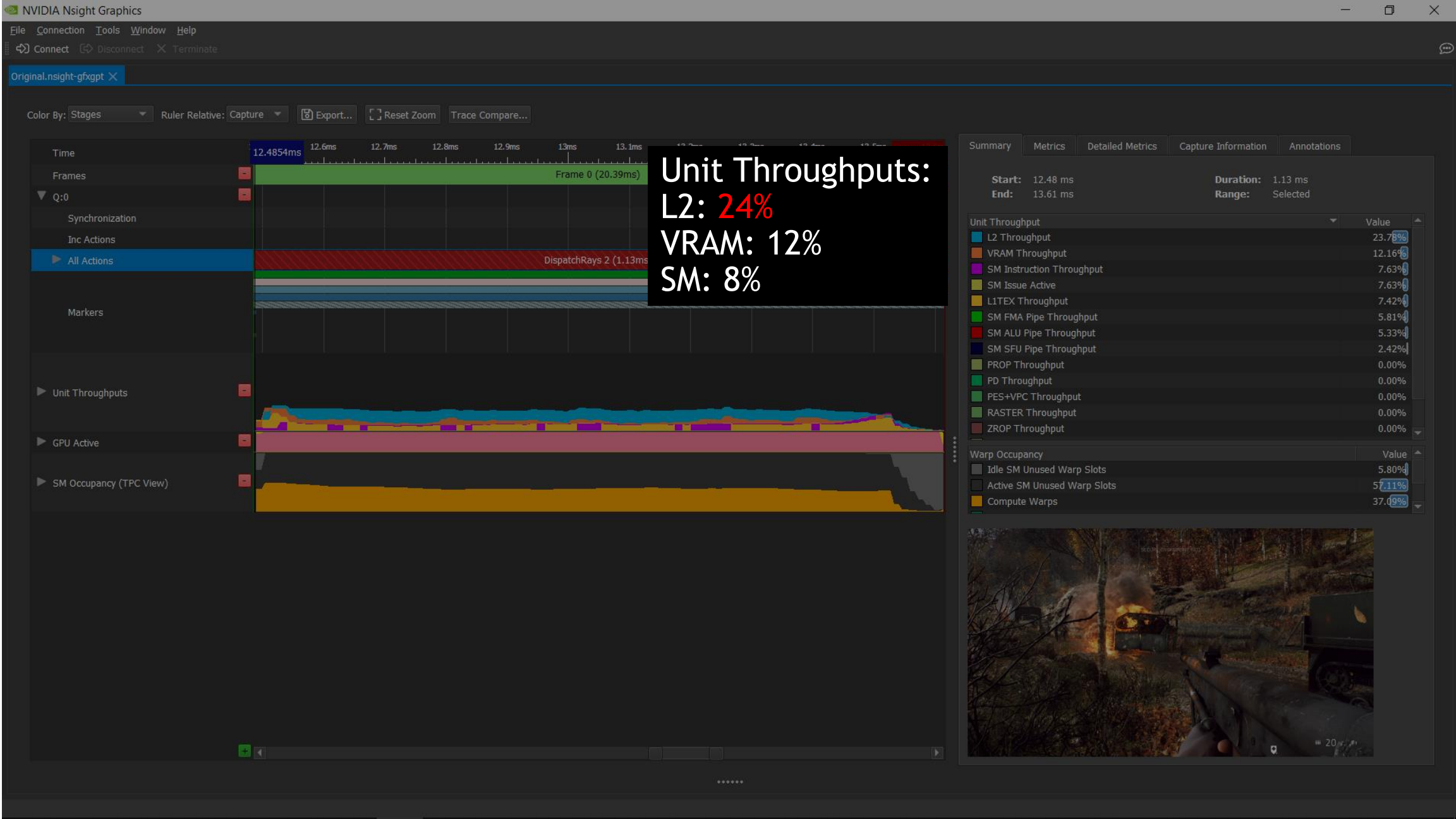
Case Study #6:
Thread-Divergence-Limited
DispatchRays in Battlefield V

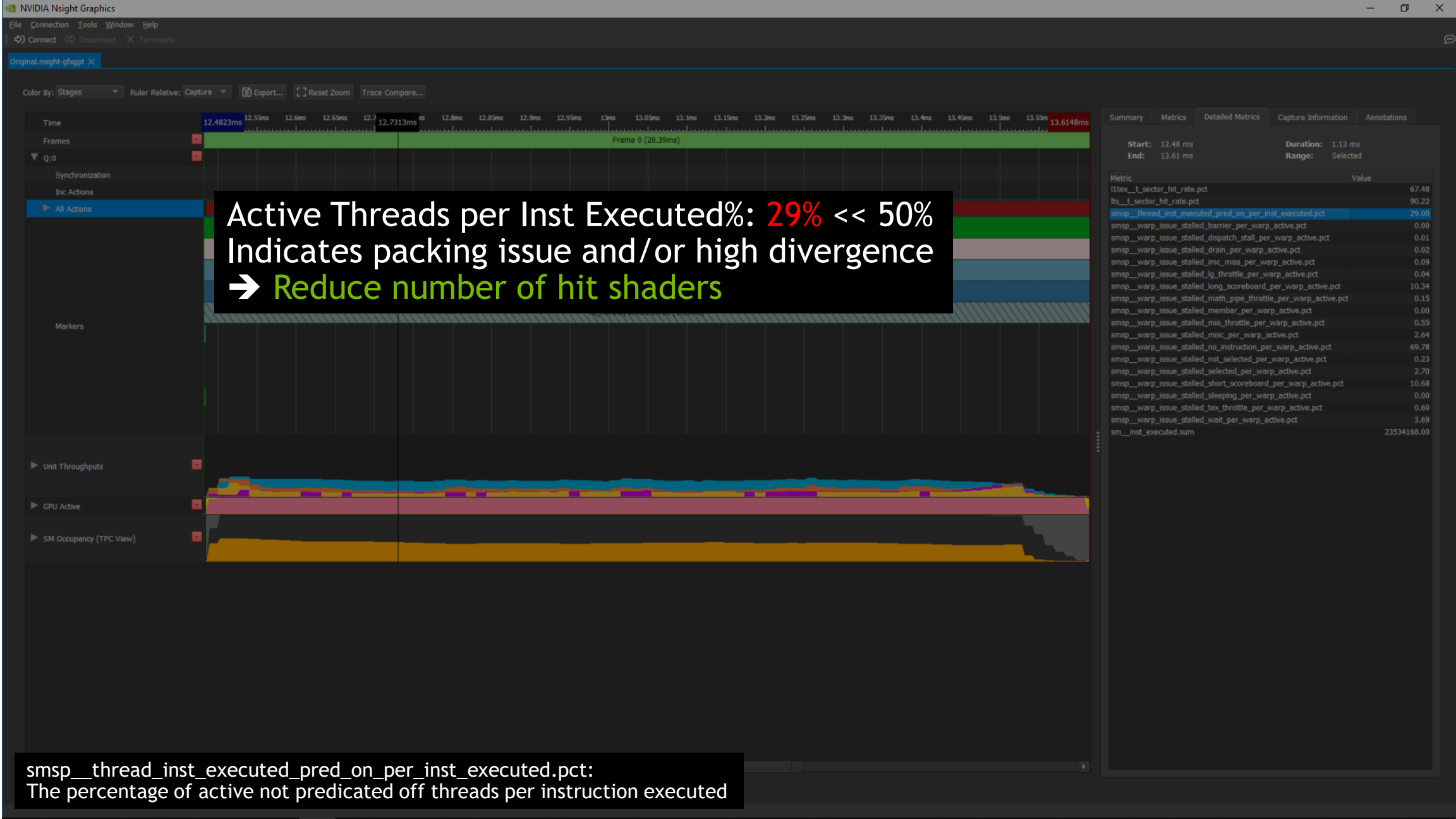
S
SECURE CHOKEPOINT KILO



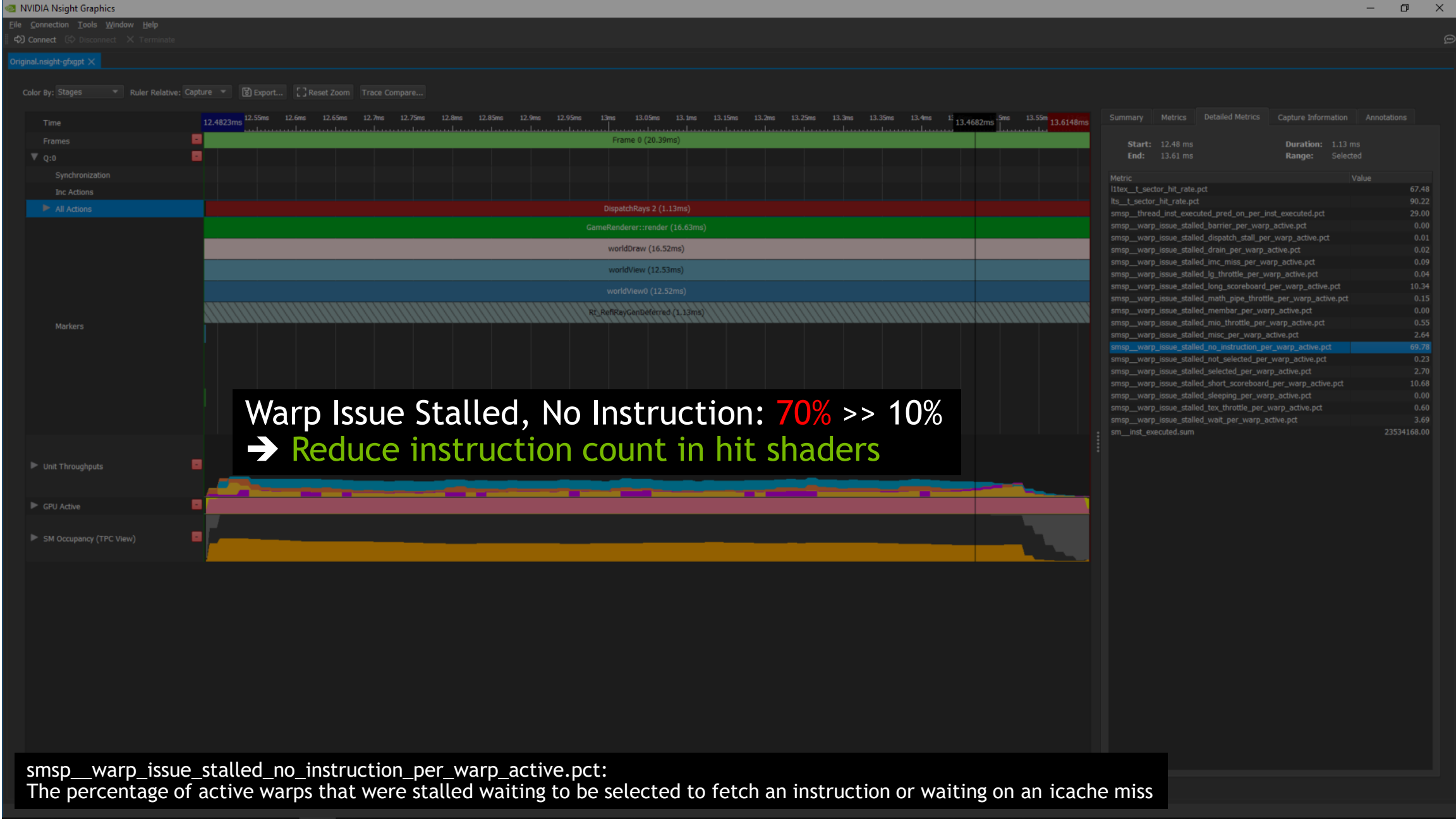
20/40 3







smsp__thread_inst_executed_pred_on_per_inst_executed.pct:
 The percentage of active not predicated off threads per instruction executed



smsp__warp_issue_stalled_no_instruction_per_warp_active.pct:
 The percentage of active warps that were stalled waiting to be selected to fetch an instruction or waiting on an icache miss

THE OPAQUE FLAGS

- ▶ D3D12_RAYTRACING_GEOMETRY_FLAG_OPAQUE
 - ▶ D3D12_RAYTRACING_INSTANCE_FLAG_FORCE_OPAQUE
 - ▶ RAY_FLAG_FORCE_OPAQUE
- ▶ Indicates whether triangles are fully opaque or not
 - ▶ Stop rays earlier
 - ▶ Helps reduce the number of evaluated materials

Battlefield V DXR
Original OPAQUE flags

DispatchRays: 1.13 ms

GeForce RTX 2080 +
SetStablePowerState



Battlefield V DXR
+OPAQUE flag for ALL geometries

DispatchRays: 0.58 ms (1.95x)

GeForce RTX 2080 +
SetStablePowerState



SECURE CHOKEPOINT KILO

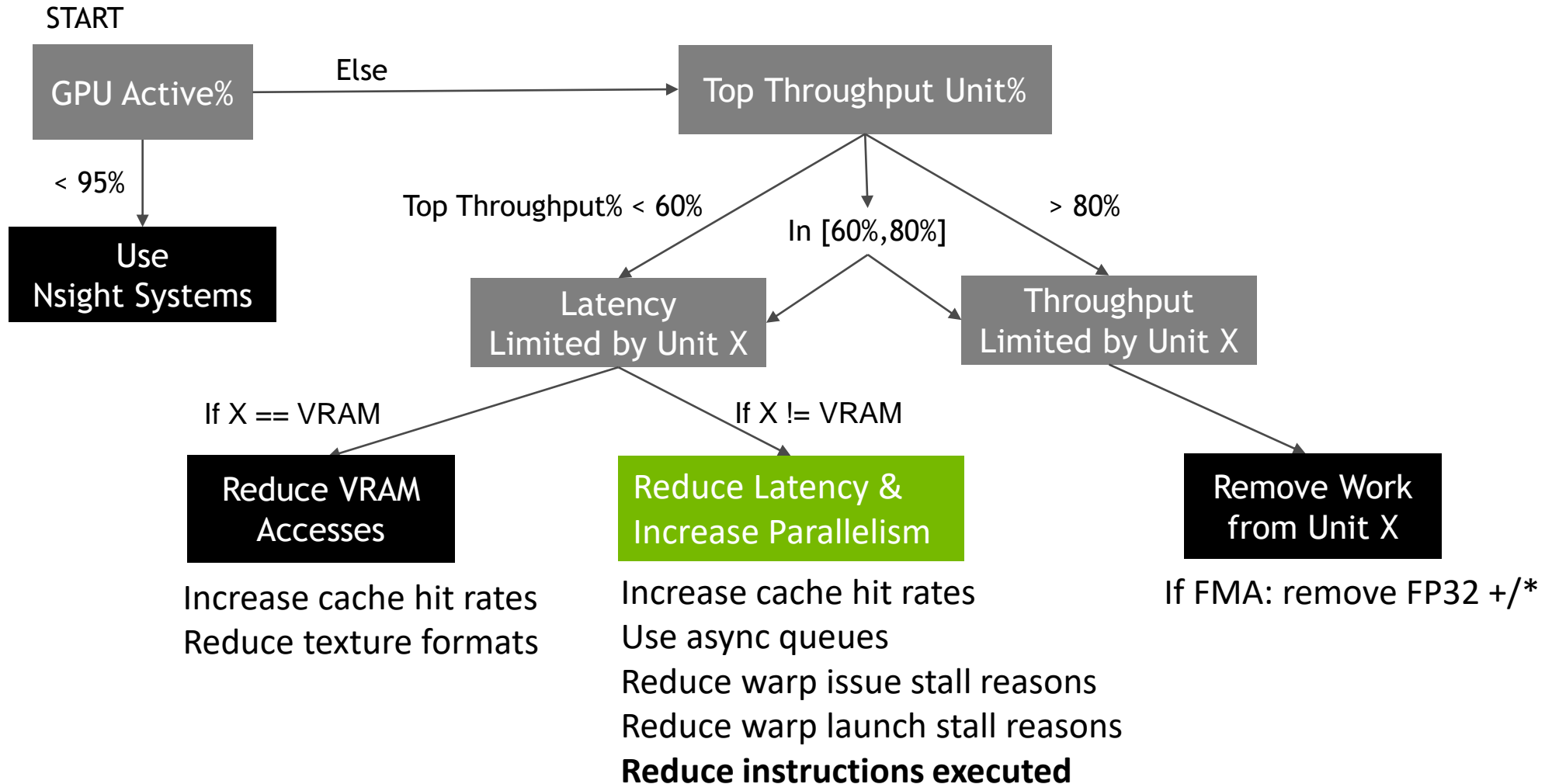
20 40 3

FORCE_OPAQUE EXPERIMENT

	BEFORE	AFTER	RATIO
GPU Elapsed Time	1.13 ms	0.58 ms	1.95x Gain
Top Throughput: L2	23.8%	24.0%	1.01x
SM Active Threads Per Instruction Executed	29.0%	40.5%	1.40x
SM Instructions Executed Per Warp	725.1	425.2	0.59x
SM Warp Issue Stalls, No Instruction	69.8%	68.2%	0.98x

On RTX 2080 with SetStablePowerState(TRUE)

THE P3 (PEAK-PERF%) METHOD



NVIDIA TOOLS FOR GPU PROFILING

Two Nsight: Graphics modules

	GPU TRACE METRICS	RANGE PROFILER
Graphs over time	Y	N
Async queue support	Unbiased	Serialized
APIs	DX12 only so far	DX12, DX11 and Vk
GPUs	>= Turing	>= Kepler

CONCLUSION

The P3 Method (Peak-Perf-Percentage)

- ▶ A **method** to triage the performance of **any** GPU workload:
 - ▶ Start from the « Top Throughput% » Metrics (aka SOL% or Peak-Perf%)
 - ▶ Do NOT start from SM Warp Occupancy
 - ▶ We have been keeping [a blog post](#) up-to-date as tools and GPUs evolve.
- ▶ Async Compute rule of thumb:
 - ▶ Do not overlap 2 VRAM-latency-limited workloads



QUESTIONS?

Louis Bavoil | lbavoil@nvidia.com

www.nvidia.com/GDC