# Occupancy explained through the AMD RDNA™ architecture

**François Guthmann**

**Graphics Programming Conference**

# GPUOPEN

AMD GPUOpen

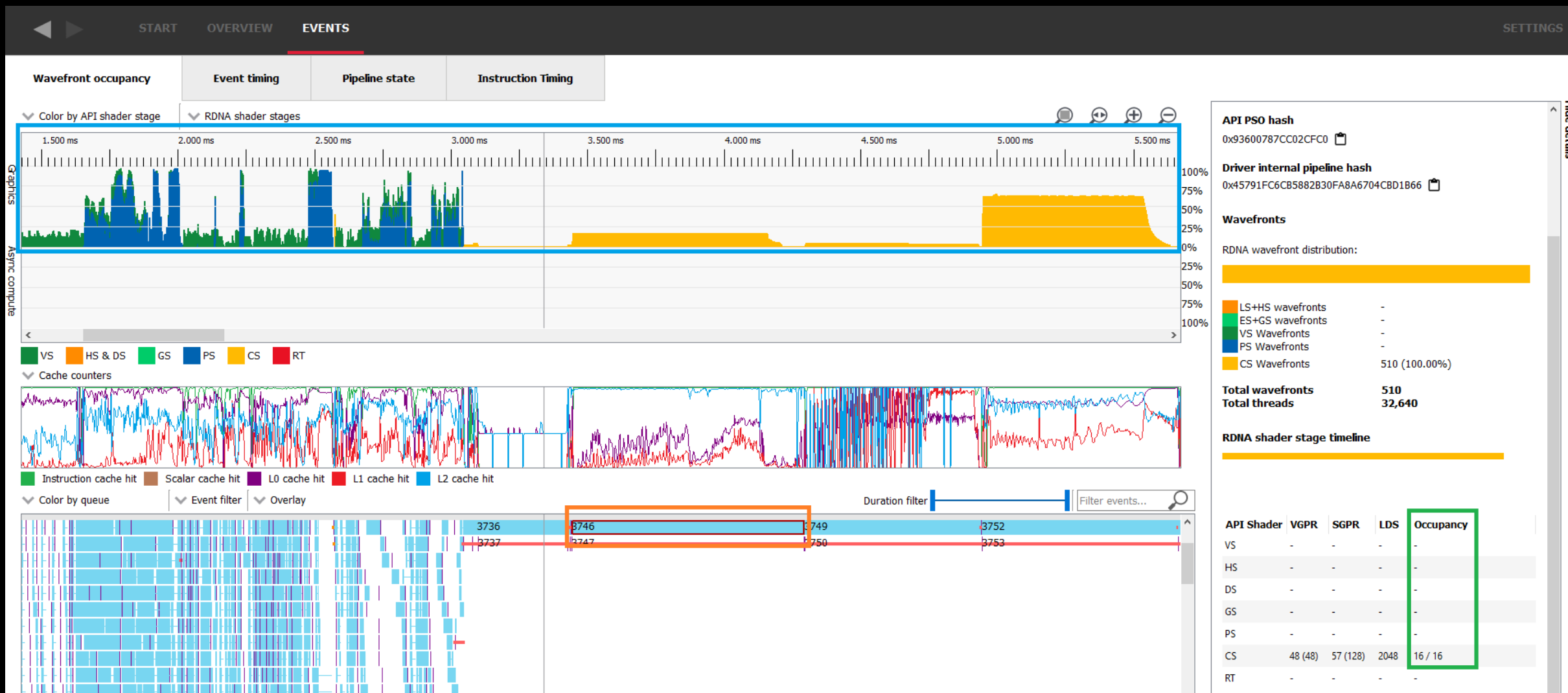## Occupancy explained

**François Guthmann**
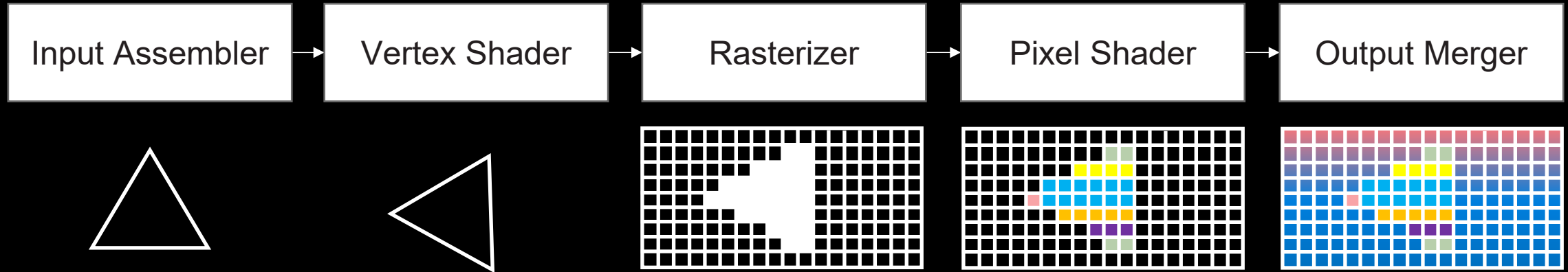
Originally posted December 20, 2023

---

If you're working with GPUs, chances are you've heard the term *occupancy* thrown around in the context of shader performance. You might have heard it helps hiding memory latency but are not sure exactly what that means. If that's the case, then you are exactly where you should be! In this blog post we will try to demystify what exactly this metric is. We will first talk a bit about the hardware architecture to understand where this metric is coming from. We will then explain the factors that can limit occupancy both statically at compile time and dynamically at run time. We will also help you identify occupancy-limited workloads using tools like the Radeon™ GPU Profiler and offer potential leads to alleviate the issues. Finally, the last section will try to summarize all the concepts touched upon in this post and offer practical solutions to practical problems.

This article however assumes you have a basic understanding of how to work with a GPU. Mainly, we expect you to know how to use the GPU from a graphics API perspective (draws, dispatches, barriers etc.) and that the workloads are executed in groups of threads on the GPU. We also expect you to know about the basic resources a shader uses like the scalar registers, vector registers, and shared memory.
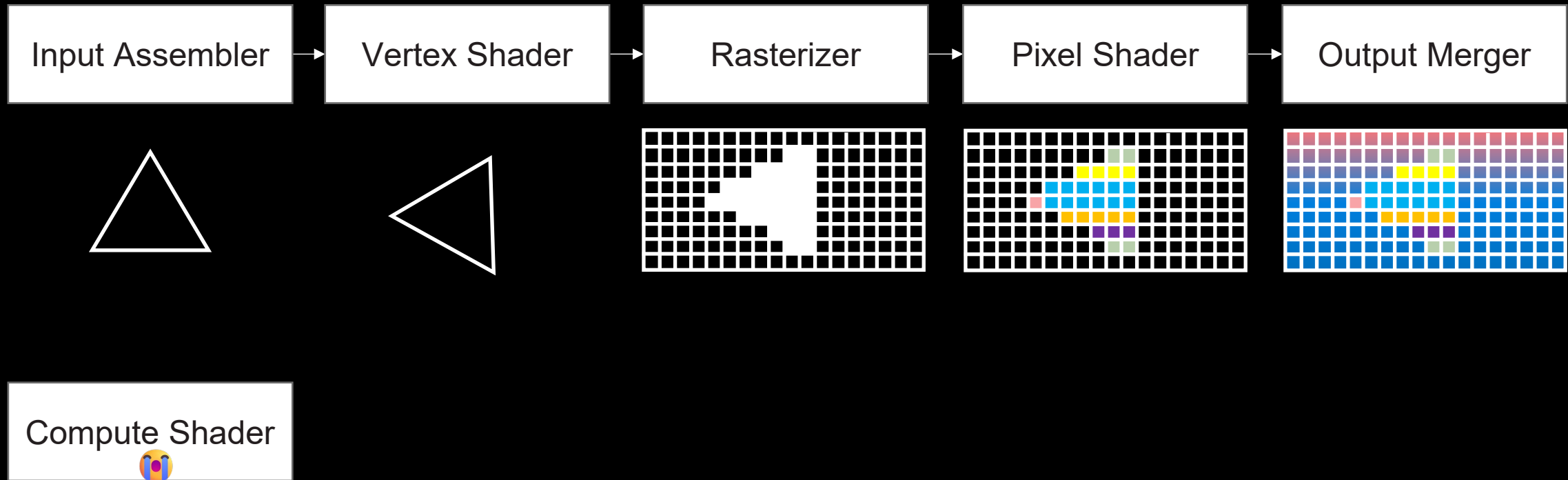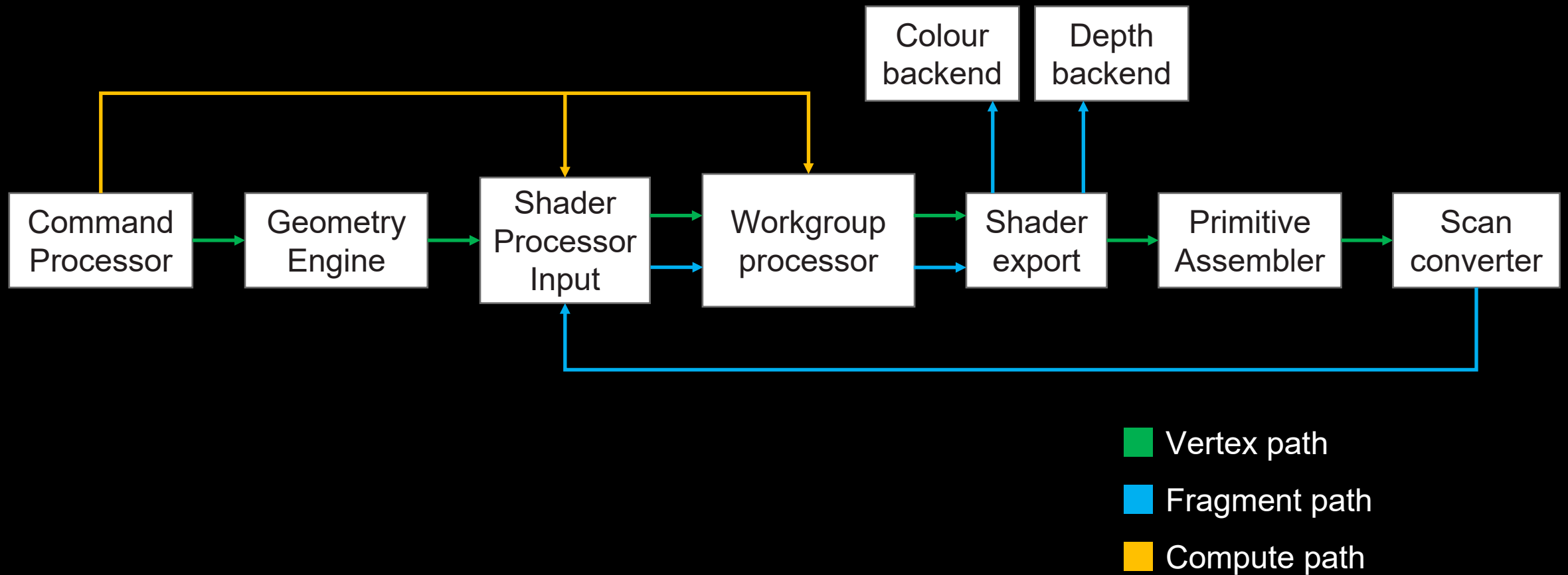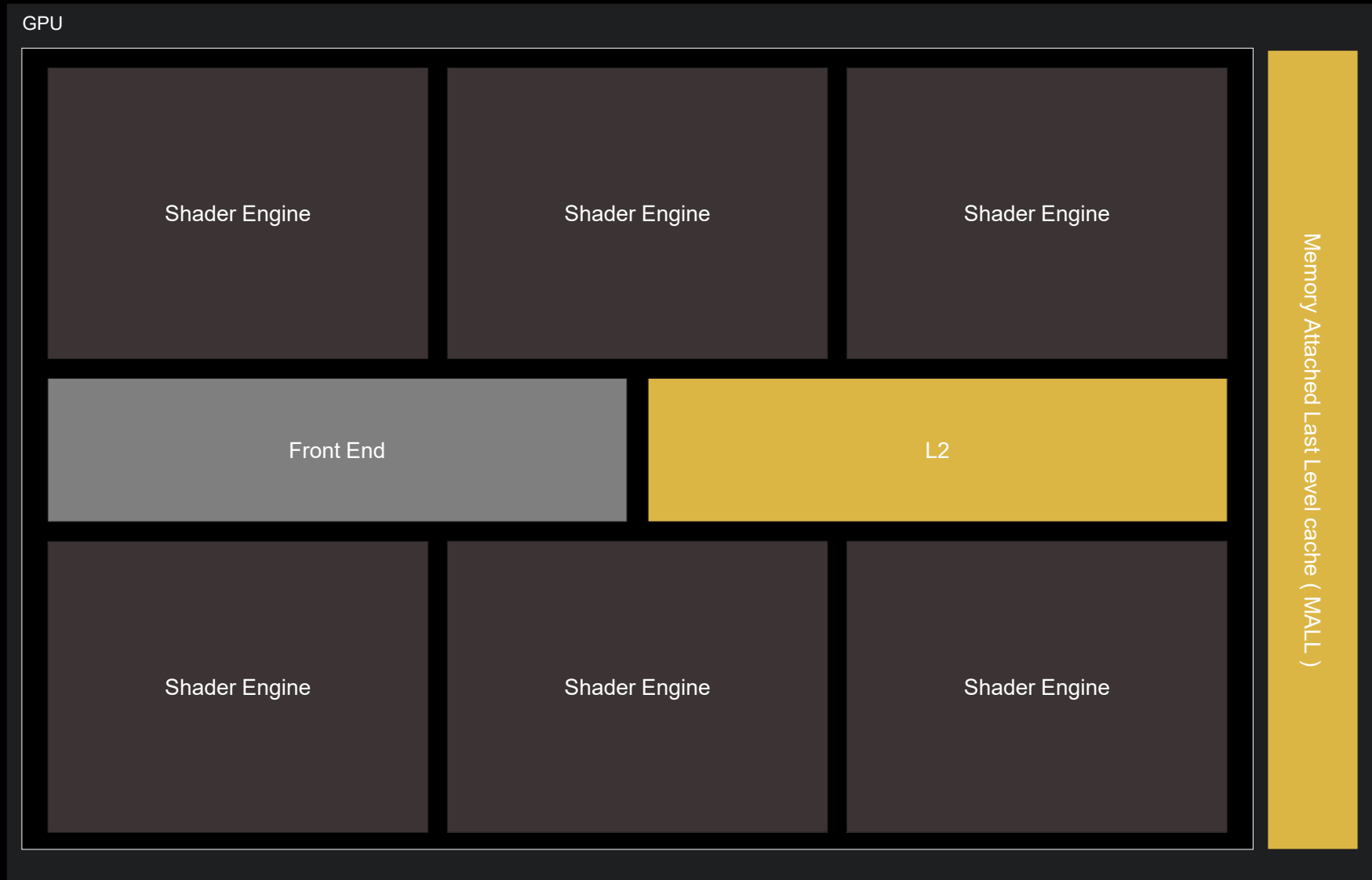
# OCCUPANCY?

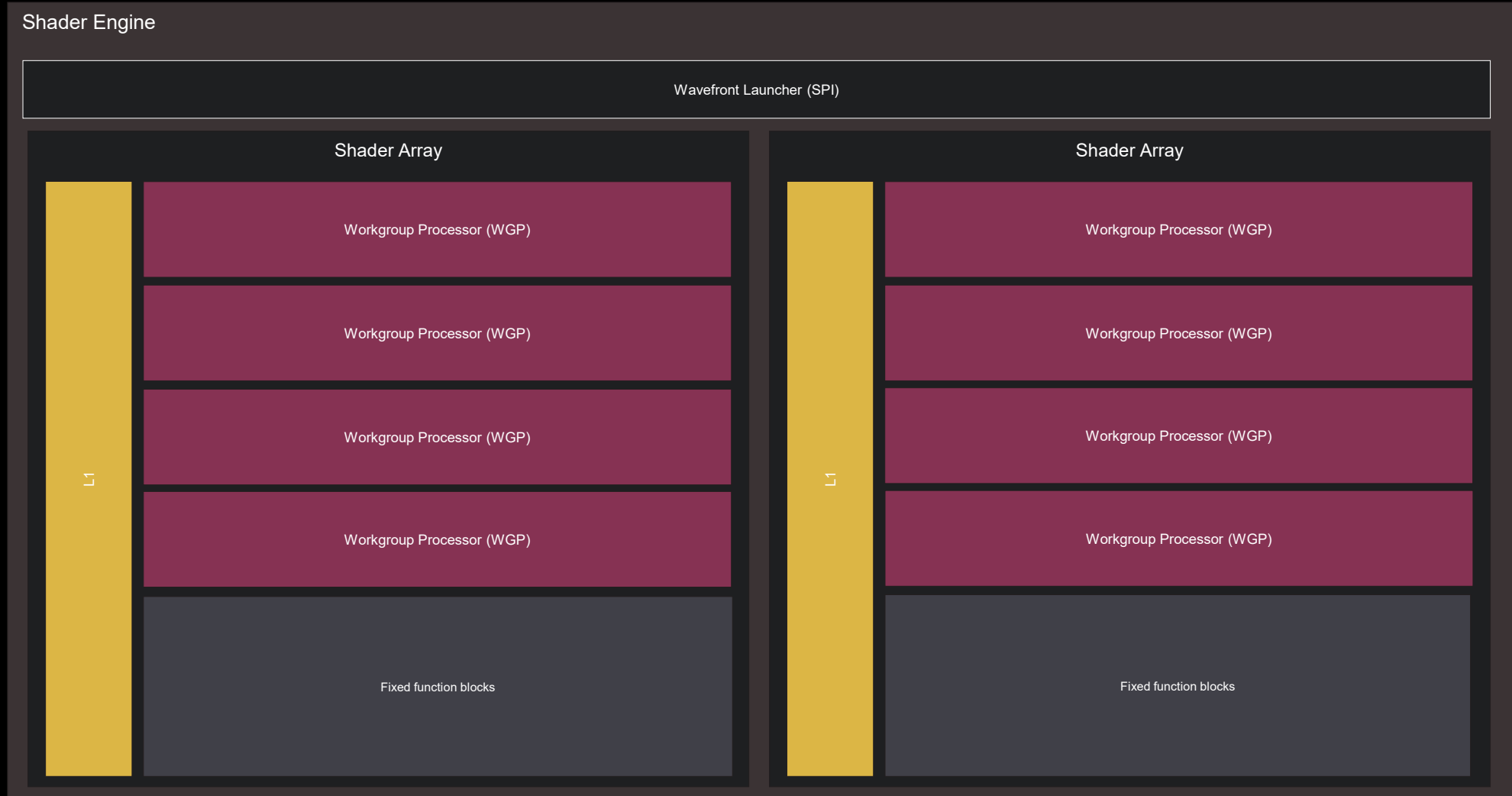# LOGICAL GRAPHICS PIPELINE

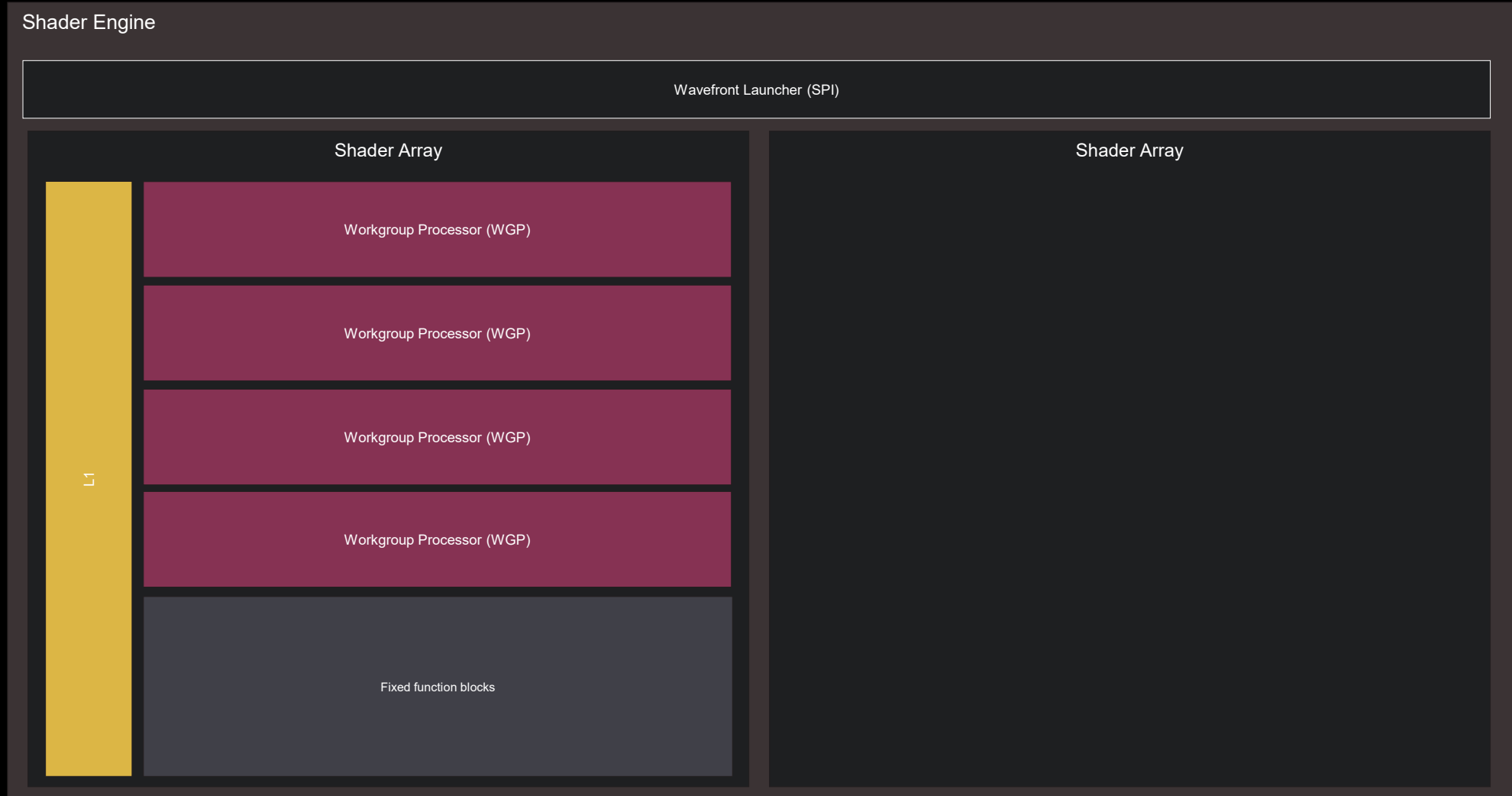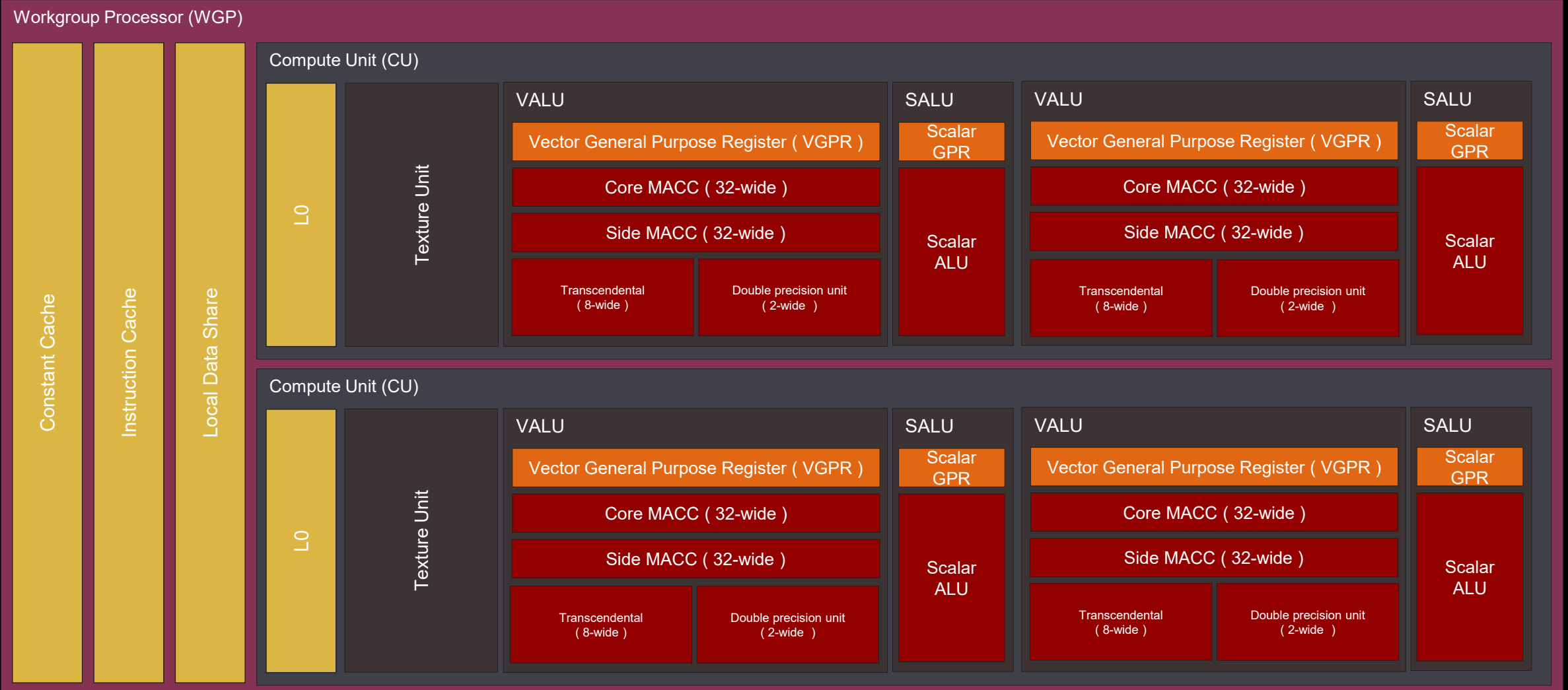# LOGICAL GRAPHICS PIPELINE

# HARDWARE GRAPHICS PIPELINE
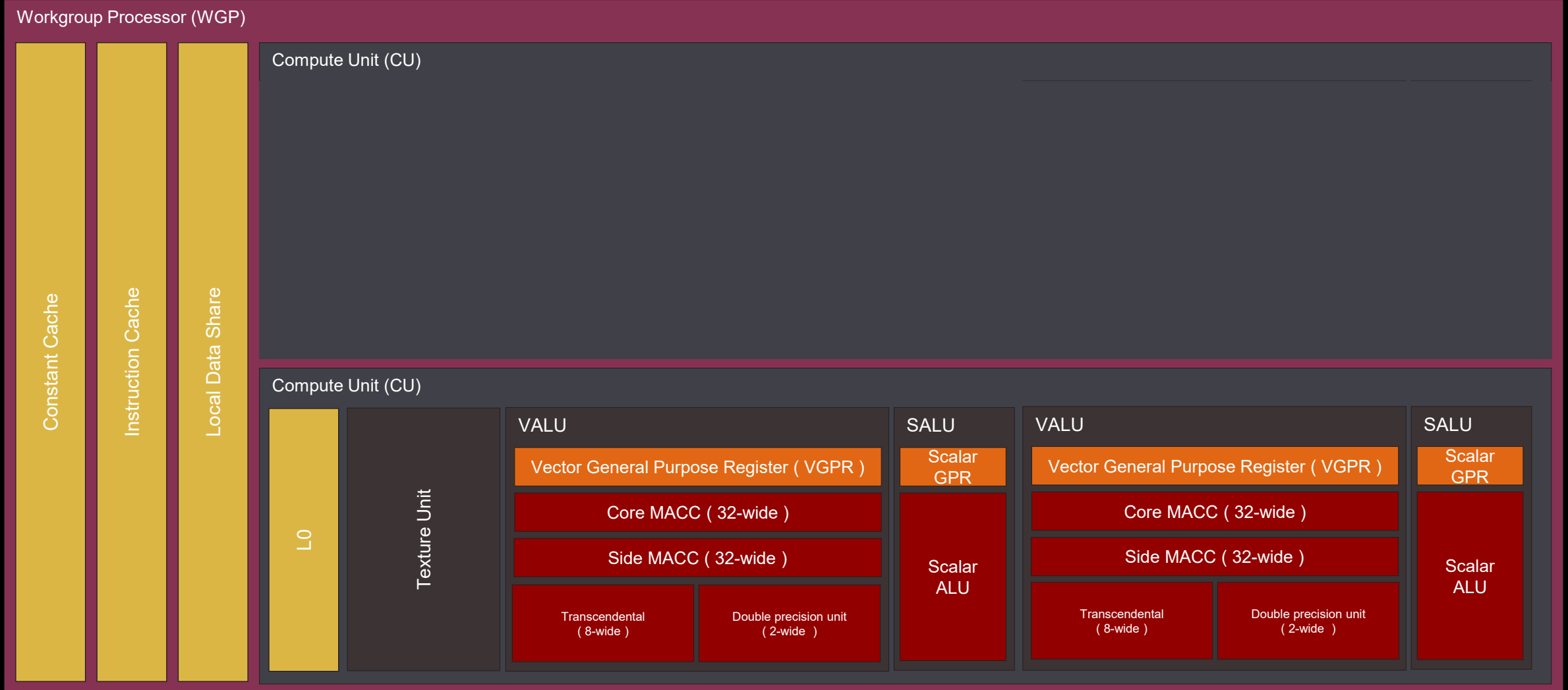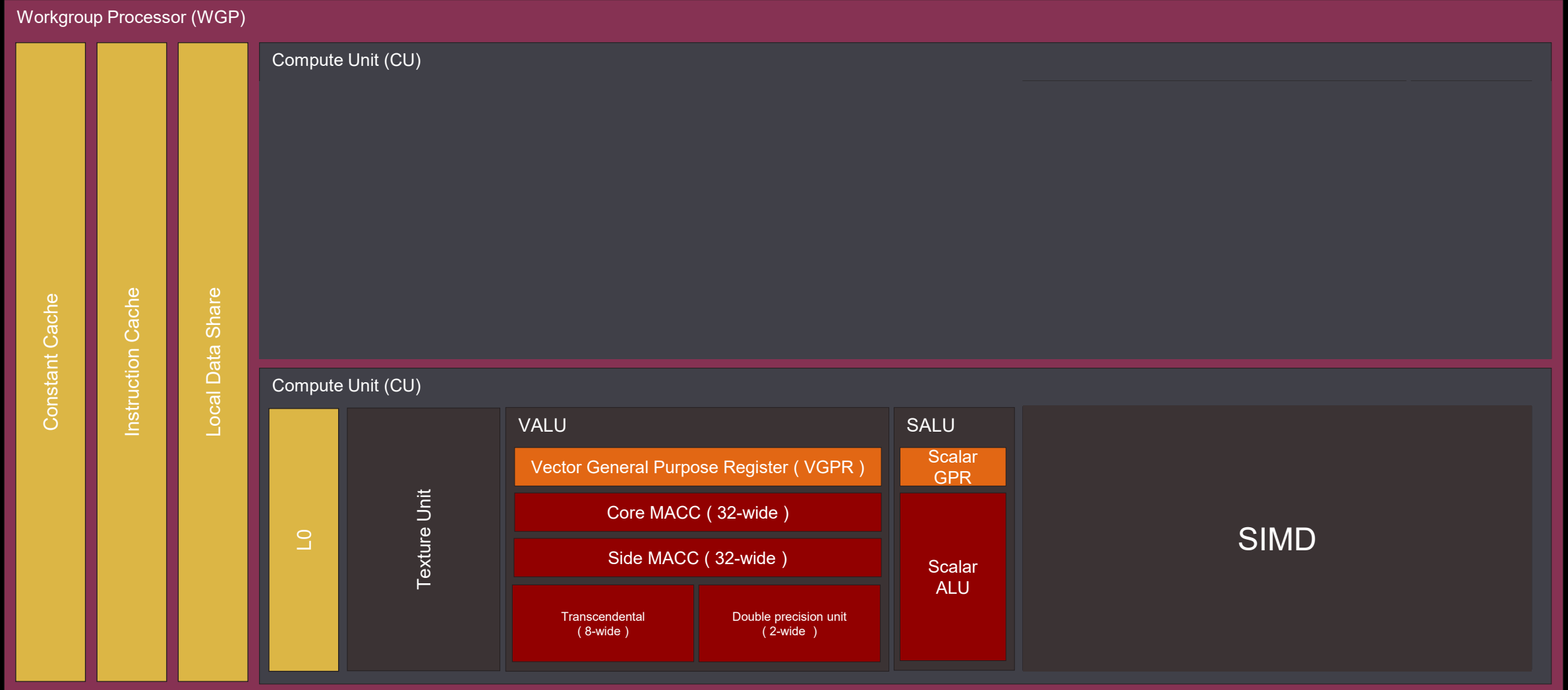
# HIGH LEVEL OVERVIEW
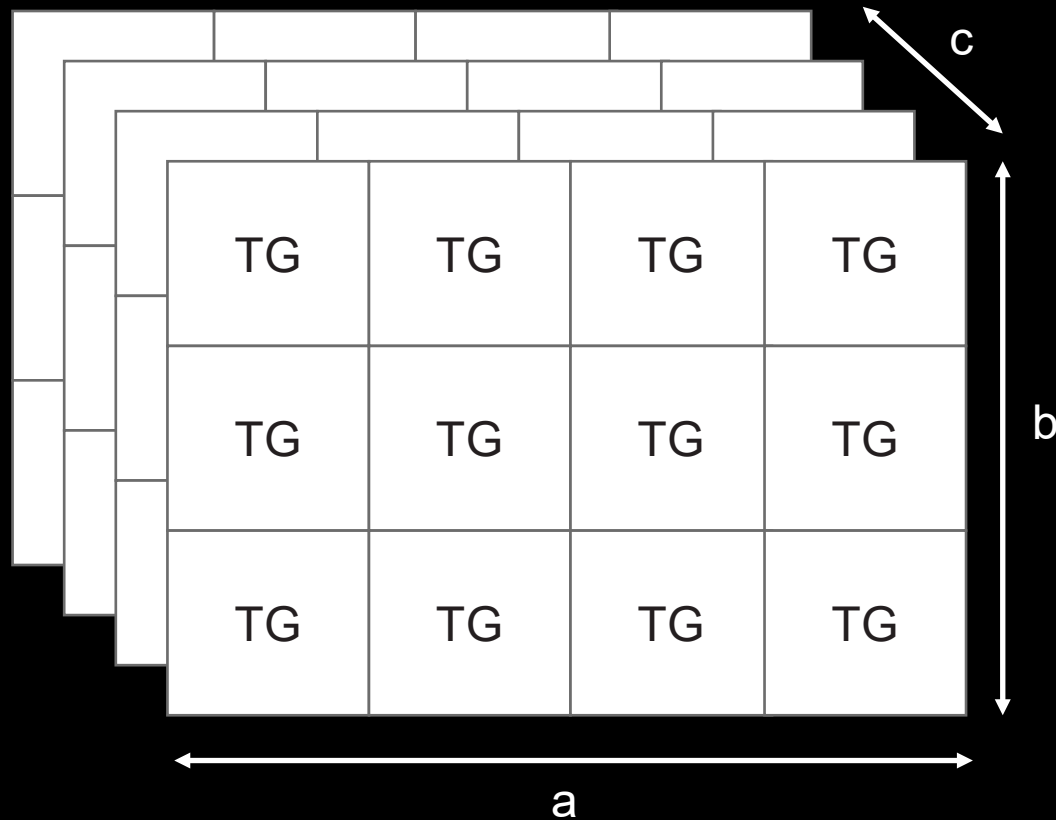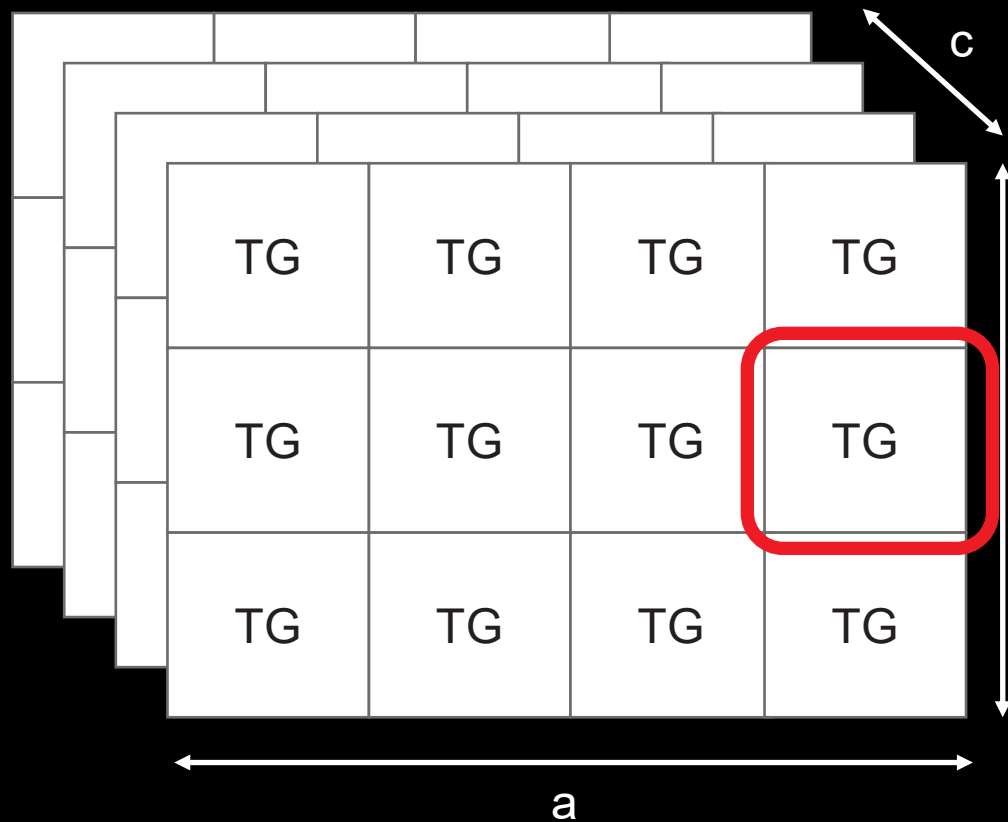
# SHADER ENGINE

# SHADER ENGINE

# WGP

# WGP

# WGP

# COMPUTE & THREADGROUPS

## dispatch(a, b, c)
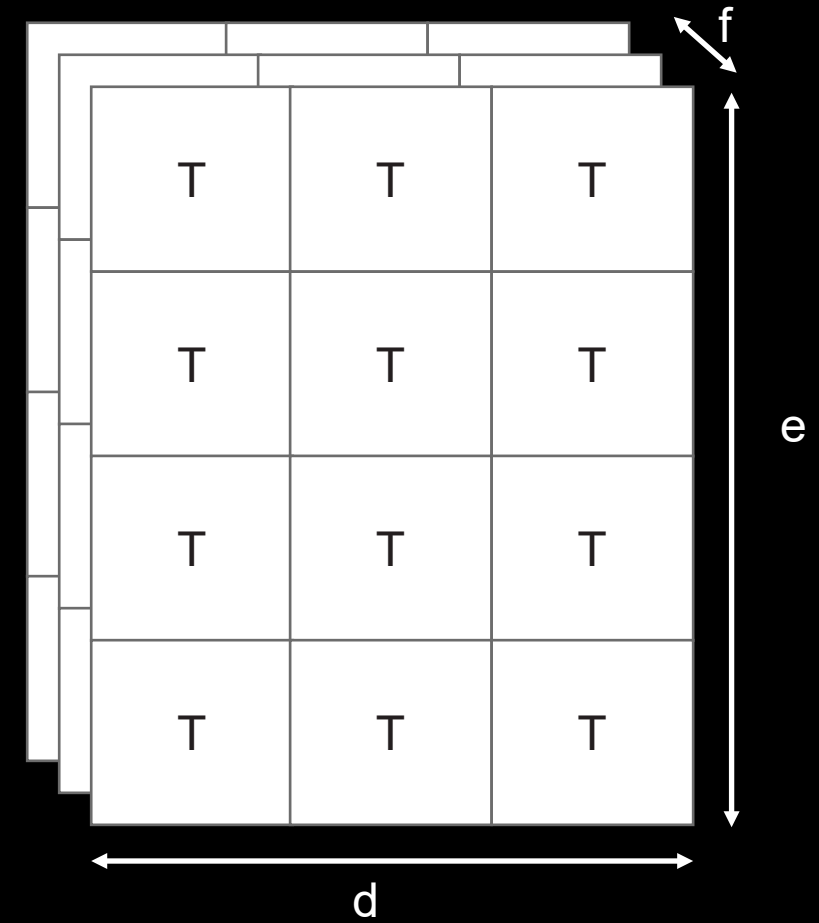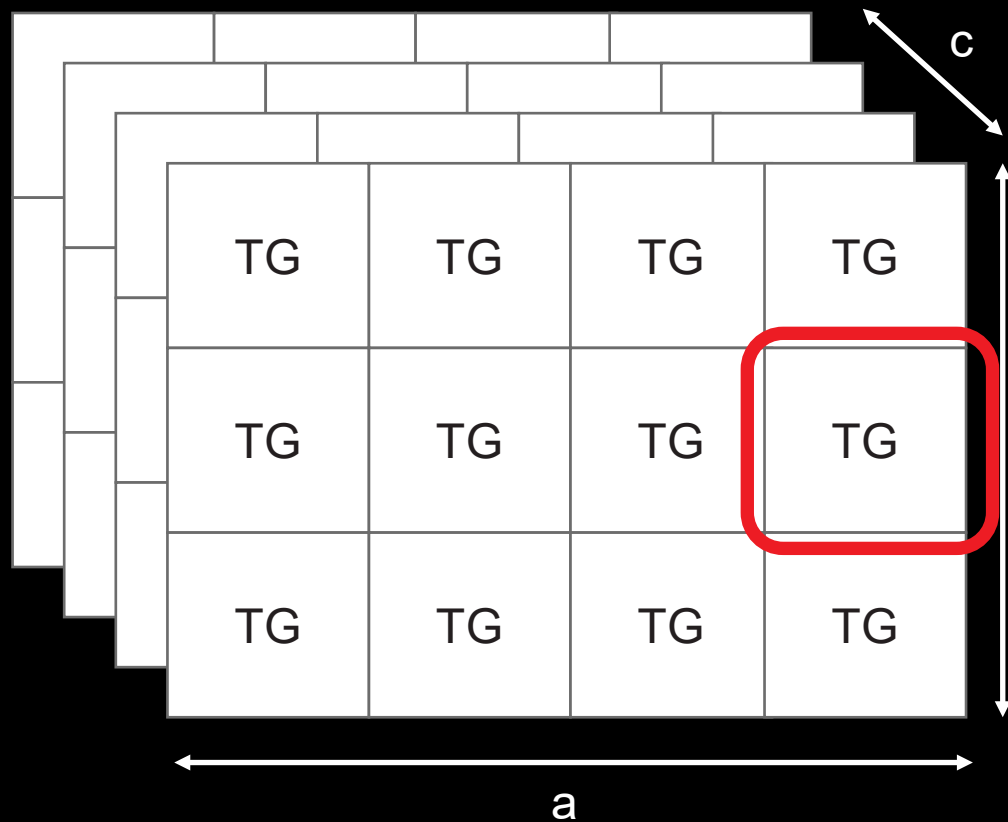
# COMPUTE & THREADGROUPS
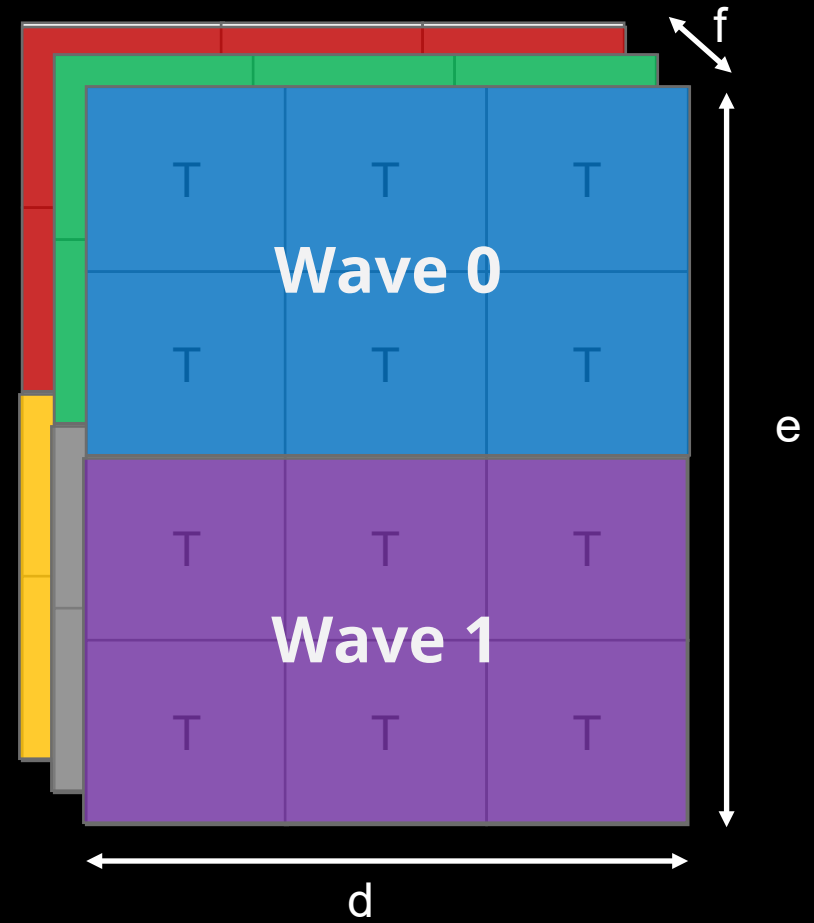


dispatch(a, b, c)

numthreads(d, e, f)

# WAVEFRONTS



dispatch(a, b, c) → numthreads(d, e, f)

# LOCKSTEP EXECUTION

```
[numthreads(32, 1, 1)]
void CSMain( uint threadIndex :SV_DispatchThreadID )
{
    int sum = 0;
    if(threadIndex < 16)
    {
        sum += 1;
    }
    else
    {
        sum += 2;
    }

    data[threadIndex] = sum;
}
```
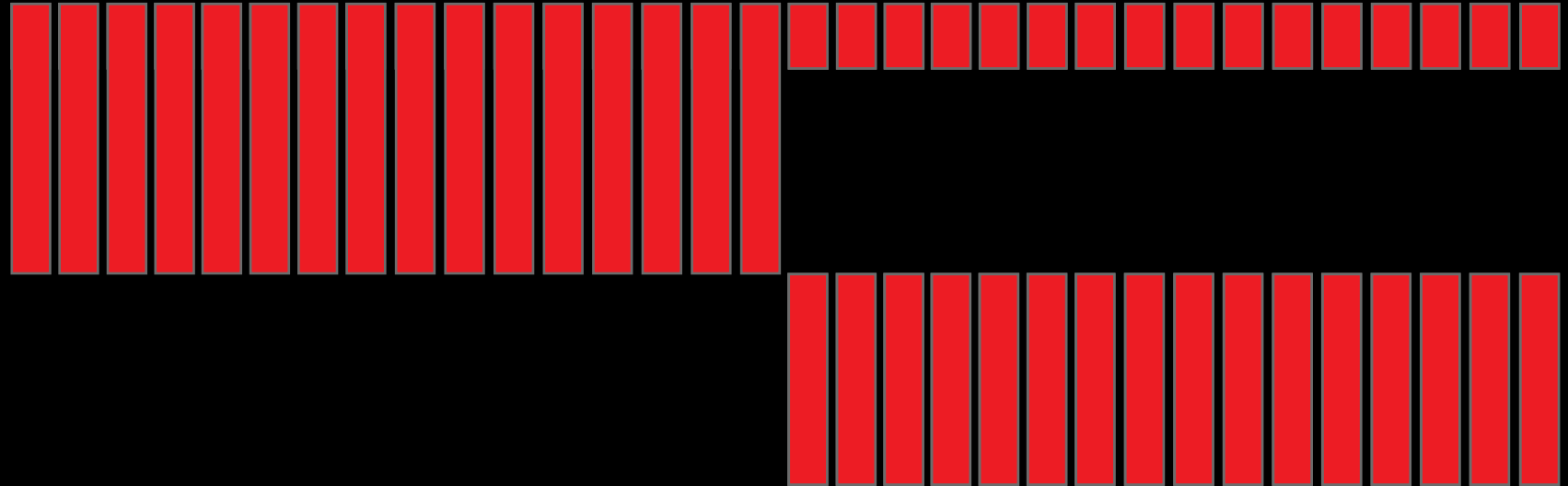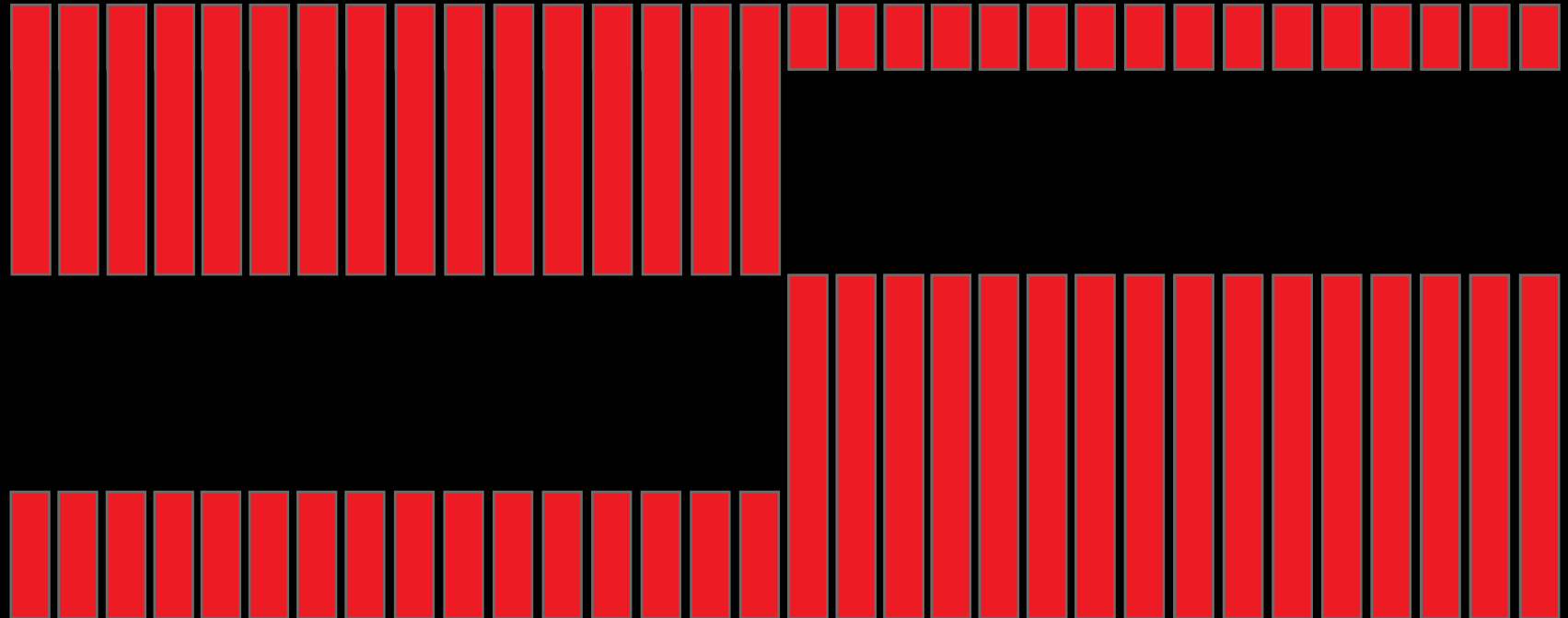
# LOCKSTEP EXECUTION

```
[numthreads(32, 1, 1)]
void CSMain( uint threadIndex :SV_DispatchThreadID )
{
    int sum = 0;
    if(threadIndex < 16)
    {
        sum += 1;
    }
    else
    {
        sum += 2;
    }

    data[threadIndex] = sum;
}
```

# LOCKSTEP EXECUTION

```
[numthreads(32, 1, 1)]
void CSMain( uint threadIndex :SV_DispatchThreadID )
{
    int sum = 0;
    if(threadIndex < 16)
    {
        sum += 1;
    }
    else
    {
        sum += 2;
    }

    data[threadIndex] = sum;
}
```

AMD
GPUOpen

# LOCKSTEP EXECUTION

```
[numthreads(32, 1, 1)]
void CSMain( uint threadIndex :SV_DispatchThreadID )
{
    int sum = 0;
    if(threadIndex < 16)
    {
        sum += 1;
    }
    else
    {
        sum += 2;
    }

    data[threadIndex] = sum;
}
```
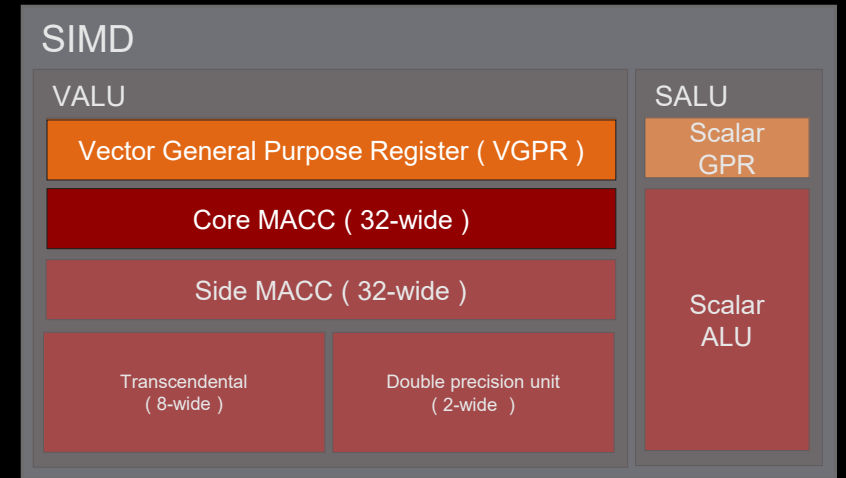
AMD
GPUOpen

# SIMD & LOCKSTEP

```
cbuffer input : register(b0){ int data; };
RWBuffer<int> output : register(u0);

[numthreads(32, 1, 1)]
void CSMain( uint threadIndex : SV_DispatchThreadID )
{
    int sum = 0;
    sum += threadIndex;
    sum += data;
    output[threadIndex] = sum;
}
```
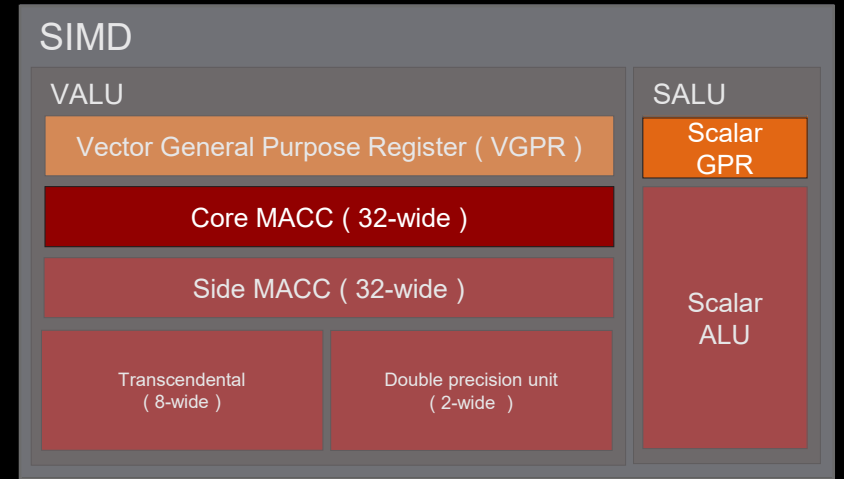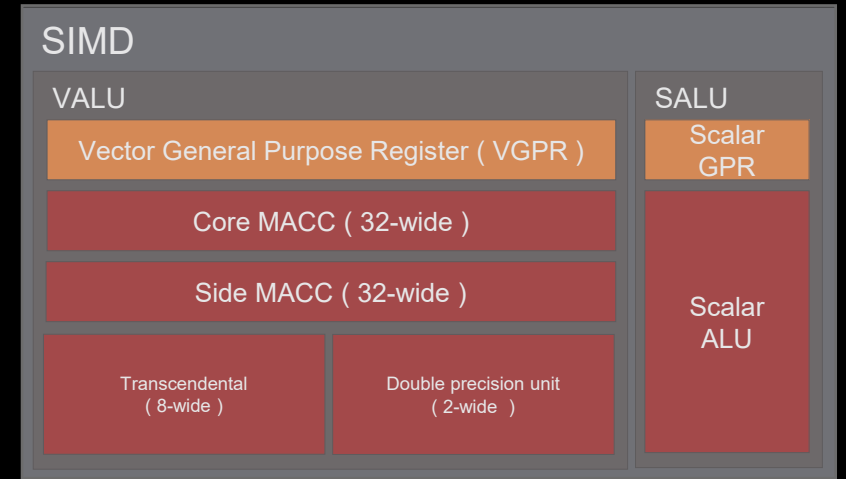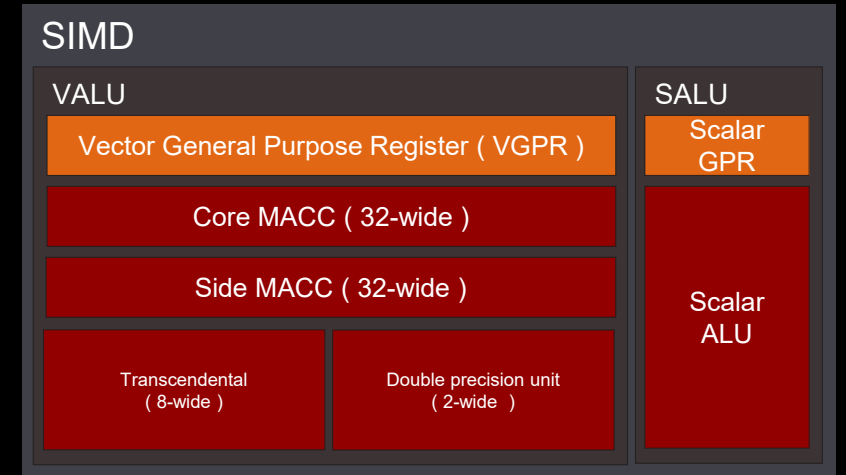
# SIMD & LOCKSTEP

```hlsl
cbuffer input : register(b0){ int data; };
RWBuffer<int> output : register(u0);

[numthreads(32, 1, 1)]
void CSMain( uint threadIndex : SV_DispatchThreadID )
{
    int sum = 0;
    sum += threadIndex;
    sum += data;
    output[threadIndex] = sum;
}
```

| sum: 0 | sum: 0 | sum: 0 | sum: 0 | … |

# SIMD & LOCKSTEP

```
cbuffer input : register(b0){ int data; };
RWBuffer<int> output : register(u0);

[numthreads(32, 1, 1)]
void CSMain( uint threadIndex : SV_DispatchThreadID )
{
    int sum = 0;
    sum += threadIndex;
    sum += data;
    output[threadIndex] = sum;
}
```

**SIMD**

**VALU**

Vector General Purpose Register ( VGPR )

Core MACC ( 32-wide )

Side MACC ( 32-wide )

Transcendental ( 8-wide )

Double precision unit ( 2-wide )

**SALU**

Scalar GPR

Scalar ALU

**VGPR**

| sum: 0 | sum: 0 | sum: 0 | sum: 0 | … |

AMD GPUOpen

# SIMD & LOCKSTEP

```
cbuffer input : register(b0){ int data; };
RWBuffer<int> output : register(u0);

[numthreads(32, 1, 1)]
void CSMain( uint threadIndex : SV_DispatchThreadID )
{
    int sum = 0;
→   sum += threadIndex;
    sum += data;
    output[threadIndex] = sum;
}
```

**SIMD**

**VALU**

Vector General Purpose Register ( VGPR )

Core MACC ( 32-wide )

Side MACC ( 32-wide )

| Transcendental ( 8-wide ) | Double precision unit ( 2-wide ) |

**SALU**

Scalar GPR

Scalar ALU

## VGPR

| sum: 0 | sum: 0 | sum: 0 | sum: 0 | … |
| sum: 0 | sum: 1 | sum: 2 | sum: 3 | … |

# SIMD & LOCKSTEP

```
cbuffer input : register(b0){ int data; };
RWBuffer<int> output : register(u0);

[numthreads(32, 1, 1)]
void CSMain( uint threadIndex : SV_DispatchThreadID )
{
    int sum = 0;
    sum += threadIndex;
→   sum += data;
    output[threadIndex] = sum;
}
```



SIMD

VALU

Vector General Purpose Register ( VGPR )

Core MACC ( 32-wide )

Side MACC ( 32-wide )

Transcendental ( 8-wide )

Double precision unit ( 2-wide )

SALU

Scalar GPR

Scalar ALU

VGPR

| sum: 0 | sum: 0 | sum: 0 | sum: 0 | … |
| sum: 0 | sum: 1 | sum: 2 | sum: 3 | … |
| sum: 5 | sum: 6 | sum: 7 | sum: 8 | … |

SGPR

| data: 5 |
| data: 5 |
| data: 5 |

AMD GPUOpen

# SIMD & LOCKSTEP

```hlsl
cbuffer input : register(b0){ int data; };
RWBuffer<int> output : register(u0);

[numthreads(32, 1, 1)]
void CSMain( uint threadIndex : SV_DispatchThreadID )
{
    int sum = 0;
    sum += threadIndex;
    sum += data;
    output[threadIndex] = sum;
}
```

**SIMD**

**VALU**

Vector General Purpose Register ( VGPR )

Core MACC ( 32-wide )

Side MACC ( 32-wide )

Transcendental ( 8-wide )

Double precision unit ( 2-wide )

**SALU**

Scalar GPR

Scalar ALU

**VGPR**

| sum: 0 | sum: 0 | sum: 0 | sum: 0 | … |
| sum: 0 | sum: 1 | sum: 2 | sum: 3 | … |
| sum: 5 | sum: 6 | sum: 7 | sum: 8 | … |
| sum: 5 | sum: 6 | sum: 7 | sum: 8 | … |

**SGPR**

| data: 5 |
| data: 5 |
| data: 5 |
| data: 5 |

AMD GPUOpen

# SIMD & LOCKSTEP

```
cbuffer input : register(b0){ int data; };
RWBuffer<int> output : register(u0);

[numthreads(32, 1, 1)]
void CSMain( uint threadIndex : SV_DispatchThreadID )
{
    int sum = 0;
    sum += threadIndex;
    sum += data;
    output[threadIndex] = sum;
}
```

### SIMD

| VALU | | SALU |
|---|---|---|
| Vector General Purpose Register ( VGPR ) | | Scalar GPR |
| Core MACC ( 32-wide ) | | Scalar ALU |
| Side MACC ( 32-wide ) | | |
| Transcendental ( 8-wide ) | Double precision unit ( 2-wide ) | |

### VGPR

| | | | | | |
|---|---|---|---|---|---|
| sum: 0 | sum: 0 | sum: 0 | sum: 0 | … | data: 5 |
| sum: 0 | sum: 1 | sum: 2 | sum: 3 | … | data: 5 |
| sum: 5 | sum: 6 | sum: 7 | sum: 8 | … | data: 5 |
| sum: 5 | sum: 6 | sum: 7 | sum: 8 | … | data: 5 |

AMD GPUOpen

# ASSIGNED WAVEFRONTS

# ASSIGNED WAVEFRONTS

## Wavefronts don't have to be executed in order

## Wavefronts execution can be interrupted and resumed at any time

AMD
GPUOpen

# WAVEFRONT SCHEDULING & LATENCY HIDING



Wavefront A

Wavefront B

Slot 0

Slot 1

Time

# WAVEFRONT SCHEDULING & LATENCY HIDING



Wavefront A

Wavefront B

Slot 0    ALU

Slot 1

Time

# WAVEFRONT SCHEDULING & LATENCY HIDING

# WAVEFRONT SCHEDULING & LATENCY HIDING



Wavefront A

Wavefront B

Slot 0 — ALU ALU MEM

Slot 1

Time

AMD GPUOpen

# WAVEFRONT SCHEDULING & LATENCY HIDING

# WAVEFRONT SCHEDULING & LATENCY HIDING

# WAVEFRONT SCHEDULING & LATENCY HIDING

# WAVEFRONT SCHEDULING & LATENCY HIDING

# WAVEFRONT SCHEDULING & LATENCY HIDING

# WAVEFRONT SCHEDULING & LATENCY HIDING

# WAVEFRONT SCHEDULING & LATENCY HIDING

# LATENCY HIDING IN RGP

# Occupancy is the ratio of assigned wavefronts to the maximum available slots



SIMD – 25% occupancy

# Occupancy is the ratio of assigned wavefronts to the maximum available slots



SIMD – 25% occupancy

SIMD - 75% occupancy

AMD
GPUOpen

# Better occupancy doesn't mean better performance !

# Latency bound workloads *might* benefit from increased occupancy

# In memory bound scenarios, increasing occupancy might thrash the caches

# THEORETICAL OCCUPANCY – GPRS

```hlsl
[numthreads(32, 1, 1)]
void CSMain( uint threadIndex :
        SV_DispatchThreadID )
{
    int sum = 0;
    if(threadIndex < 16)
    {
        sum += 1;
    }
    else
    {
        sum += 2;
    }

    data[threadIndex] = sum;
}
```

AMD
GPUOpen

# THEORETICAL OCCUPANCY – GPRS

```cpp
[numthreads(32, 1, 1)]
void CSMain( uint threadIndex :
        SV_DispatchThreadID )
{
    int sum = 0;
    if(threadIndex < 16)
    {
        sum += 1;
    }
    else
    {
        sum += 2;
    }

    data[threadIndex] = sum;
}
```

```
shader main
    asic(GFX10_3)
    type(CS)
    sgpr_count(6)
    vgpr_count(8)
    wave_size(32)
    s_version            UC_VERSION_GFX10 | UC_VERSION_W32_BIT
    s_inst_prefetch 0x0003
    s_getpc_b64       s[0:1]
    s_mov_b32         s0, s2
    s_load_dwordx4    s[4:7], s[0:1], null
    v_mad_u32_u24     v1, s3, 32, v0
    v_cmp_gt_u32      vcc_lo, 16, v1
    v_cndmask_b32     v2, 2, 1, vcc_lo
    v_mov_b32         v3, v2
    v_mov_b32         v4, v2
    v_mov_b32         v5, v2
    s_waitcnt         lgkmcnt(0)
    buffer_store_format_xyzw  v[2:5], v1, s[4:7], 0 idxen glc
    s_endpgm
```

# THEORETICAL OCCUPANCY – GPRS

```hlsl
[numthreads(32, 1, 1)]
void CSMain( uint threadIndex :
        SV_DispatchThreadID )
{

    int sum = 0;
    if(threadIndex < 16)
    {
        sum += 1;
    }
    else
    {
        sum += 2;
    }

    data[threadIndex] = sum;
}
```

```
shader main
  asic(GFX10_3)
  type(CS)
  sgpr_count(6)
  vgpr_count(8)
  wave_size(32)
  s_version          UC_VERSION_GFX10 | UC_VERSION_W32_BIT
  s_inst_prefetch 0x0003
  s_getpc_b64        s[0:1]
  s_mov_b32          s0, s2
  s_load_dwordx4     s[4:7], s[0:1], null
  v_mad_u32_u24      v1, s3, 32, v0
  v_cmp_gt_u32       vcc_lo, 16, v1
  v_cndmask_b32      v2, 2, 1, vcc_lo
  v_mov_b32          v3, v2
  v_mov_b32          v4, v2
  v_mov_b32          v5, v2
  s_waitcnt          lgkmcnt(0)
  buffer_store_format_xyzw  v[2:5], v1, s[4:7], 0 idxen glc
  s_endpgm
```

# THEORETICAL OCCUPANCY – GPRS

```hlsl
[numthreads(32, 1, 1)]
void CSMain( uint threadIndex :
        SV_DispatchThreadID )
{
    int sum = 0;
    if(threadIndex < 16)
    {
        sum += 1;
    }
    else
    {
        sum += 2;
    }

    data[threadIndex] = sum;
}
```

```
shader main
  asic(GFX10_3)
  type(CS)
  sgpr_count(6)
  vgpr_count(8)
  wave_size(32)
```

```
s_version         UC_VERSION_GFX10 | UC_VERSION_W32_BIT
s_inst_prefetch  0x0003
s_getpc_b64       s[0:1]
s_mov_b32         s0, s2
s_load_dwordx4    s[4:7], s[0:1], null
v_mad_u32_u24     v1, s3, 32, v0
v_cmp_gt_u32      vcc_lo, 16, v1
v_cndmask_b32     v2, 2, 1, vcc_lo
v_mov_b32         v3, v2
v_mov_b32         v4, v2
v_mov_b32         v5, v2
s_waitcnt         lgkmcnt(0)
buffer_store_format_xyzw  v[2:5], v1, s[4:7], 0 idxen glc
s_endpgm
```

AMD
GPUOpen

# RADEON™ GPU PROFILER TO THE RESCUE



**For max occupancy**
- **Wave32 - 1536 / 16 = 96 VGPR per wave**
- **Wave64 - 1536 / 2 / 16 = 48 VGPR per wave**

# RADEON™ GPU PROFILER TO THE RESCUE

# RADEON™ GPU ANALYZER TO THE RESCUE
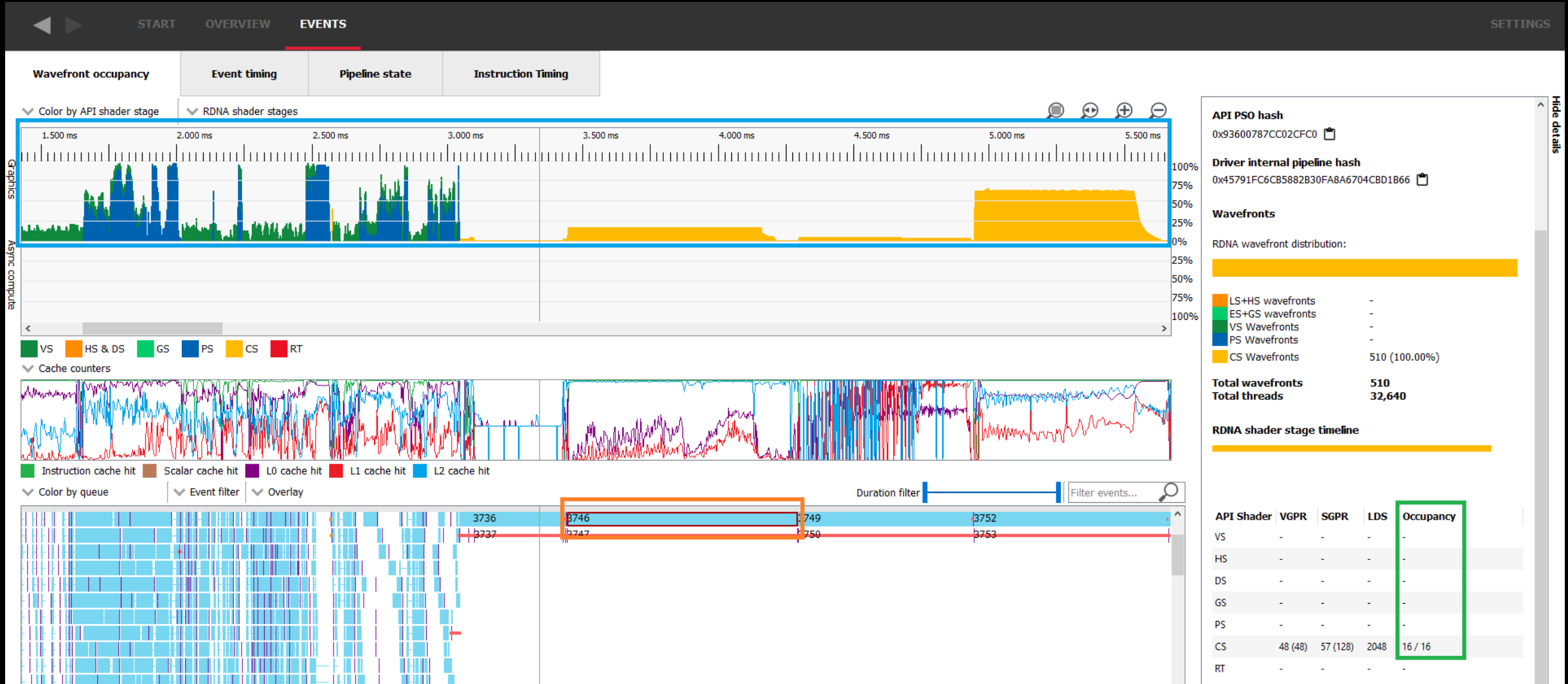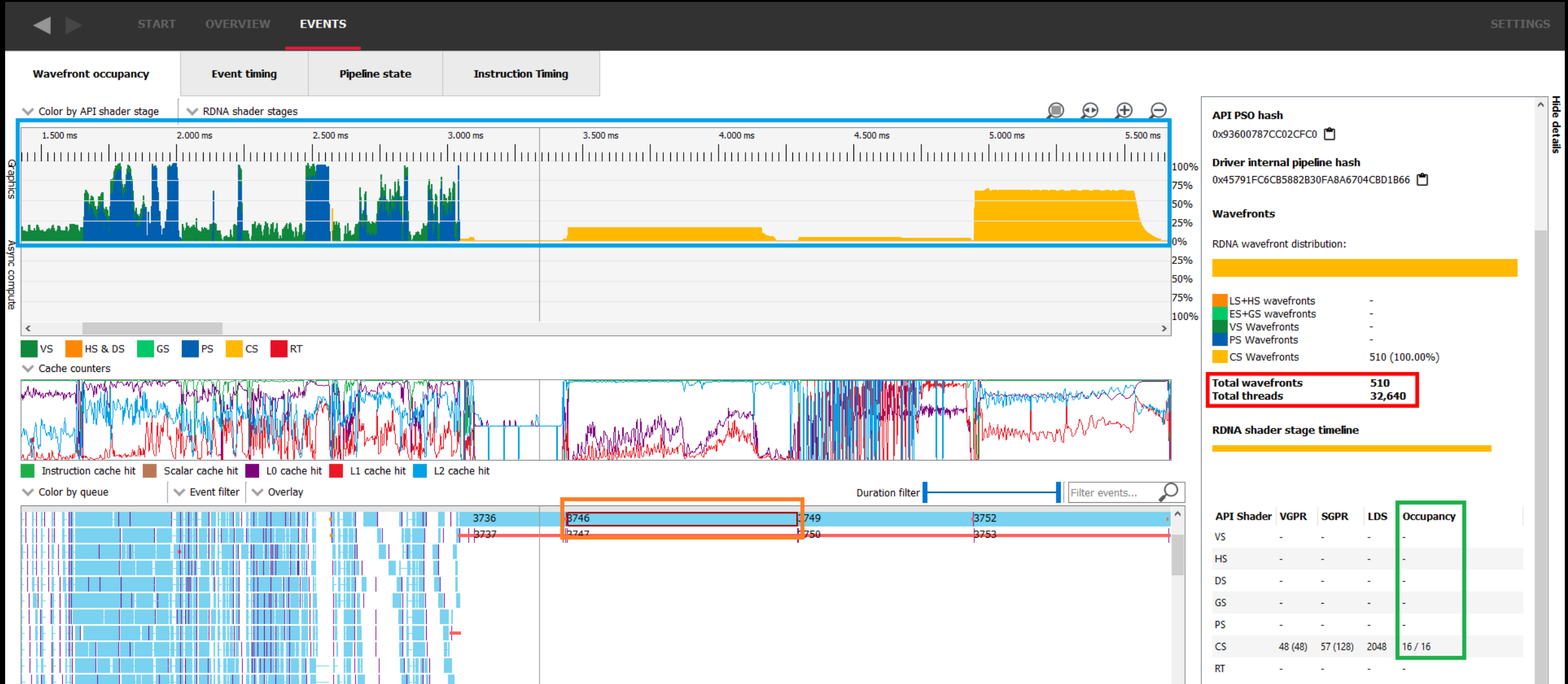
# RADEON™ GPU ANALYZER TO THE RESCUE

# THEORETICAL OCCUPANCY – LDS & THREADGROUP_SIZE

# MEASURED OCCUPANCY

# MEASURED OCCUPANCY

# LACK OF WORK LIMITED OCCUPANCY



**Shader Engines (SE): 6**
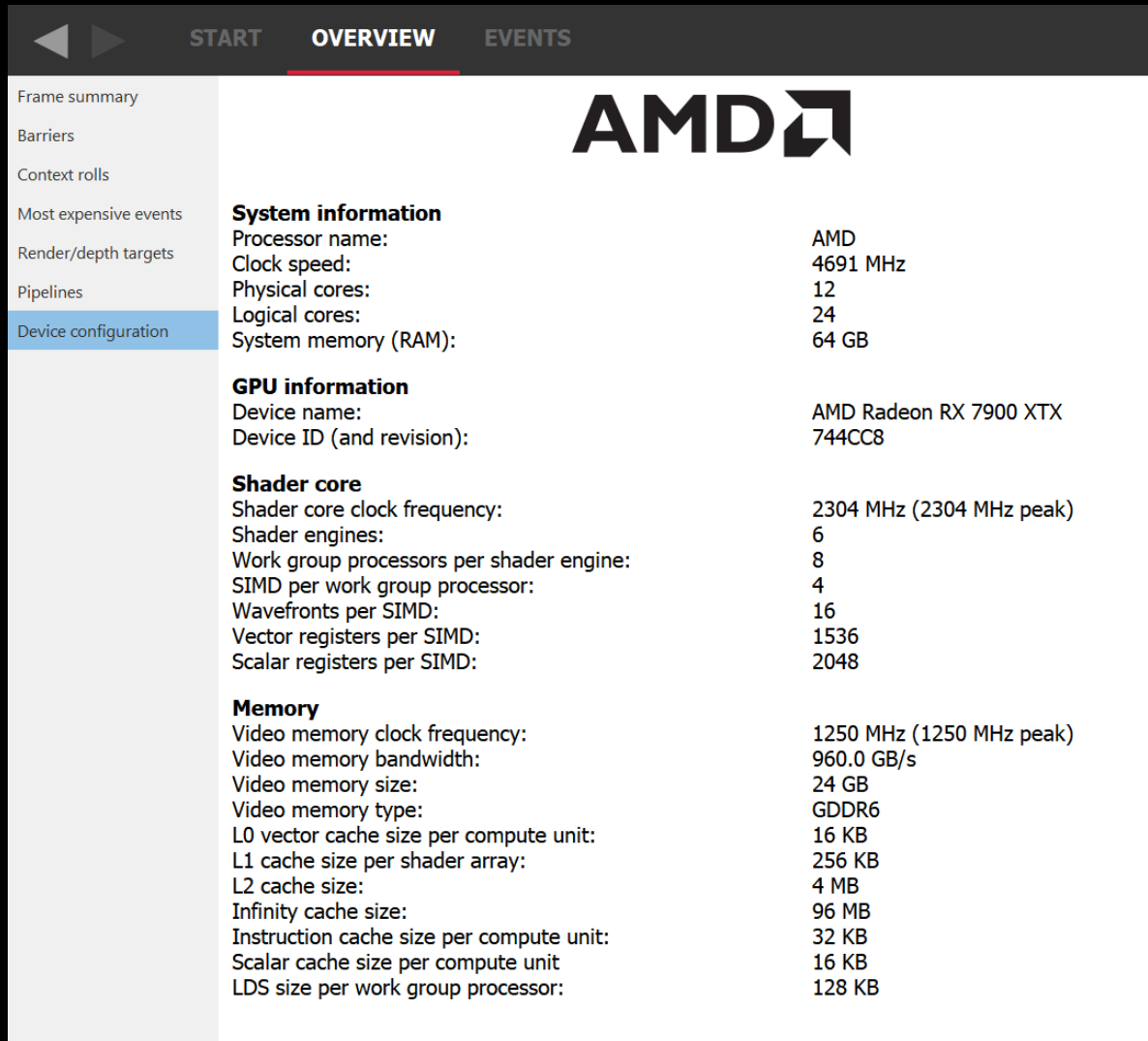
# LACK OF WORK LIMITED OCCUPANCY



**Shader Engines (SE): 6**
**WorkGroup Processors (WGP) / SE: 8**

# LACK OF WORK LIMITED OCCUPANCY



**System information**
| | |
|---|---|
| Processor name: | AMD |
| Clock speed: | 4691 MHz |
| Physical cores: | 12 |
| Logical cores: | 24 |
| System memory (RAM): | 64 GB |

**GPU information**
| | |
|---|---|
| Device name: | AMD Radeon RX 7900 XTX |
| Device ID (and revision): | 744CC8 |

**Shader core**
| | |
|---|---|
| Shader core clock frequency: | 2304 MHz (2304 MHz peak) |
| Shader engines: | 6 |
| Work group processors per shader engine: | 8 |
| SIMD per work group processor: | 4 |
| Wavefronts per SIMD: | 16 |
| Vector registers per SIMD: | 1536 |
| Scalar registers per SIMD: | 2048 |

**Memory**
| | |
|---|---|
| Video memory clock frequency: | 1250 MHz (1250 MHz peak) |
| Video memory bandwidth: | 960.0 GB/s |
| Video memory size: | 24 GB |
| Video memory type: | GDDR6 |
| L0 vector cache size per compute unit: | 16 KB |
| L1 cache size per shader array: | 256 KB |
| L2 cache size: | 4 MB |
| Infinity cache size: | 96 MB |
| Instruction cache size per compute unit: | 32 KB |
| Scalar cache size per compute unit: | 16 KB |
| LDS size per work group processor: | 128 KB |

**Shader Engines (SE): 6**
**WorkGroup Processors (WGP) / SE: 8**
**Total WGP: 6 * 8 = 48**
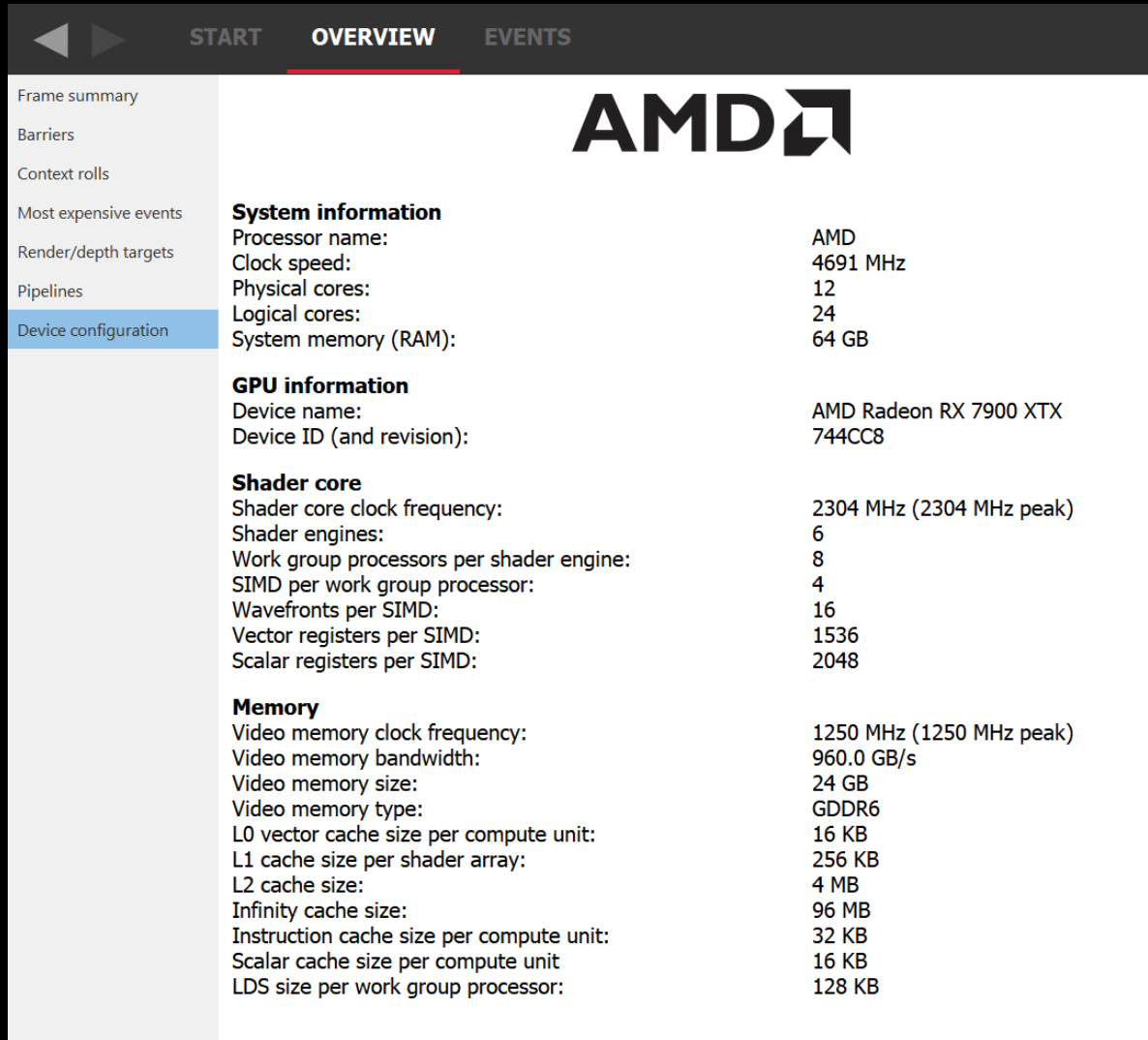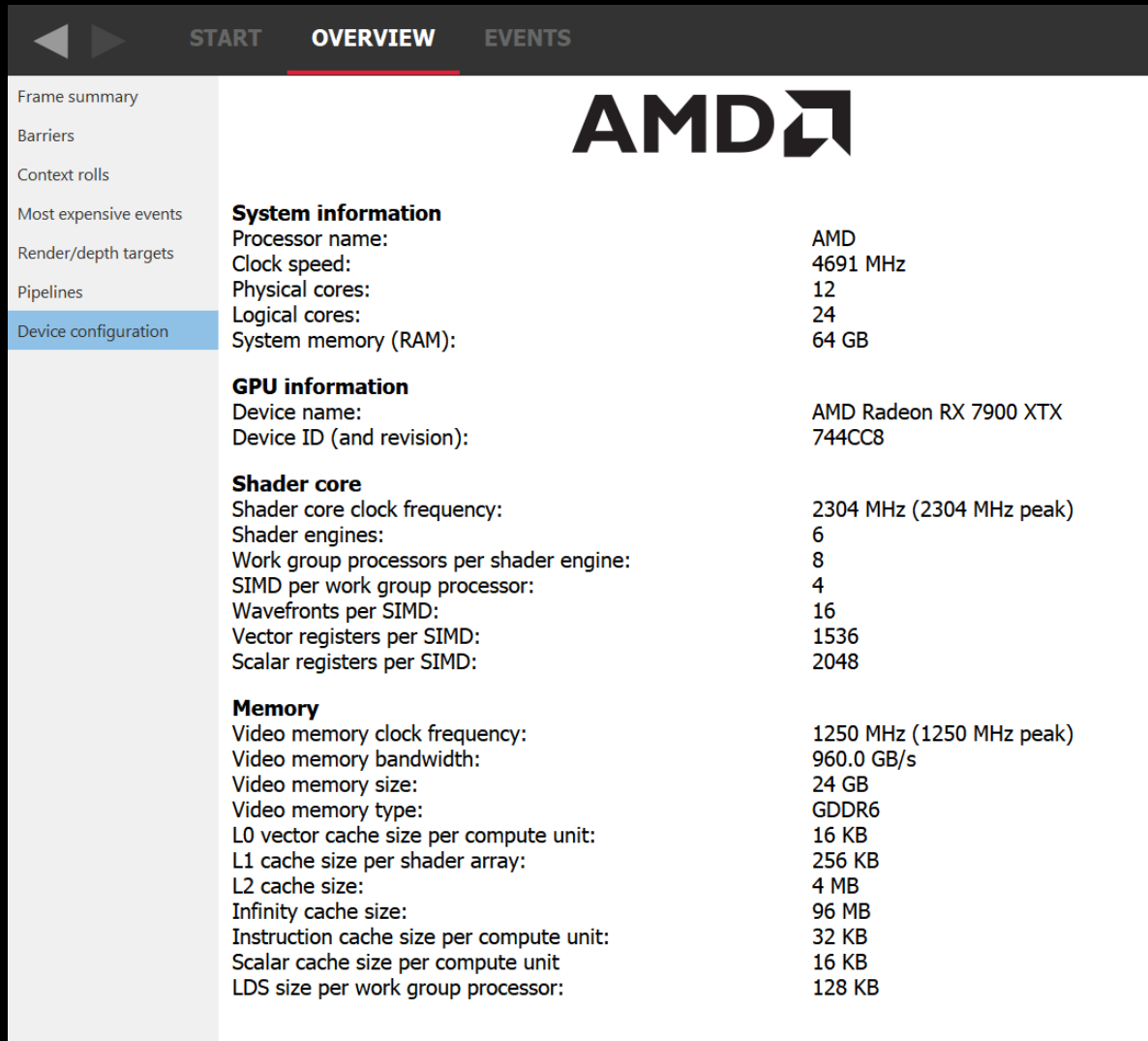
# LACK OF WORK LIMITED OCCUPANCY



**Shader Engines (SE): 6**
**WorkGroup Processors (WGP) / SE: 8**
**Total WGP: 6 * 8 = 48**
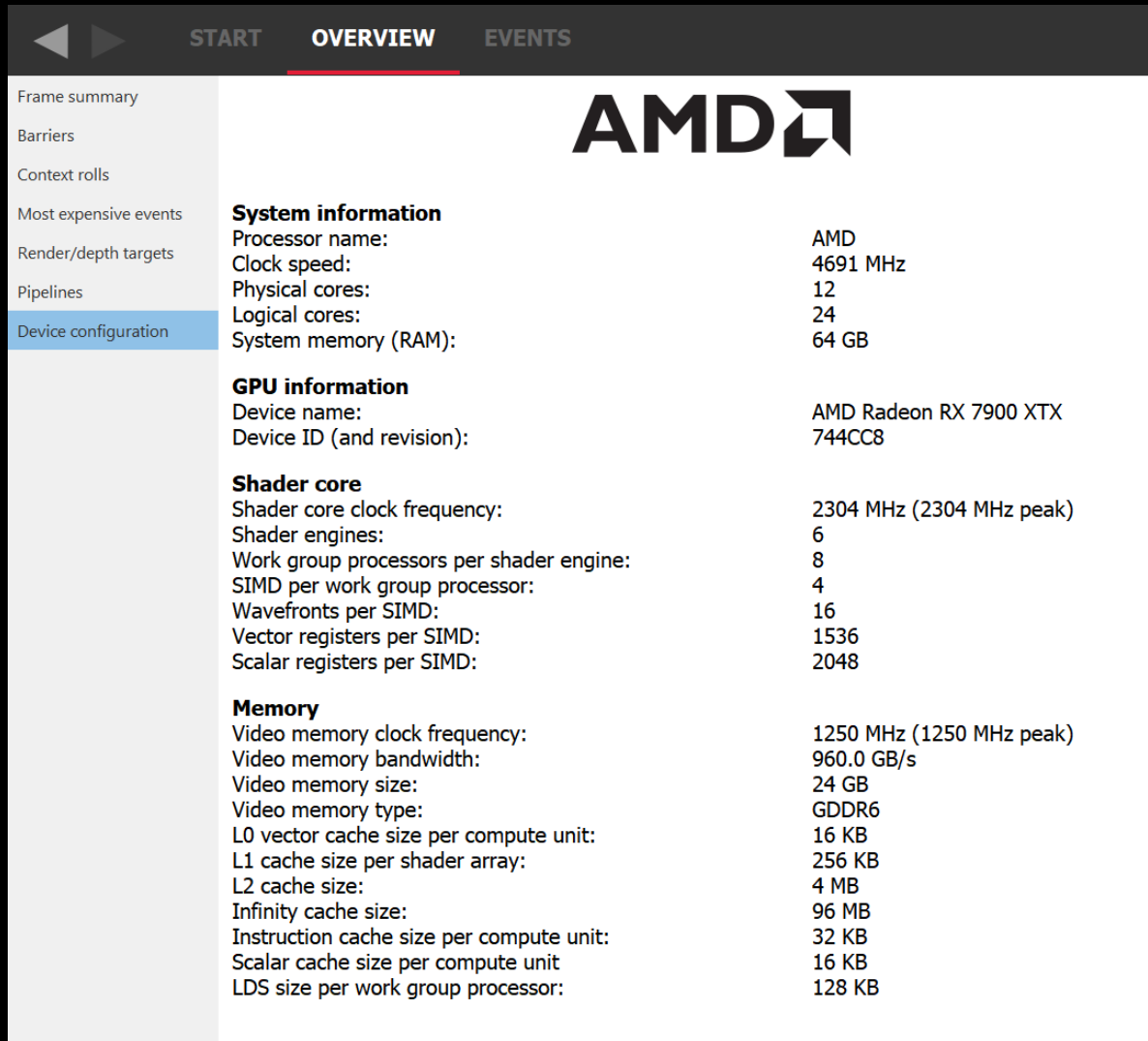**SIMD per WGP: 4**

# LACK OF WORK LIMITED OCCUPANCY



**System information**

| | |
|---|---|
| Processor name: | AMD |
| Clock speed: | 4691 MHz |
| Physical cores: | 12 |
| Logical cores: | 24 |
| System memory (RAM): | 64 GB |

**GPU information**

| | |
|---|---|
| Device name: | AMD Radeon RX 7900 XTX |
| Device ID (and revision): | 744CC8 |

**Shader core**

| | |
|---|---|
| Shader core clock frequency: | 2304 MHz (2304 MHz peak) |
| Shader engines: | 6 |
| Work group processors per shader engine: | 8 |
| SIMD per work group processor: | 4 |
| Wavefronts per SIMD: | 16 |
| Vector registers per SIMD: | 1536 |
| Scalar registers per SIMD: | 2048 |

**Memory**

| | |
|---|---|
| Video memory clock frequency: | 1250 MHz (1250 MHz peak) |
| Video memory bandwidth: | 960.0 GB/s |
| Video memory size: | 24 GB |
| Video memory type: | GDDR6 |
| L0 vector cache size per compute unit: | 16 KB |
| L1 cache size per shader array: | 256 KB |
| L2 cache size: | 4 MB |
| Infinity cache size: | 96 MB |
| Instruction cache size per compute unit: | 32 KB |
| Scalar cache size per compute unit: | 16 KB |
| LDS size per work group processor: | 128 KB |

**Shader Engines (SE): 6**
**WorkGroup Processors (WGP) / SE: 8**
**Total WGP: 6 * 8 = 48**
**SIMD per WGP: 4**
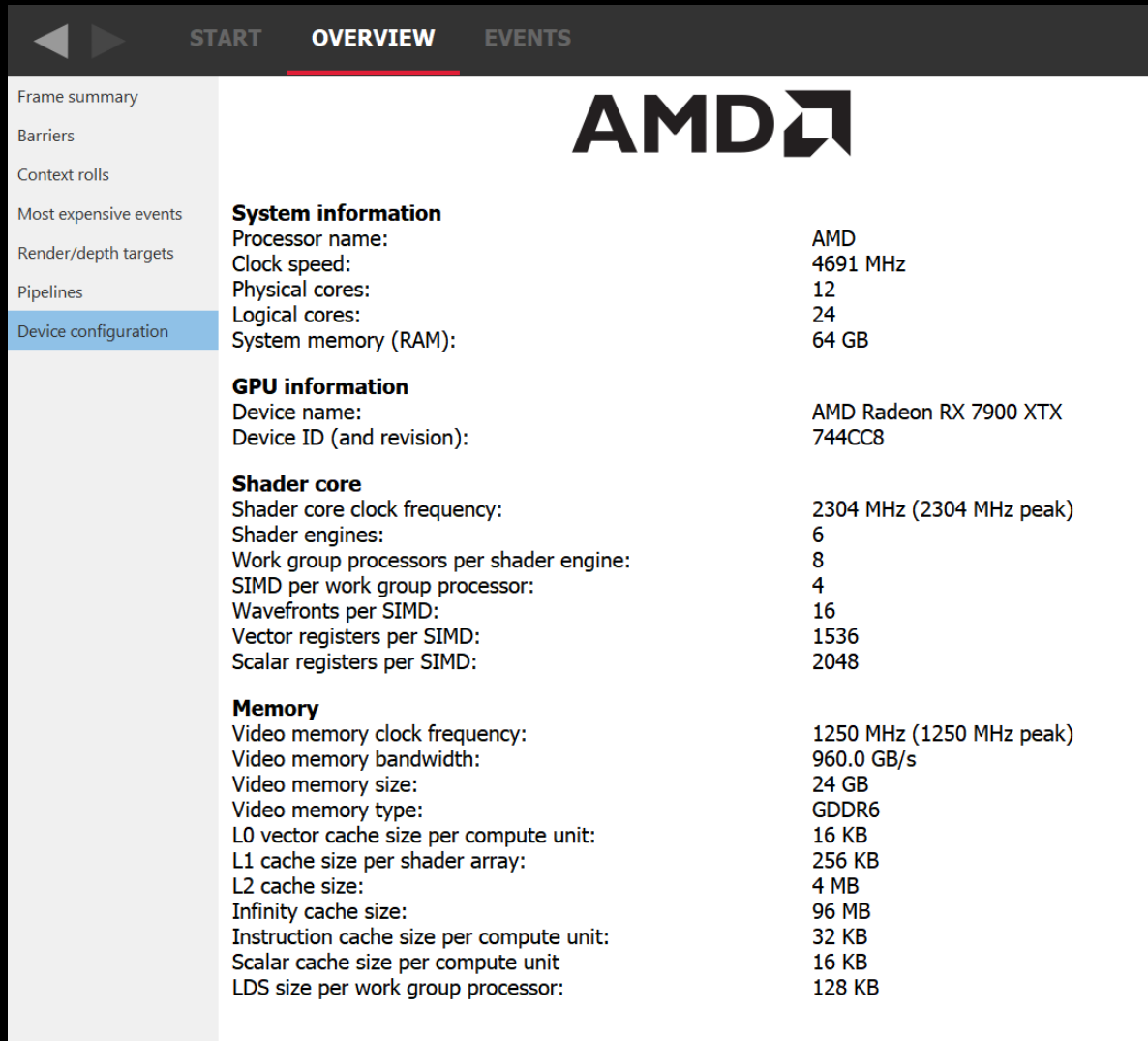**Total SIMD: 4 * 48 = 192**

# LACK OF WORK LIMITED OCCUPANCY



**Shader Engines (SE): 6**
**WorkGroup Processors (WGP) / SE: 8**
**Total WGP: 6 * 8 = 48**
**SIMD per WGP: 4**
**Total SIMD: 4 * 48 = 192**
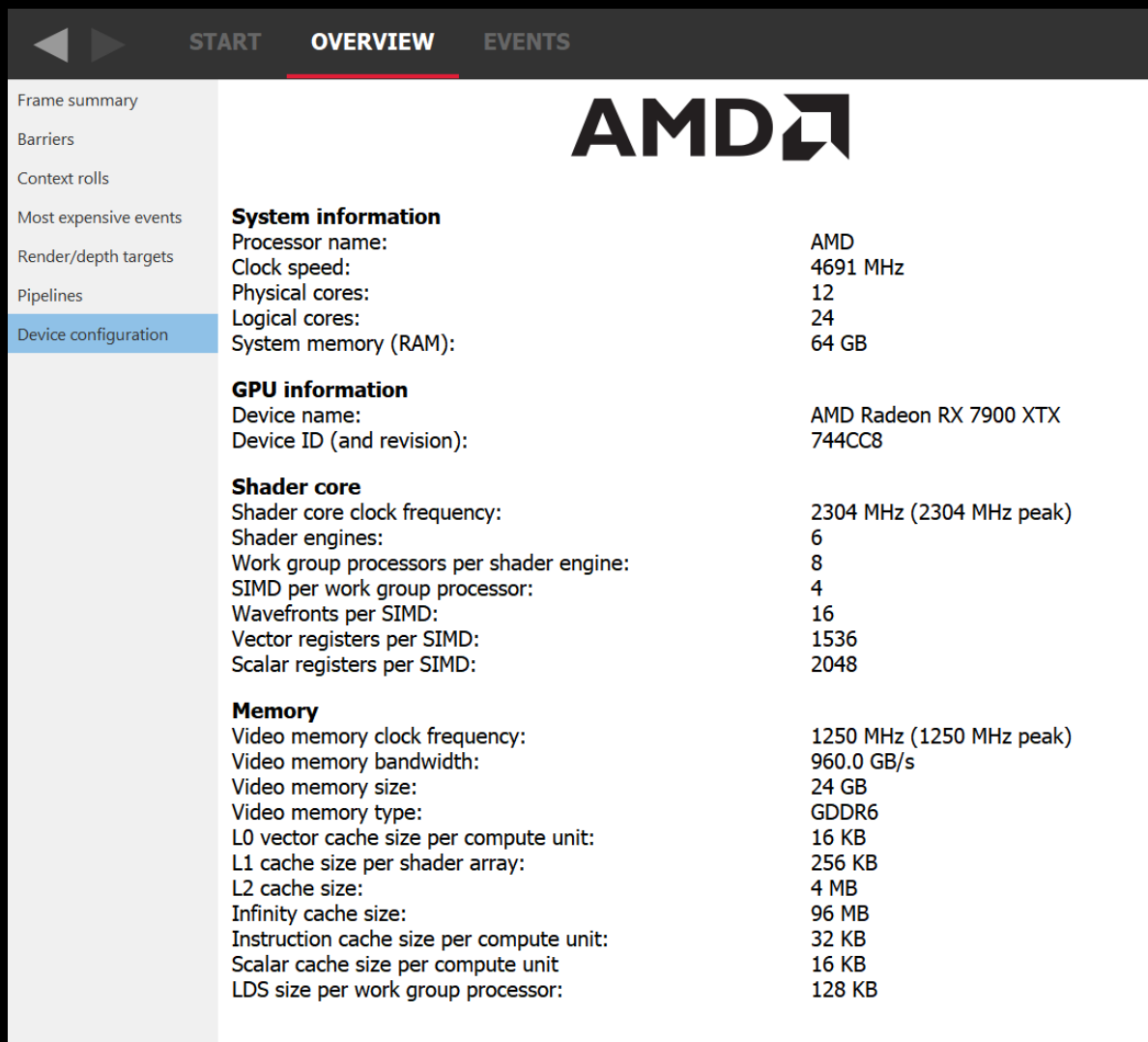**Wave slots per SIMD: 16**

# LACK OF WORK LIMITED OCCUPANCY



**Shader Engines (SE): 6**
**WorkGroup Processors (WGP) / SE: 8**
**Total WGP: 6 * 8 = 48**
**SIMD per WGP: 4**
**Total SIMD: 4 * 48 = 192**
**Wave slots per SIMD: 16**
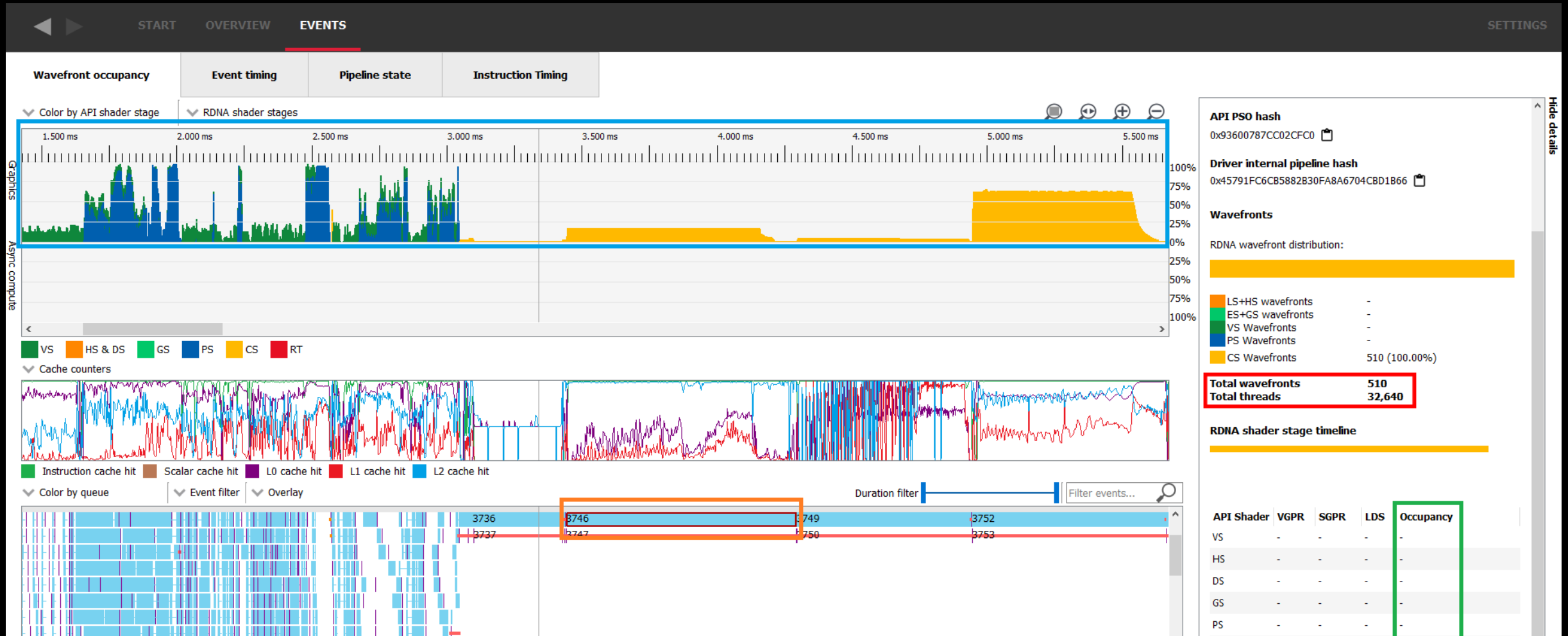**Wave slots: 16 * 192 = 3072**

# LACK OF WORK LIMITED OCCUPANCY



**System information**
Processor name: AMD
Clock speed: 4691 MHz
Physical cores: 12
Logical cores: 24
System memory (RAM): 64 GB

**GPU information**
Device name: AMD Radeon RX 7900 XTX
Device ID (and revision): 744CC8

**Shader core**
Shader core clock frequency: 2304 MHz (2304 MHz peak)
Shader engines: 6
Work group processors per shader engine: 8
SIMD per work group processor: 4
Wavefronts per SIMD: 16
Vector registers per SIMD: 1536
Scalar registers per SIMD: 2048

**Memory**
Video memory clock frequency: 1250 MHz (1250 MHz peak)
Video memory bandwidth: 960.0 GB/s
Video memory size: 24 GB
Video memory type: GDDR6
L0 vector cache size per compute unit: 16 KB
L1 cache size per shader array: 256 KB
L2 cache size: 4 MB
Infinity cache size: 96 MB
Instruction cache size per compute unit: 32 KB
Scalar cache size per compute unit: 16 KB
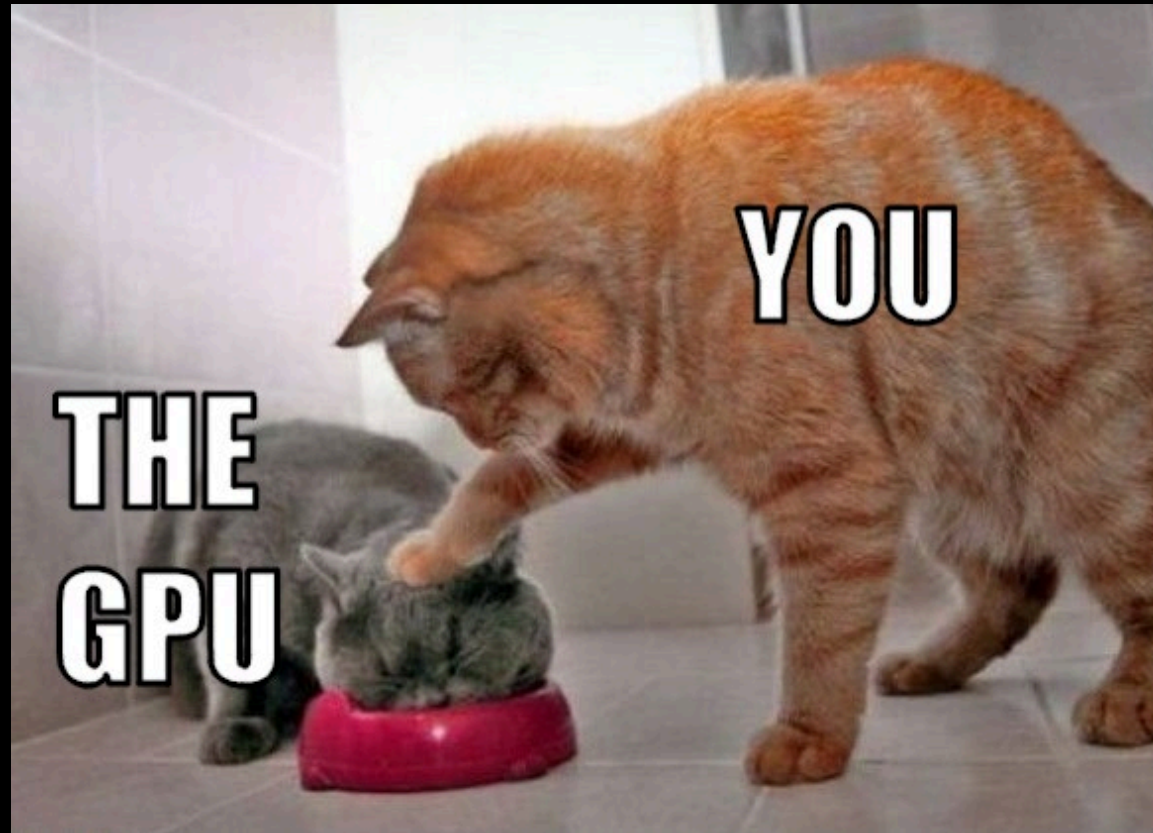LDS size per work group processor: 128 KB

**Shader Engines (SE): 6**
**WorkGroup Processors (WGP) / SE: 8**
**Total WGP: 6 * 8 = 48**
**SIMD per WGP: 4**
**Total SIMD: 4 * 48 = 192**
**Wave slots per SIMD: 16**
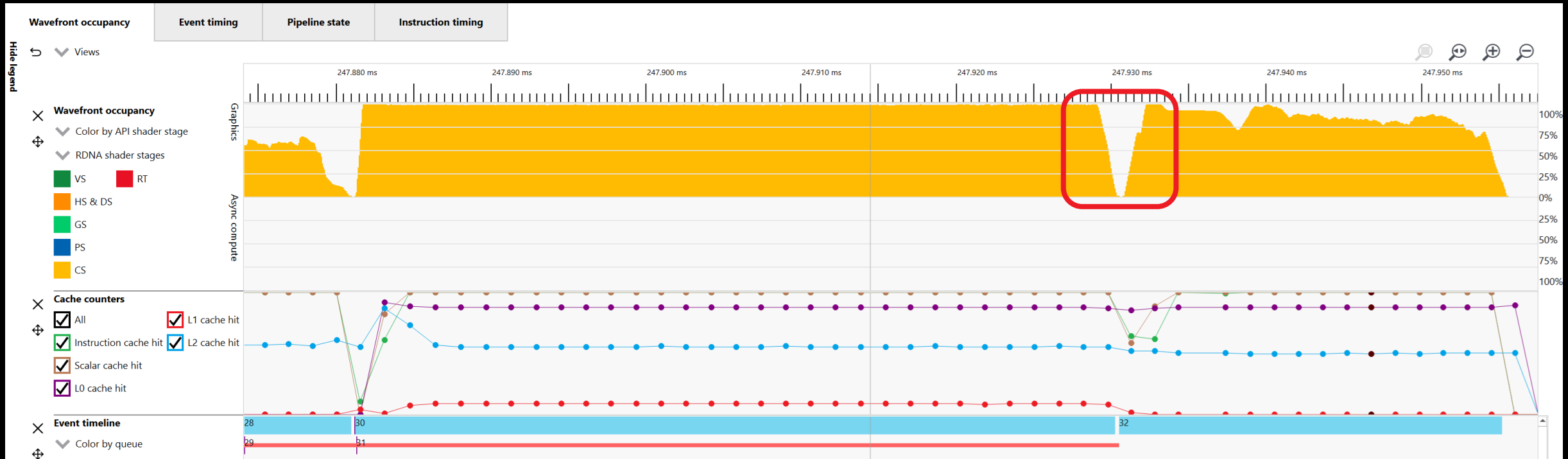**Wave slots: 16 * 192 = 3072**

# Max occupancy
# 510 / 3072 = 16.6%
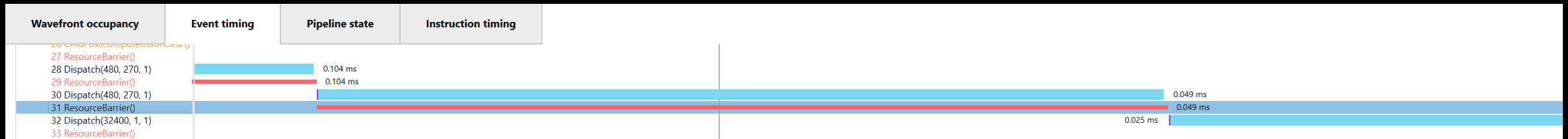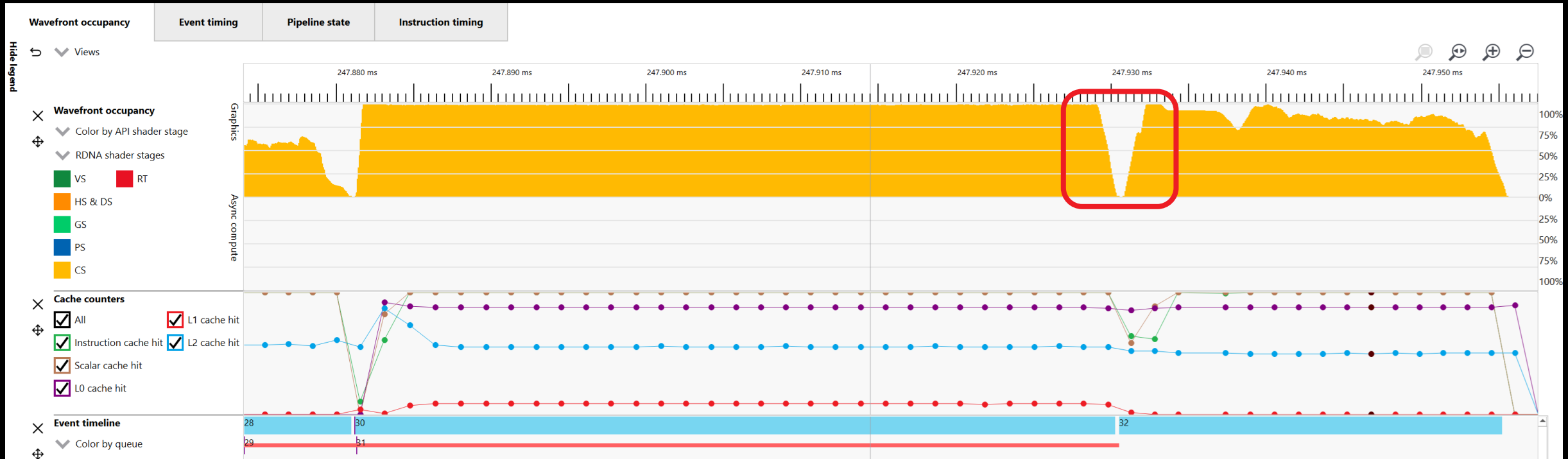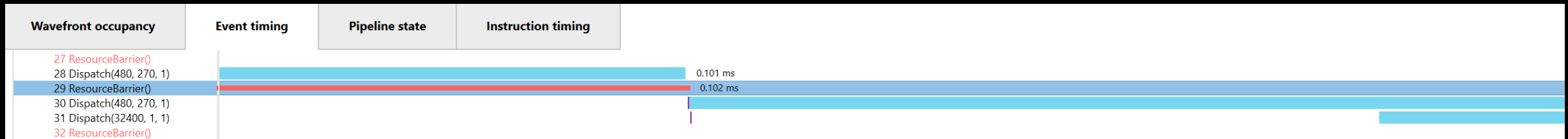
# FILL THE GPU WITH ENOUGH WORK
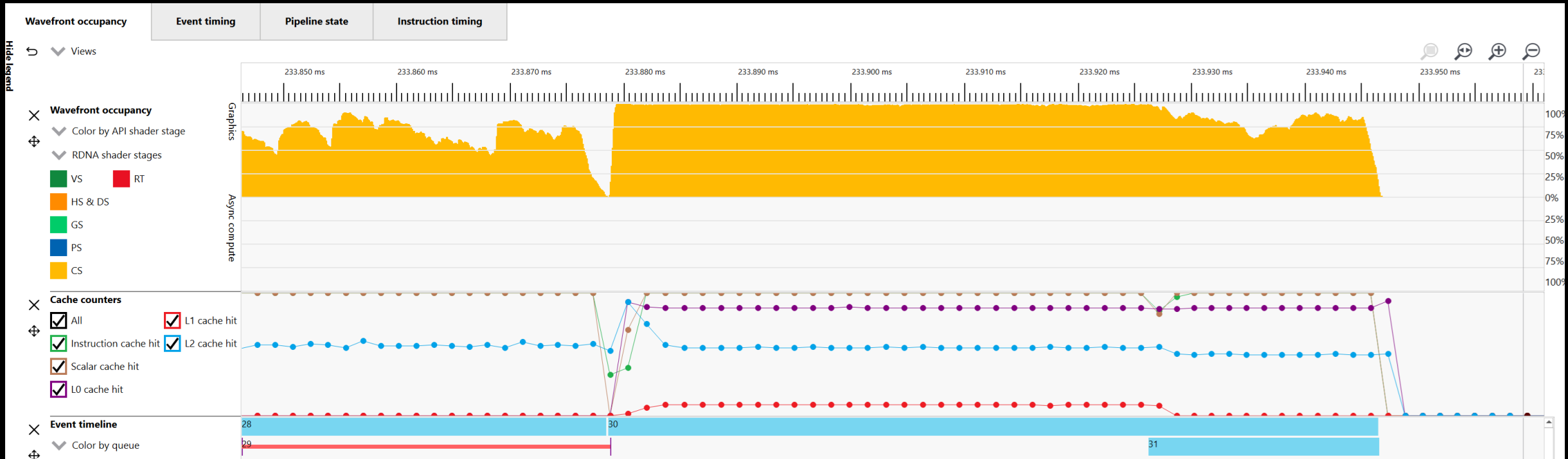
# FEED DEM GPUs

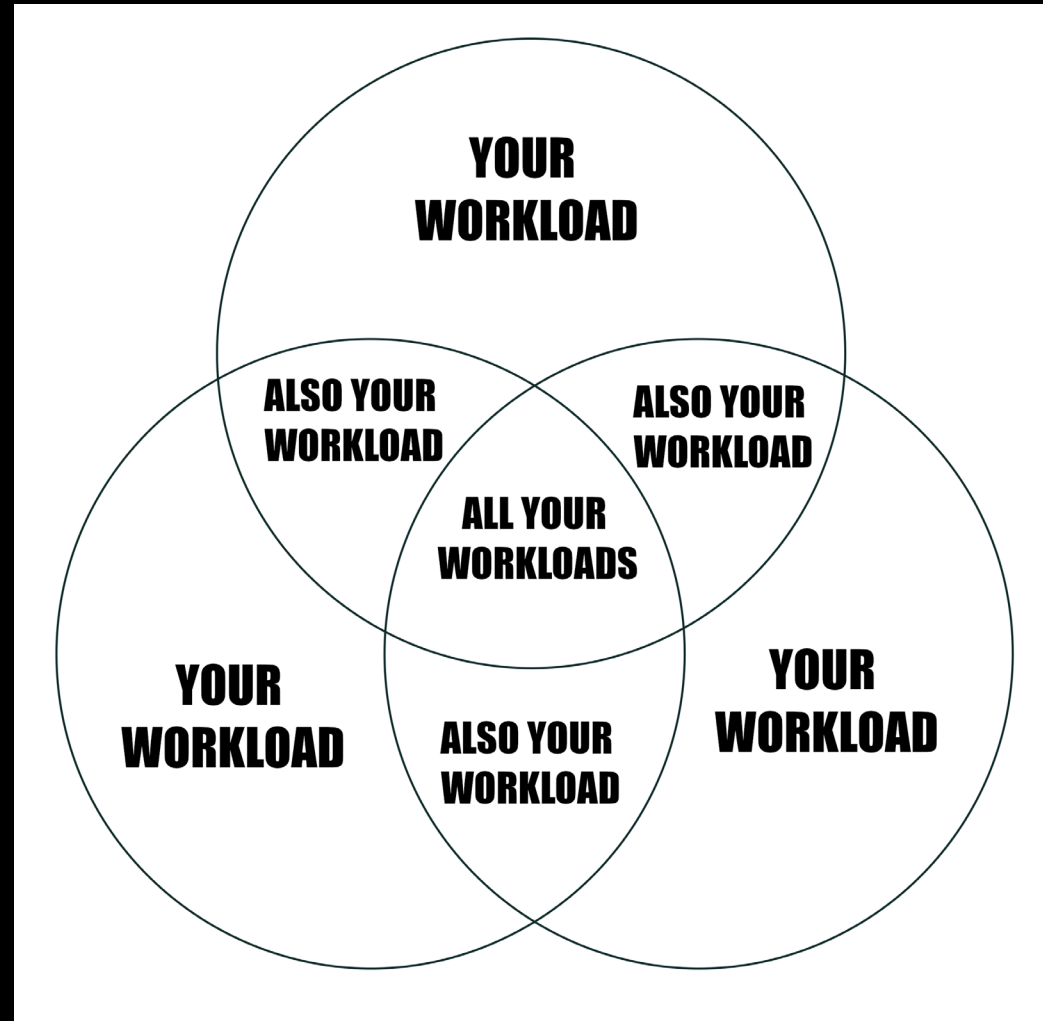# OCCUPANCY GAP

# OCCUPANCY GAP

# LET YOUR WORKLOADS OVERLAP
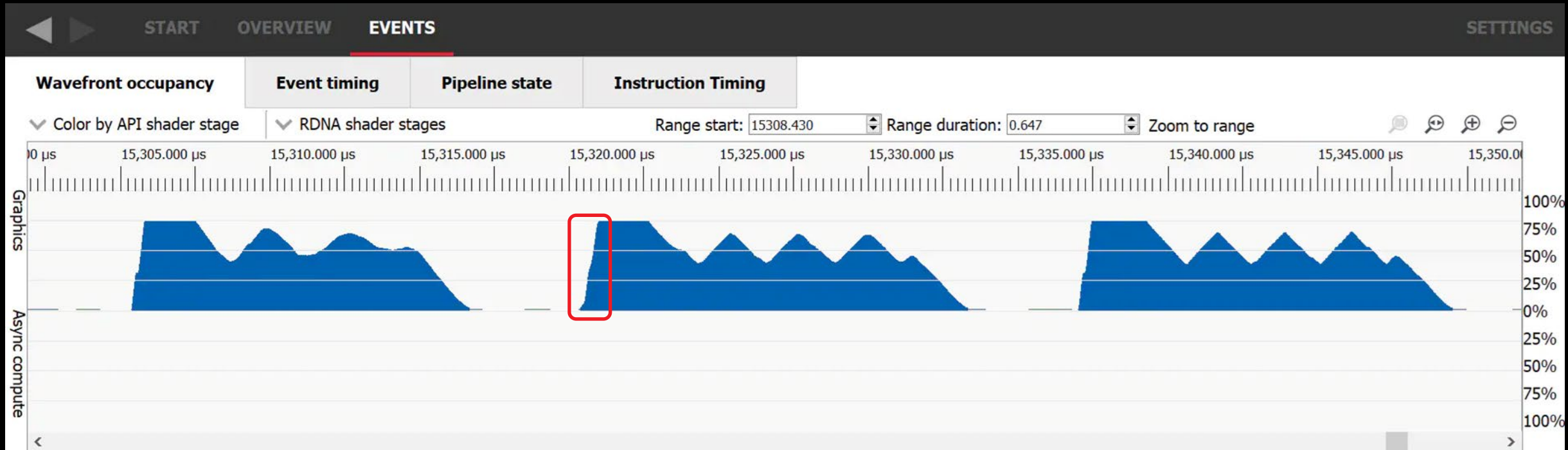
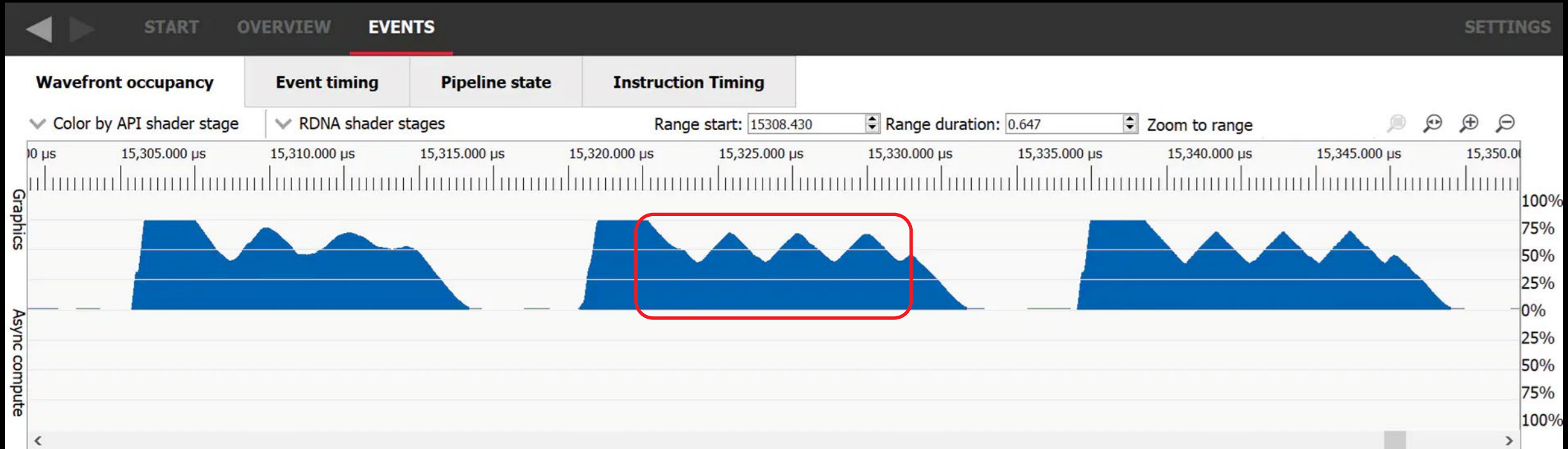# LET YOUR WORKLOADS OVERLAP

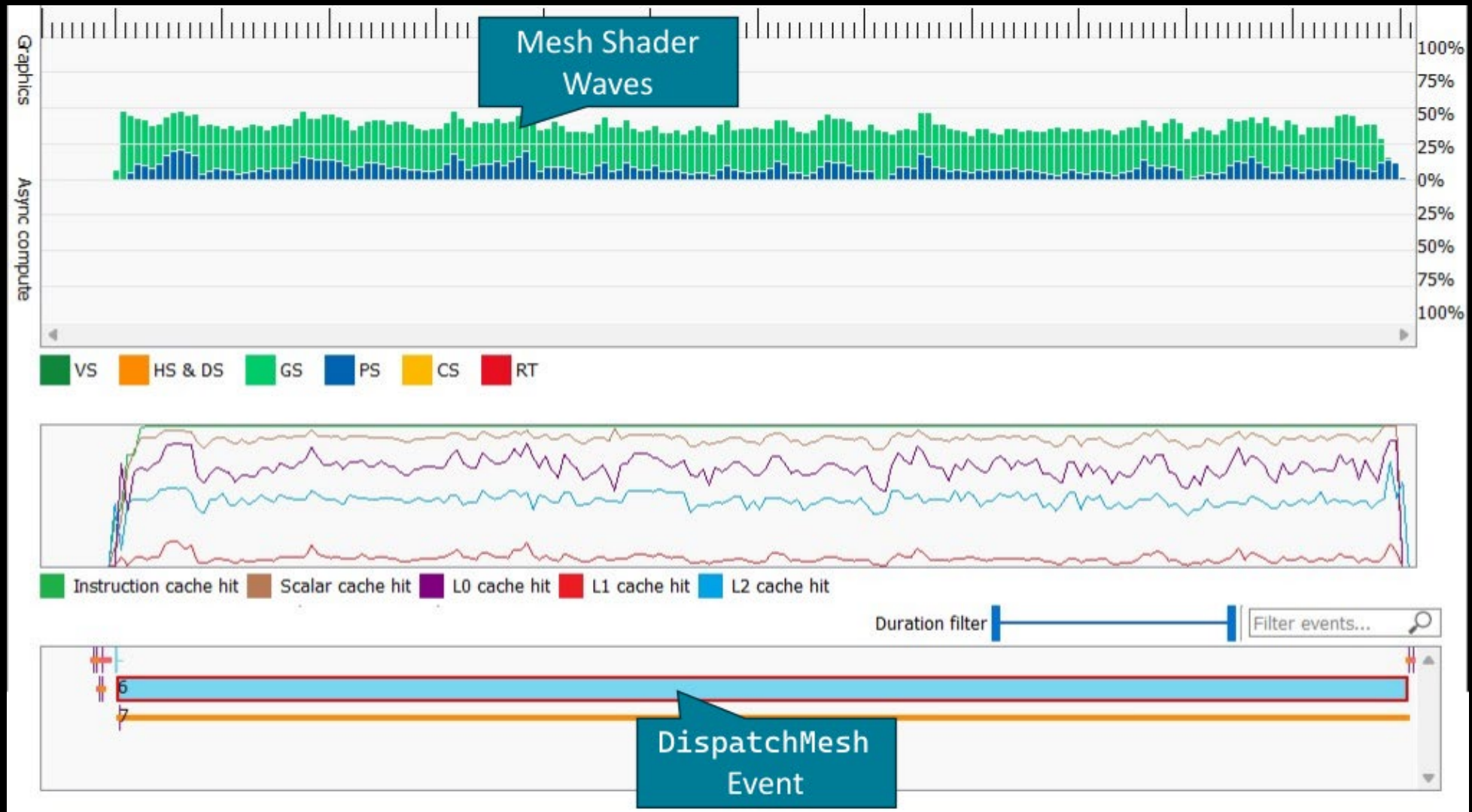# LET YOUR WORKLOADS OVERLAP

# LET YOUR WORKLOADS OVERLAP

# LAUNCH RATE LIMITED WORKLOAD

# LAUNCH RATE LIMITED WORKLOAD

# GEOMETRY WORKLOADS

# GEOMETRY WORKLOADS



**Mesh shaders on AMD RDNA™ graphics cards**
- From vertex shader to mesh shader
- Optimization and best practices
- Font- and vector-art rendering with mesh shaders
- Procedural grass rendering

## Mesh shaders on AMD RDNA™ graphics cards 🔗

Despite the flexibility and performance mesh shading can add to the geometry stage, we find that the technology has not been widely adopted in rendering engines so far. The purpose of this article series is to revisit mesh shading five years after its initial rollout between 2018-2019.

As a result, this blog series aims to demystify mesh shading by providing more detailed explanations, analysis, use-case examples, tutorials, and general advice.

- Part 1: From vertex shader to mesh shader
- Part 2: Optimization and best practices
- Part 3: Font- and Vector-Art Rendering with Mesh Shaders
- Part 4: Procedural grass rendering

---

### Max Oberberger

*Max is part of AMD's GPU Architecture and Software Technologies Team. His current focus is GPU work graphs and mesh shader research.*

### Bastian Kuth

*Bastian is a PhD candidate at Coburg University and University of Erlangen-Nuremberg. His research focuses on real-time geometry processing on GPUs.*
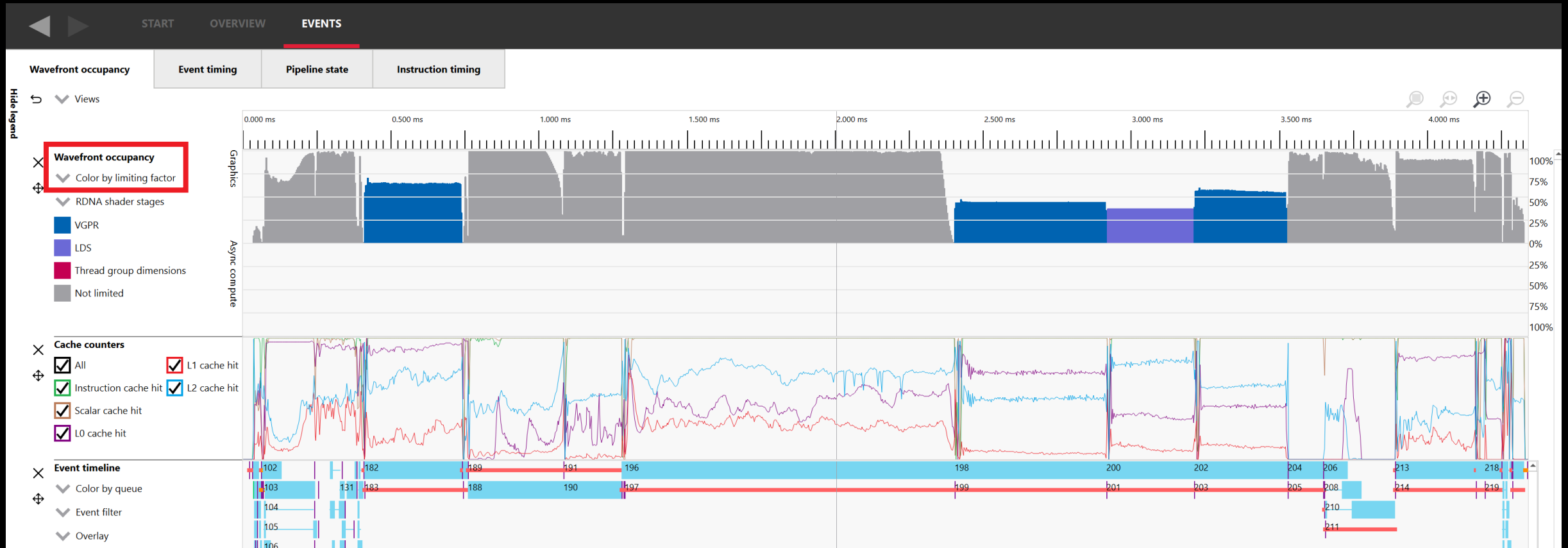
### Quirin Meyer

*Before becoming a computer graphics professor at Coburg University, Quirin Meyer obtained a Ph.D. in graphics and worked as a software engineer in the industry. His research focuses on real-time geometry processing primarily on GPUs.*
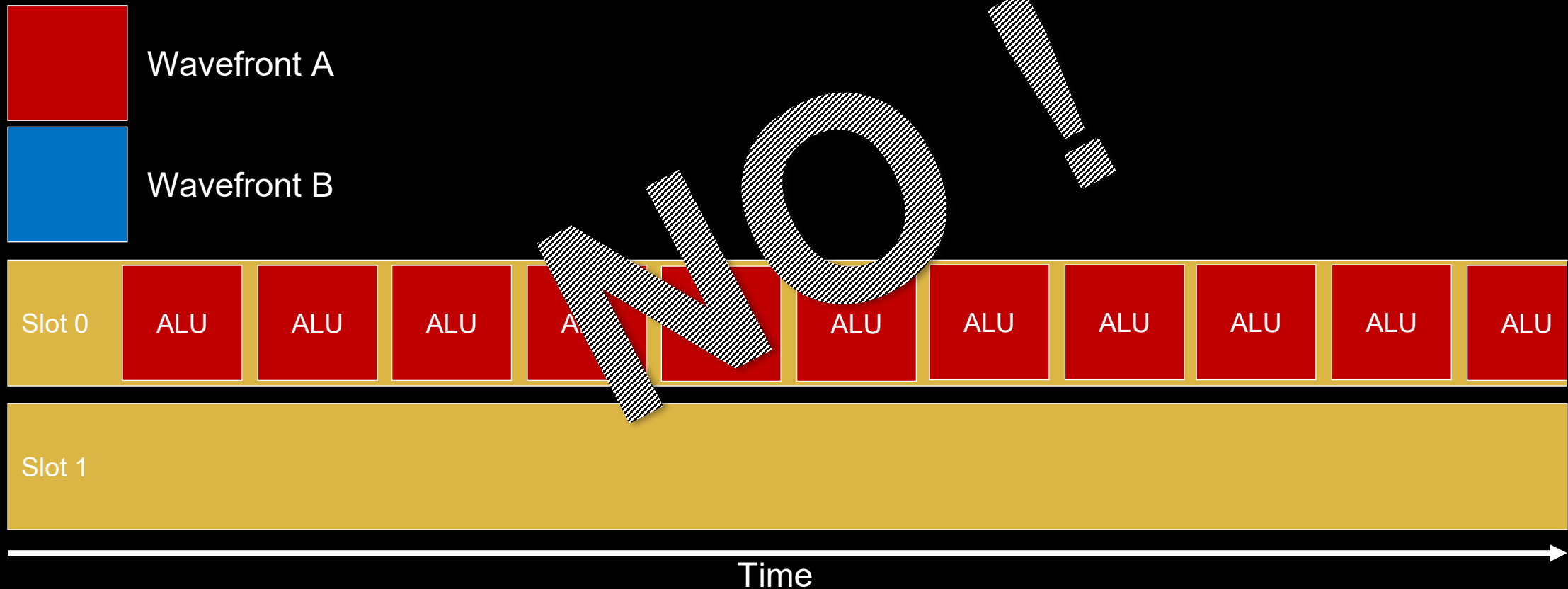
# OCCUPANCY LIMITERS

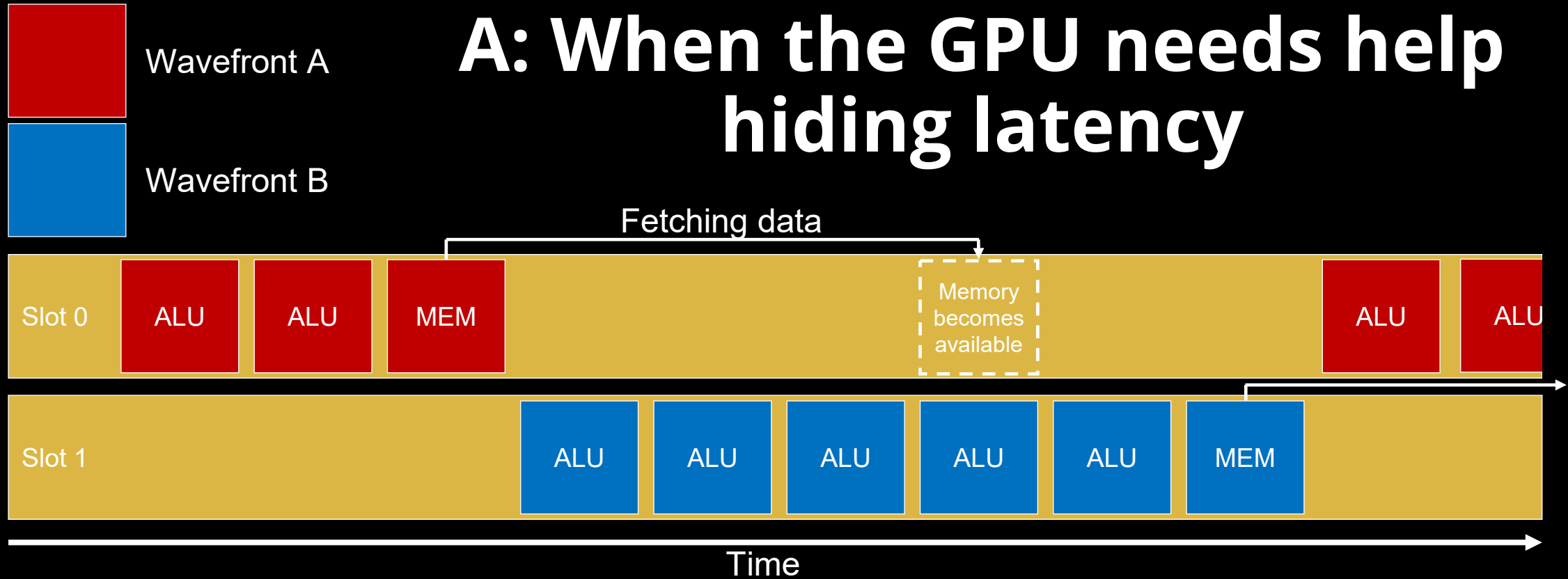# Q: Does better occupancy necessarily mean better performance?

# Q: When should I care about occupancy?

# Q: Does maximum occupancy mean that all the memory access latency from my shader is hidden?

# Q: Does maximum occupancy mean that all the memory access latency from my shader is hidden?

# Q: Is lower theoretical occupancy always bad for performance?

# Q: Is lower theoretical occupancy always bad for performance?

NO !

# Q: Is lower theoretical occupancy always bad for performance?

NO !
BUT ALSO,
YES?

# REGISTER SPILLING



| IA | VS | HS | DS | GS | RS | PS | OM | CS |

**Information**   ISA

**Dispatch properties**

| Total thread groups | Thread group dimensions | Ordered append |
|---|---|---|
| {480, 270, 1} | {8, 8, 1} | ⬤ OFF |

Strict shader processor interpolator (SPI) ordering
⬤ OFF

**Wavefronts and threads**

| Total wavefronts | Total threads | Average wavefront duration | Average threads per wavefront |
|---|---|---|---|
| 129,600 | 8,294,400 | 0.026 ms | 64 |

Wavefront mode
wave64

**Per-wavefront resources**

| Vector registers | Scalar registers | Registers spilled to scratch memory | Local data share per thread group |
|---|---|---|---|
| 135 (144 allocated) | 88 (128 allocated) | ⬤ ON | - |

**Theoretical wavefront occupancy**

The occupancy of this shader is limited by its vector register usage.
This shader could potentially run 5 wavefronts out of 16 wavefronts per SIMD.

However, if you reduce vector register usage by 16 you could run another wavefront.

# REGISTER SPILLING



| IA | VS | HS | DS | GS | RS | PS | OM | CS |

**Information**  ISA

**Dispatch properties**

Total thread groups
{480, 270, 1}

Thread group dimensions
{8, 8, 1}

Ordered append
◯ OFF

Strict shader processor interpolator (SPI) ordering
◯ OFF

**Wavefronts and threads**

Total wavefronts
129,600

Total threads
8,294,400

Average wavefront duration
0.026 ms

Average threads per wavefront
64

Wavefront mode
wave64

**Per-wavefront resources**

Vector registers
135 (144 allocated)

Scalar registers
88 (128 allocated)

Registers spilled to scratch memory
◯ ON

Local data share per thread group
-

**Theoretical wavefront occupancy**

The occupancy of this shader is limited by its vector register usage.
This shader could potentially run 5 wavefronts out of 16 wavefronts per SIMD.

However, if you reduce vector register usage by 16 you could run another wavefront.

AMD GPUOpen

# COPYRIGHT AND DISCLAIMER

# THANK YOU !

✉ - francois.guthmann@amd.com
🦋 / X - @frguthmann

AMD↗
GPUOpen