

Apple Metal

optimisation

George Ostrobrod



2. Рисуем остаток совы

Why do we need optimisation?

Example task

very naive implementation

- 256 particles
- each sprite - an image with transparent parts
- each sprite's transform is updated each frame
- target framebuffer size: 1242 x 1242px



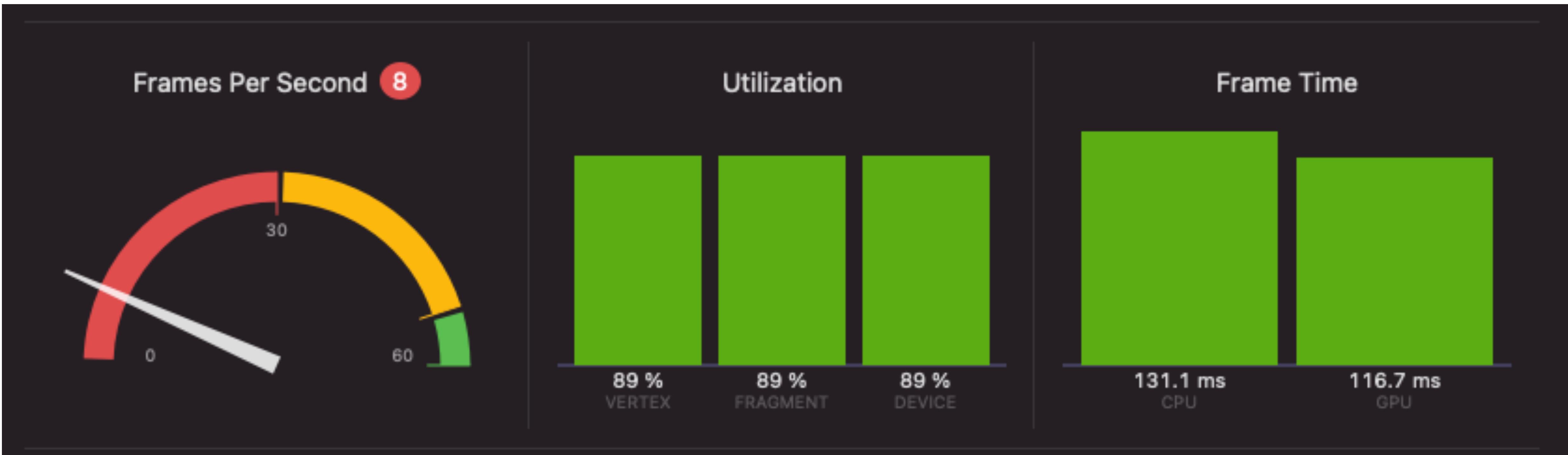
Example task

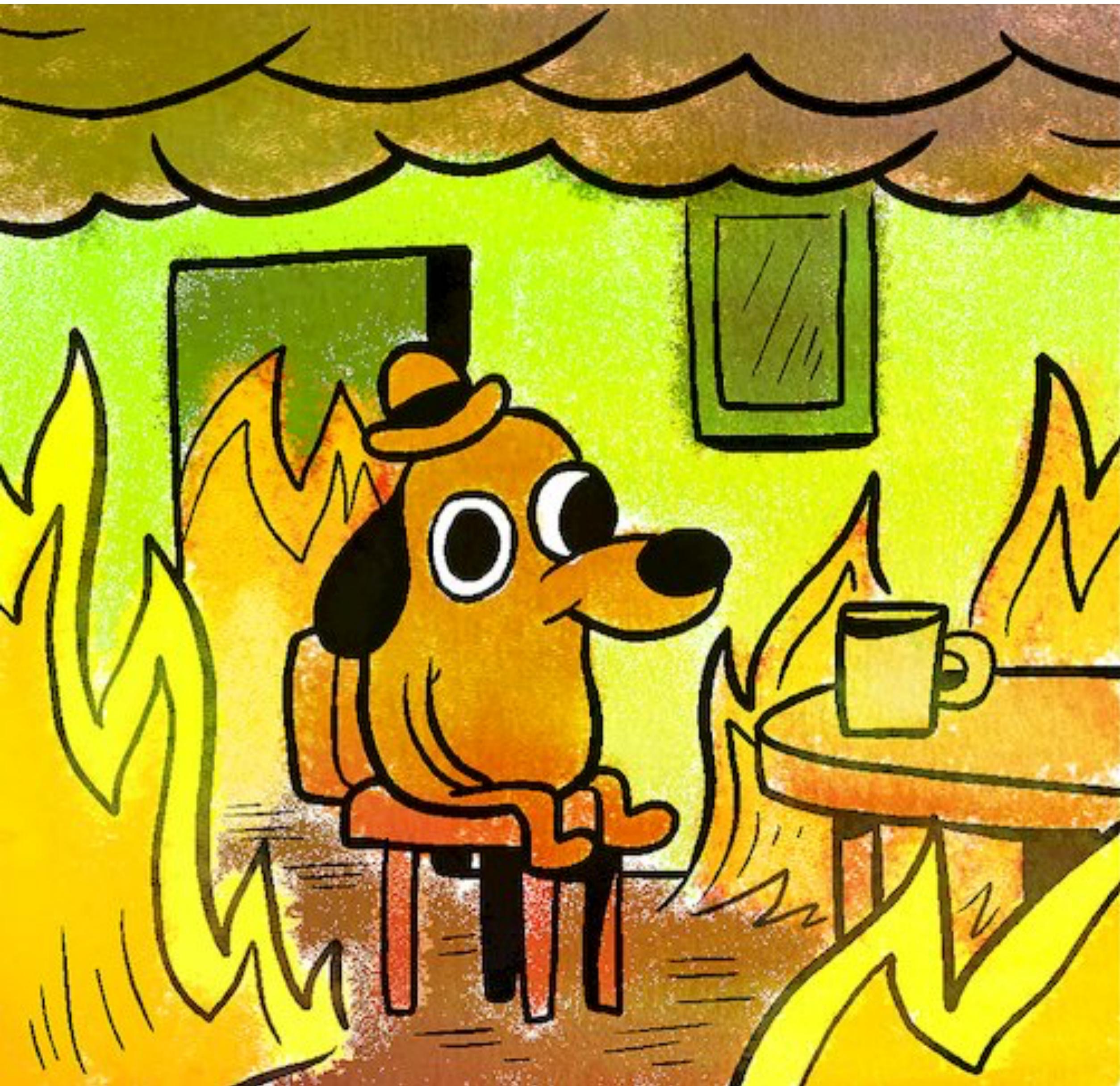
very naive implementation



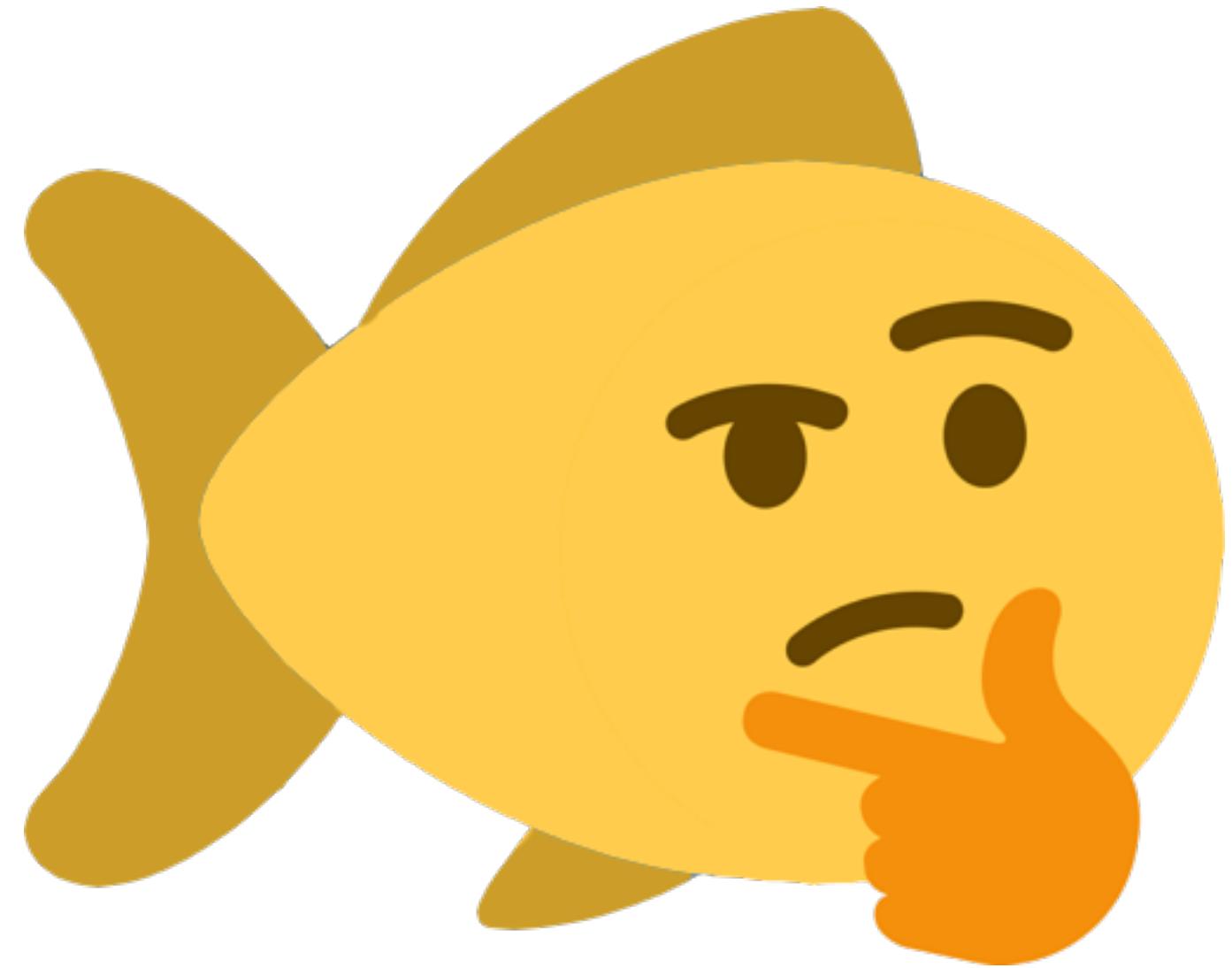
Example task

very naive implementation





Common sense vs Tools

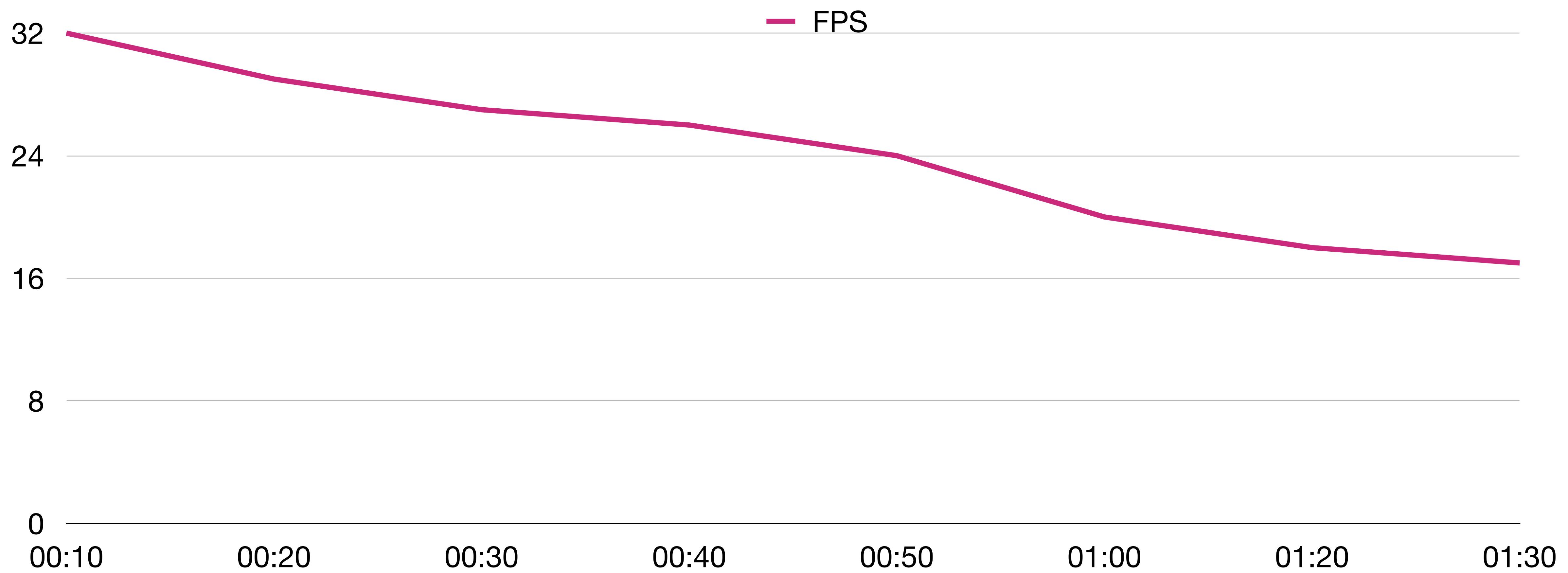


do not even try to optimise **sin** and **cos** in

$$\sin^2(t) + \cos^2(t) = 1$$

Throttling

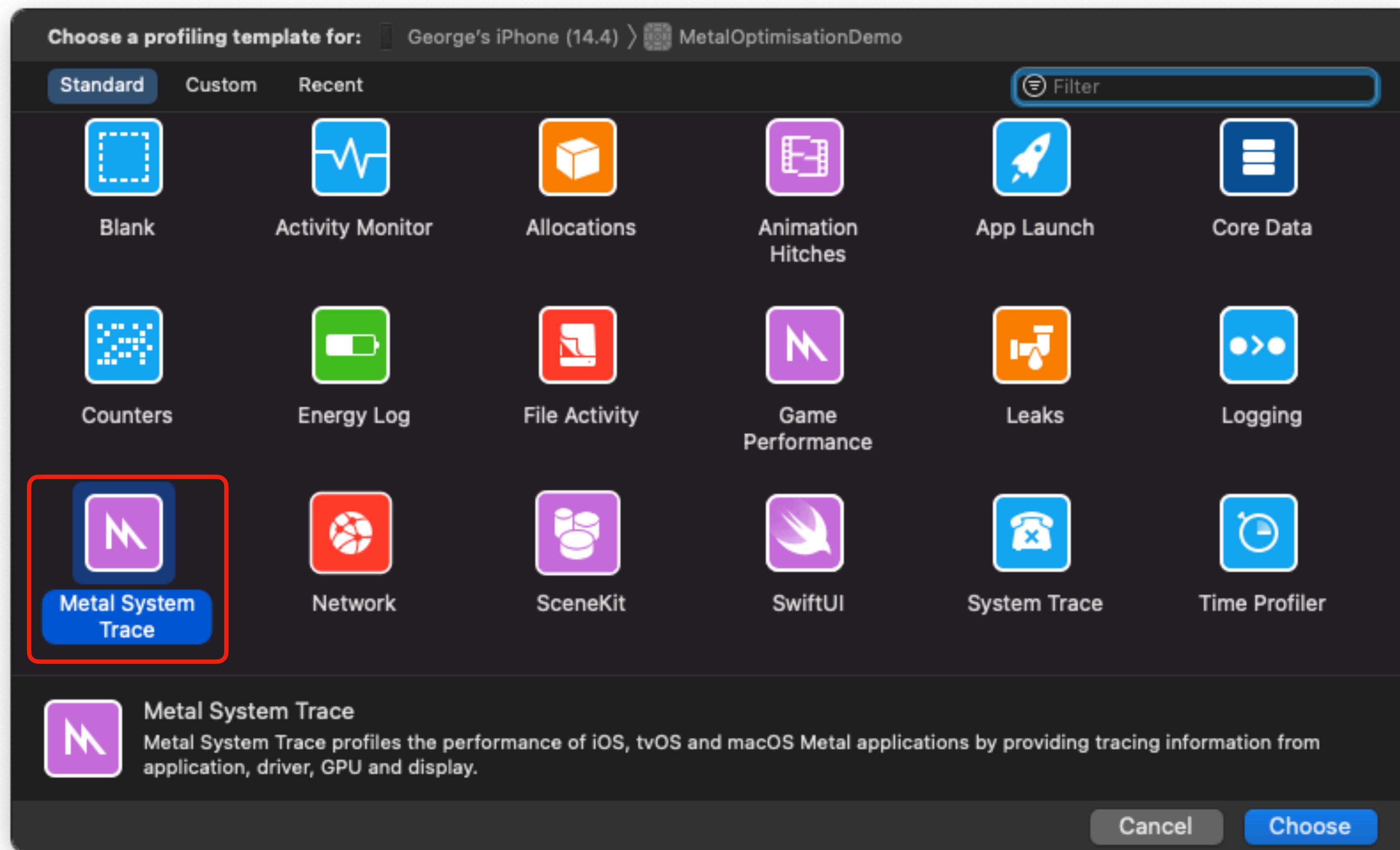
iOS is smart, too smart



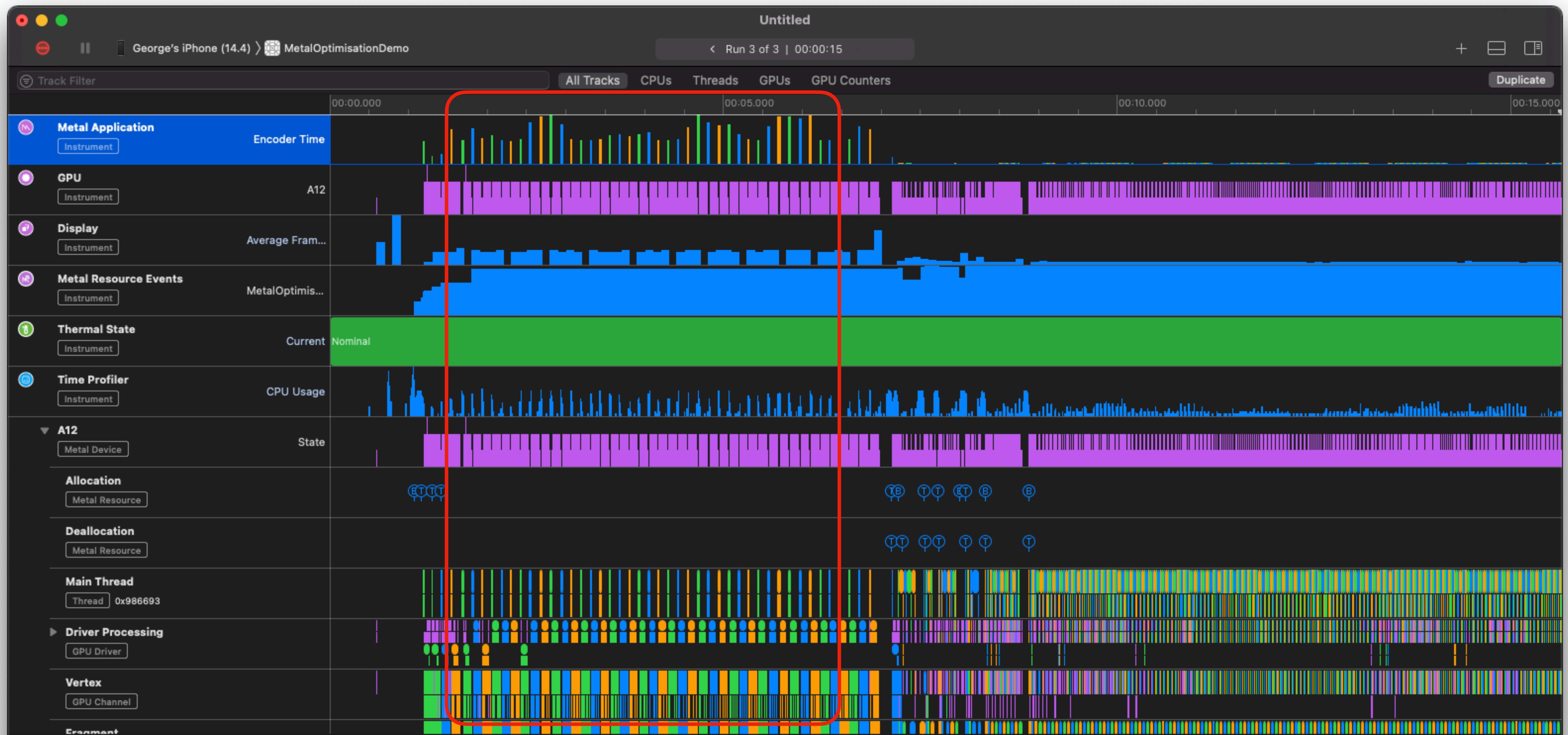
Profiling

Instruments

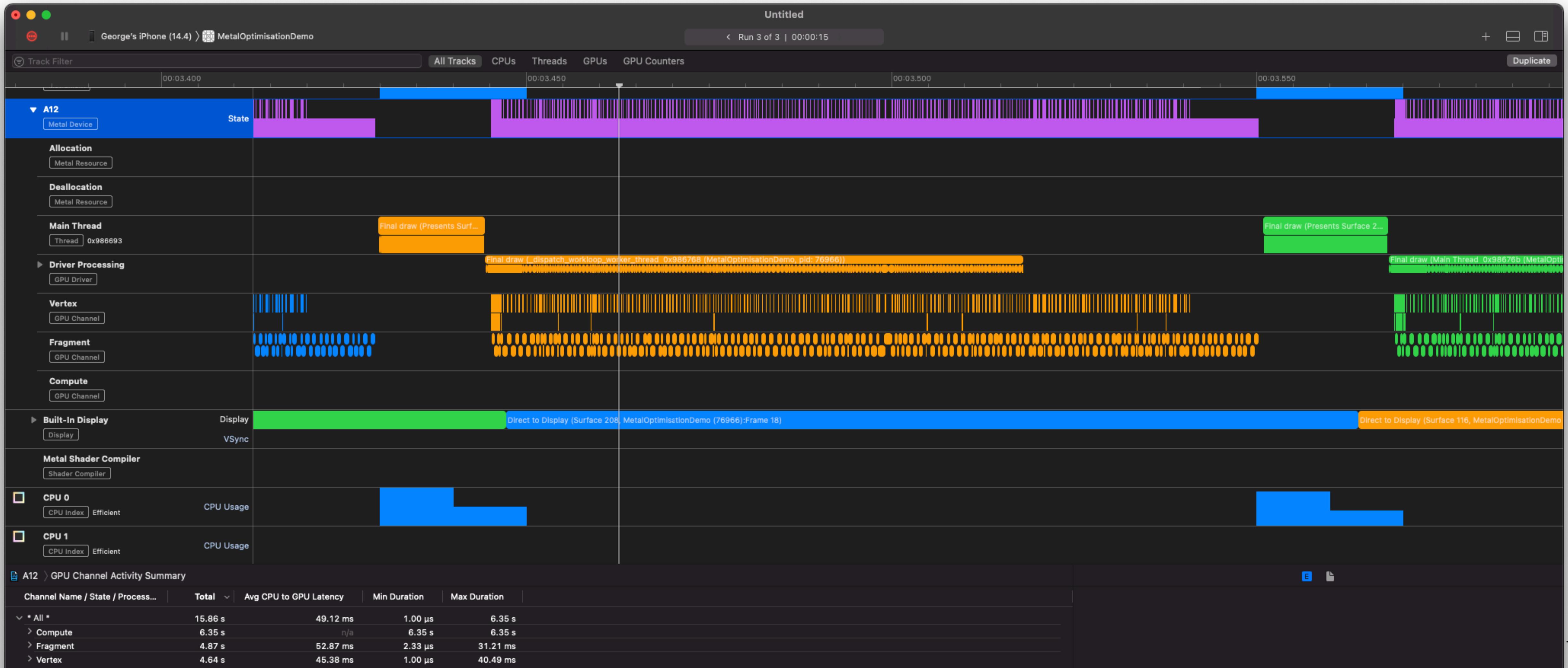
Metal System Trace



Metal System Trace



Metal System Trace





Rendering

Useless synchronisation

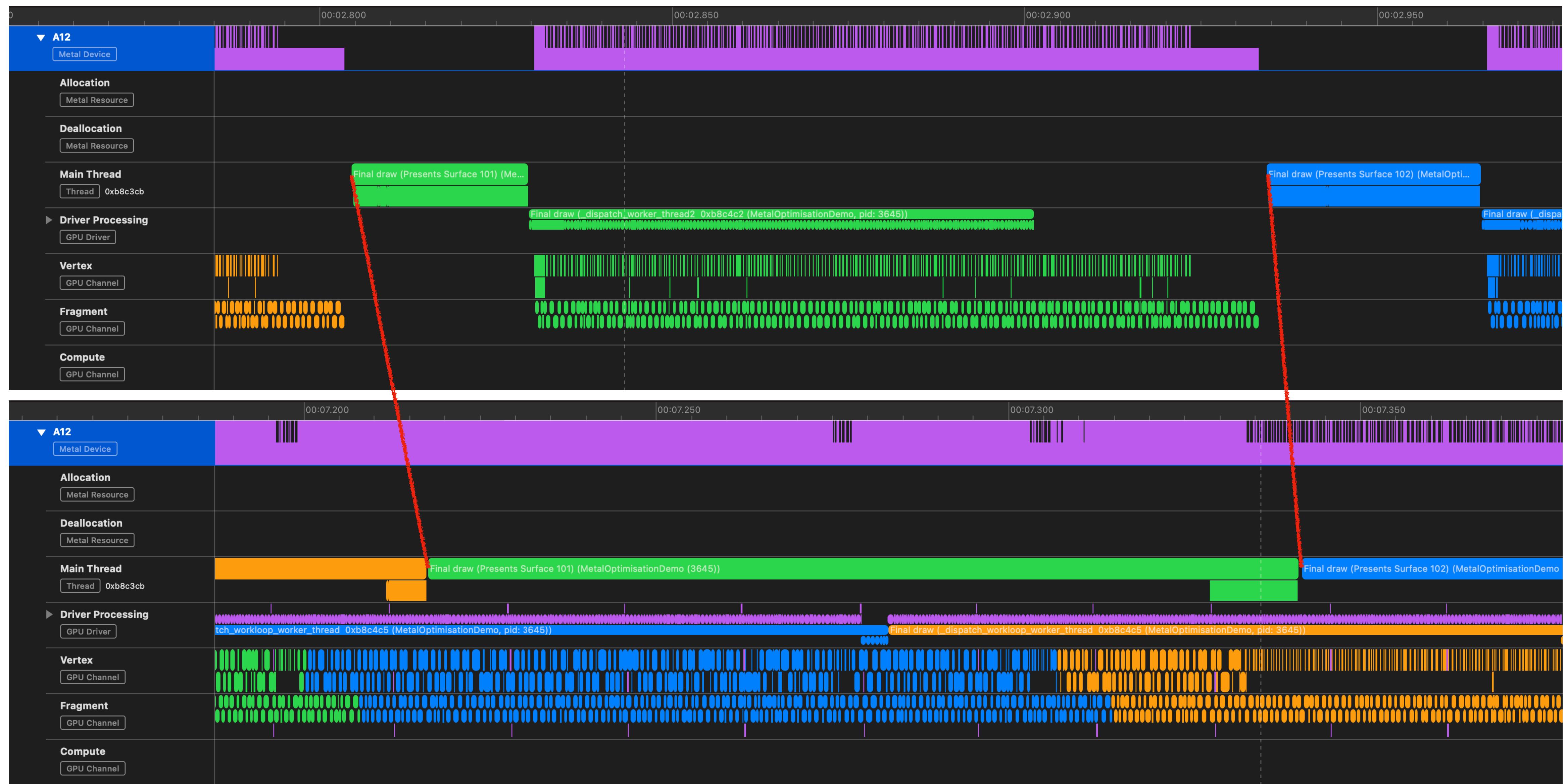
render and wait until buffer started running

```
// ...
commandBuffer.waitUntilScheduled()
// ...
```

render and wait until whole command buffer is completed

```
// ...
commandBuffer.waitUntilCompleted()
// ...
```

Useless synchronisation



Number of encoders

```
MTLRenderPassDescriptor* descriptor =
    self.mtkView.currentRenderPassDescriptor;
descriptor.colorAttachments[0].clearColor =
    (MTLClearColor){0, 0.2, 0.2, 1};
descriptor.colorAttachments[0].loadAction =
    MTLLoadActionClear;

for (int i = 0; i < MaxInstances; ++i) {
    id <MTLRenderCommandEncoder> encoder =
        [commandBuffer renderCommandEncoderWithDescriptor:
            descriptor];

    // . . .

    [encoder drawPrimitives:MTLPrimitiveTypeTriangleStrip
        vertexStart:0
        vertexCount:4];
    [encoder endEncoding];

    descriptor.colorAttachments[0].loadAction =
        MTLLoadActionLoad;
}
```

```
MTLRenderPassDescriptor* descriptor =
    self.mtkView.currentRenderPassDescriptor;
descriptor.colorAttachments[0].clearColor =
    (MTLClearColor){0, 0.2, 0.2, 1};
descriptor.colorAttachments[0].loadAction =
    MTLLoadActionClear;

id <MTLRenderCommandEncoder> encoder =
    [commandBuffer renderCommandEncoderWithDescriptor:
        descriptor];

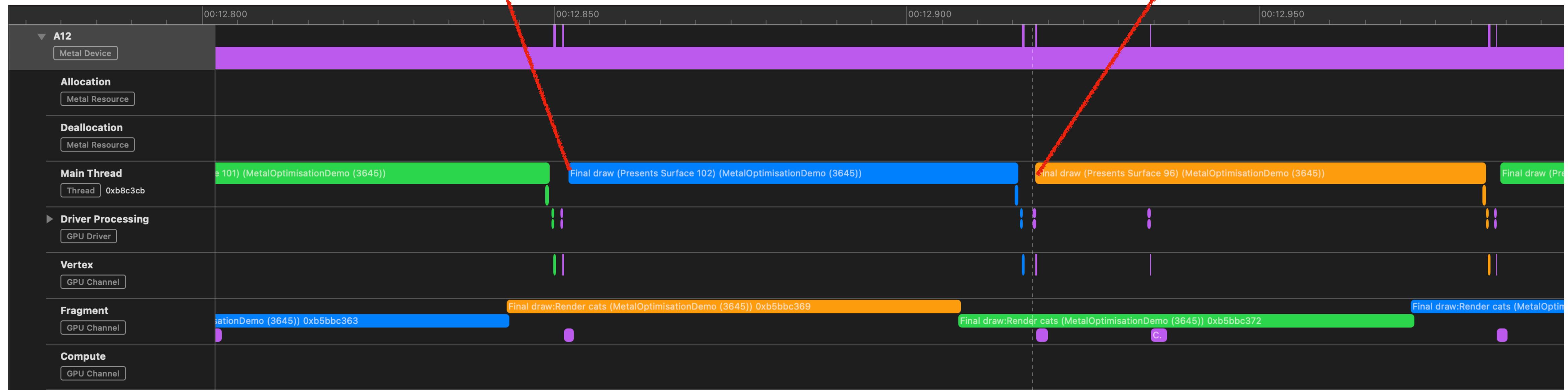
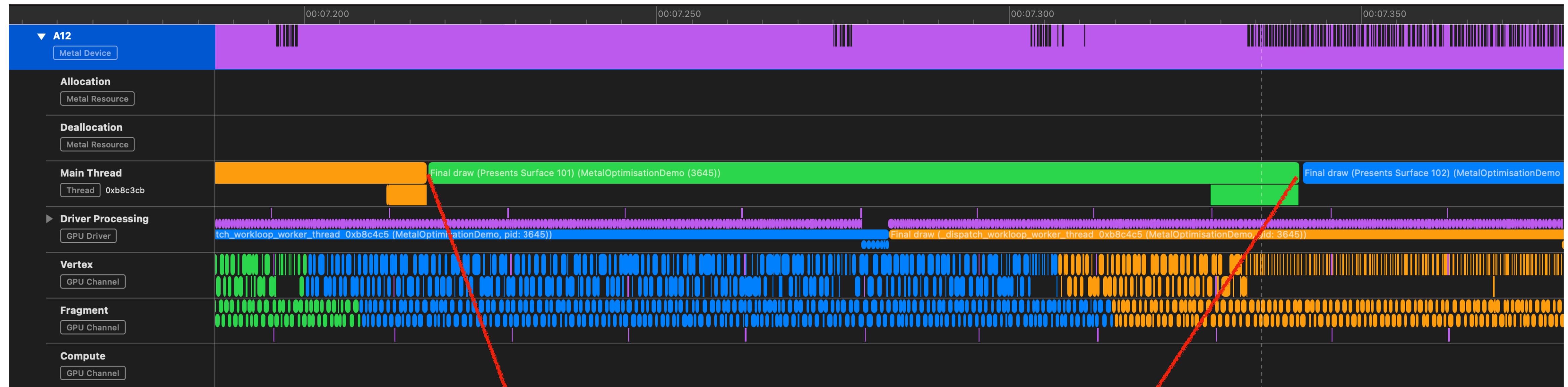
// . . .

for (int i = 0; i < MaxInstances; ++i) {
    // . . .

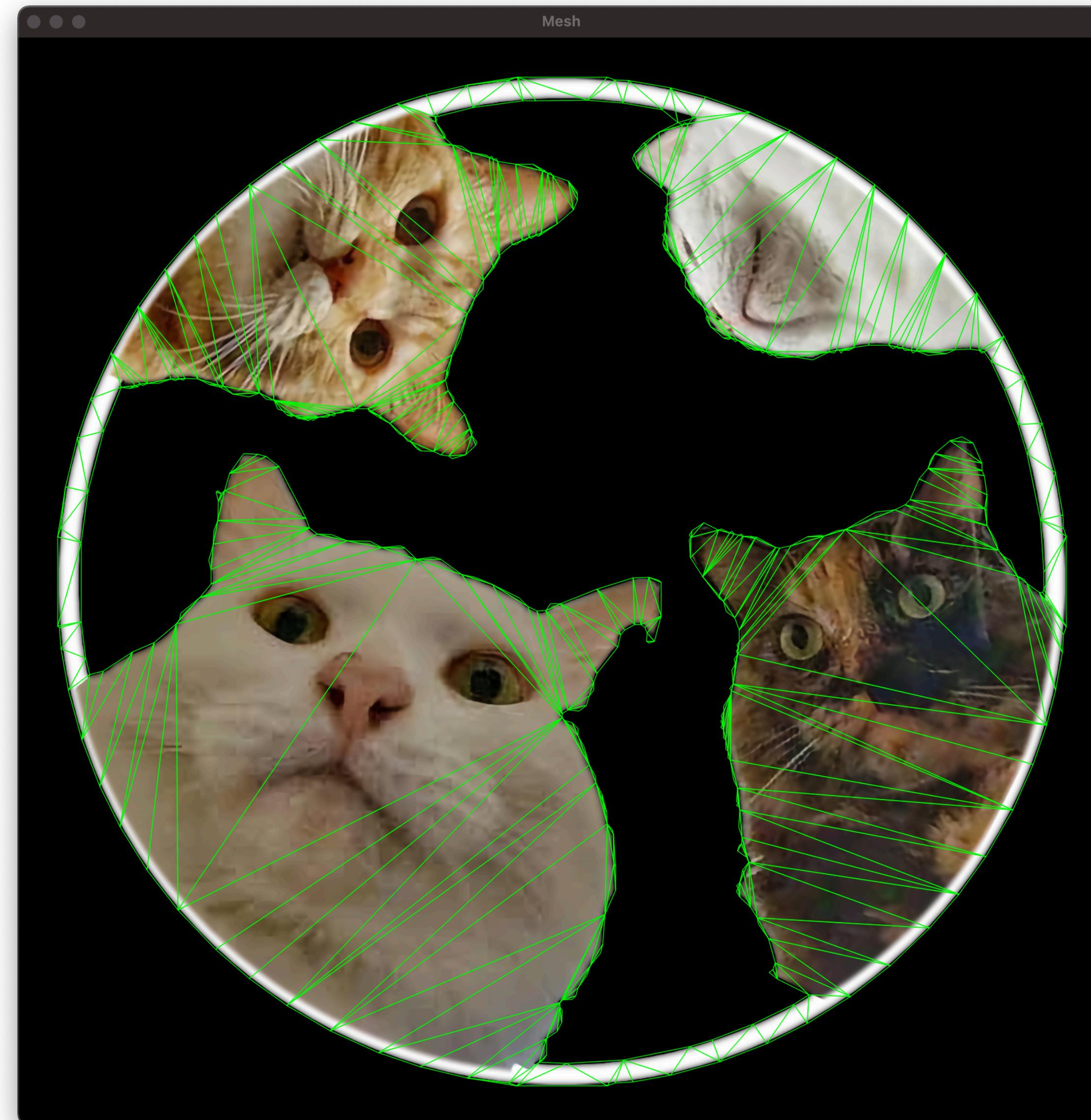
    [encoder
        drawPrimitives:MTLPrimitiveTypeTriangleStrip
        vertexStart:0
        vertexCount:4];
}

[encoder endEncoding];
```

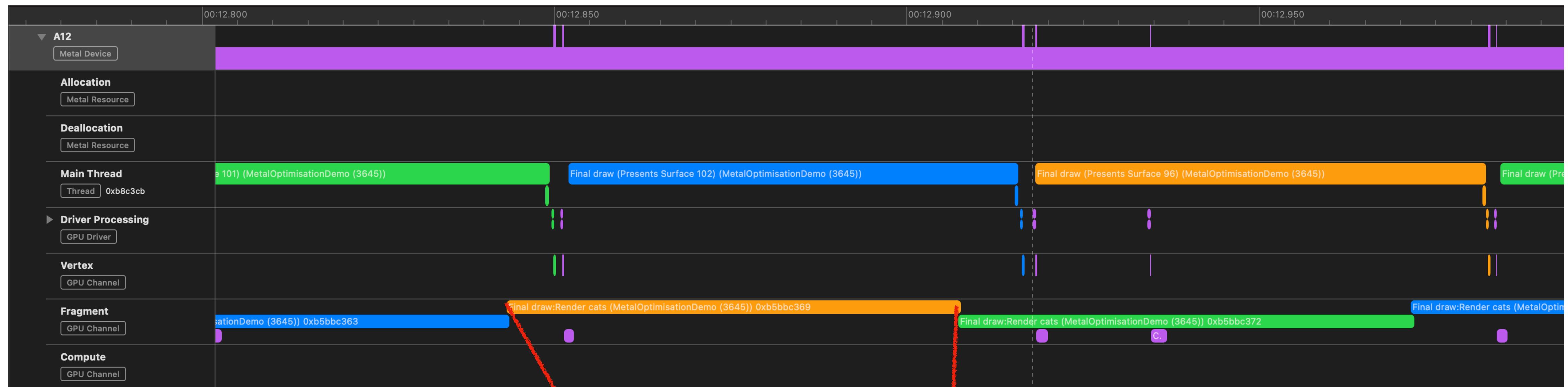
Number of encoders



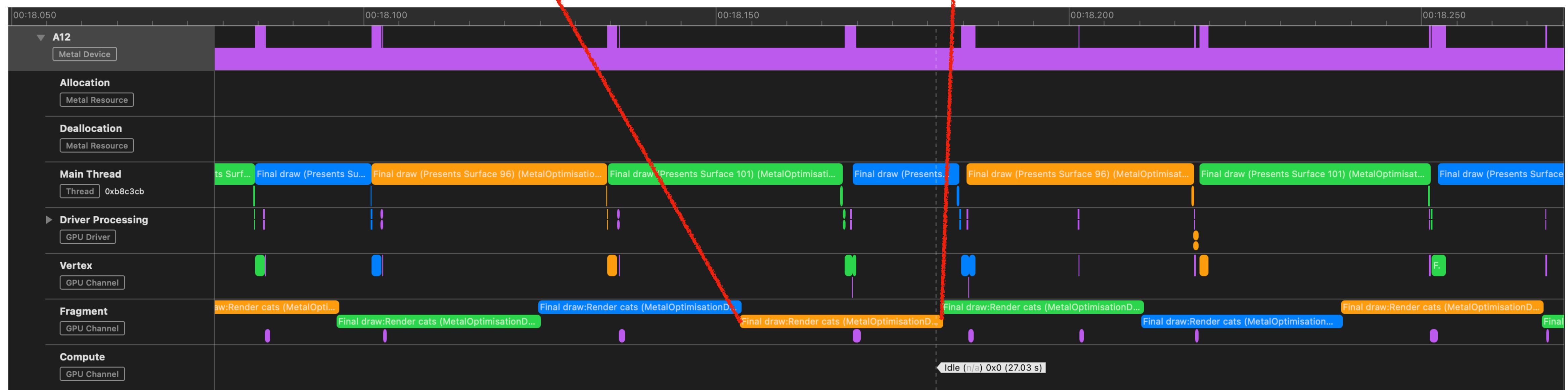
Fillrate



Fillrate



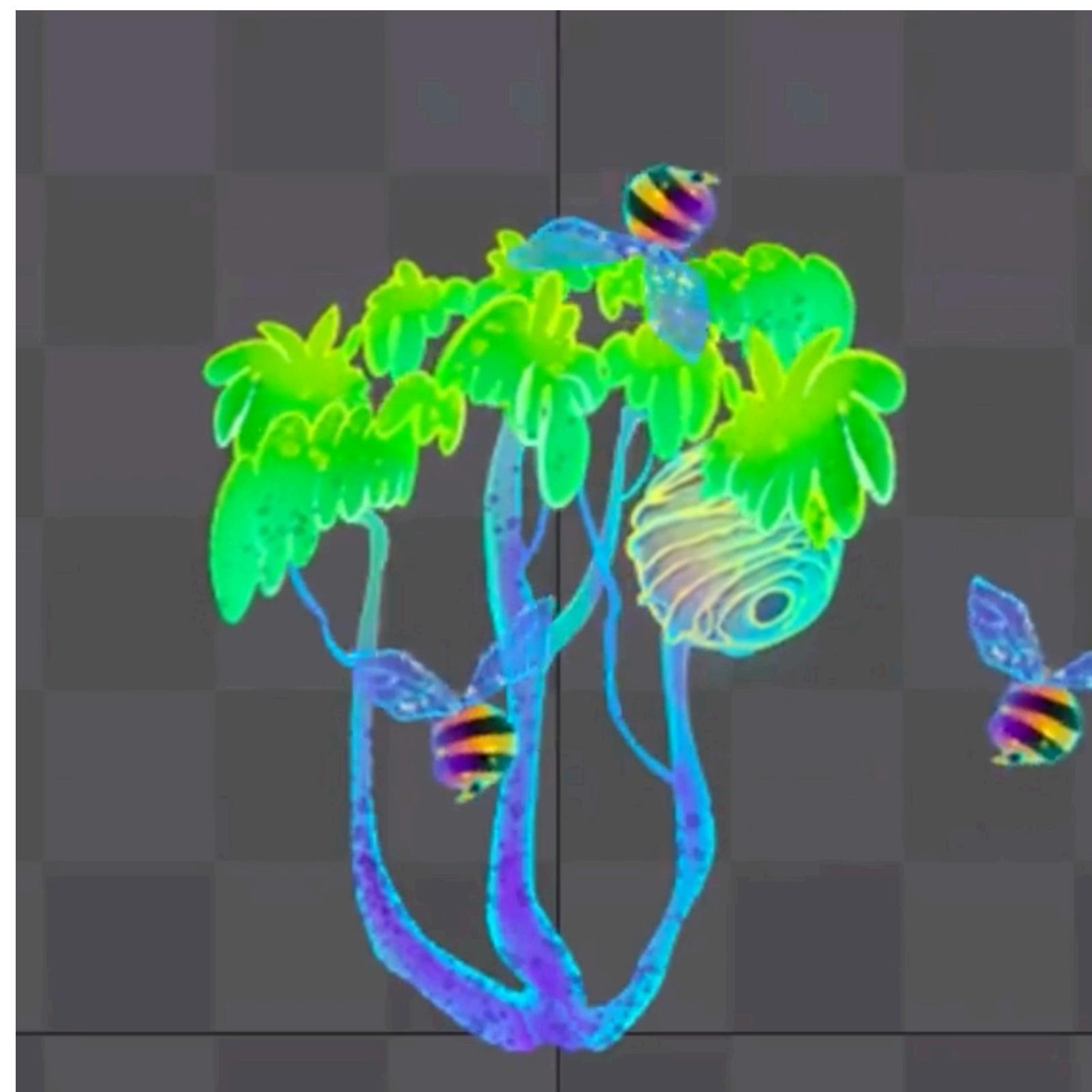
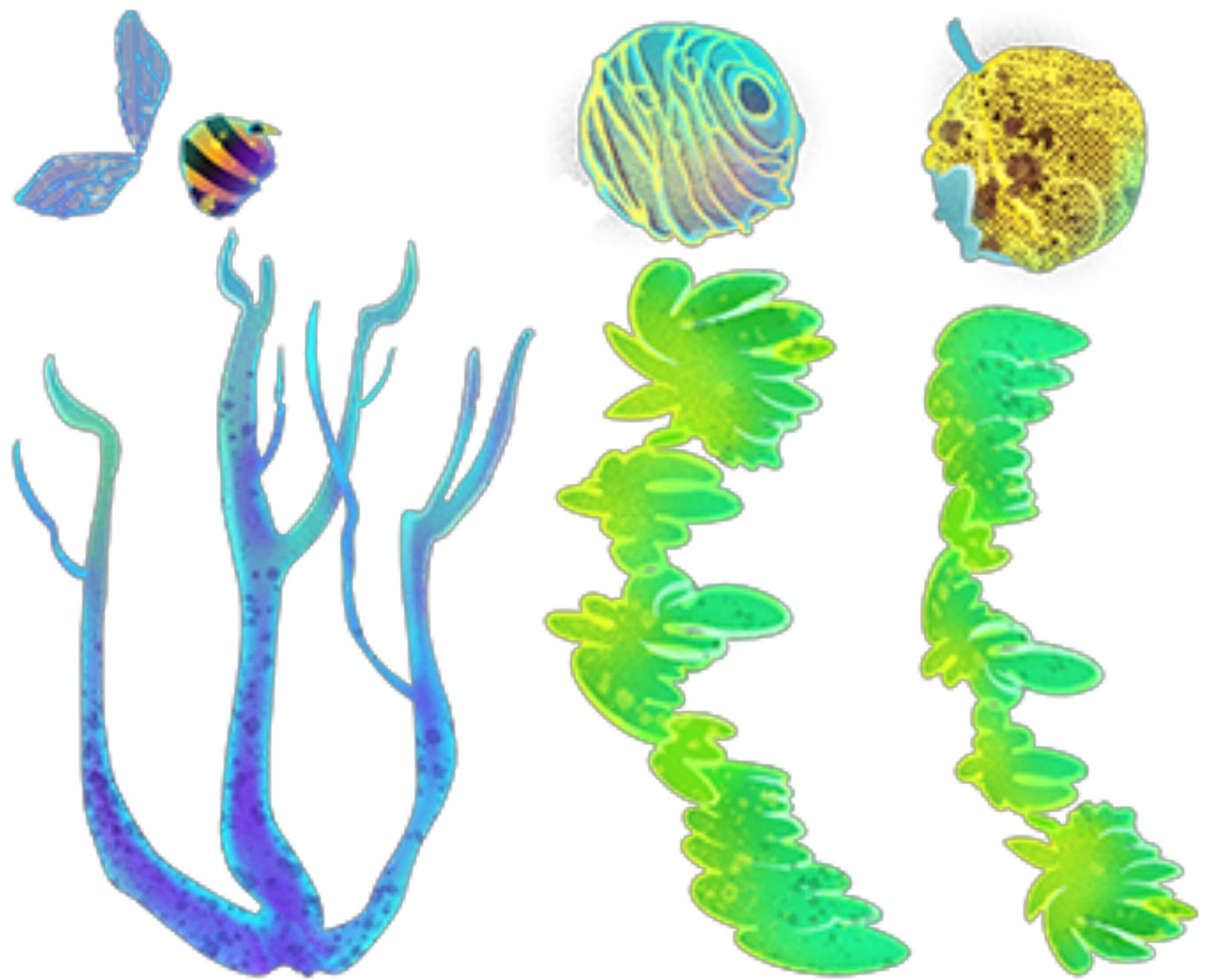
14 FPS



37 FPS

Number of draw calls

Texture atlases



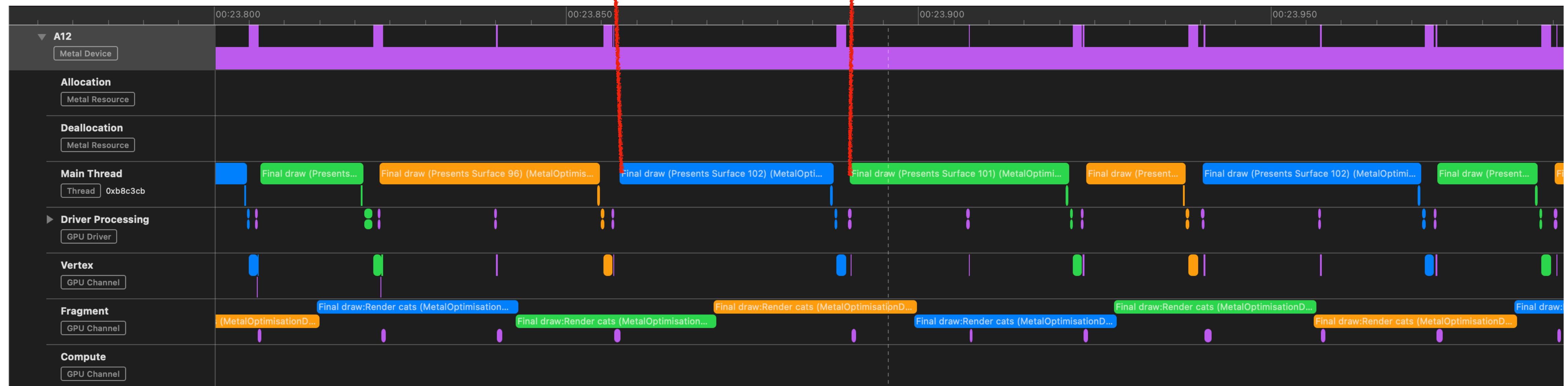
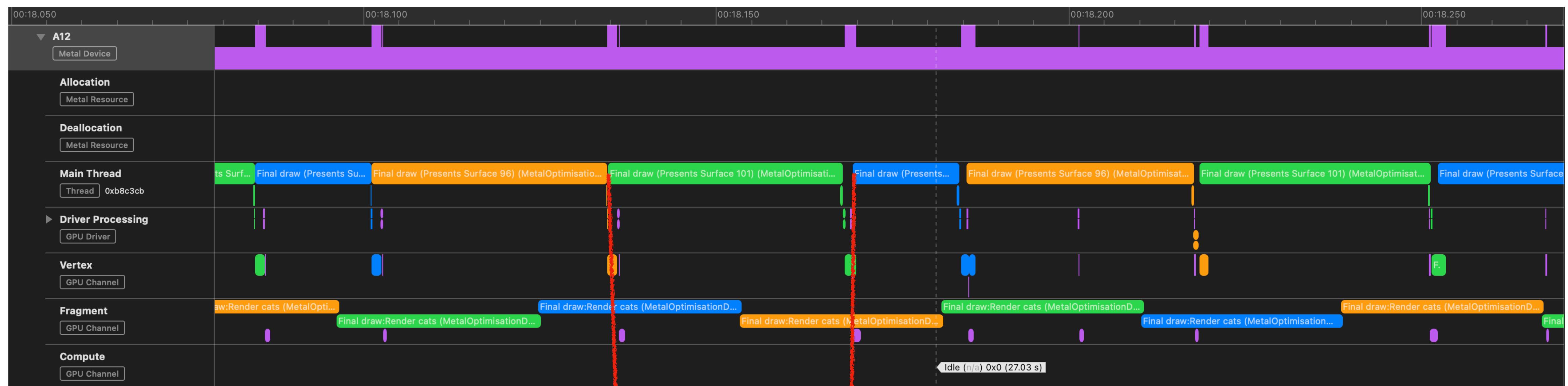
Number of draw calls

Instances

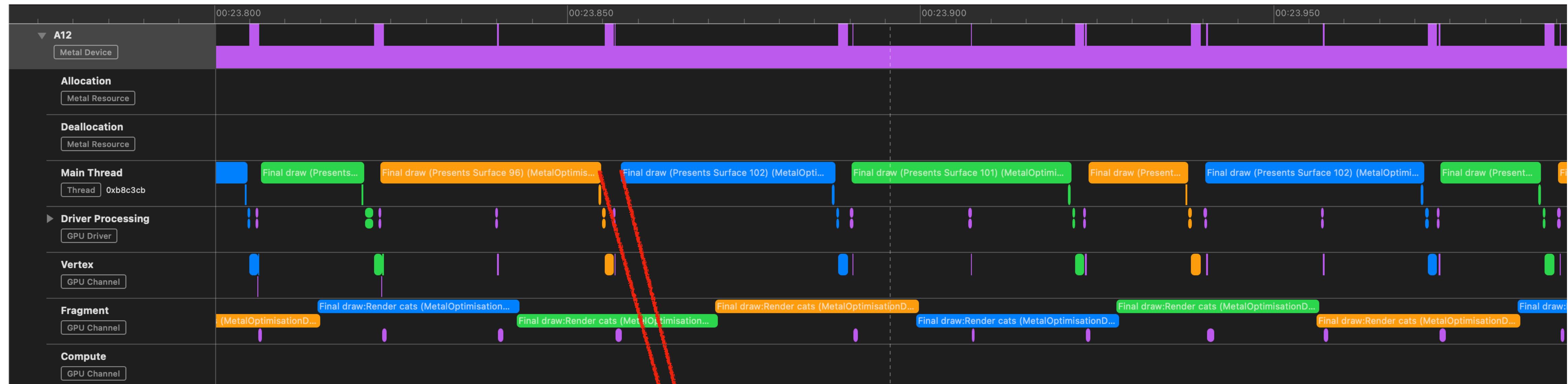
```
for (int i = 0; i < MaxInstances; ++i) {
    [encoder setVertexBytes:&i length:sizeof(int) atIndex:1];
    [encoder drawIndexedPrimitives:MTLPrimitiveTypeTriangle
        indexCount:_mesh.indices_count
        indexType:MTLIndexTypeUInt16
        indexBuffer:_bufferIndices
        indexBufferOffset:0];
}
```

```
[encoder drawIndexedPrimitives:MTLPrimitiveTypeTriangle
    indexCount:_mesh.indices_count
    indexType:MTLIndexTypeUInt16
    indexBuffer:_bufferIndices
    indexBufferOffset:0
    instanceCount:_instances];
```

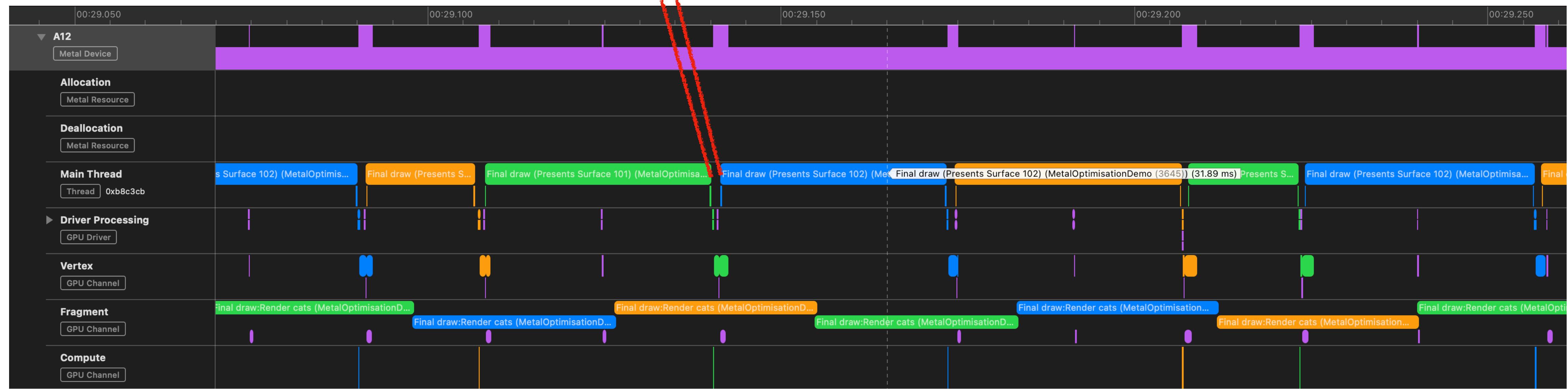
Number of draw calls



Update state on GPU



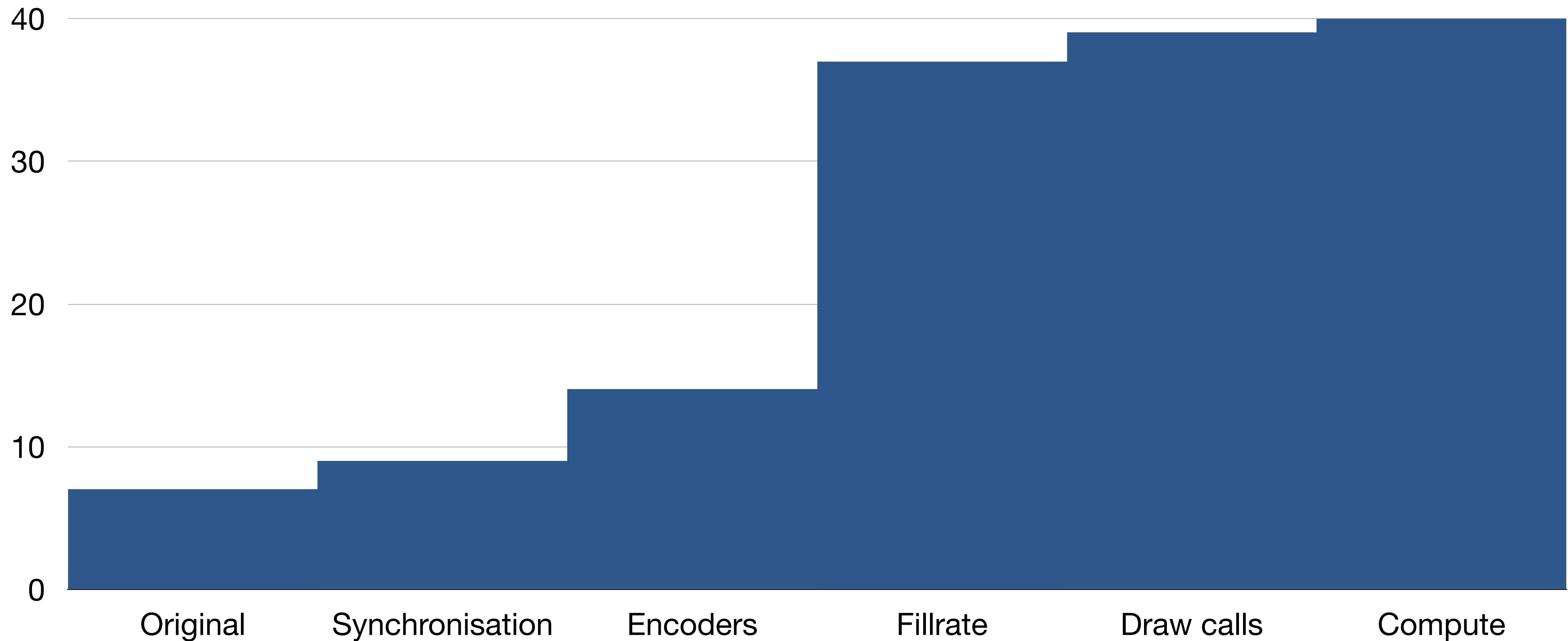
39 FPS



40 FPS

Summary (FPS)

Task 1: Render



Shaders

Sampling: `nearest` vs `linear`

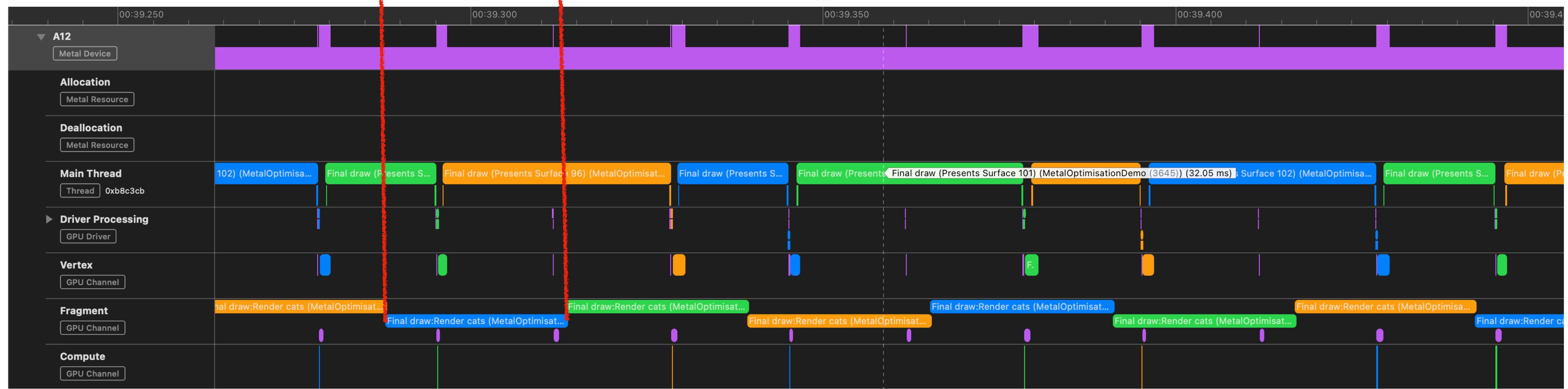
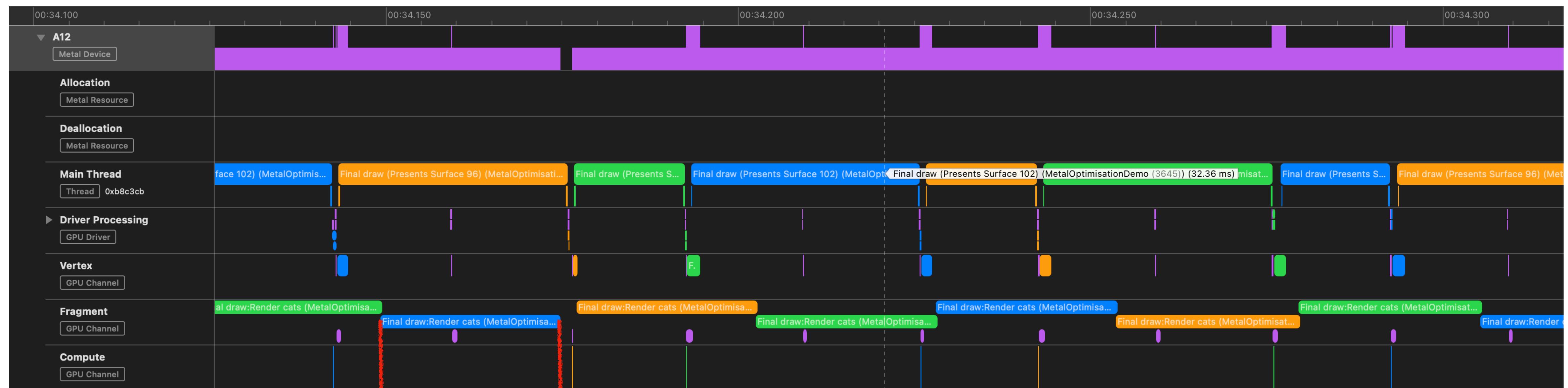
```
// ...
constexpr sampler imageSampler(address::clamp_to_zero, filter::linear);
// ...
```

```
// ...
constexpr sampler imageSampler(address::clamp_to_zero, filter::nearest);
// ...
```

Sampling: `nearest` vs `linear`



`half` vs `float`



`half` vs `float`

```
fragment float4 fshRenderer07(ColorInOut in [[stage_in]],
    float4 back [[ color(0) ]],
    texture2d<float> image [[ texture(0) ]])
{
    constexpr sampler imageSampler(address::clamp_to_zero, filter::linear);

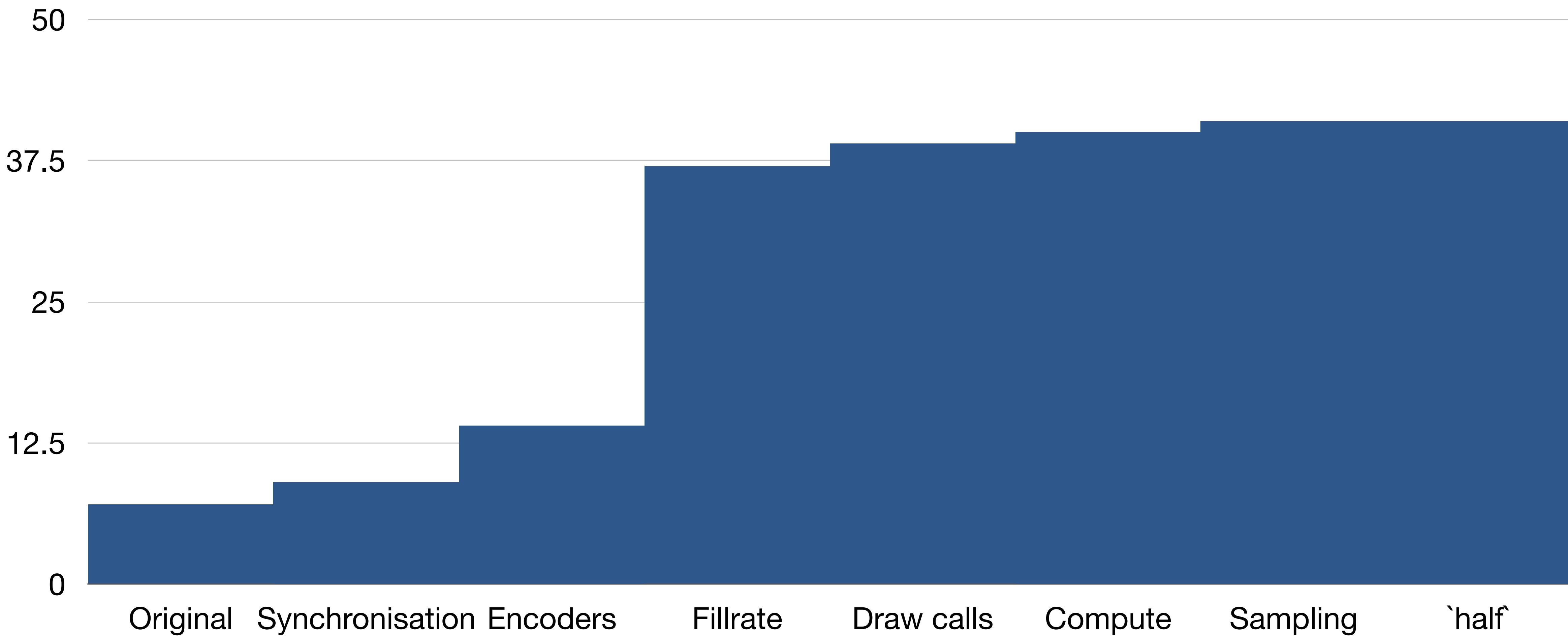
    float4 res = image.sample(imageSampler, in.texCoord);

    res.rgb = mix(back.rgb, res.rgb, res.a);

    return res;
}
```

Summary (FPS)

Task 1: Render + Shader



Summary

Task 1: Render + Shader



Another example

heavy shader

- two input textures
 - with transparent parts
 - without transparency
- sharpen
- color transform (RGB-HSV-RGB)
- displacement transform
- target texture size: 8192 x 8192 px



Another example

heavy shader

10 FPS



`half` vs `float`

```
static float3 rgb2hsv(float3 c) {
    float M = max(c.x, max(c.y, c.z));
    float m = min(c.x, min(c.y, c.z));
    float C = M - m;

    float h = 0;
    if (C > 0) {
        if (c.r == M) {
            h = fmod((c.g - c.b) / C, 6.0);
        } else if (c.g == M) {
            h = (c.b - c.r) / C + 2;
        } else {
            h = (c.r - c.g) / C + 4;
        }
    }
    h /= 6.0;

    float v = M;

    float s = 0;
    if (v > 0) {
        s = C / v;
    }

    return clamp(float3(h, s, v), 0.0, 1.0);
}
```

```
static half3 rgb2hsv(half3 c) {
    half M = max(c.x, max(c.y, c.z));
    half m = min(c.x, min(c.y, c.z));
    half C = M - m;

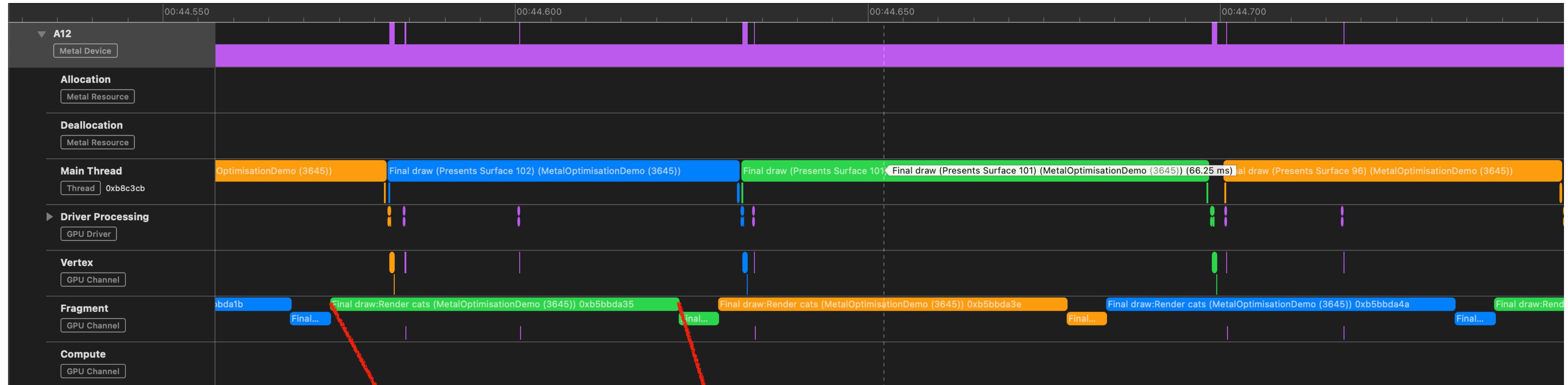
    half h = 0.0h;
    if (C > 0.0h) {
        if (c.r == M) {
            h = fmod((c.g - c.b) / C, 6.0h);
        } else if (c.g == M) {
            h = (c.b - c.r) / C + 2.0h;
        } else {
            h = (c.r - c.g) / C + 4.0h;
        }
    }
    h /= 6.0h;

    half v = M;

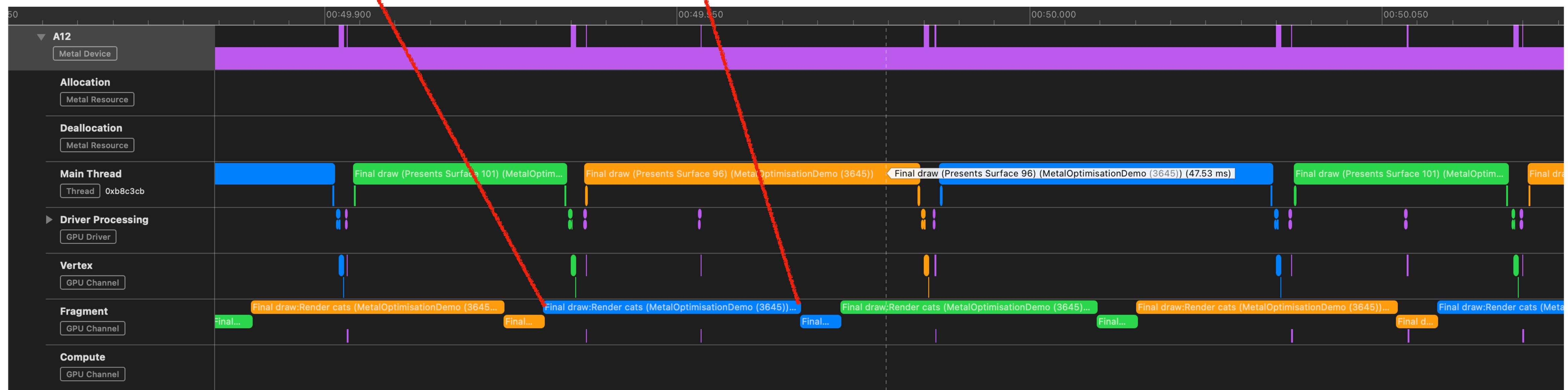
    half s = 0.0h;
    if (v > 0.0h) {
        s = C / v;
    }

    return clamp(half3(h, s, v), 0.0h, 1.0h);
}
```

`half` vs `float`



10 FPS



13 FPS

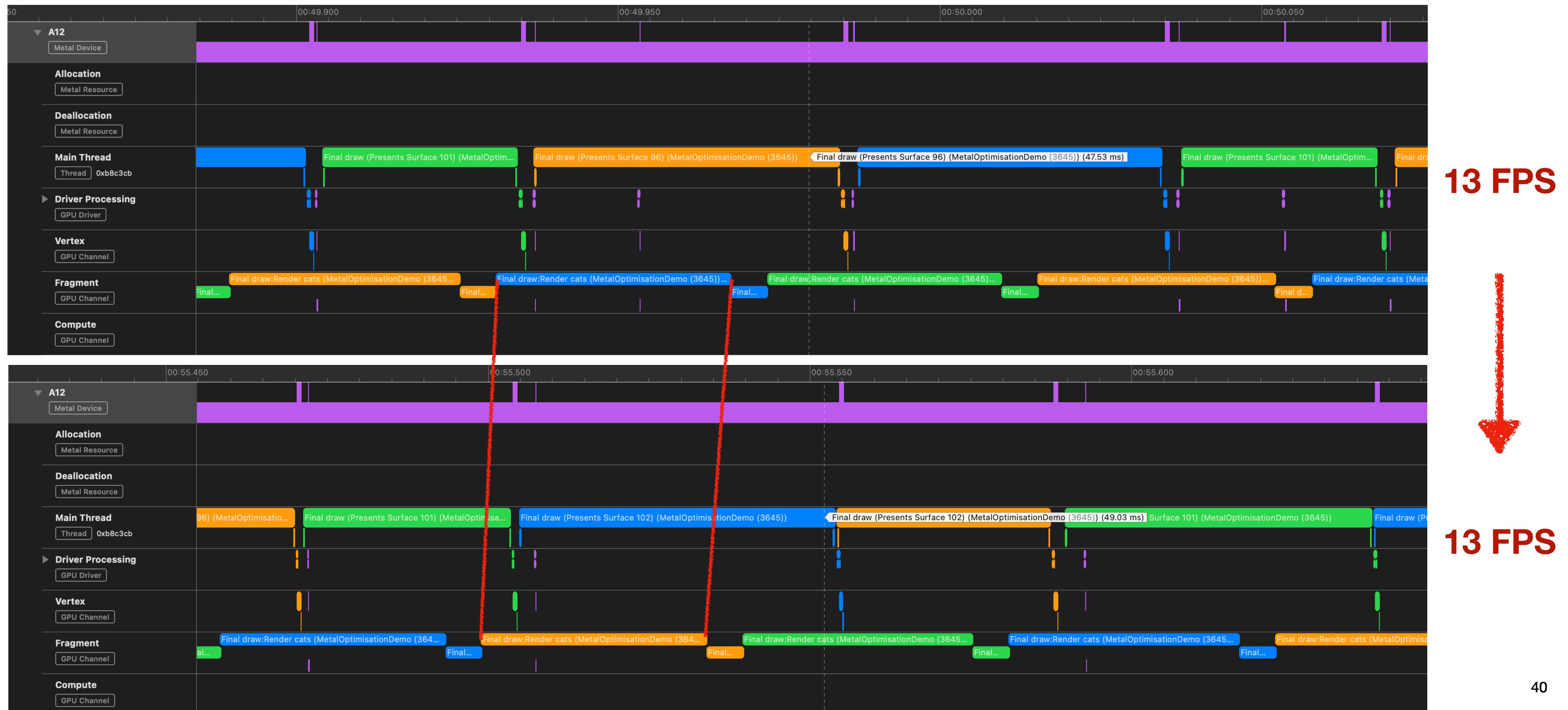
Reduce number of operations

```
half4 dst_pixels[9];
for (char i = 0; i < 9; ++i) {
    float2 dp = pixel * float2(i % 3 - 1, i / 3 - 1);
    dst_pixels[i] = image1.sample(image1Sampler, uv1 + dp);
}
half4 outline = 8 * dst_pixels[4];
for (char i = 0; i < 9; ++i) {
    if (i == 4) {
        continue;
    }
    outline -= dst_pixels[i];
}
```

```
half4 dst_pixels[9] = {
    image1.sample(image1Sampler, uv1 + pixel * float2(-1.0, -1.0)),
    image1.sample(image1Sampler, uv1 + pixel * float2(00.0, -1.0)),
    image1.sample(image1Sampler, uv1 + pixel * float2(+1.0, -1.0)),
    image1.sample(image1Sampler, uv1 + pixel * float2(-1.0, 00.0)),
    image1.sample(image1Sampler, uv1 + pixel * float2(00.0, 00.0)),
    image1.sample(image1Sampler, uv1 + pixel * float2(+1.0, 00.0)),
    image1.sample(image1Sampler, uv1 + pixel * float2(-1.0, +1.0)),
    image1.sample(image1Sampler, uv1 + pixel * float2(00.0, +1.0)),
    image1.sample(image1Sampler, uv1 + pixel * float2(+1.0, +1.0)),
};

half4 outline = 8.0h * dst_pixels[4] - (dst_pixels[0] + dst_pixels[1] +
dst_pixels[2] + dst_pixels[3] + dst_pixels[5] + dst_pixels[6] + dst_pixels[7] +
dst_pixels[8]);
```

Reduce number of operations



Use vectorisation and standard methods

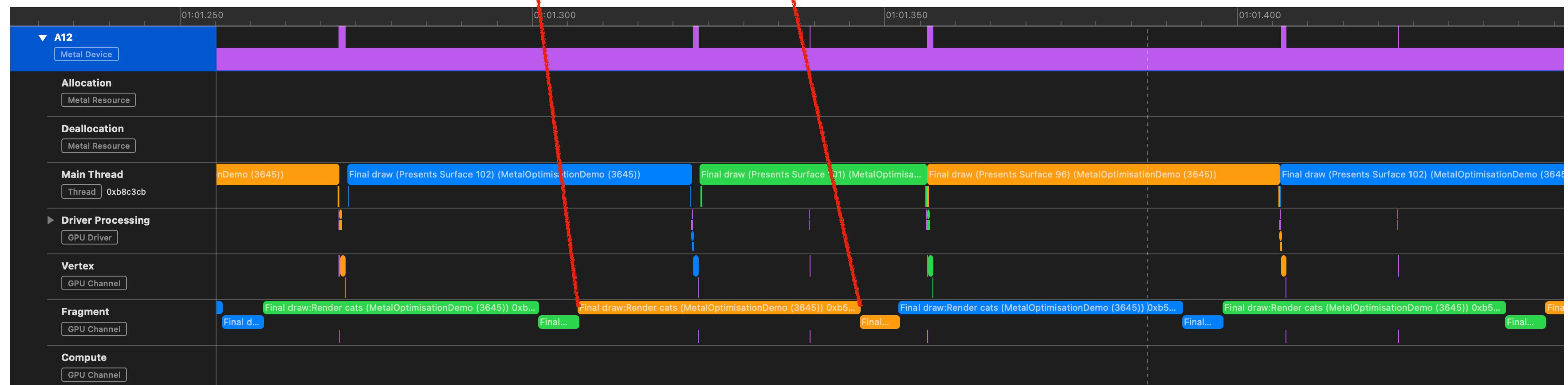
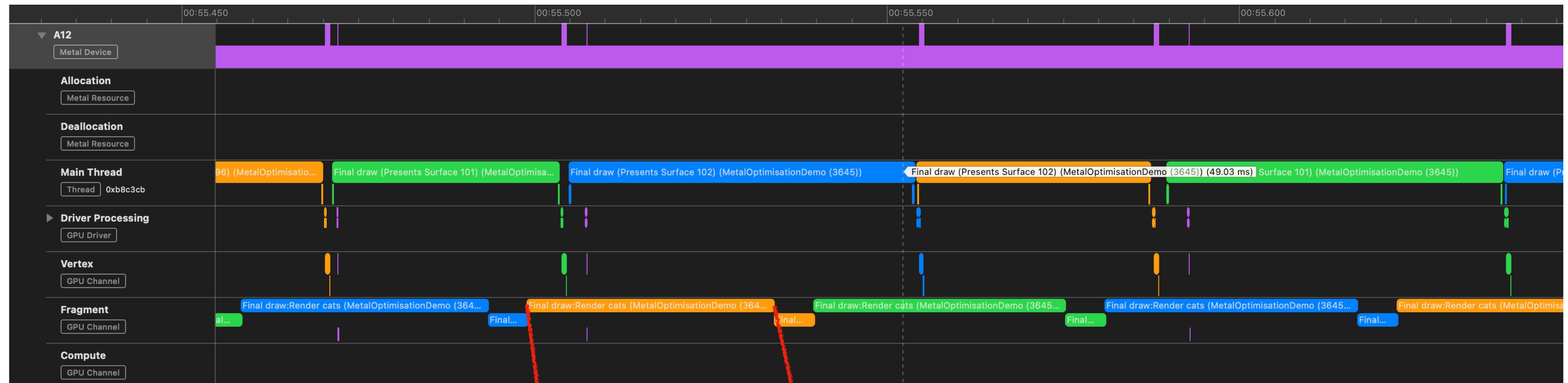
```
static half4 blend(half4 top, half4 btm) {
    half4 res;

    res.rgb = top.rgb + btm.rgb * (1.0h - top.a);
    res.a = top.a + btm.a * (1.0h - top.a);

    return res;
}
```

```
static half4 blend(half4 top, half4 btm) {
    return fma(1.0 - top.a, btm, top);
}
```

Use vectorisation and standard methods



Use alternative algorithms

```
static half3 rgb2hsv(half3 c) {
    half M = max(c.x, max(c.y, c.z));
    half m = min(c.x, min(c.y, c.z));
    half C = M - m;

    half h = 0.0h;
    if (C > 0.0h) {
        if (c.r == M) {
            h = fmod((c.g - c.b) / C, 6.0h);
        } else if (c.g == M) {
            h = (c.b - c.r) / C + 2.0h;
        } else {
            h = (c.r - c.g) / C + 4.0h;
        }
    }
    h /= 6.0h;

    half v = M;

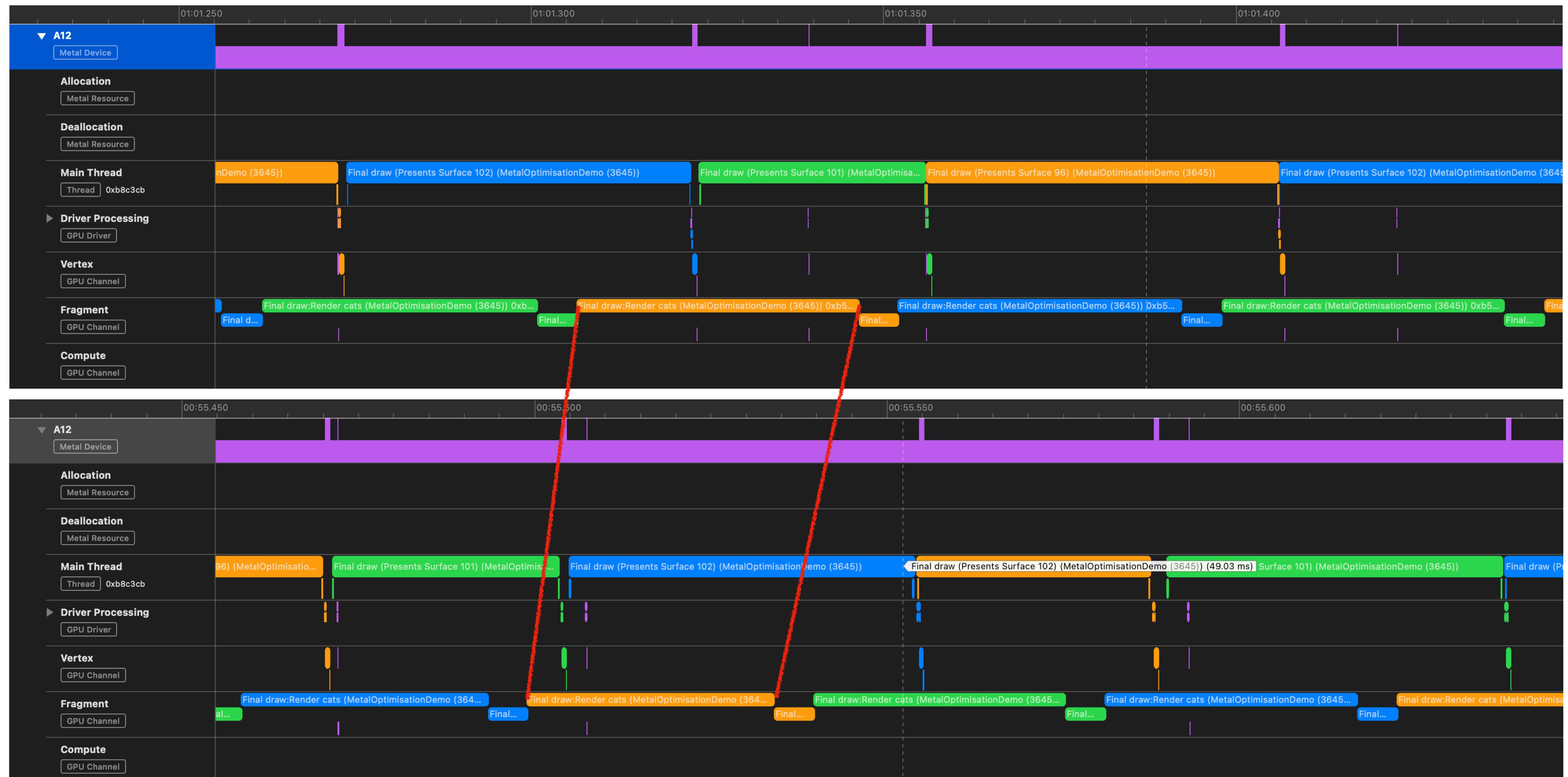
    half s = 0.0h;
    if (v > 0.0h) {
        s = C / v;
    }

    return clamp(half3(h, s, v), 0.0h, 1.0h);
}
```

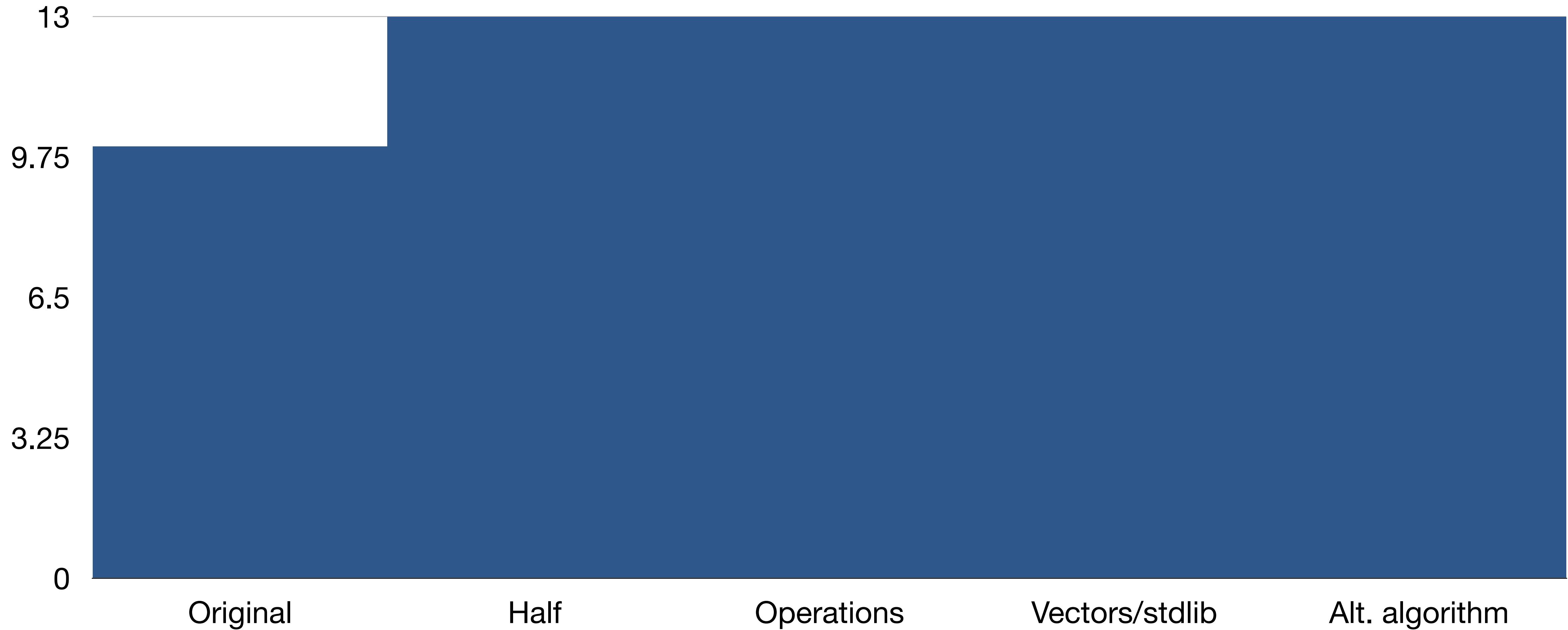
```
static half3 rgb2hsv(half3 c)
{
    constexpr half4 K = half4(0.0h, -1.0h / 3.0h, 2.0h / 3.0h, -1.0h);
    half4 p = select(half4(c.bg, K.wz), half4(c.gb, K.xy), c.b < c.g);
    half4 q = select(half4(p.xyw, c.r), half4(c.r, p.yzx), p.x < c.r);

    half d = q.x - min(q.w, q.y);
    half e = 6.0e-8h;
    return clamp(half3(abs(q.z + (q.w - q.y) / (6.0h * d + e)),
                      d / (q.x + e), q.x),
                0.0h,
                1.0h);
}
```

Use alternative algorithms



Fiasco? (FPS)



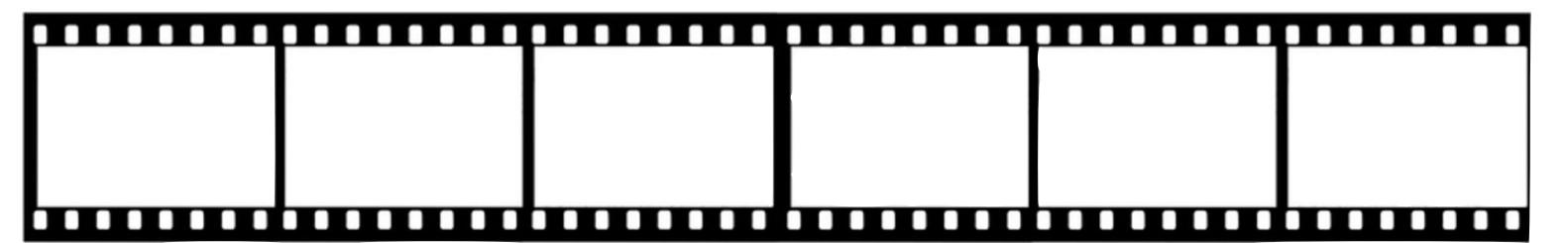
Use right profiling tool



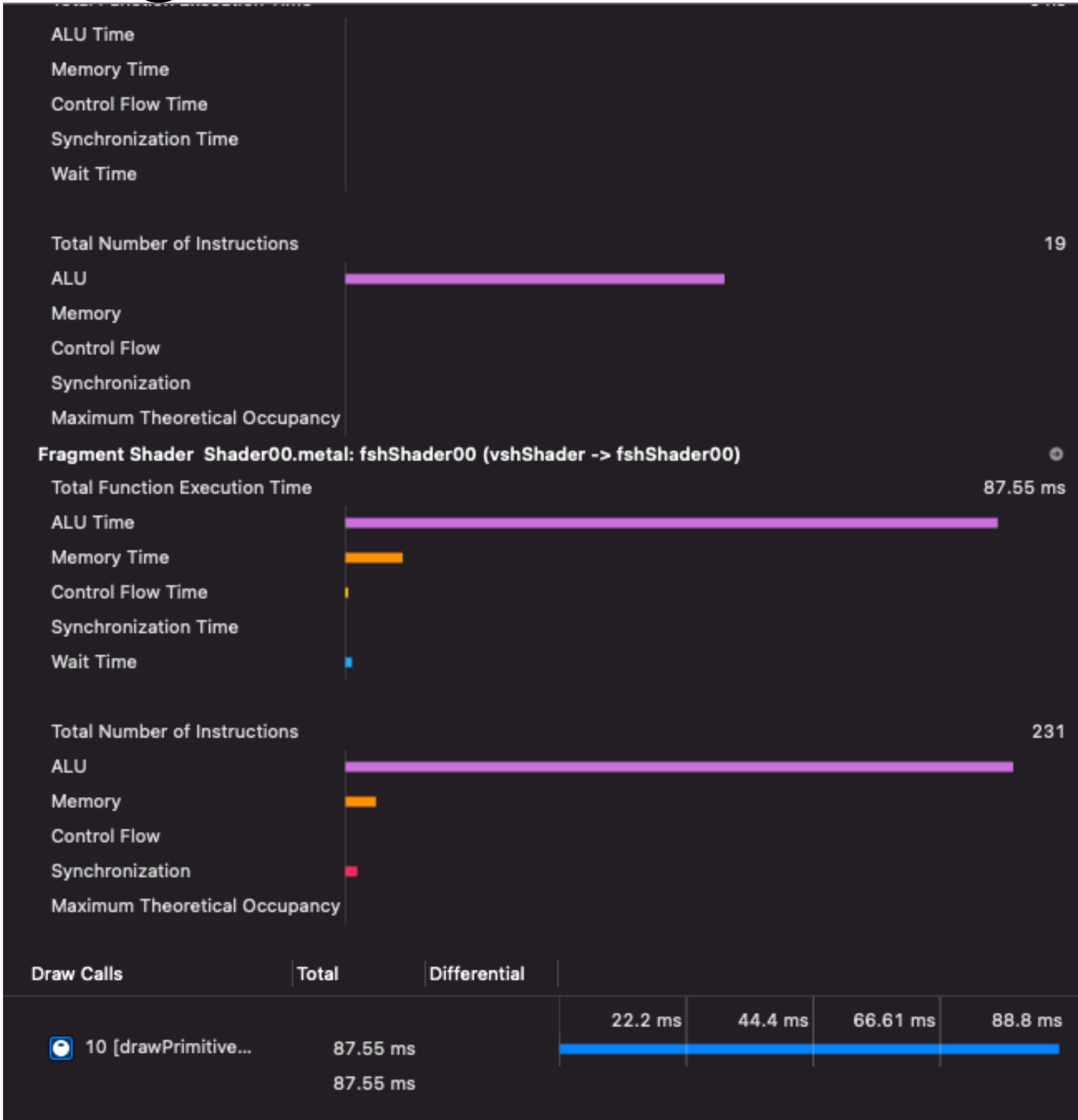
Instruments: Metal System Trace



Xcode: Capture GPU frame

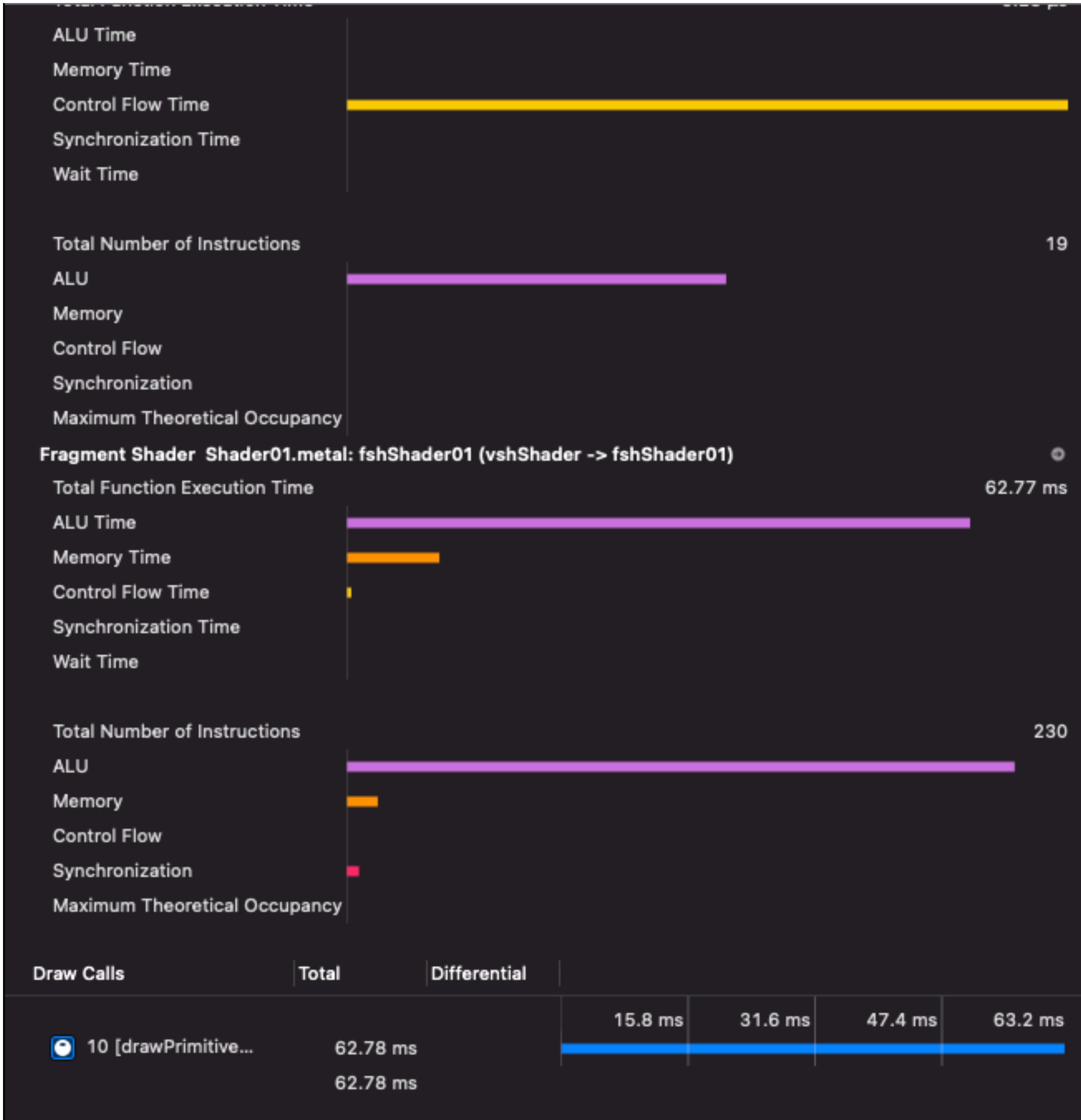


Original shader



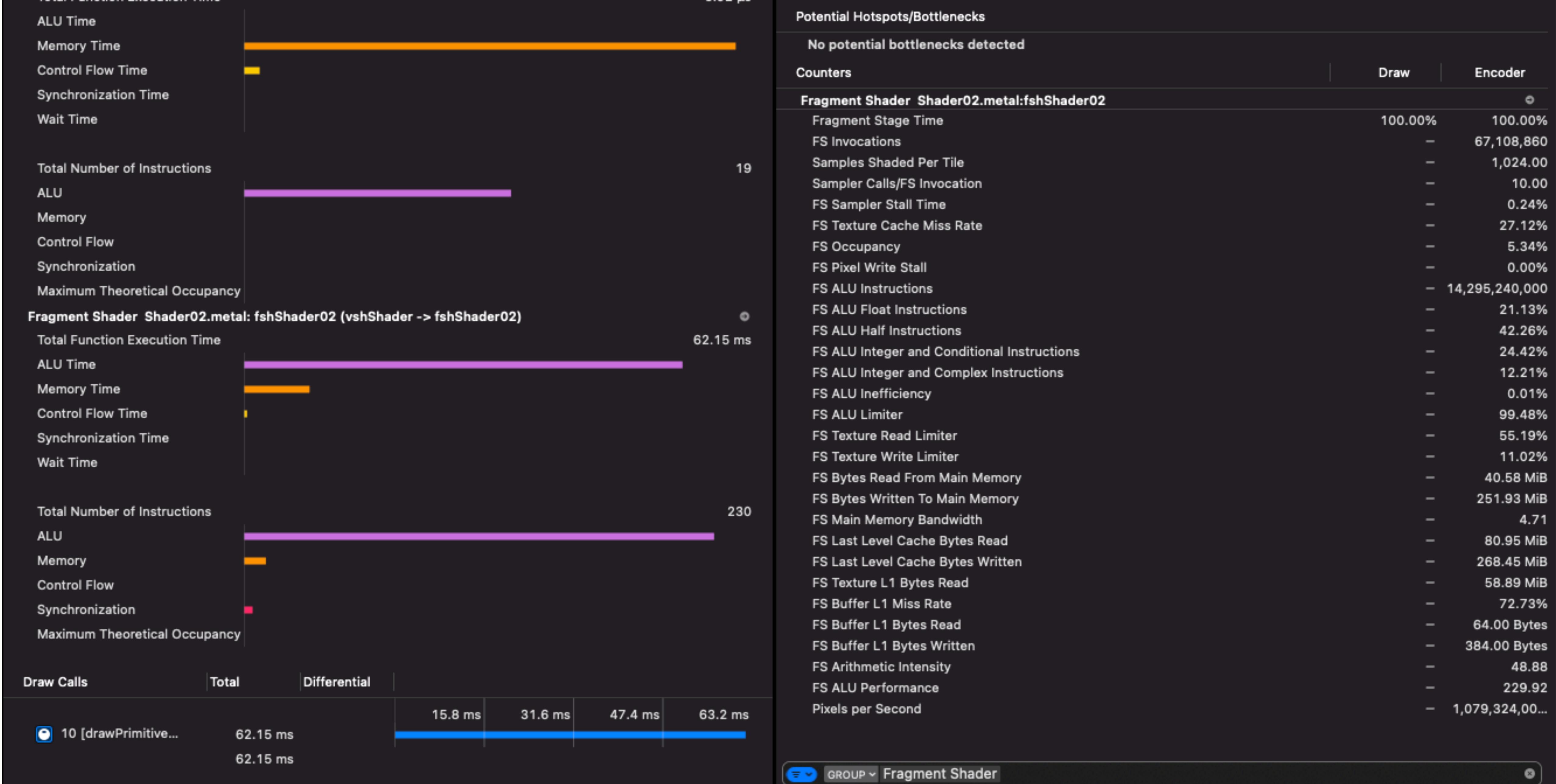
Potential Hotspots/Bottlenecks		
No potential bottlenecks detected		
Counters		
Fragment Shader Shader00.metal:fshShader00		
Fragment Stage Time	100.00%	100.00%
FS Invocations	67,108,860	67,108,860
Samples Shaded Per Tile	N/A	1,024.00
Sampler Calls/FS Invocation	10.00	10.00
FS Sampler Stall Time	N/A	1.63%
FS Texture Cache Miss Rate	N/A	27.93%
FS Occupancy	N/A	1.48%
FS Pixel Write Stall	N/A	0.00%
FS ALU Instructions	14,362,350,000	14,362,350,000
FS ALU Float Instructions	64.02%	64.02%
FS ALU Integer and Conditional Instructions	23.84%	23.84%
FS ALU Integer and Complex Instructions	12.15%	12.15%
FS ALU Inefficiency	0.01%	0.01%
FS ALU Limiter	N/A	91.40%
FS Texture Read Limiter	N/A	40.78%
FS Texture Write Limiter	N/A	7.83%
FS Bytes Read From Main Memory	N/A	40.56 MiB
FS Bytes Written To Main Memory	N/A	251.93 MiB
FS Main Memory Bandwidth	N/A	3.34
FS Last Level Cache Bytes Read	N/A	80.80 MiB
FS Last Level Cache Bytes Written	N/A	268.46 MiB
FS Texture L1 Bytes Read	57.30 MiB	56.86 MiB
FS Buffer L1 Miss Rate	N/A	44.45%
FS Buffer L1 Bytes Read	N/A	64.00 Bytes
FS Buffer L1 Bytes Written	N/A	384.00 Bytes
FS Arithmetic Intensity	N/A	49.11
FS ALU Performance	164.06	163.95
Pixels per Second	766,576,000.00	766,064,800.00

`float` to `half`

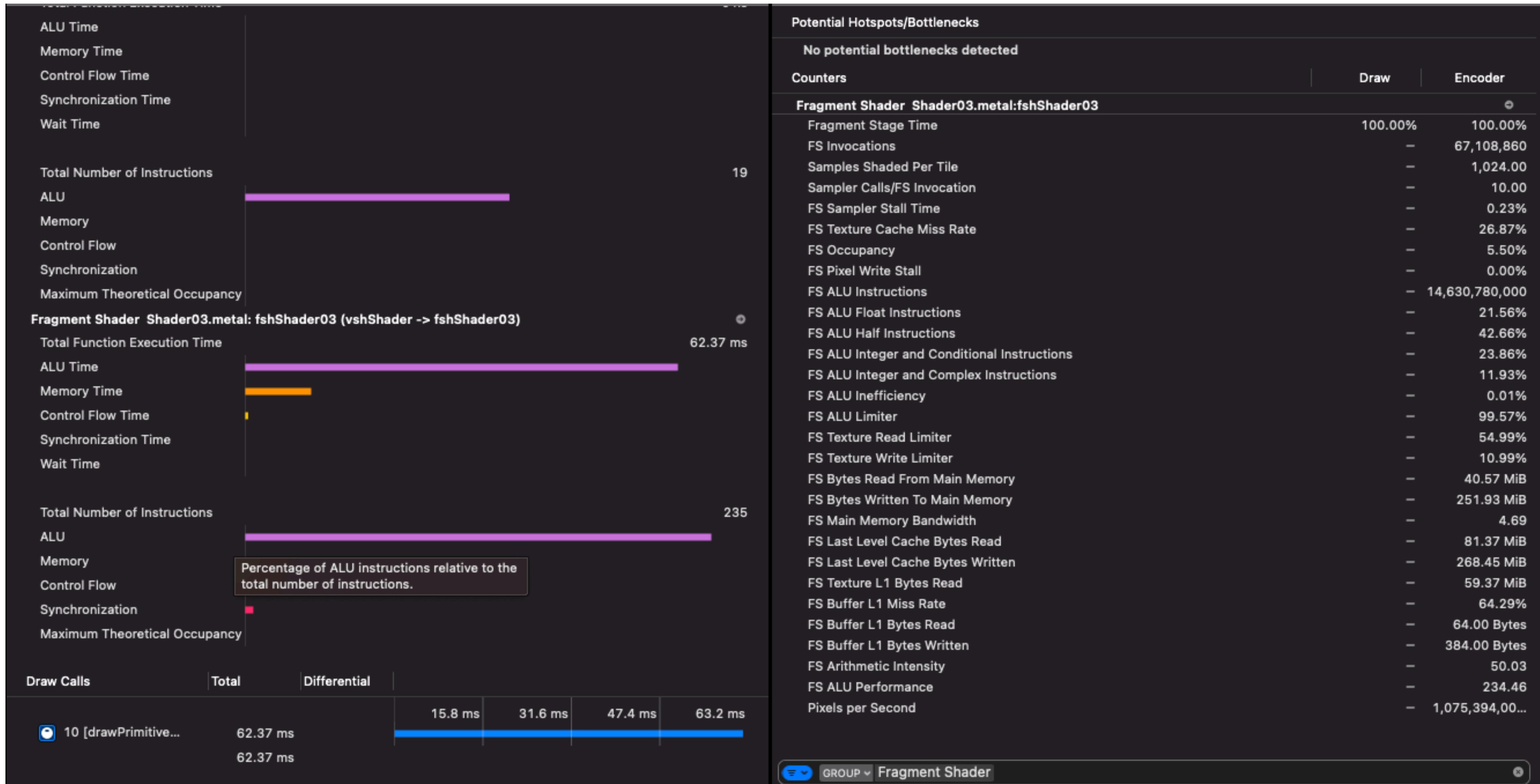


Potential Hotspots/Bottlenecks		
No potential bottlenecks detected		
Counters		
Fragment Shader Shader01.metal:fshShader01	Draw	Encoder
Fragment Stage Time	100.00%	100.00%
FS Invocations	67,108,860	67,108,860
Samples Shaded Per Tile	N/A	1,024.00
Sampler Calls/FS Invocation	10.00	10.00
FS Sampler Stall Time	N/A	0.24%
FS Texture Cache Miss Rate	N/A	27.12%
FS Occupancy	N/A	5.79%
FS Pixel Write Stall	N/A	0.00%
FS ALU Instructions	14,295,240,000	14,295,240,000
FS ALU Float Instructions	23.48%	23.48%
FS ALU Half Instructions	39.91%	39.91%
FS ALU Integer and Conditional Instructions	24.42%	24.42%
FS ALU Integer and Complex Instructions	12.21%	12.21%
FS ALU Inefficiency	0.01%	0.01%
FS ALU Limiter	N/A	99.51%
FS Texture Read Limiter	N/A	54.63%
FS Texture Write Limiter	N/A	10.91%
FS Bytes Read From Main Memory	N/A	40.58 MiB
FS Bytes Written To Main Memory	N/A	251.93 MiB
FS Main Memory Bandwidth	N/A	4.66
FS Last Level Cache Bytes Read	N/A	81.32 MiB
FS Last Level Cache Bytes Written	N/A	268.46 MiB
FS Texture L1 Bytes Read	59.76 MiB	58.89 MiB
FS Buffer L1 Miss Rate	N/A	44.45%
FS Buffer L1 Bytes Read	N/A	64.00 Bytes
FS Buffer L1 Bytes Written	N/A	384.00 Bytes
FS Arithmetic Intensity	N/A	48.88
FS ALU Performance	227.75	227.62
Pixels per Second	1,069,147,00...	1,068,526,00...

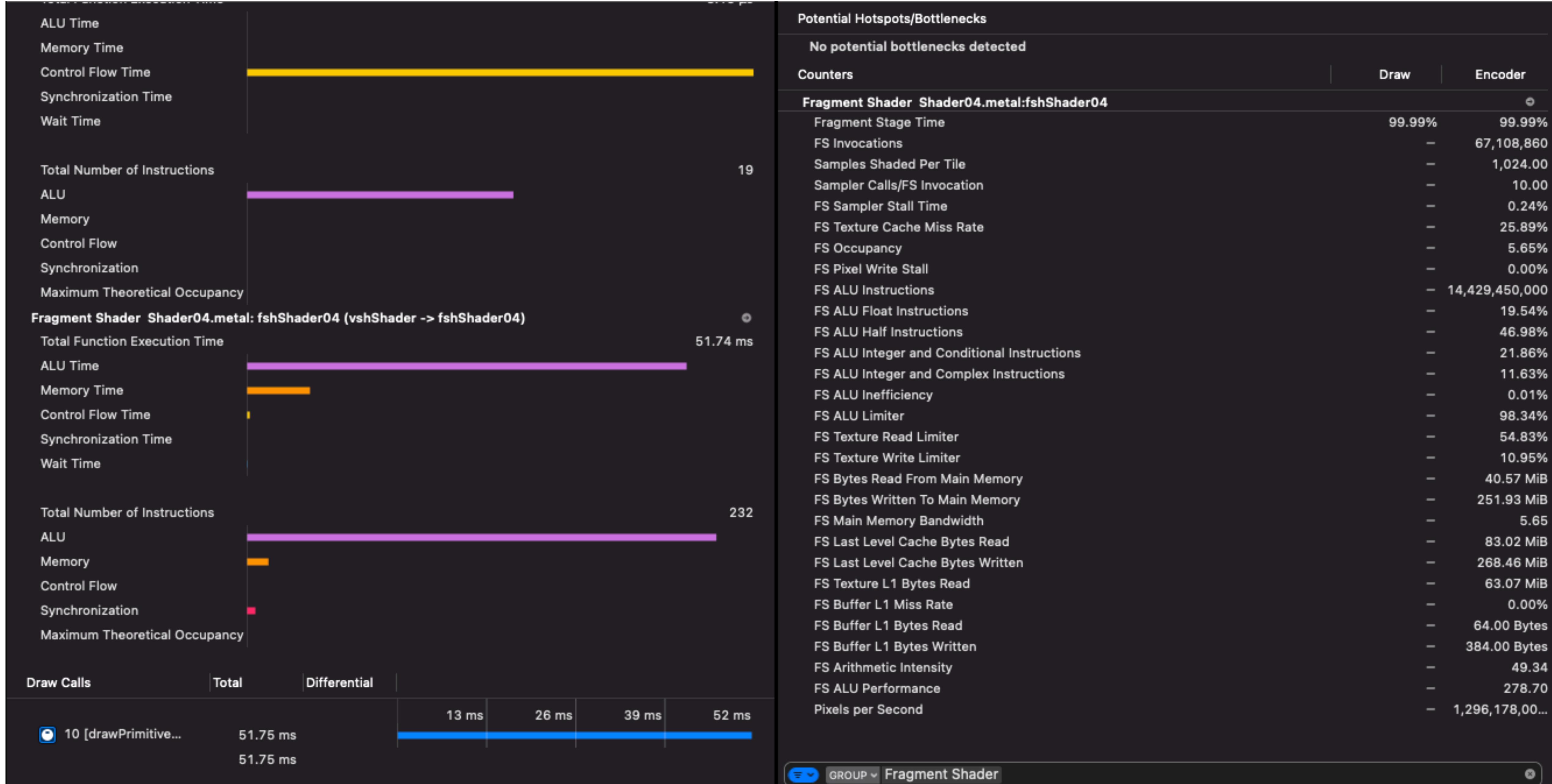
Reduce number of operations



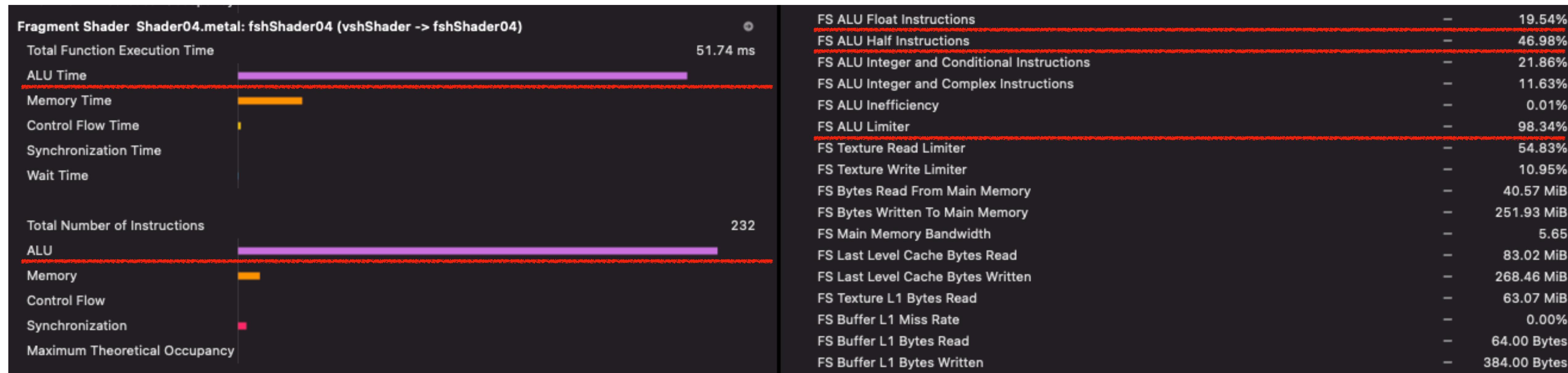
Vectorisation + stdlib



Alternative algorithm



Recognize bottleneck



Recognize bottleneck

half4 dst = dst_pixels[4] * smoothstep(0.8h, 1.0h, 1.0h - length(outline));	9.91%	●
// - HSV	3.86%	●
half3 hsv = rgb2hsv(dst.rgb);	16.95%	●
half3 col = time_color(half2(uv1), time);	10.77%	●
hsv.x = fract(fma(half(angle), 0.5h / M_PI_H, 0.5h) - fma(time, 0.25h, -rgb2hsv(col).x));	2.89%	●
half k = smoothstep(0.0h, 1.0h, length(half2(shift)));	1.21%	●
hsv.y = fma(hsv.y, 1.0h - k, k);	18.04%	●
dst.rgb = hsv2rgb(hsv);		

Reducing weights

Use your own functions

```
template <typename T>
inline static T pos_fract(T x) {
    return x - floor(x);
}
```

Reducing weights

Reduce access to vector items

```
static half3 hsv2rgb(half3 c) {
    constexpr half4 K = half4(1.0, 2.0 / 3.0, 1.0 / 3.0, 3.0);
    half3 p = abs(fract(c.xxx + K.xyz) * 6.0h - K.www);
    return clamp(c.z * mix(K.www, clamp(p - K.www, 0.0h, 1.0h), c.y), 0.0h, 1.0h);
}
```

```
static half3 hsv2rgb(half3 c) {
    constexpr half3 K = half3(1.0h, 2.0h / 3.0h, 1.0h / 3.0h);
    half3 p = abs(pos_fract(c.x + K) * 6.0h - 3.0h);
    return c.z * fma(clamp(p - 1.0h, 0.0h, 1.0h), c.y, 1.0h - c.y);
}
```

Reducing weights

Use school math :)

```
static half3 time_color(half2 uv, half time) {
    return fma(abs(half3(
        sin(fma(cos(fma(3.0h, uv.y, time)), 2.0h * uv.x, time)),
        cos(fma(sin(fma(2.0h, uv.x, time)), 3.0h * uv.y, time)),
        0.4h / 0.9h)), 0.9h, 0.1h);
}
```

$$\sin\left(\frac{\pi}{2} + \theta\right) = \cos \theta$$

Reducing weights

Use school math :)

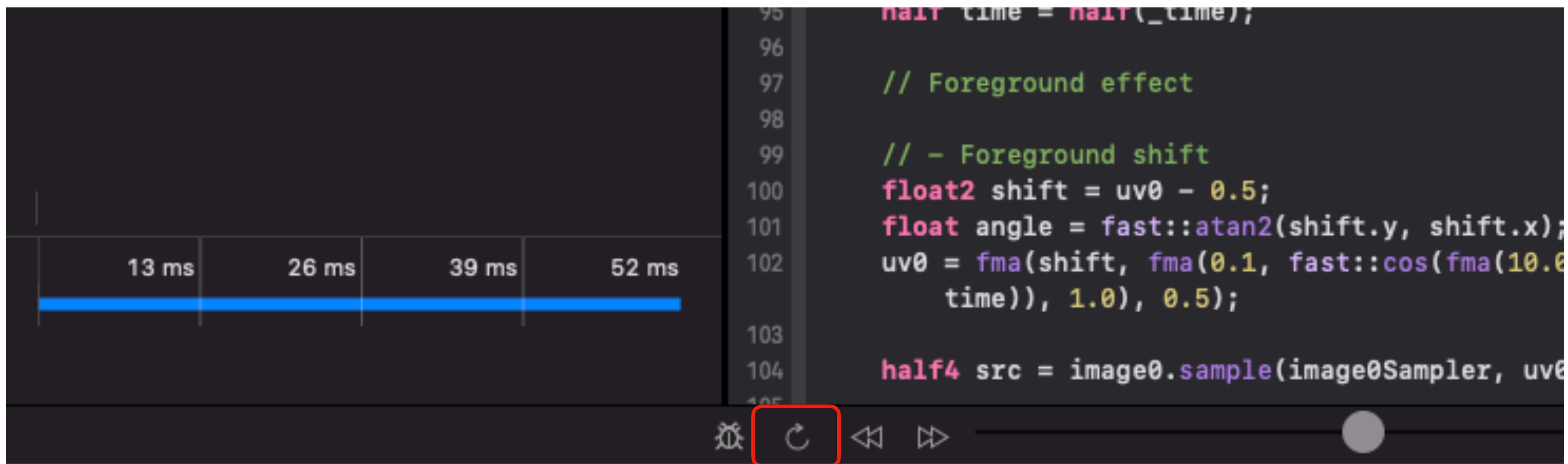
```
static half3 time_color(half2 uv, half time) {
    return fma(abs(half3(
        sin(fma(cos(fma(3.0h, uv.y, time)), 2.0h * uv.x, time)),
        cos(fma(sin(fma(2.0h, uv.x, time)), 3.0h * uv.y, time)),
        0.4h / 0.9h)), 0.9h, 0.1h);
}
```

```
static half3 time_color(half2 uv, half time) {
    half2 uv_l = uv * half2(3.0h, 2.0h);
    half2 tcs = time + half2(M_PI_2_H, 0.0h);
    half2 tsc = time + half2(0.0h, M_PI_2_H);

    half2 a = uv_l + tcs;
    half2 a_cs = sin(a);
    half2 b = fma(uv_l, a_cs, tsc);
    half2 b_sc = sin(b);

    return fma(abs(half3(b_sc, 0.4h / 0.9h)), 0.9h, 0.1h);
}
```

Update shader on-fly



The screenshot shows a code editor with a dark theme. At the top, there is a horizontal timeline bar divided into four segments with labels: "13 ms", "26 ms", "39 ms", and "52 ms". A blue horizontal bar spans across the first three segments. Below the timeline, there is a line of code with line numbers on the left:

```
95     half time = half(_time);  
96  
97     // Foreground effect  
98  
99     // - Foreground shift  
100    float2 shift = uv0 - 0.5;  
101    float angle = fast::atan2(shift.y, shift.x);  
102    uv0 = fma(shift, fma(0.1, fast::cos(fma(10.6  
           time)), 1.0), 0.5);  
103  
104    half4 src = image0.sample(image0Sampler, uv0);  
105
```

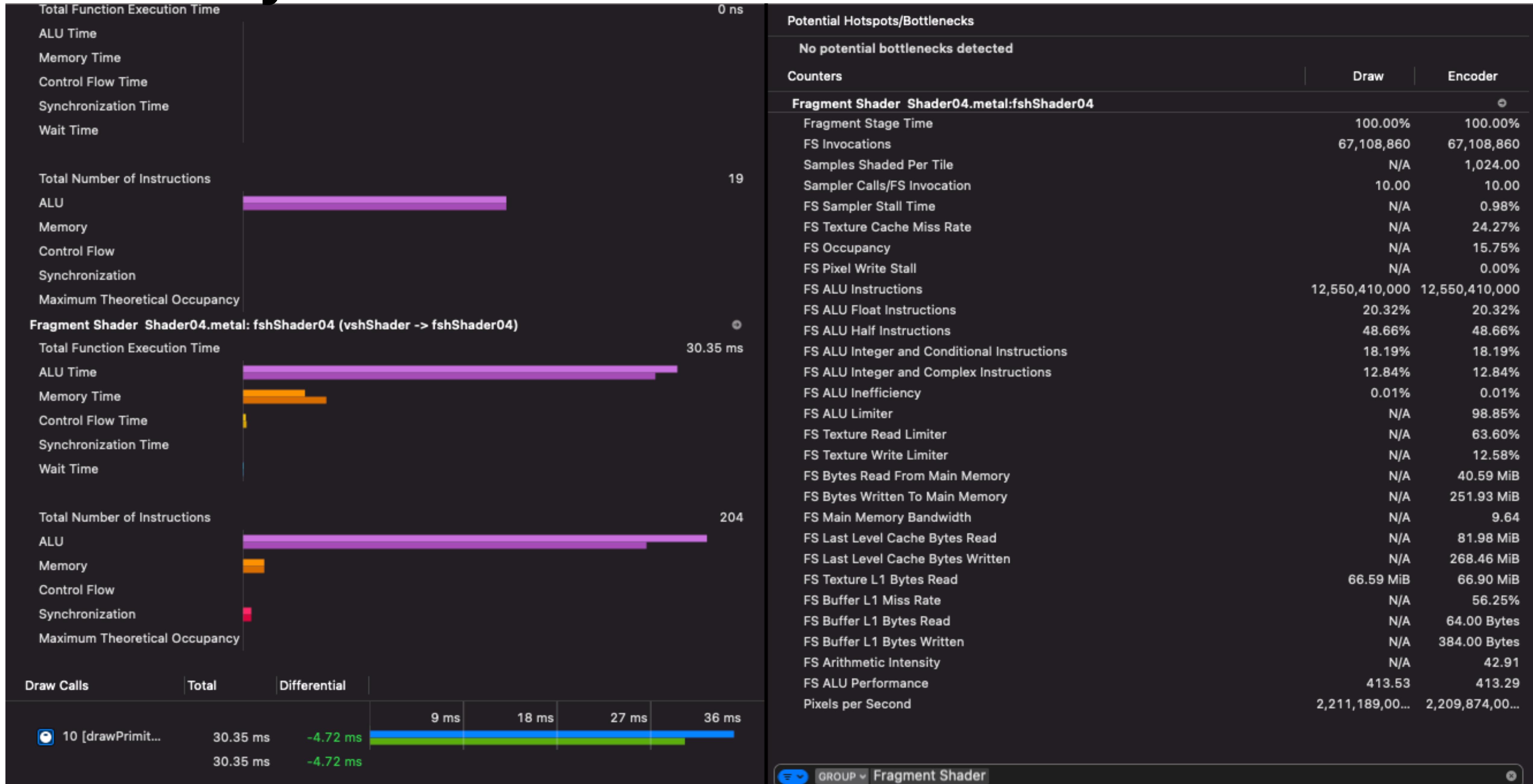
At the bottom of the editor window, there is a toolbar with several icons. One icon, which is a circular arrow symbol, is highlighted with a red rectangular box.

Reducing weights

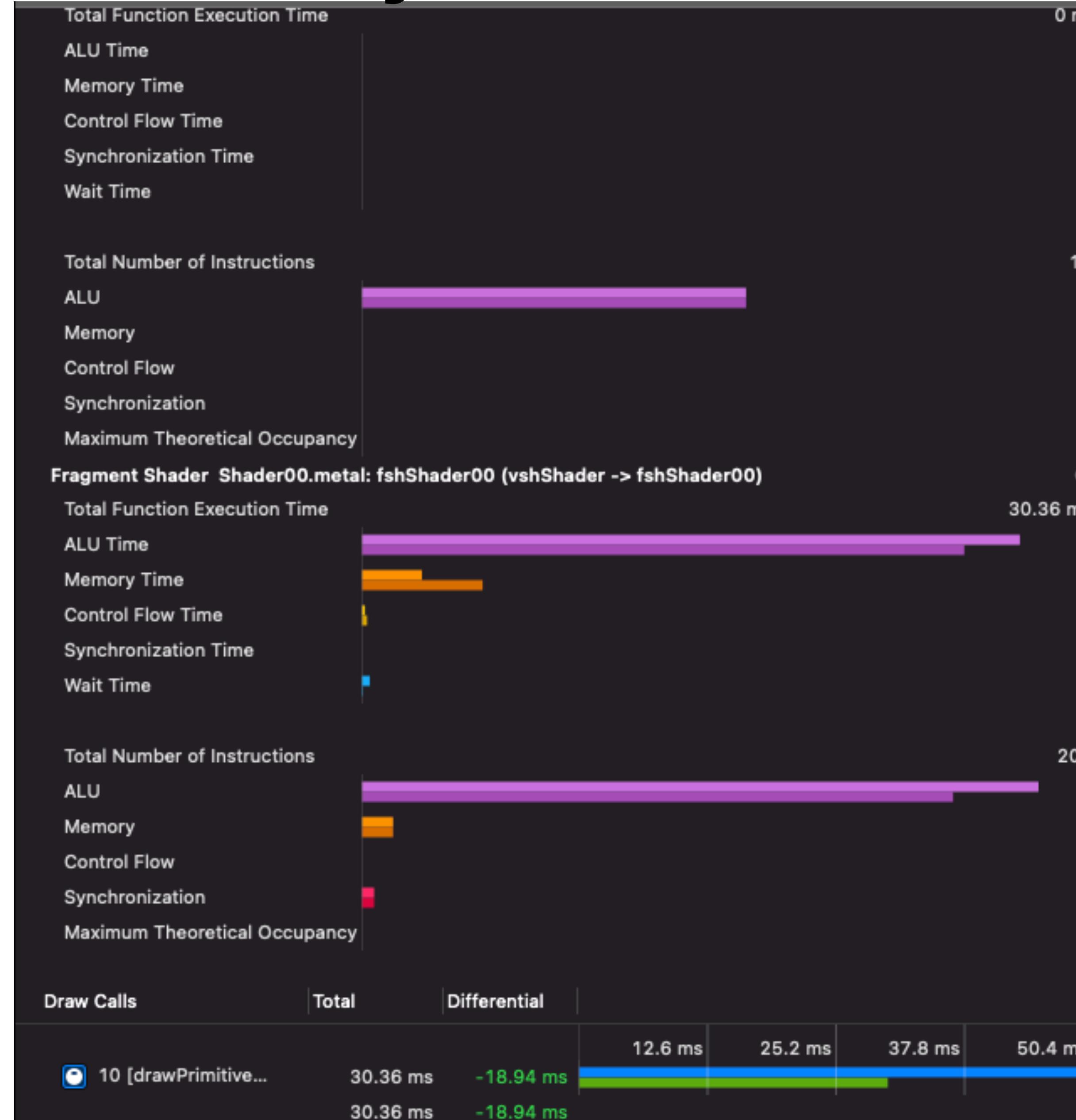
Summary

<pre>half4 dst = dst_pixels[4] * smoothstep(0.8h, 1.0h, 1.0h - length(outline));</pre>	9.91%	
<pre>// - HSV half3 hsv = rgb2hsv(dst.rgb);</pre>	3.86%	
<pre>half3 col = time_color(half2(uv1), time);</pre>	16.95%	
<pre>hsv.x = fract(fma(half(angle), 0.5h / M_PI_H, 0.5h) - fma(time, 0.25h, -rgb2hsv(col).x));</pre>	10.77%	
<pre>half k = smoothstep(0.0h, 1.0h, length(half2(shift)));</pre>	2.89%	
<pre>hsv.y = fma(hsv.y, 1.0h - k, k);</pre>	1.21%	
<pre>dst.rgb = hsv2rgb(hsv);</pre>	18.04%	
<pre>half4 dst = dst_pixels[4] * smoothstep(0.8h, 1.0h, 1.0h - border);</pre>	4.28%	
<pre>// - HSV half3 hsv = rgb2hsv(dst.rgb);</pre>	4.90%	
<pre>half3 col = time_color(half2(uv1), time);</pre>	13.96%	
<pre>hsv.x = pos_fract(fma(half(angle), 0.5h / M_PI_H, 0.5h) - fma(time, 0.25h, -rgb2hsv(col).x));</pre>	9.31%	
<pre>half k = smoothstep(0.0h, 1.0h, length(half2(shift)));</pre>	4.08%	
<pre>hsv.y = fma(hsv.y, 1.0h - k, k);</pre>	1.21%	
<pre>dst.rgb = hsv2rgb(max(hsv, 0.0h));</pre>	16.10%	

Summary



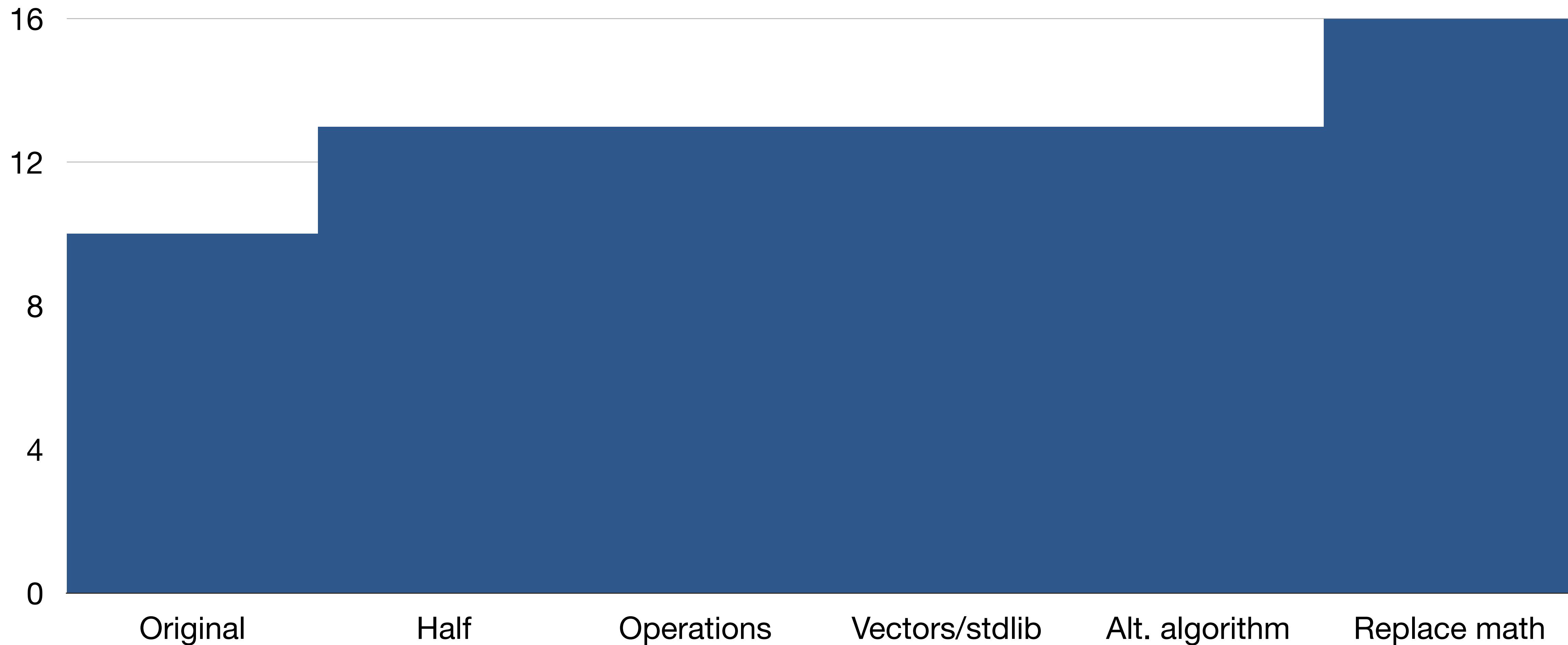
Summary



Potential Hotspots/Bottlenecks		Draw	Encoder
No potential bottlenecks detected			
Counters			
Fragment Shader Shader00.metal:fshShader00			
Fragment Stage Time	100.00%	100.00%	
FS Invocations	-	67,108,860	
Samples Shaded Per Tile	-	1,024.00	
Sampler Calls/FS Invocation	-	10.00	
FS Sampler Stall Time	-	0.89%	
FS Texture Cache Miss Rate	-	24.38%	
FS Occupancy	-	15.78%	
FS Pixel Write Stall	-	0.00%	
FS ALU Instructions	-	12,550,410,000	
FS ALU Float Instructions	-	20.32%	
FS ALU Half Instructions	-	48.66%	
FS ALU Integer and Conditional Instructions	-	18.19%	
FS ALU Integer and Complex Instructions	-	12.84%	
FS ALU Inefficiency	-	0.01%	
FS ALU Limiter	-	98.86%	
FS Texture Read Limiter	-	63.60%	
FS Texture Write Limiter	-	12.58%	
FS Bytes Read From Main Memory	-	40.59 MiB	
FS Bytes Written To Main Memory	-	251.93 MiB	
FS Main Memory Bandwidth	-	9.64	
FS Last Level Cache Bytes Read	-	81.90 MiB	
FS Last Level Cache Bytes Written	-	268.46 MiB	
FS Texture L1 Bytes Read	-	66.59 MiB	
FS Buffer L1 Miss Rate	-	44.45%	
FS Buffer L1 Bytes Read	-	64.00 Bytes	
FS Buffer L1 Bytes Written	-	384.00 Bytes	
FS Arithmetic Intensity	-	42.91	
FS ALU Performance	-	413.28	
Pixels per Second	-	2,209,862,00...	

GROUP ▾ Fragment Shader

Summary (FPS)



Classes in shaders

- Encapsulate context
- Cache computation results
- Reduce passing objects
- More understandable code

Textures

Another example

displacing

- one input textures
- displacement transforms
 - random tile
 - **sin** offset
- target texture size: 8192 x 8192 px



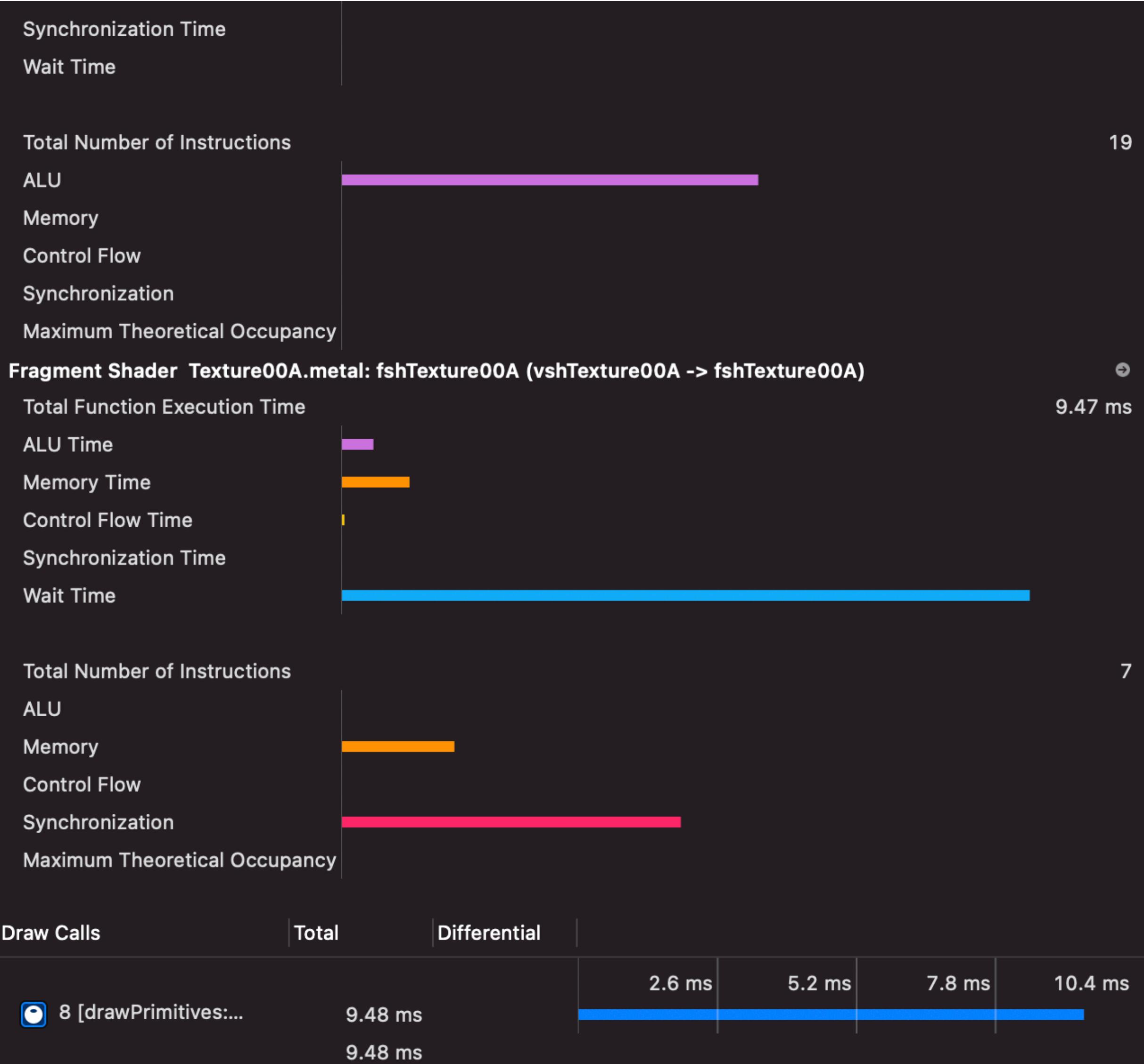
Localisation (0)



~ 10.46 ms



Localisation (0)



Potential Hotspots/Bottlenecks

No potential bottlenecks detected

Counters

Fragment Shader Texture00A.metal:fshTexture00A

	Draw	Encoder
FS ALU Limiter	N/A	0.02%
FS Texture Read Limiter	N/A	18.94%
FS Texture Write Limiter	N/A	98.86%
FS Bytes Read From Main Memory	N/A	25.90 MiB
FS Bytes Written To Main Memory	N/A	163.75 MiB
FS Main Memory Bandwidth	N/A	19.29
FS Last Level Cache Bytes Read	N/A	47.29 MiB
FS Last Level Cache Bytes Written	N/A	268.45 MiB
FS Texture L1 Bytes Read	26.61 MiB	26.68 MiB
FS Buffer L1 Miss Rate	N/A	80.00%
FS Buffer L1 Bytes Written	N/A	384.00 Bytes
FS Arithmetic Intensity	N/A	0.01
FS ALU Performance	0.12	0.11
Pixels per Second	7,093,166,00...	6,823,602,00...

Post-Fragment Stage 1 attachment, 8192 x 8192

Pixels Stored	67,108,860	67,108,860
Pixel Write Stall	N/A	61.09%
Predicated Texture Thread Writes	0.00%	0.00%
Compression Ratio of Texture Memory Written	N/A	0.00
Texture Cache Write Miss Rate	N/A	59.70%
Texture		
Texture Cache Miss Rate	N/A	28.02%
Texture L1 Bytes Read	26.61 MiB	26.68 MiB
Texture Sampler Calls	67,108,860	67,108,860
Lossless Compressed Texture Samples	N/A	0.00%
Uncompressed Texture Samples	N/A	100.00%
Compression Ratio of Texture Memory Read	N/A	0.00

Filter

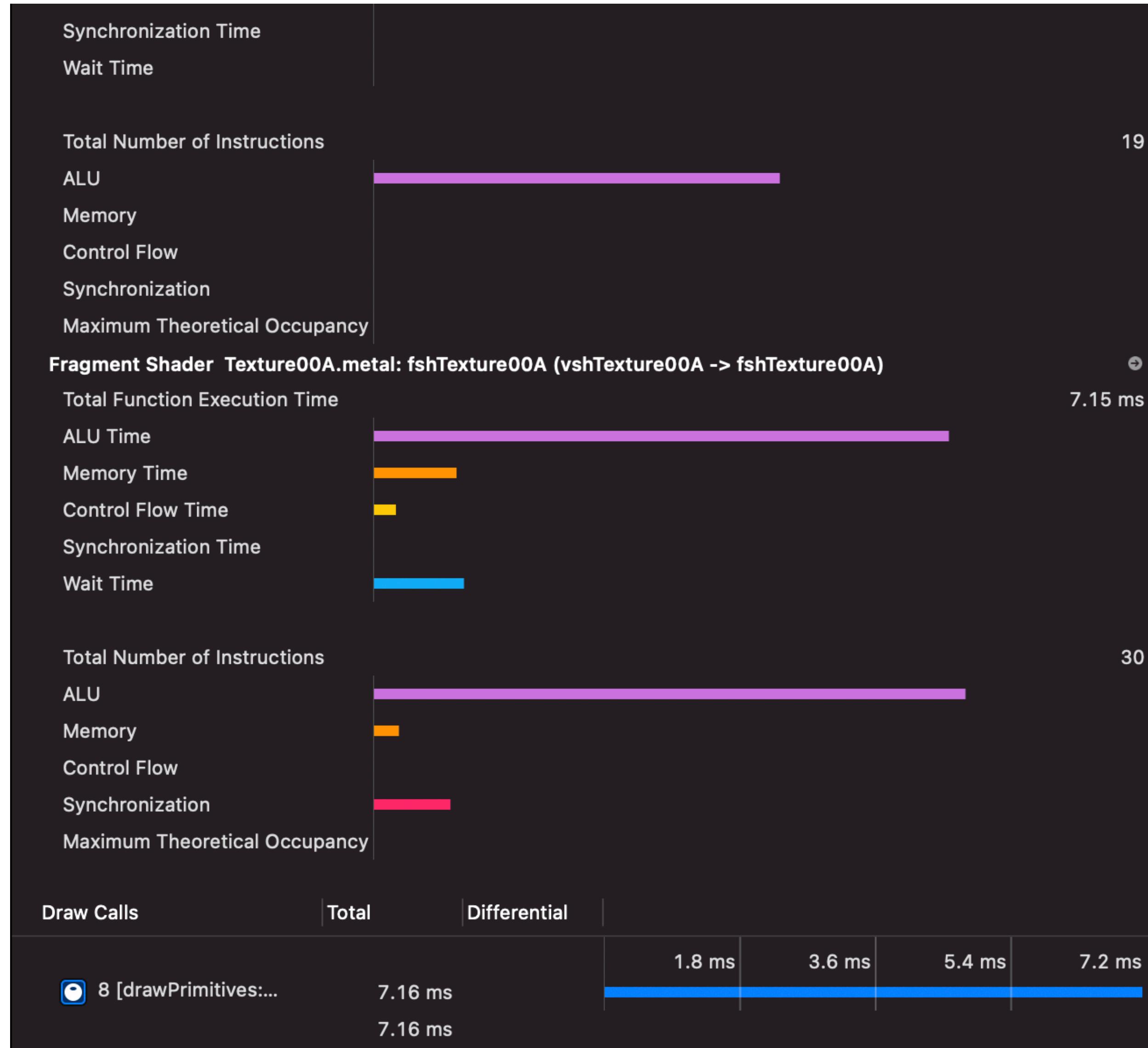
Localisation (1)



~ 10.24 ms



Localisation (1)



Potential Hotspots/Bottlenecks			
No potential bottlenecks detected			
Counters			
Fragment Shader Texture00A.metal:fshTexture00A			
FS ALU Limiter	N/A	61.32%	
FS Texture Read Limiter	N/A	18.87%	
FS Texture Write Limiter	N/A	98.90%	
FS Bytes Read From Main Memory	N/A	26.00 MiB	
FS Bytes Written To Main Memory	N/A	162.57 MiB	
FS Main Memory Bandwidth	N/A	19.33	
FS Last Level Cache Bytes Read	N/A	47.40 MiB	
FS Last Level Cache Bytes Written	N/A	268.46 MiB	
FS Texture L1 Bytes Read	24.71 MiB	24.60 MiB	
FS Buffer L1 Miss Rate	N/A	87.50%	
FS Buffer L1 Bytes Written	N/A	384.00 Bytes	
FS Arithmetic Intensity	N/A	8.20	
FS ALU Performance	216.15	158.34	
Pixels per Second	9,391,196,00...	6,879,651,00...	
Post-Fragment Stage 1 attachment, 8192 x 8192			
Pixels Stored	67,108,860	67,108,860	
Pixel Write Stall	N/A	64.10%	
Predicated Texture Thread Writes	0.00%	0.00%	
Compression Ratio of Texture Memory Written	N/A	0.00	
Texture Cache Write Miss Rate	N/A	62.82%	
Texture			
Texture Cache Miss Rate	N/A	30.47%	
Texture L1 Bytes Read	24.71 MiB	24.60 MiB	
Texture Sampler Calls	67,108,860	67,108,860	
Lossless Compressed Texture Samples	N/A	0.00%	
Uncompressed Texture Samples	N/A	100.00%	
Compression Ratio of Texture Memory Read	N/A	0.00	

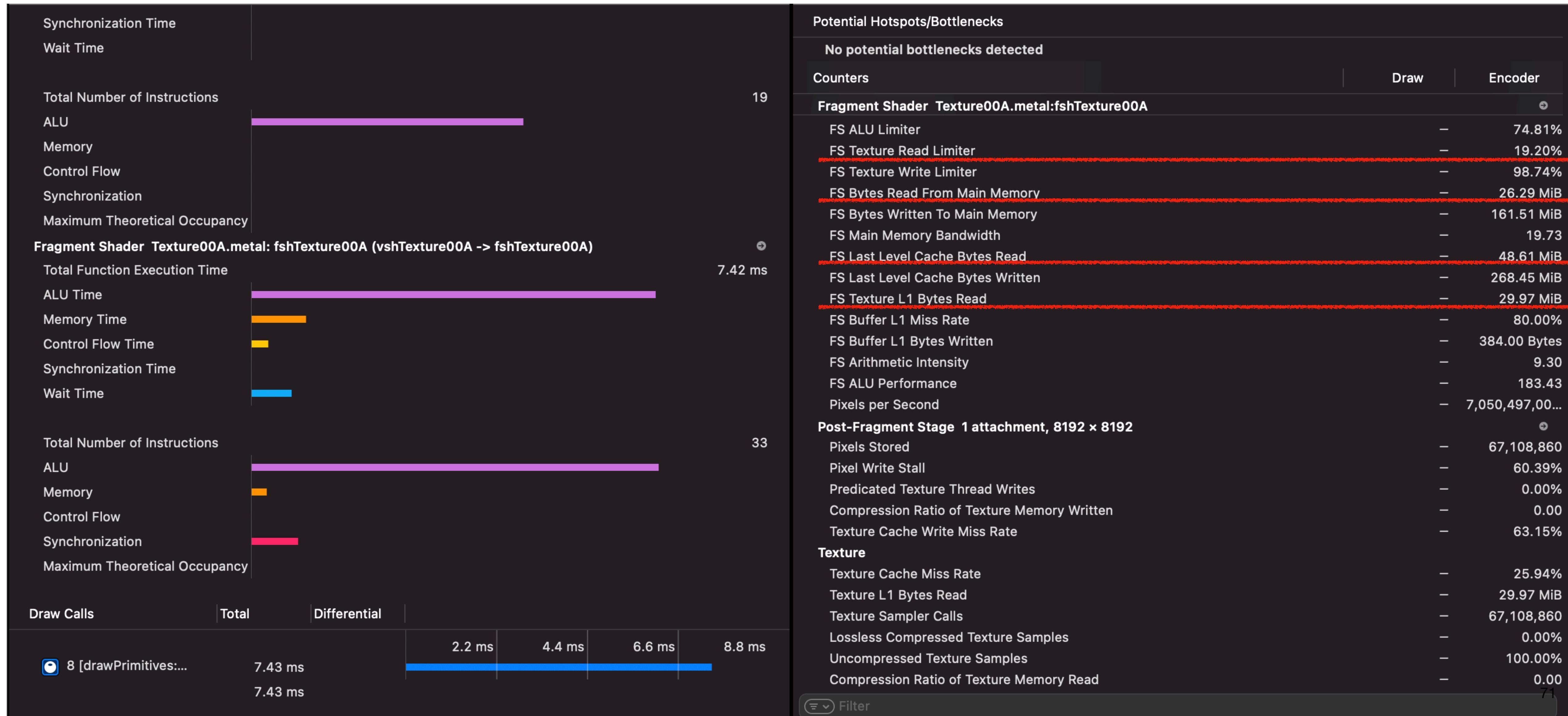
Localisation (2)



~ 10.01 ms



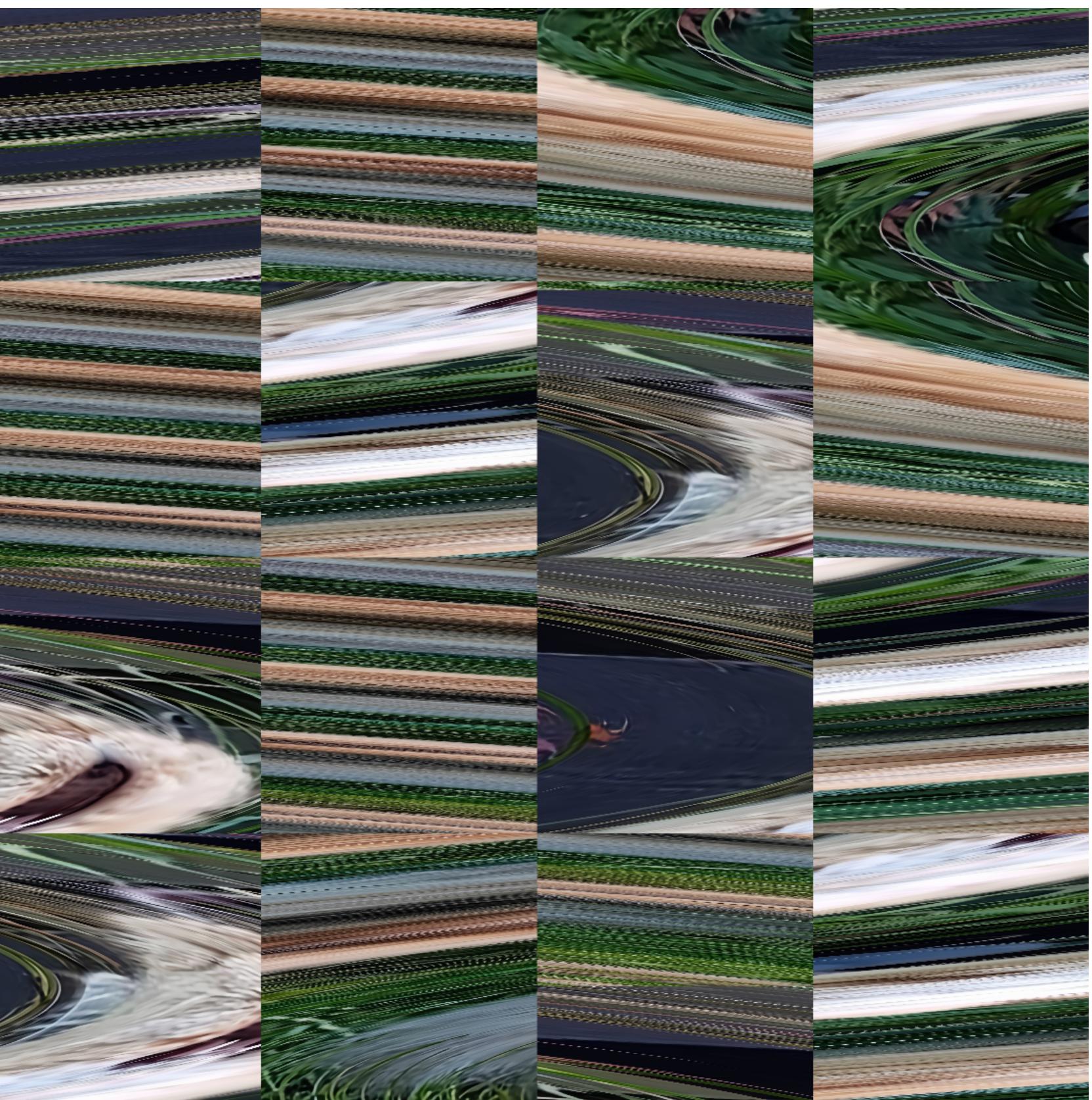
Localisation (2)



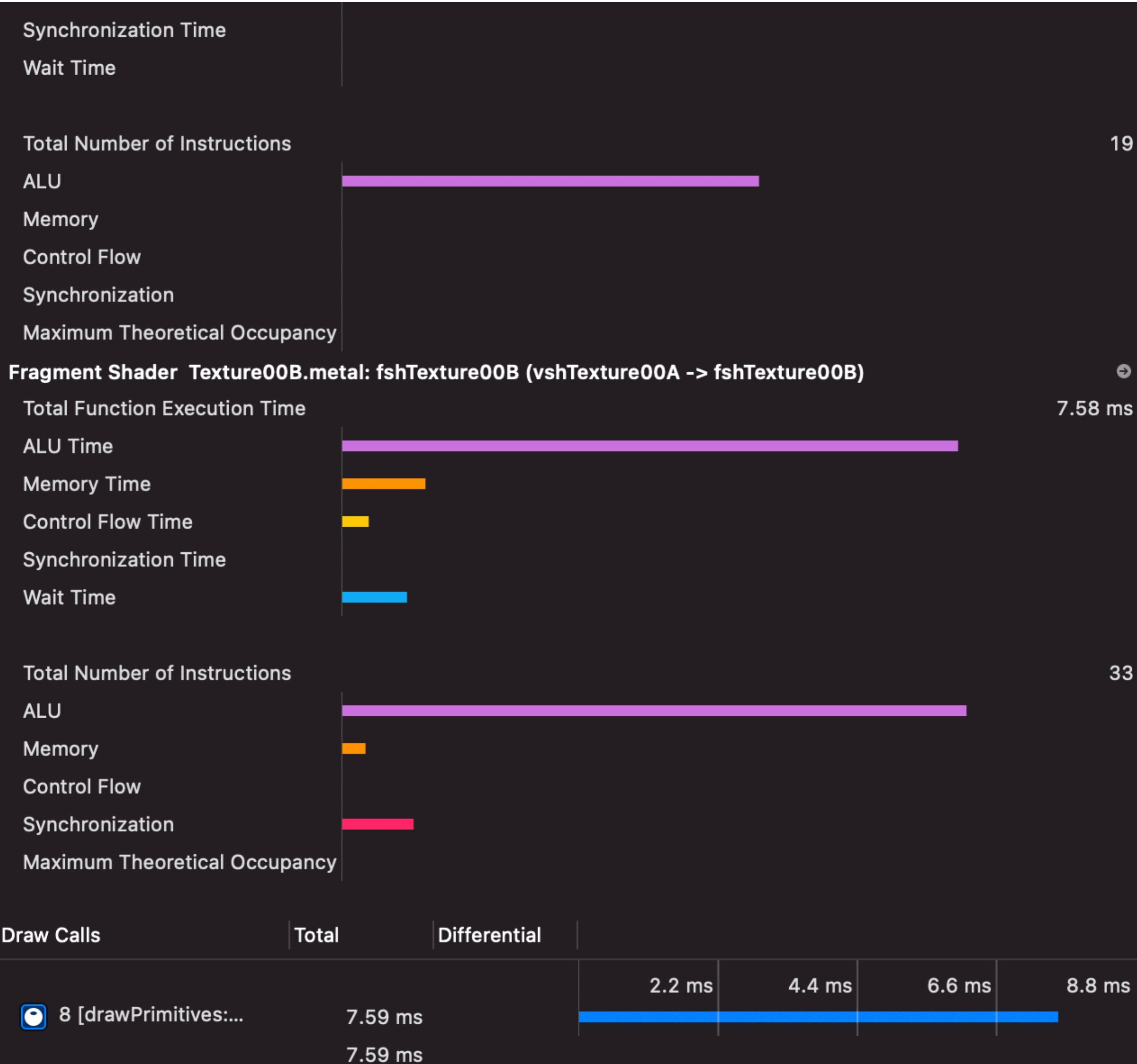
Localisation (4)



~ 10.25 ms



Localisation (4)



Potential Hotspots/Bottlenecks

No potential bottlenecks detected

Counters

Fragment Shader Texture00B.metal:fshTexture00B

	Draw	Encoder
FS ALU Limiter	N/A	73.54%
FS Texture Read Limiter	N/A	20.15%
FS Texture Write Limiter	N/A	98.41%
FS Bytes Read From Main Memory	N/A	27.21 MiB
FS Bytes Written To Main Memory	N/A	162.79 MiB
FS Main Memory Bandwidth	N/A	20.09
FS Last Level Cache Bytes Read	N/A	52.61 MiB
FS Last Level Cache Bytes Written	N/A	268.45 MiB
FS Texture L1 Bytes Read	109.49 MiB	109.77 MiB
FS Buffer L1 Miss Rate	N/A	66.67%
FS Buffer L1 Bytes Written	N/A	384.00 Bytes
FS Arithmetic Intensity	N/A	9.19
FS ALU Performance	230.45	184.56
Pixels per Second	8,857,891,00...	7,093,814,00...

Post-Fragment Stage 1 attachment, 8192 x 8192

	Draw	Encoder
Pixels Stored	67,108,860	67,108,860
Pixel Write Stall	N/A	60.80%
Predicated Texture Thread Writes	0.00%	0.00%
Compression Ratio of Texture Memory Written	N/A	0.00
Texture Cache Write Miss Rate	N/A	61.89%

Texture

	Draw	Encoder
Texture Cache Miss Rate	N/A	7.90%
Texture L1 Bytes Read	109.49 MiB	109.77 MiB
Texture Sampler Calls	67,108,860	67,108,860
Lossless Compressed Texture Samples	N/A	0.00%
Uncompressed Texture Samples	N/A	100.00%
Compression Ratio of Texture Memory Read	N/A	0.00

Filter

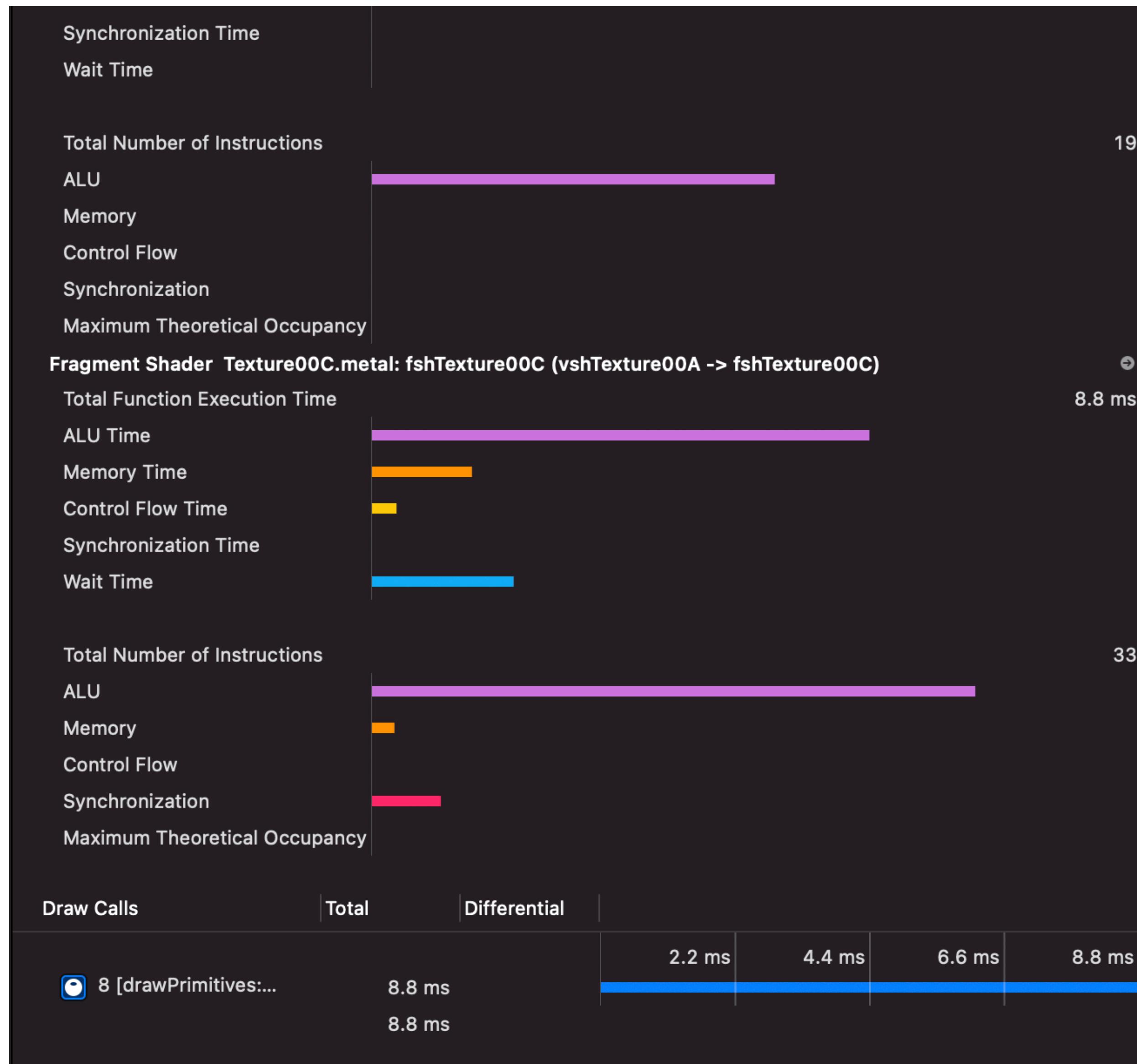
Localisation (8)



~ 10.64 ms



Localisation (8)



Potential Hotspots/Bottlenecks

No potential bottlenecks detected

Counters

Fragment Shader Texture00C.metal:fshTexture00C

	Draw	Encoder
FS ALU Limiter	—	69.41%
FS Texture Read Limiter	—	33.47%
FS Texture Write Limiter	—	85.16%
FS Bytes Read From Main Memory	—	38.43 MiB
FS Bytes Written To Main Memory	—	174.19 MiB
FS Main Memory Bandwidth	—	21.22
FS Last Level Cache Bytes Read	—	112.24 MiB
FS Last Level Cache Bytes Written	—	268.46 MiB
FS Texture L1 Bytes Read	—	288.91 MiB
FS Buffer L1 Miss Rate	—	80.00%
FS Buffer L1 Bytes Written	—	384.00 Bytes
FS Arithmetic Intensity	—	8.22
FS ALU Performance	—	174.20
Pixels per Second	—	6,695,862,00...

Post-Fragment Stage 1 attachment, 8192 x 8192

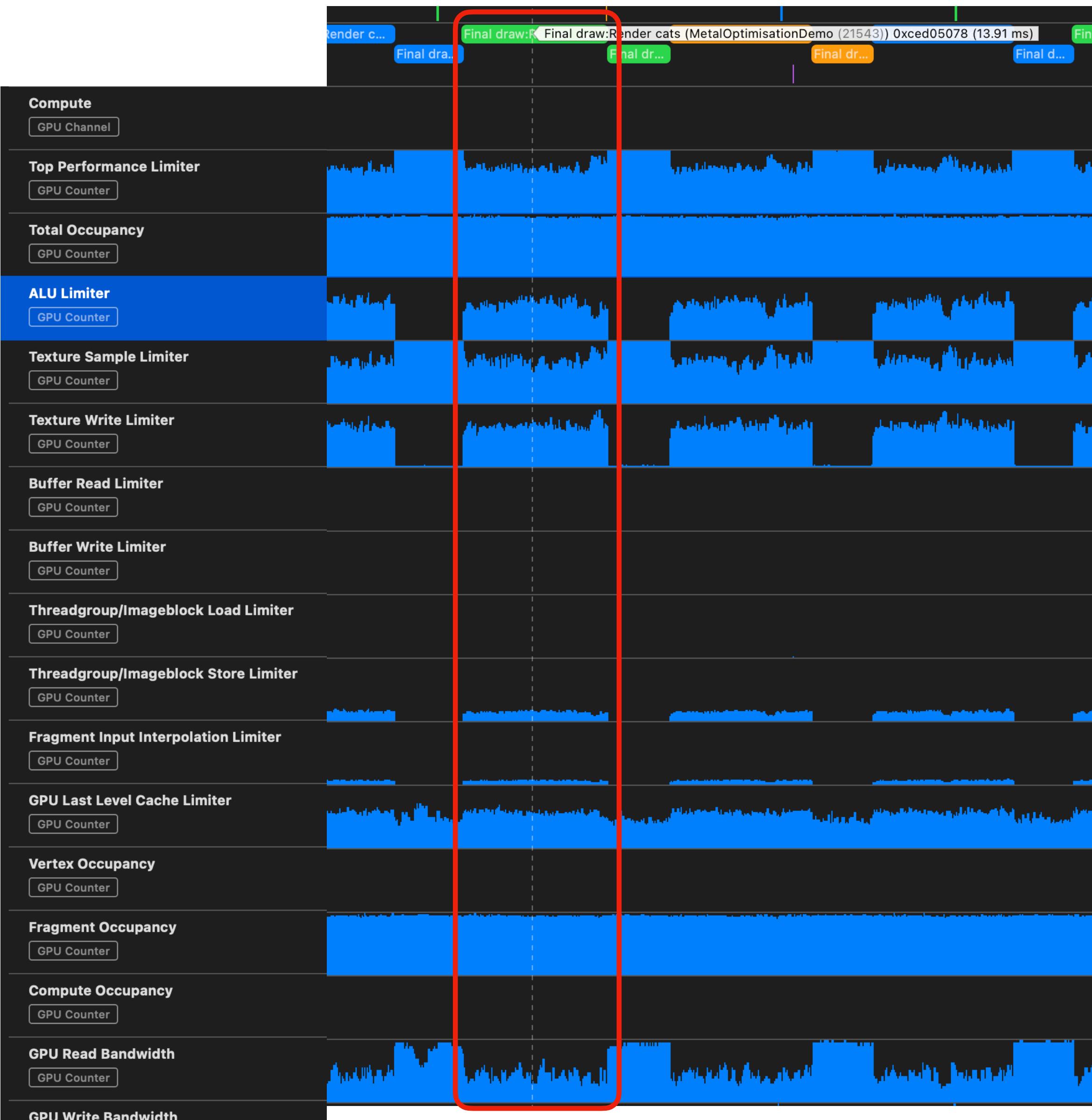
Pixels Stored	—	67,108,860
Pixel Write Stall	—	50.27%
Predicated Texture Thread Writes	—	0.00%
Compression Ratio of Texture Memory Written	—	0.00
Texture Cache Write Miss Rate	—	57.67%

Texture

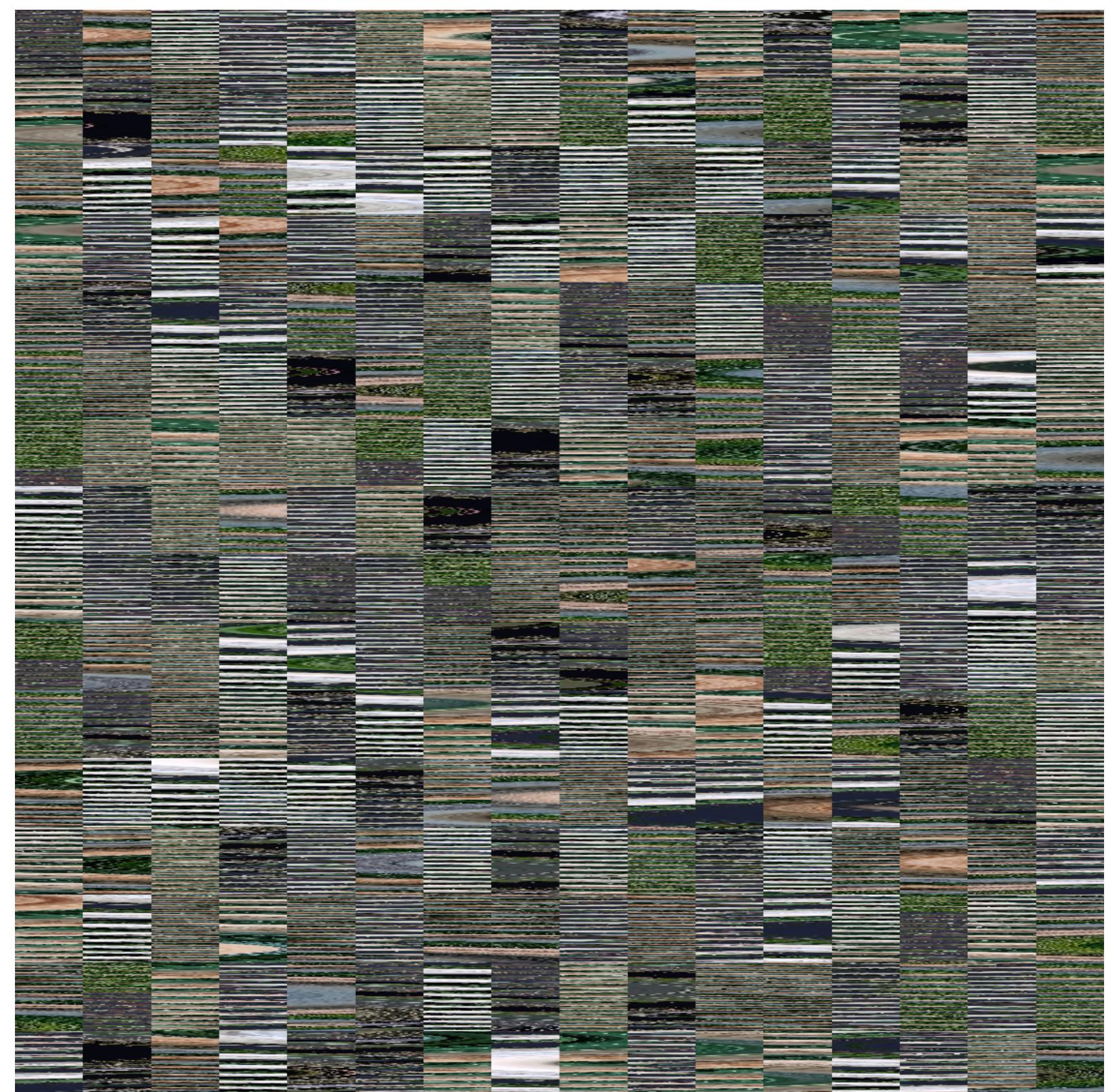
Texture Cache Miss Rate	—	8.42%
Texture L1 Bytes Read	—	288.91 MiB
Texture Sampler Calls	—	67,108,860
Lossless Compressed Texture Samples	—	0.00%
Uncompressed Texture Samples	—	100.00%
Compression Ratio of Texture Memory Read	—	0.00

Filter

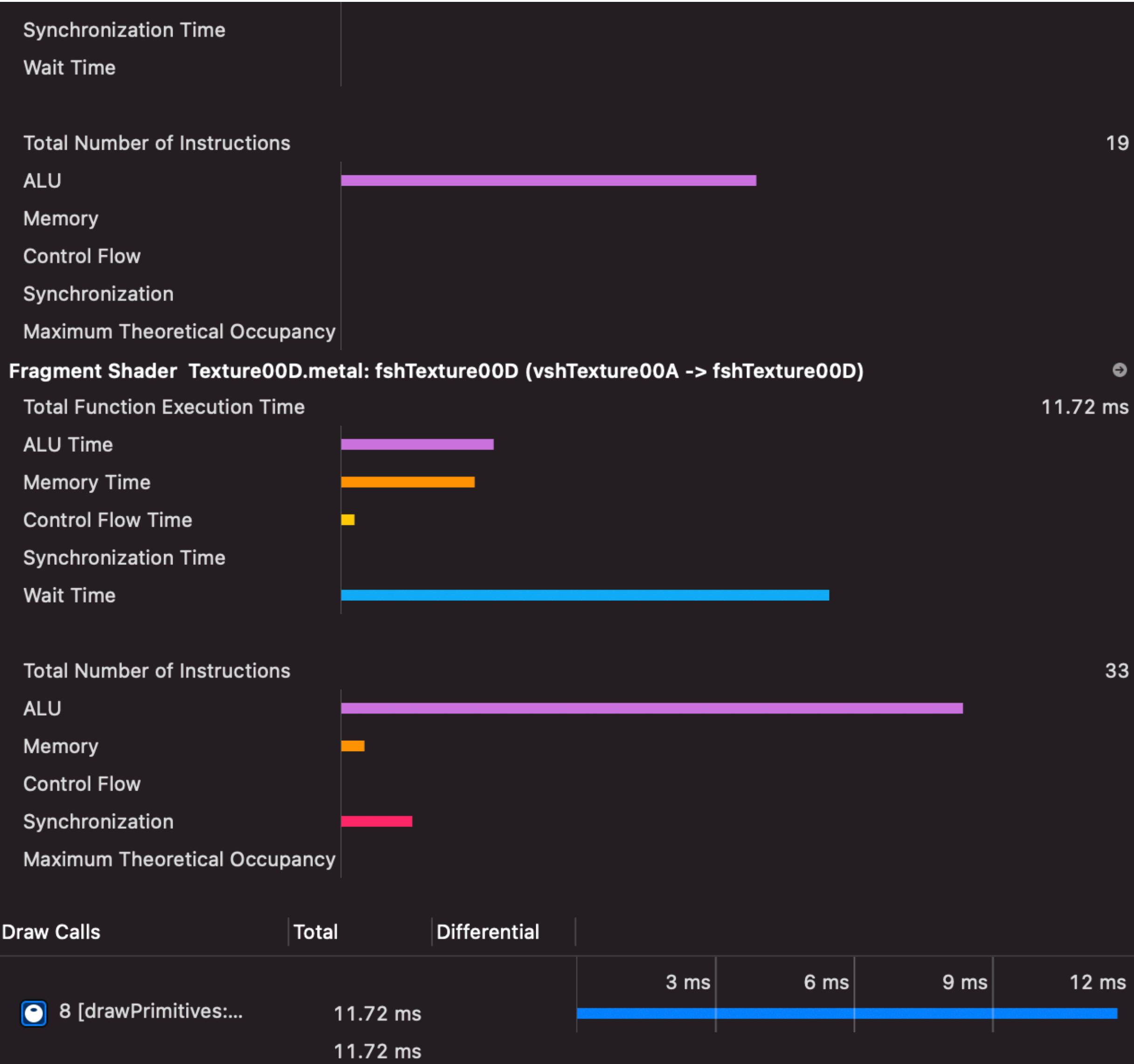
Localisation (16)



~ 13.91 ms

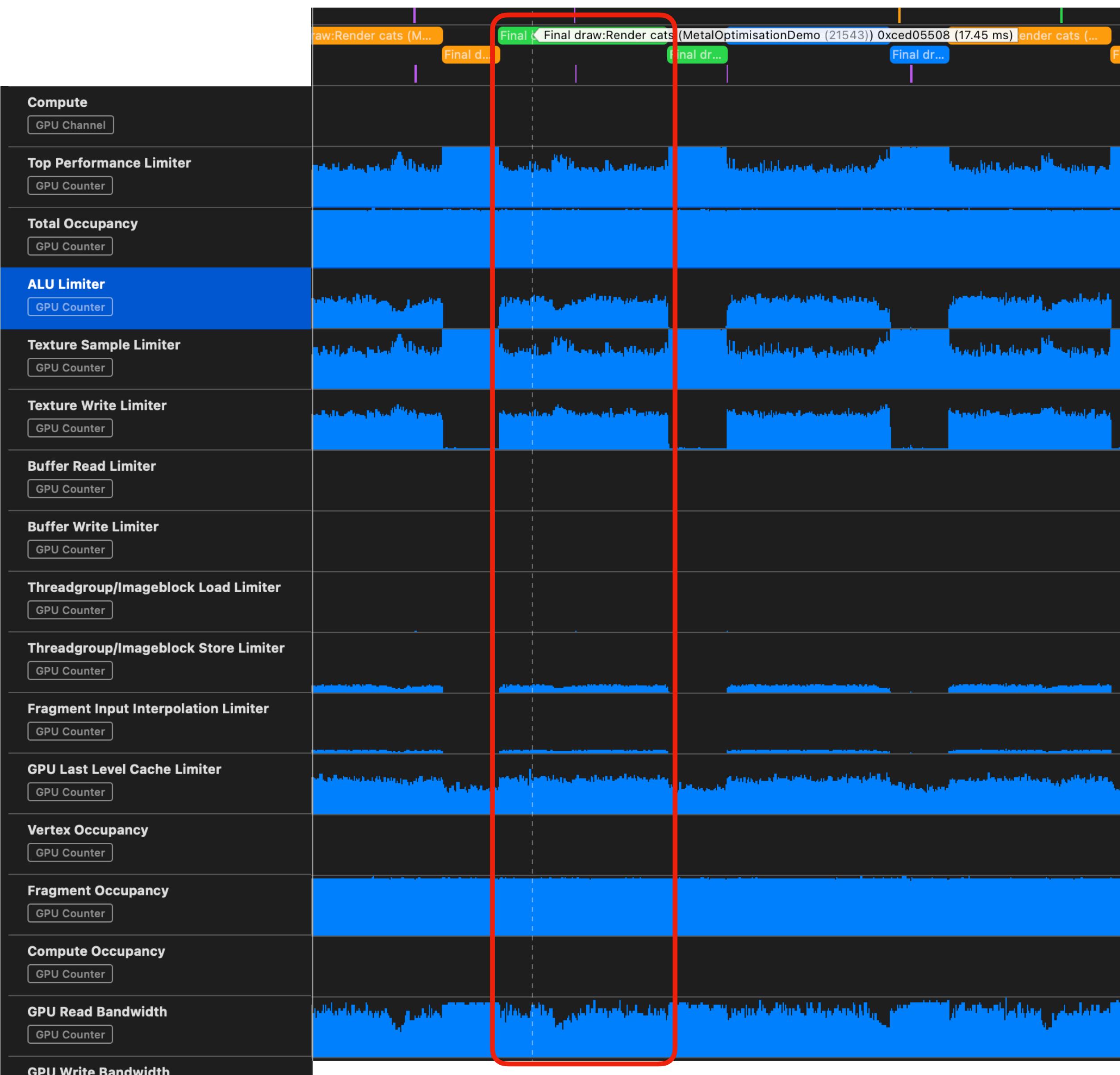


Localisation (16)

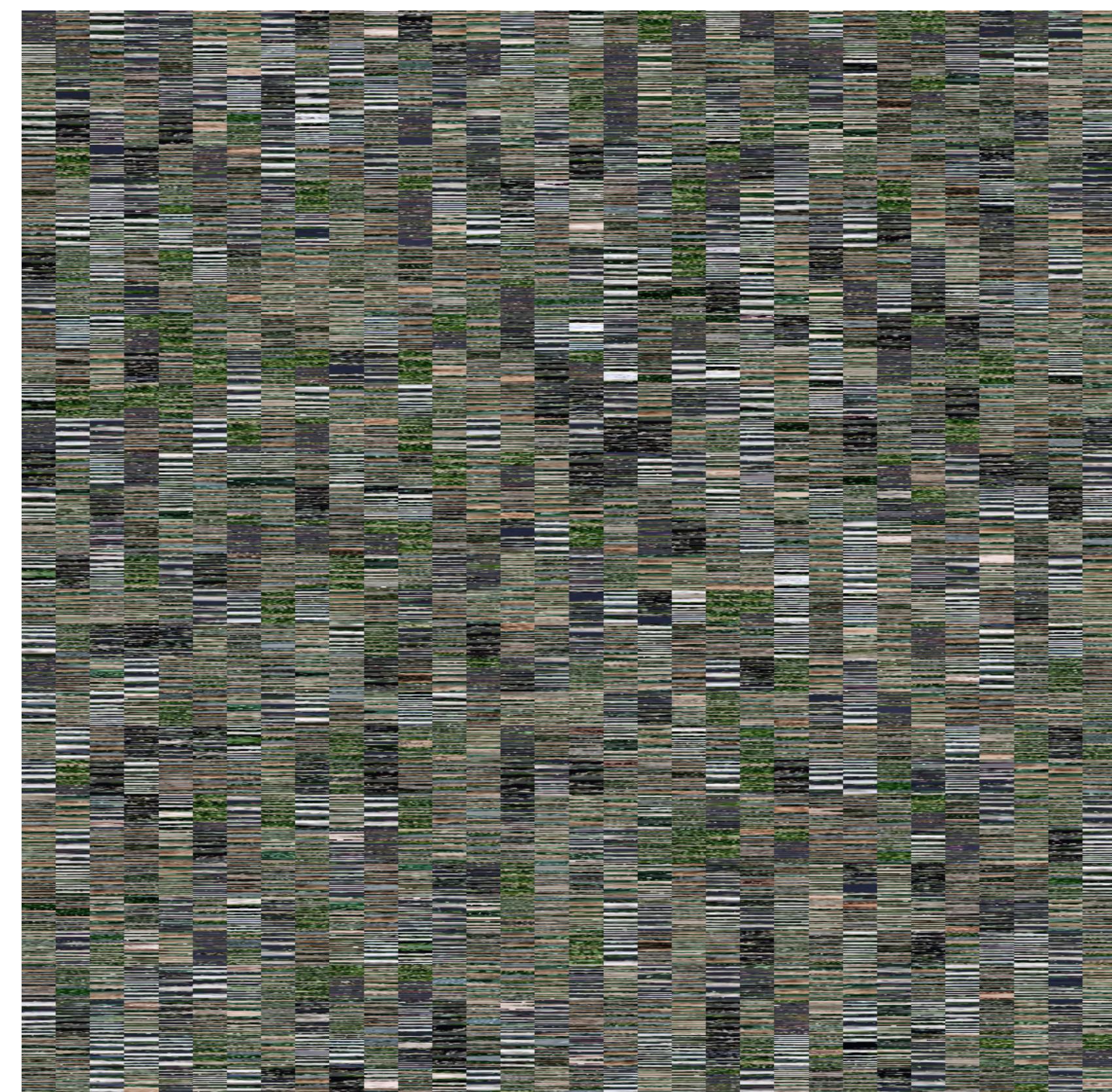


Potential Hotspots/Bottlenecks			
No potential bottlenecks detected			
Counters			
Fragment Shader Texture00D.metal:fshTexture00D			
FS ALU Limiter	N/A	60.39%	
FS Texture Read Limiter	N/A	60.36%	
FS Texture Write Limiter	N/A	60.88%	
FS Bytes Read From Main Memory	N/A	126.04 MiB	
FS Bytes Written To Main Memory	N/A	196.89 MiB	
FS Main Memory Bandwidth	N/A	27.07	
FS Last Level Cache Bytes Read	N/A	322.04 MiB	
FS Last Level Cache Bytes Written	N/A	268.46 MiB	
FS Texture L1 Bytes Read	330.74 MiB	327.83 MiB	
FS Buffer L1 Miss Rate	N/A	83.34%	
FS Buffer L1 Bytes Written	N/A	384.00 Bytes	
FS Arithmetic Intensity	N/A	5.41	
FS ALU Performance	149.04	146.35	
Pixels per Second	5,728,499,00...	5,625,432,00...	
Post-Fragment Stage 1 attachment, 8192 x 8192			
Pixels Stored	67,108,860	67,108,860	
Pixel Write Stall	N/A	29.83%	
Predicated Texture Thread Writes	0.00%	0.00%	
Compression Ratio of Texture Memory Written	N/A	0.00	
Texture Cache Write Miss Rate	N/A	51.12%	
Texture			
Texture Cache Miss Rate	N/A	23.32%	
Texture L1 Bytes Read	330.74 MiB	327.83 MiB	
Texture Sampler Calls	67,108,860	67,108,860	
Lossless Compressed Texture Samples	N/A	0.00%	
Uncompressed Texture Samples	N/A	100.00%	
Compression Ratio of Texture Memory Read	N/A	0.00	

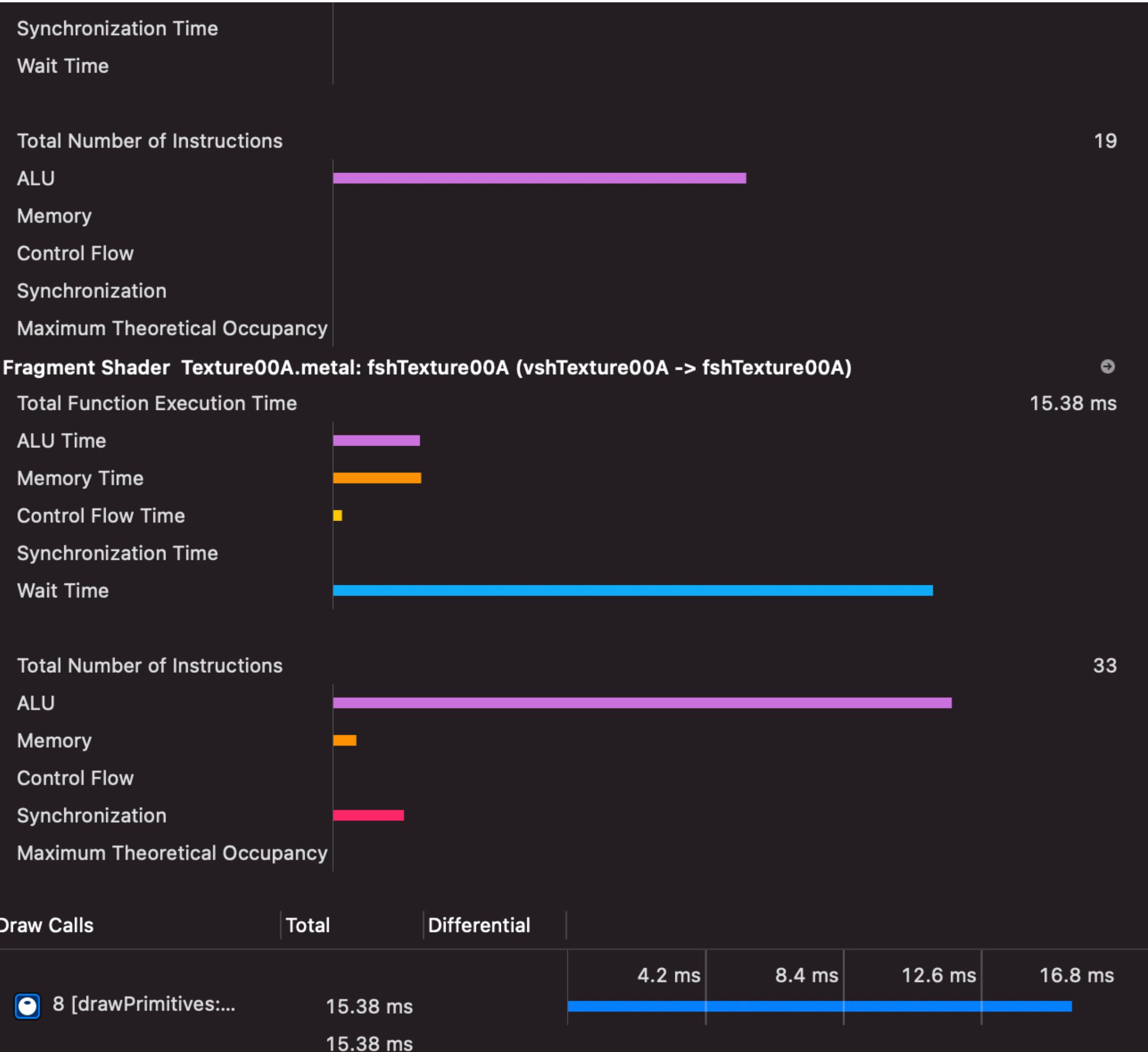
Localisation (32)



~ 17.45 ms



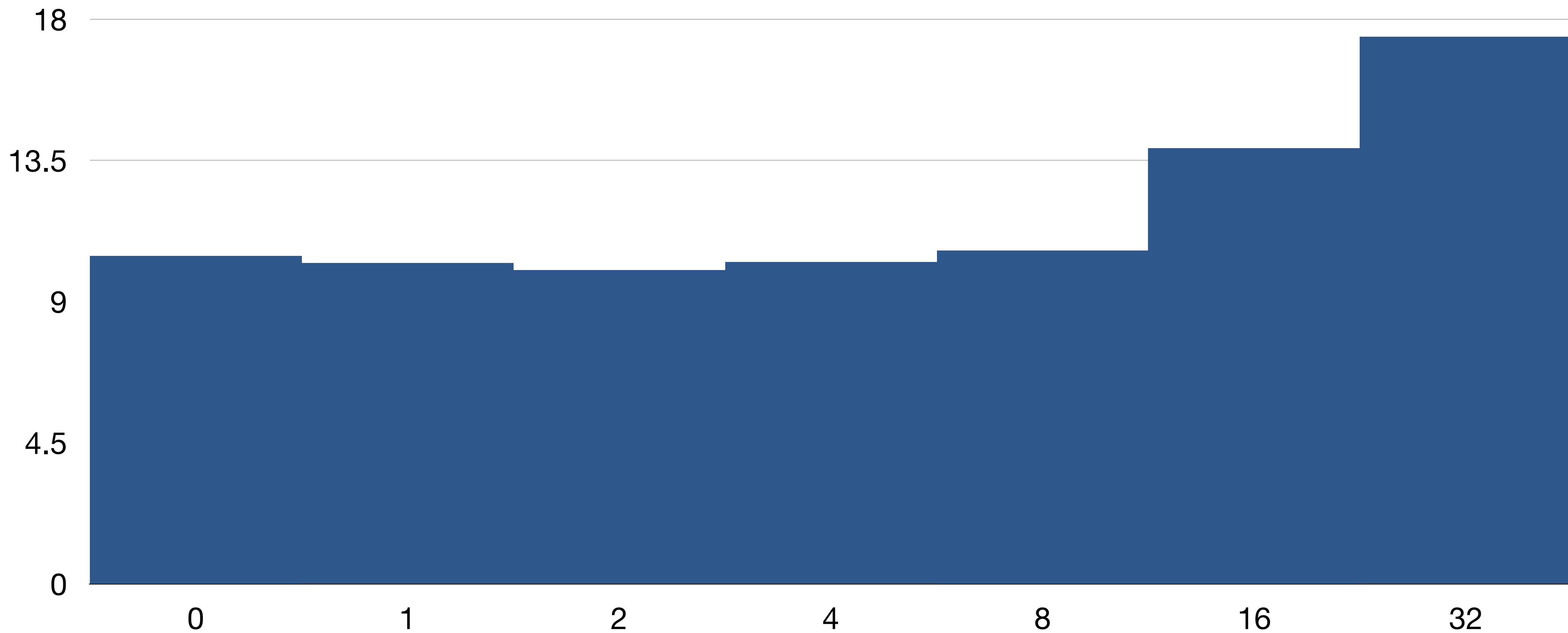
Localisation (32)



Potential Hotspots/Bottlenecks		Draw	Encoder
No potential bottlenecks detected			
Counters			
Fragment Shader Texture00A.metal:fshTexture00A			
FS ALU Limiter	N/A	45.42%	
FS Texture Read Limiter	N/A	57.81%	
FS Texture Write Limiter	N/A	56.91%	
FS Bytes Read From Main Memory	N/A	288.66 MiB	
FS Bytes Written To Main Memory	N/A	193.02 MiB	
FS Main Memory Bandwidth	N/A	30.95	
FS Last Level Cache Bytes Read	N/A	375.54 MiB	
FS Last Level Cache Bytes Written	N/A	268.46 MiB	
FS Texture L1 Bytes Read	335.49 MiB	327.75 MiB	
FS Buffer L1 Miss Rate	N/A	66.67%	
FS Buffer L1 Bytes Written	N/A	384.00 Bytes	
FS Arithmetic Intensity	N/A	3.63	
FS ALU Performance	113.52	112.17	
Pixels per Second	4,363,445,00...	4,311,385,00...	
Post-Fragment Stage 1 attachment, 8192 x 8192			
Pixels Stored	67,108,860	67,108,860	
Pixel Write Stall	N/A	33.09%	
Predicated Texture Thread Writes	0.00%	0.00%	
Compression Ratio of Texture Memory Written	N/A	0.00	
Texture Cache Write Miss Rate	N/A	57.12%	
Texture			
Texture Cache Miss Rate	N/A	27.27%	
Texture L1 Bytes Read	335.49 MiB	327.75 MiB	
Texture Sampler Calls	67,108,860	67,108,860	
Lossless Compressed Texture Samples	N/A	0.00%	
Uncompressed Texture Samples	N/A	100.00%	
Compression Ratio of Texture Memory Read	N/A	0.00	

Filter

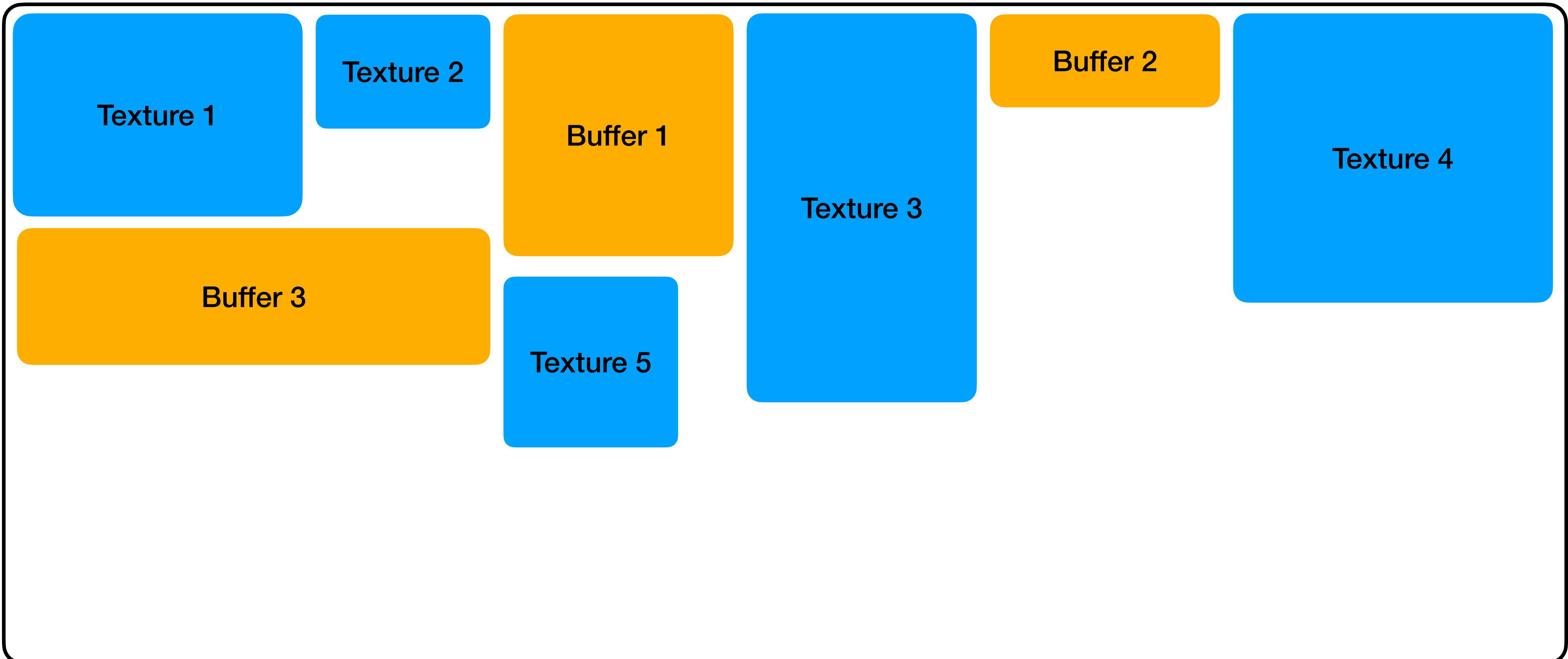
Summary (ms)



Memory

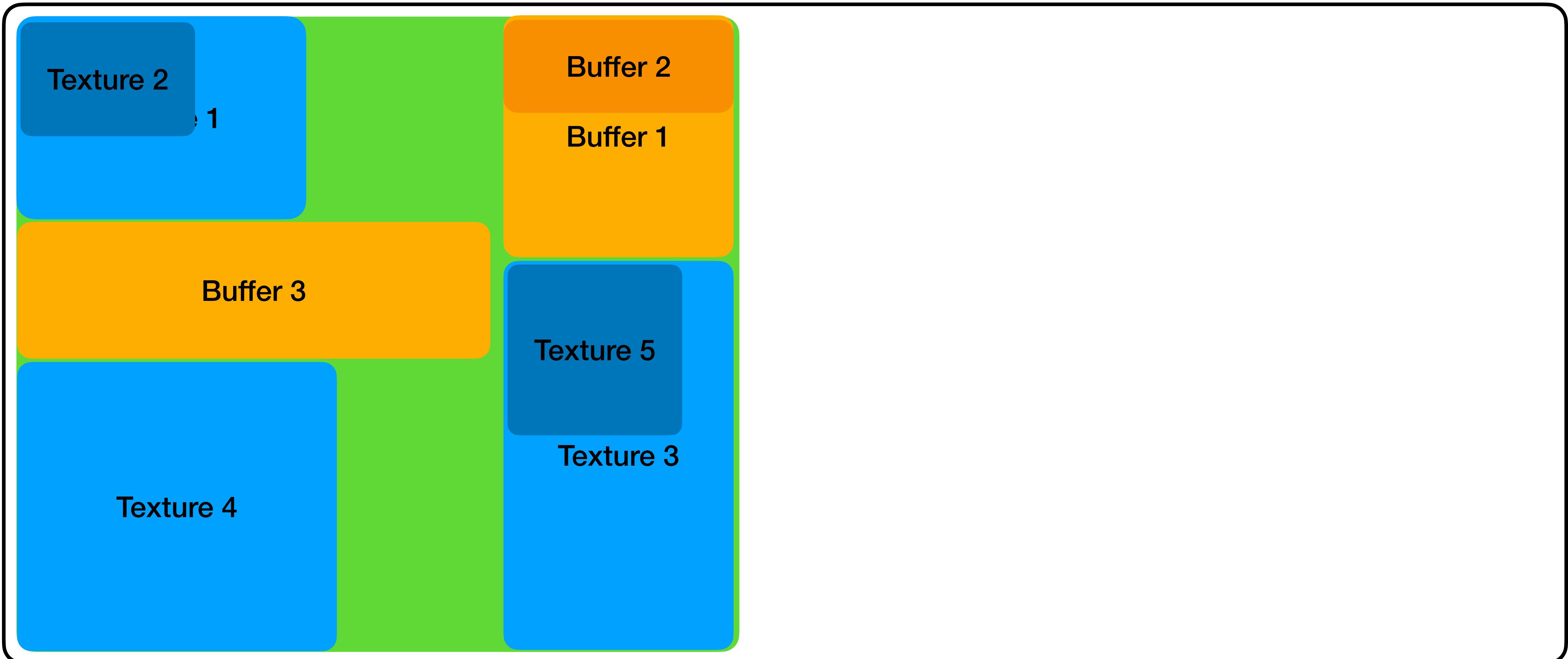
Heap and device

Device: $N \times (\text{allocate} + \text{place})$



Heap and device

Heap: allocate + N x place



Reuse buffers

Usage 1

Usage 2

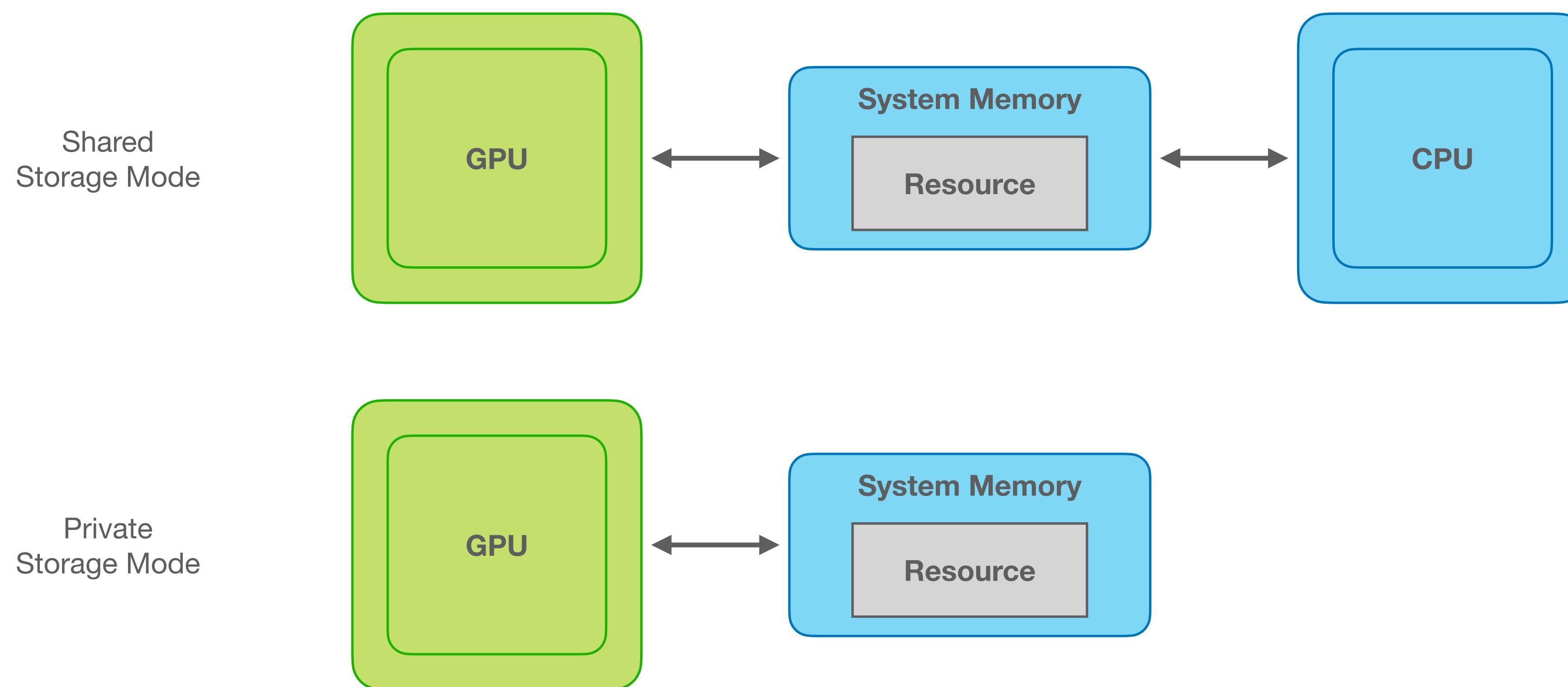
Usage 3

Usage 4

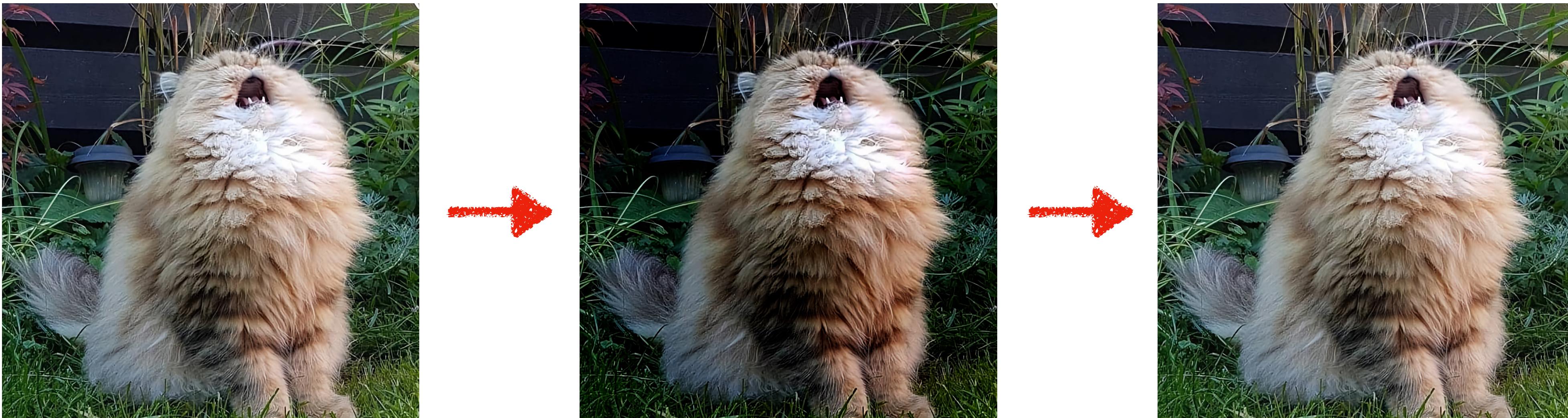
Usage 5

Usage 6

Private and shared



Several passes



Several passes

`sample`

```
// Pass 1
renderPassDescriptor.colorAttachments[0].texture = _textureTmp;
renderPassDescriptor.colorAttachments[0].loadAction = MTLLoadActionLoad;
id <MTLRenderCommandEncoder> encoder =
[commandBuffer renderCommandEncoderWithDescriptor:renderPassDescriptor];
[encoder setRenderPipelineState:_pipelineRender1];
[encoder setFragmentTexture:_texture atIndex:0];
[encoder drawPrimitives:MTLPrimitiveTypeTriangleStrip vertexStart:0 vertexCount:4];
[encoder endEncoding];

// Pass 2
renderPassDescriptor = [MTLRenderPassDescriptor new];
renderPassDescriptor.colorAttachments[0].texture = _textureOut;
renderPassDescriptor.colorAttachments[0].loadAction = MTLLoadActionLoad;
id <MTLRenderCommandEncoder> encoder =
[commandBuffer renderCommandEncoderWithDescriptor:renderPassDescriptor];
[encoder setRenderPipelineState:_pipelineRender2];
[encoder setFragmentTexture:_textureTmp atIndex:0];
[encoder drawPrimitives:MTLPrimitiveTypeTriangleStrip vertexStart:0 vertexCount:4];
[encoder endEncoding];
```

Several passes

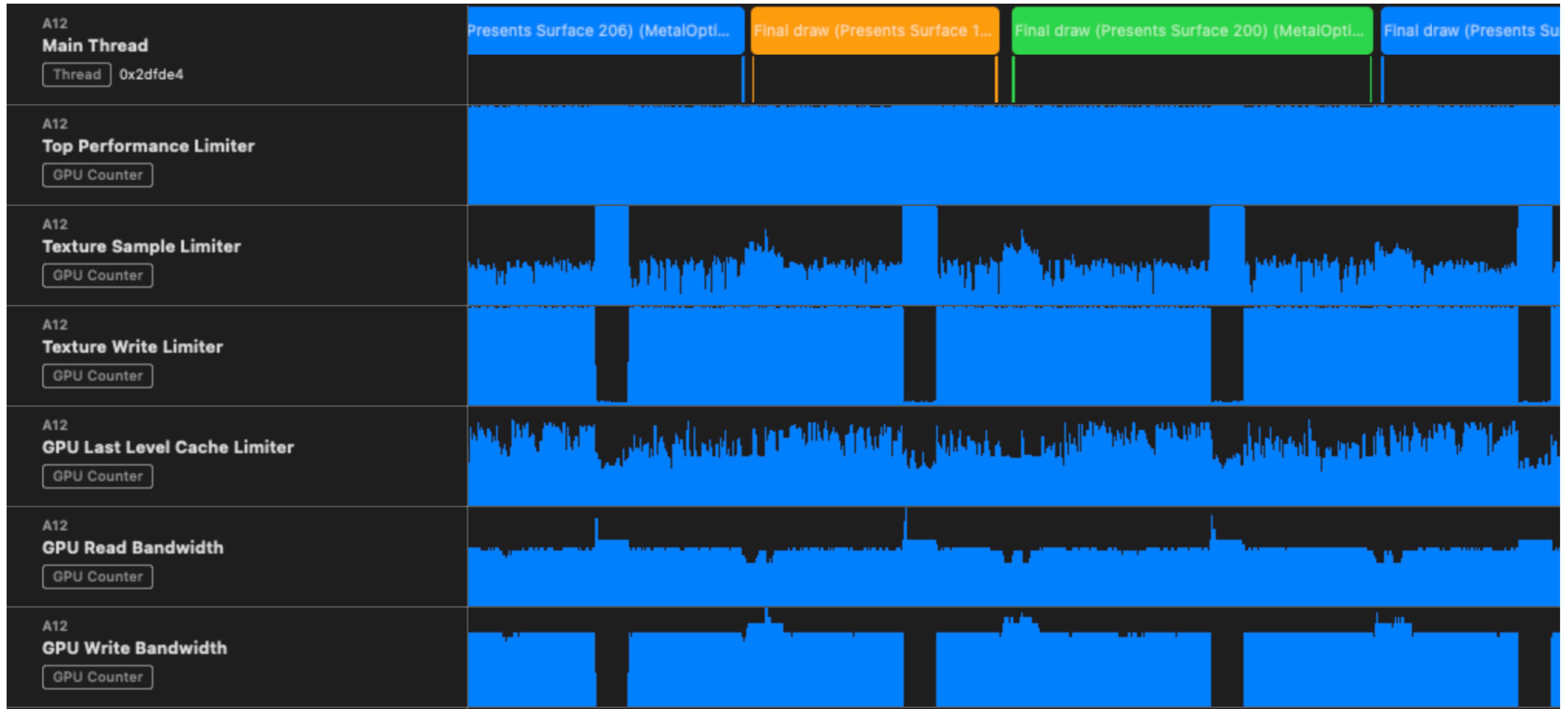
`sample`

```
fragment float4 fshTexture01A_1(ColorInOut in [[stage_in]], texture2d<float> image [[ texture(0) ]])
{
    constexpr sampler imageSampler(address::repeat, filter::nearest);
    float2 uv = in.texCoord;
    float4 src = image.sample(imageSampler, uv);
    return src * src;
}

fragment float4 fshTexture01A_2(ColorInOut in [[stage_in]], texture2d<float> image [[ texture(0) ]])
{
    constexpr sampler imageSampler(address::repeat, filter::nearest);
    float2 uv = in.texCoord;
    float4 src = image.sample(imageSampler, uv);
    return sqrt(src);
}
```

Several passes

`sample`



Several passes

attachment

```
MTLRenderPassDescriptor* renderPassDescriptor = [MTLRenderPassDescriptor new];
renderPassDescriptor.colorAttachments[0].texture = _textureOut;
renderPassDescriptor.colorAttachments[0].loadAction = MTLLoadActionLoad;
renderPassDescriptor.colorAttachments[1].texture = _textureTmp;
renderPassDescriptor.colorAttachments[1].loadAction = MTLLoadActionLoad;
if (renderPassDescriptor != nil) {
    id <MTLRenderCommandEncoder> encoder =
    [commandBuffer renderCommandEncoderWithDescriptor:renderPassDescriptor];
    encoder.label = @"Render cats";

    [encoder setRenderPipelineState:_pipelineRender1];
    [encoder setFragmentTexture:_texture atIndex:0];
    [encoder drawPrimitives:MTLPrimitiveTypeTriangleStrip vertexStart:0 vertexCount:4];

    [encoder setRenderPipelineState:_pipelineRender2];
    [encoder drawPrimitives:MTLPrimitiveTypeTriangleStrip vertexStart:0 vertexCount:4];

    [encoder endEncoding];
}
```

Several passes

attachment

```
typedef struct {
    float4 dst [[ color(1) ]];
} FragmentOut1;
```

```
typedef struct {
    float4 dst [[ color(0) ]];
} FragmentOut2;
```

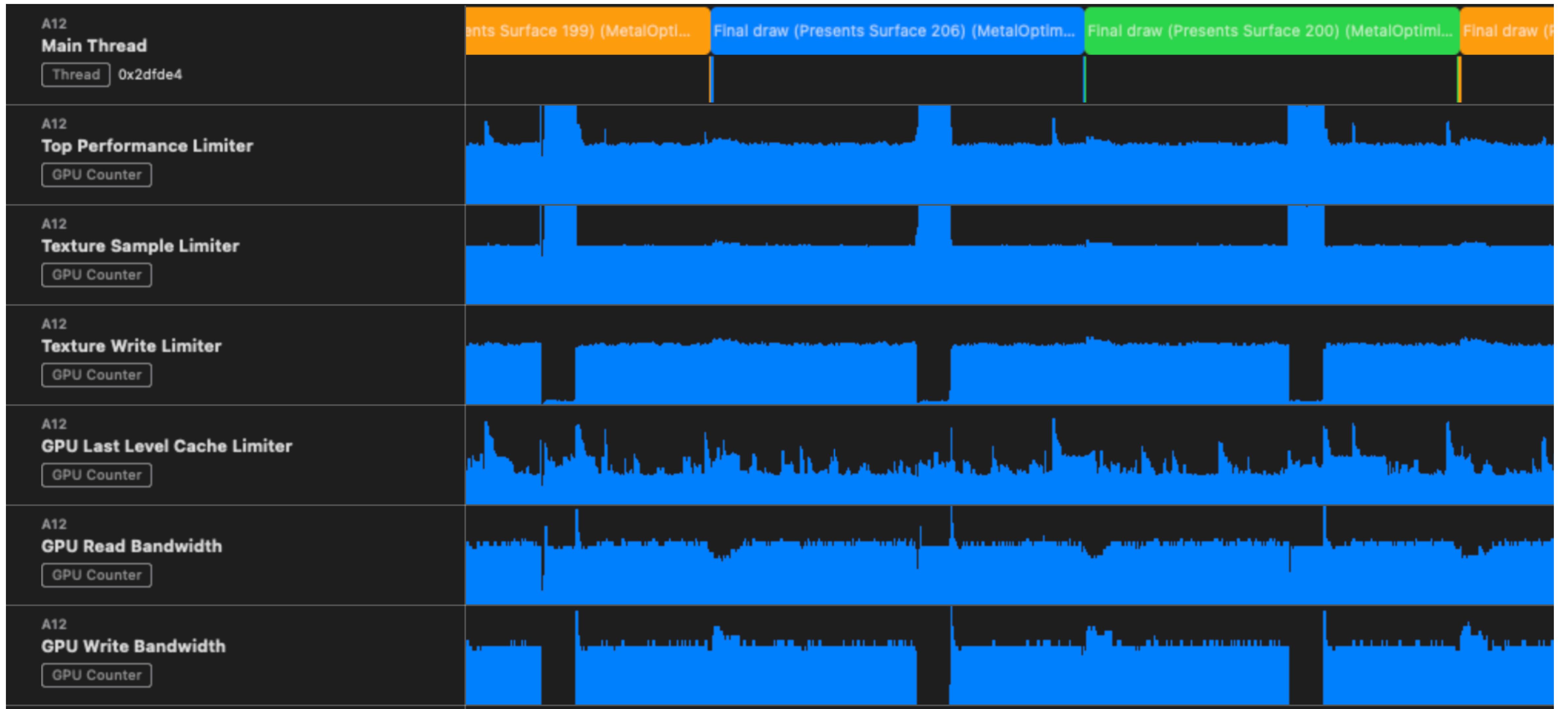
```
fragment FragmentOut1 fshTexture01B_1(
    ColorInOut in [[stage_in]],
    texture2d<float> image [[ texture(0) ]])
{
    constexpr sampler imageSampler(address::repeat, filter::nearest);

    float2 uv = in.texCoord;
    float4 src = image.sample(imageSampler, uv);
    return {src * src};
}

fragment FragmentOut2 fshTexture01B_2(
    ColorInOut in [[stage_in]],
    const float4 image [[ color(1) ]])
{
    float4 src = image;
    return {sqrt(src)};
}
```

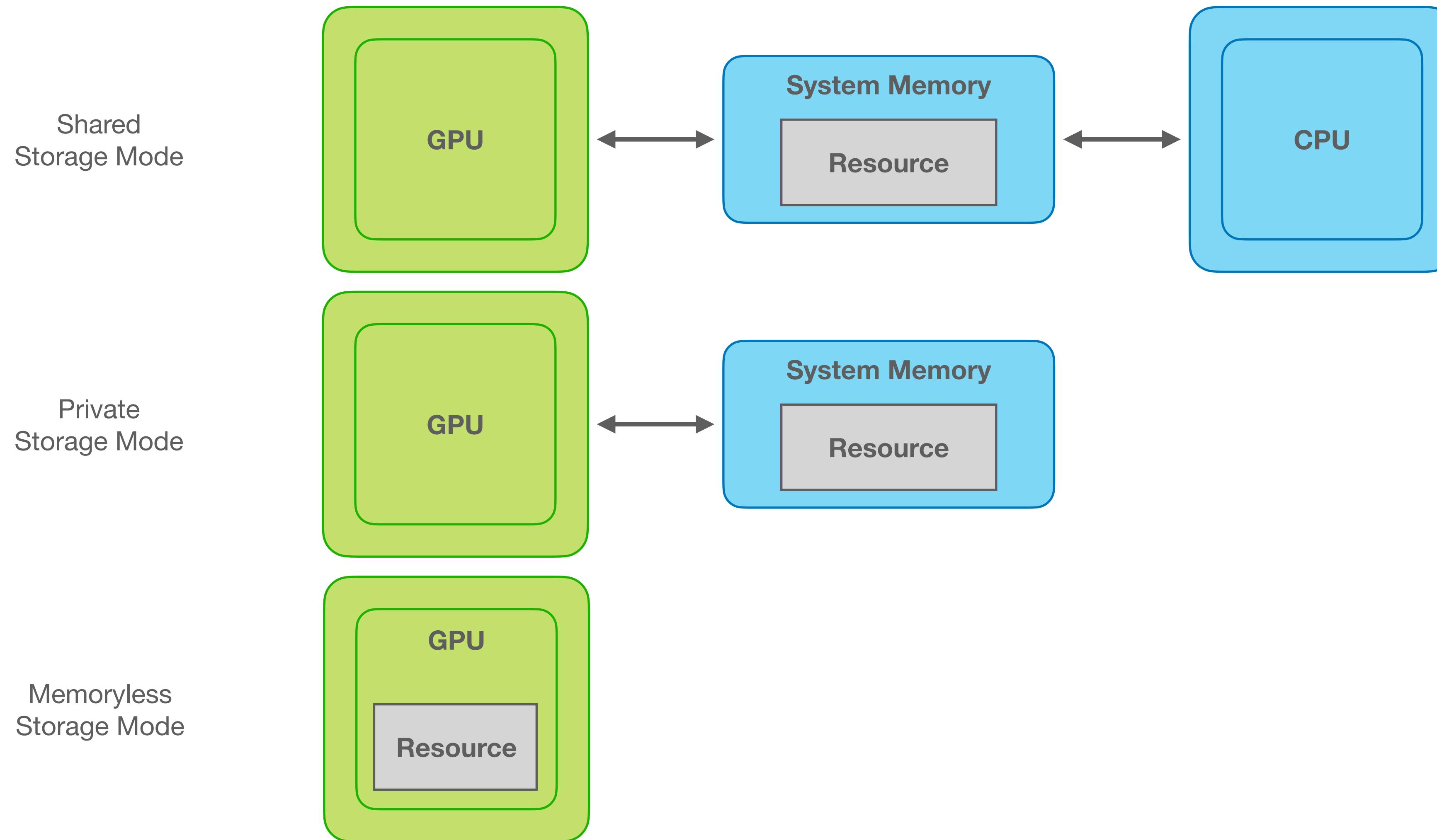
Several passes

attachment



Several passes

memoryless textures



Several passes

memoryless textures

```
MTLRenderPassDescriptor* renderPassDescriptor = [MTLRenderPassDescriptor new];
renderPassDescriptor.colorAttachments[0].texture = _textureOut;
renderPassDescriptor.colorAttachments[0].loadAction = MTLLoadActionLoad;
renderPassDescriptor.colorAttachments[1].texture = _textureTmp;
renderPassDescriptor.colorAttachments[1].loadAction = MTLLoadActionDontCare;
renderPassDescriptor.colorAttachments[1].storeAction = MTLStoreActionDontCare;
if (renderPassDescriptor != nil) {
    id <MTLRenderCommandEncoder> encoder =
    [commandBuffer renderCommandEncoderWithDescriptor:renderPassDescriptor];
    encoder.label = @"Render cats";

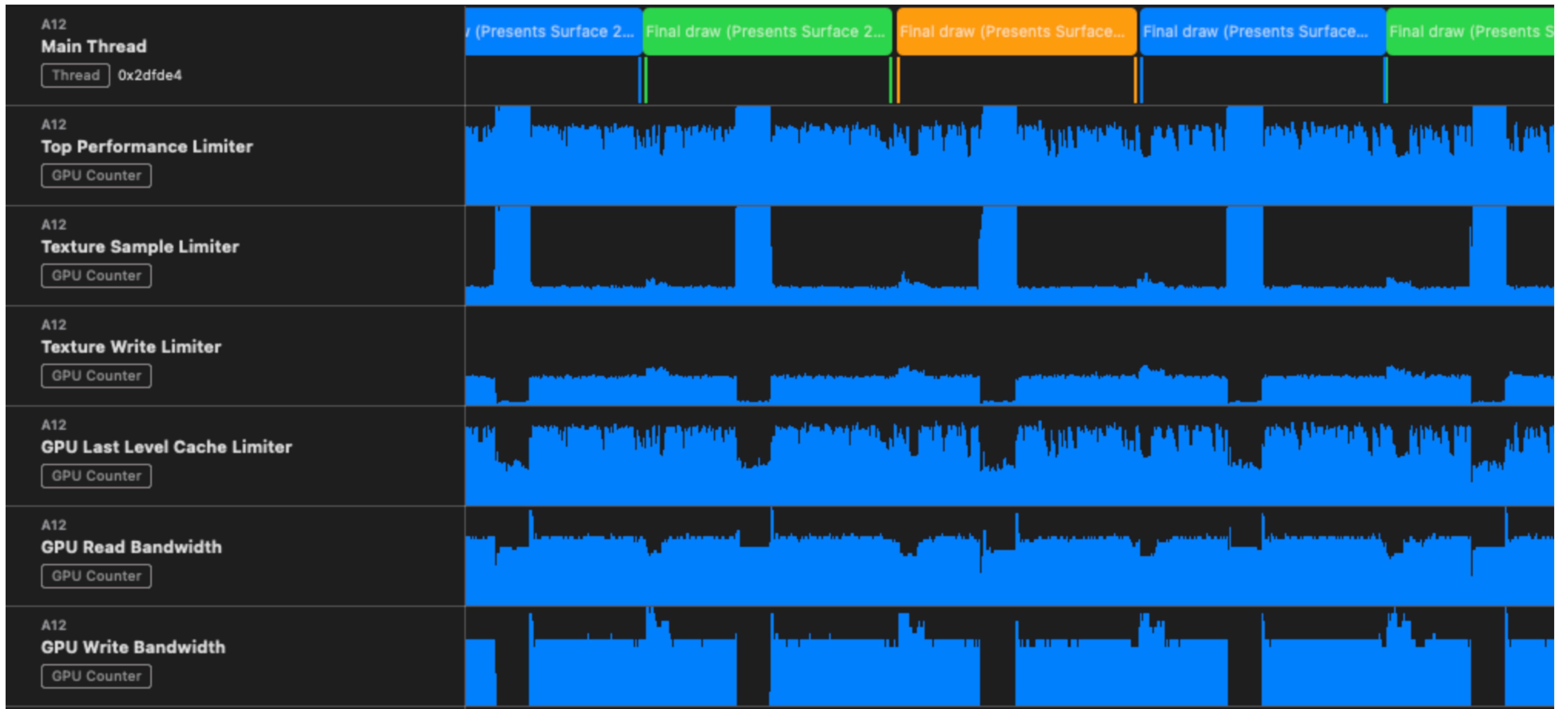
    [encoder setRenderPipelineState:_pipelineRender1];
    [encoder setFragmentTexture:_texture atIndex:0];
    [encoder drawPrimitives:MTLPrimitiveTypeTriangleStrip vertexStart:0 vertexCount:4];

    [encoder setRenderPipelineState:_pipelineRender2];
    [encoder drawPrimitives:MTLPrimitiveTypeTriangleStrip vertexStart:0 vertexCount:4];

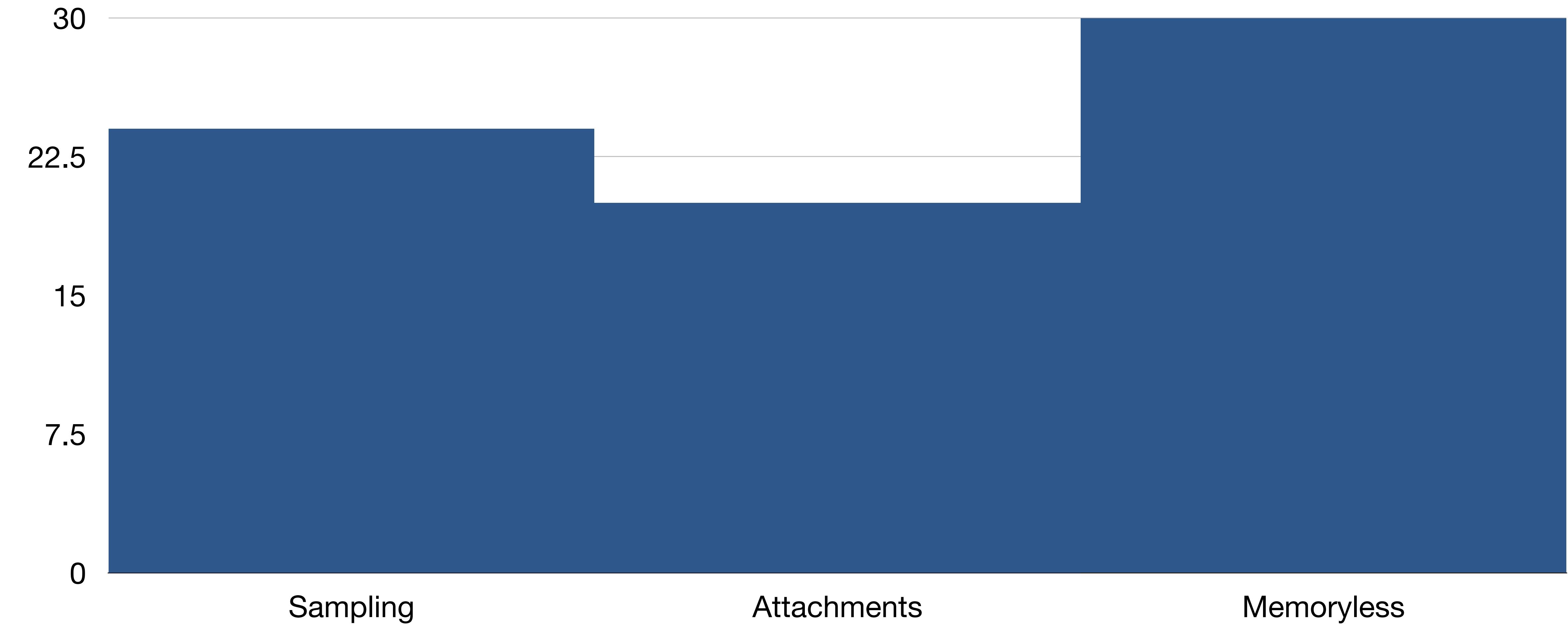
    [encoder endEncoding];
}
```

Several passes

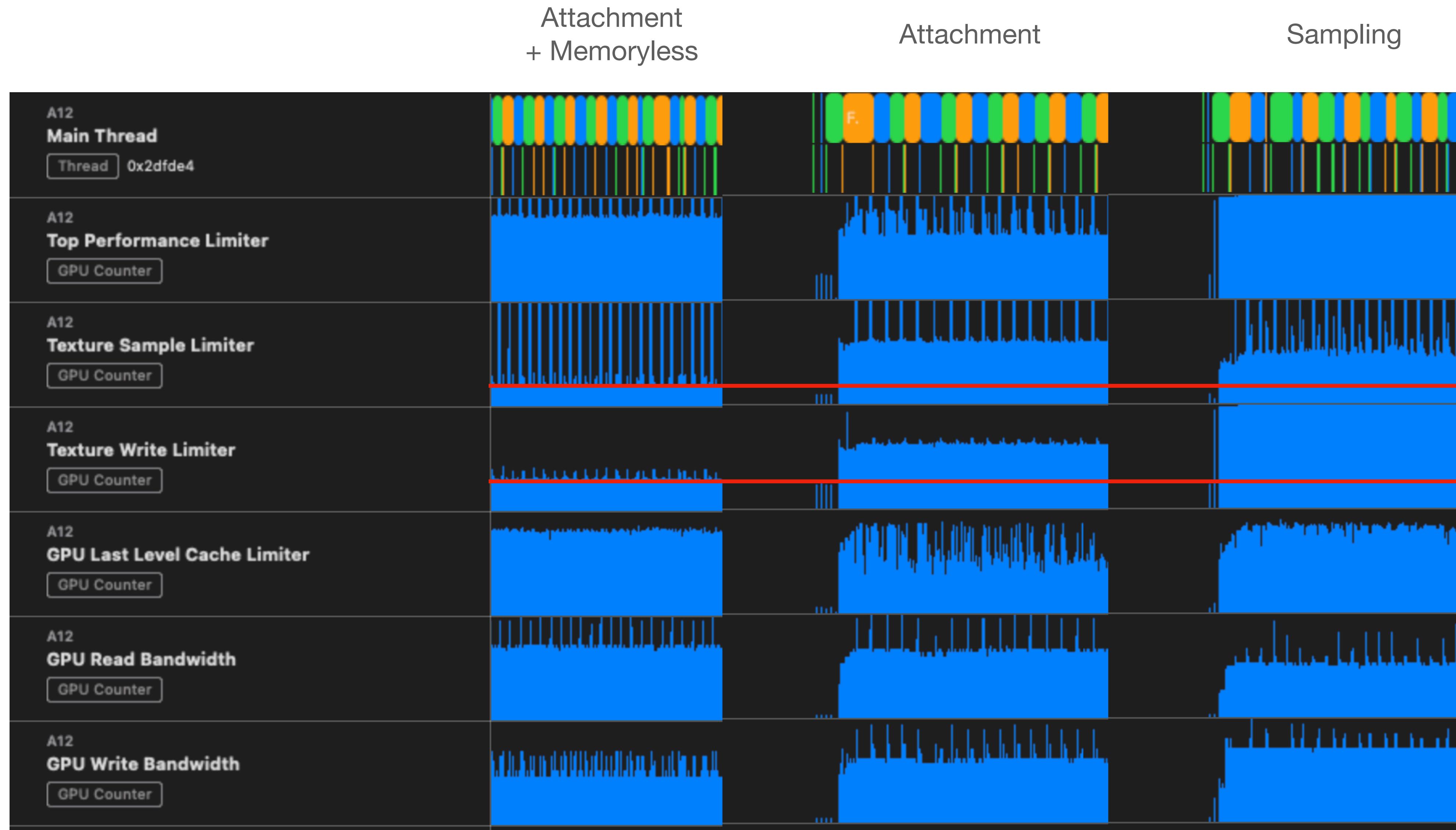
memoryless textures



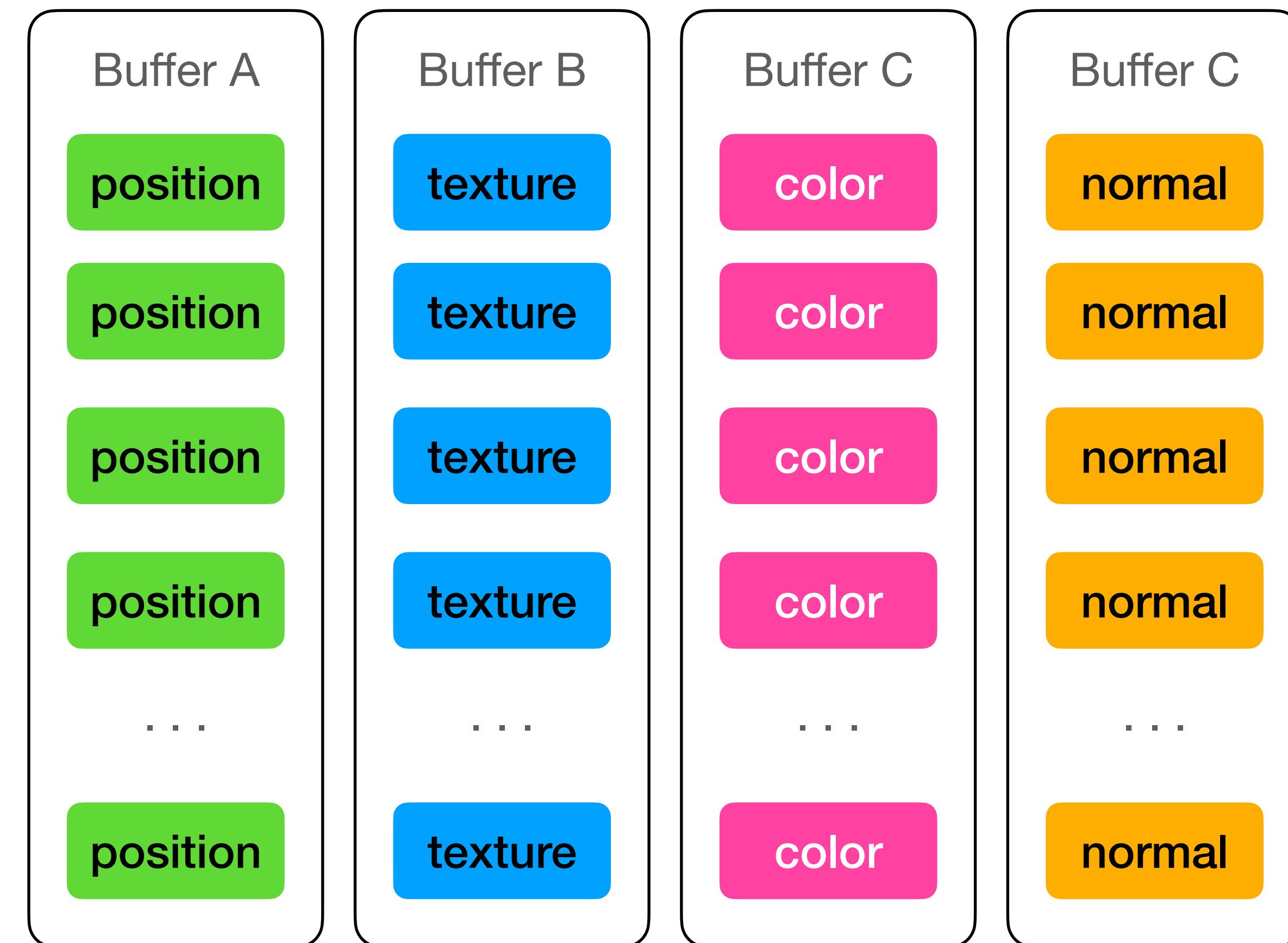
Summary (FPS)



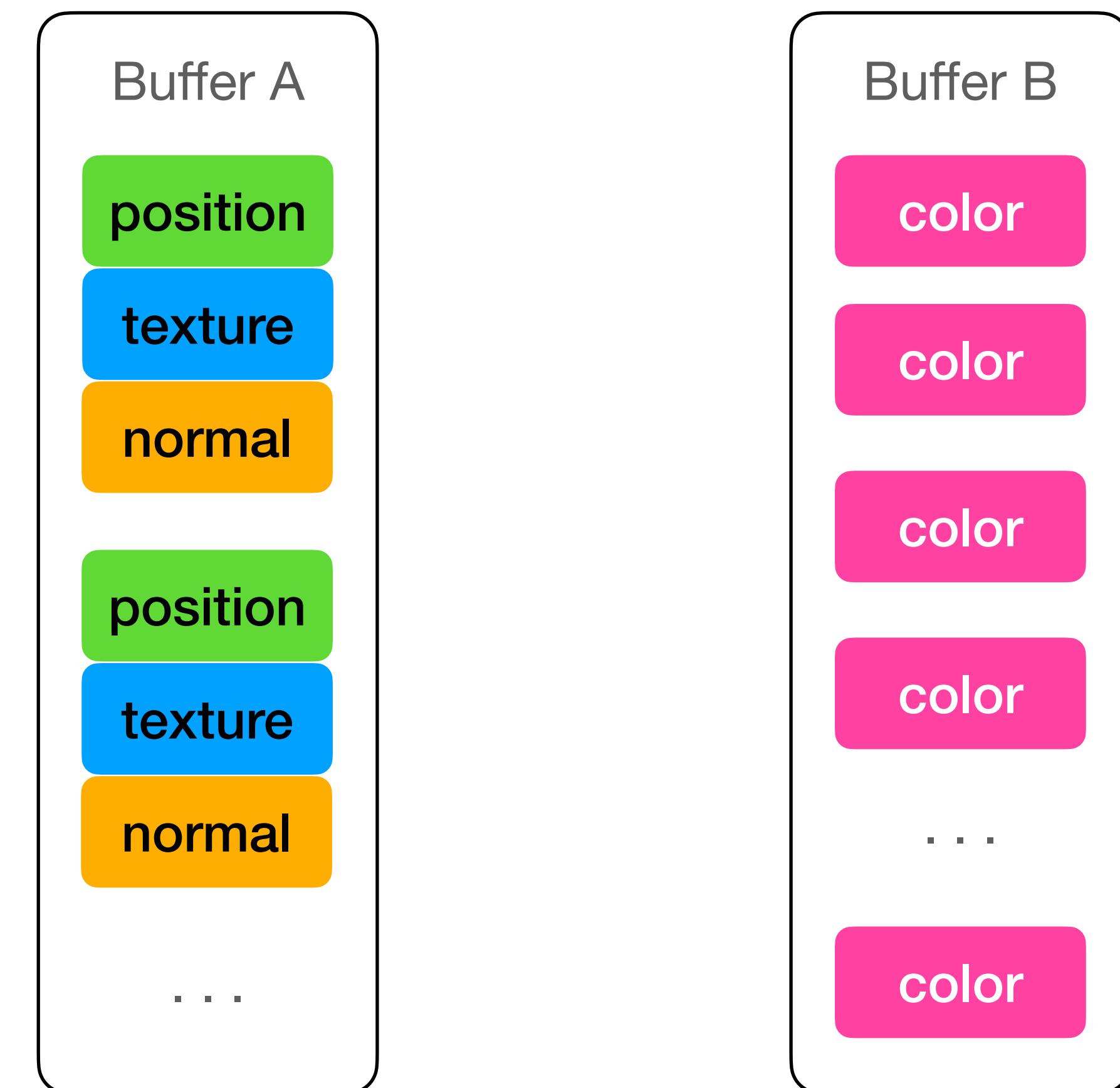
Summary



Pack your data optimally

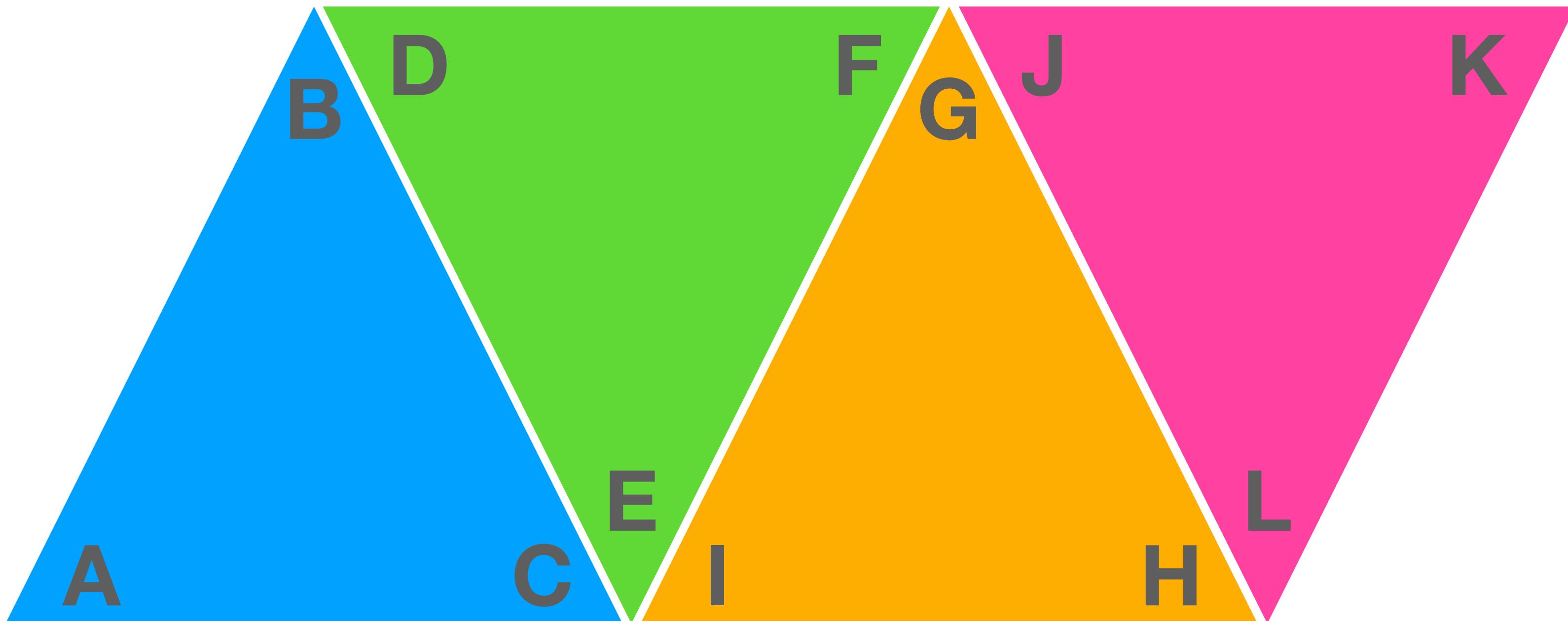


Pack your data optimally

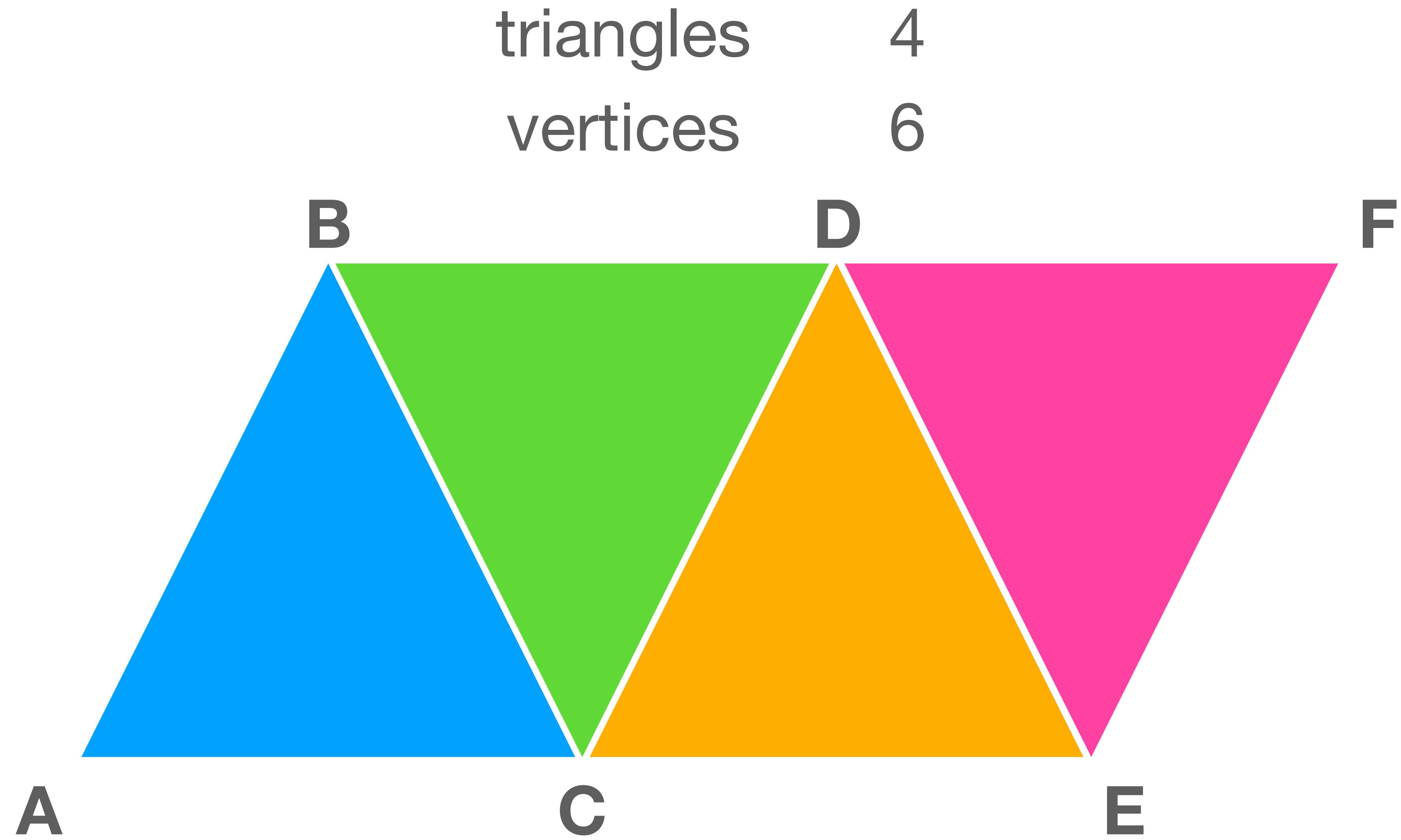


Pack your data optimally

triangles 4
vertices 12

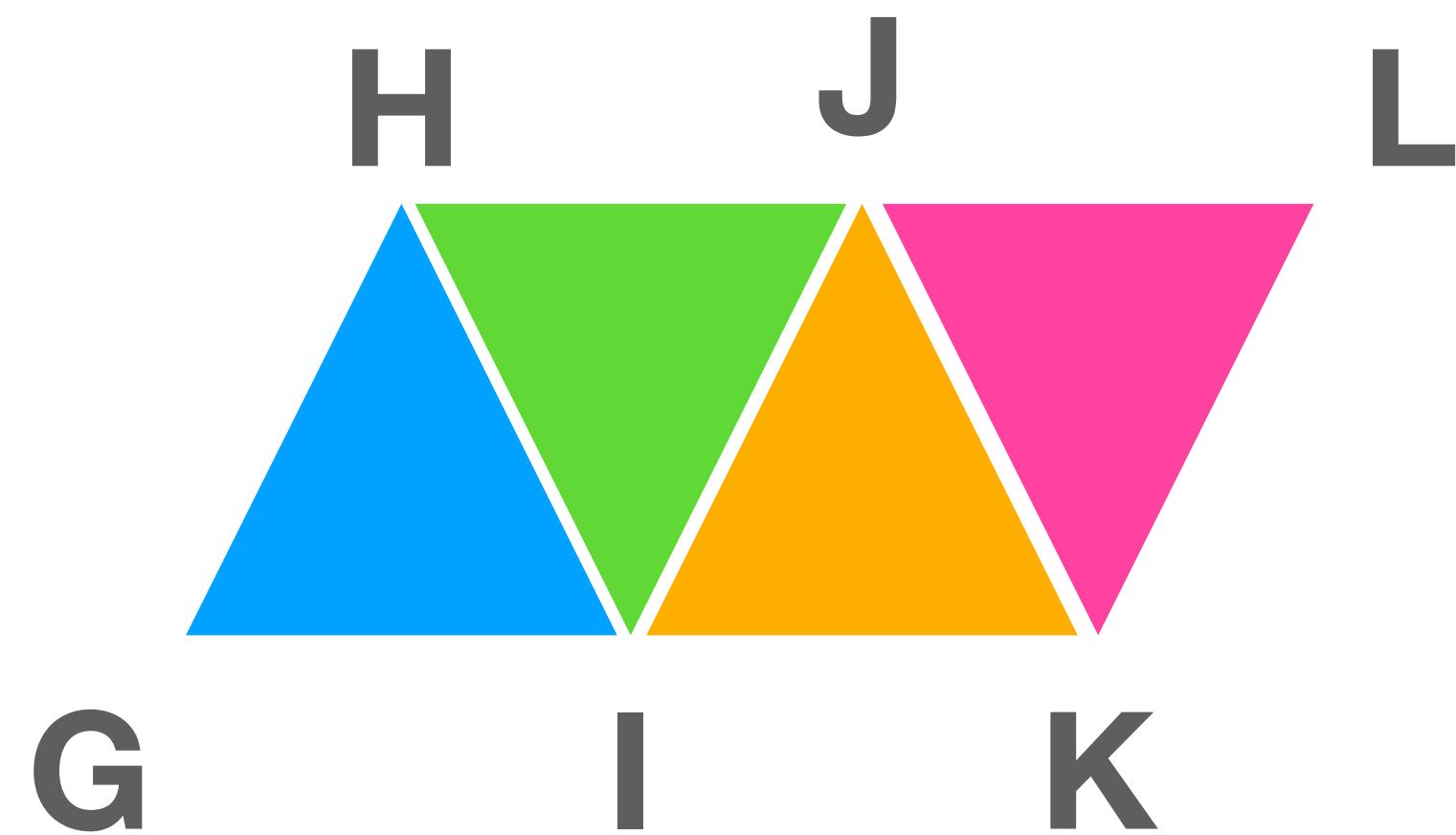
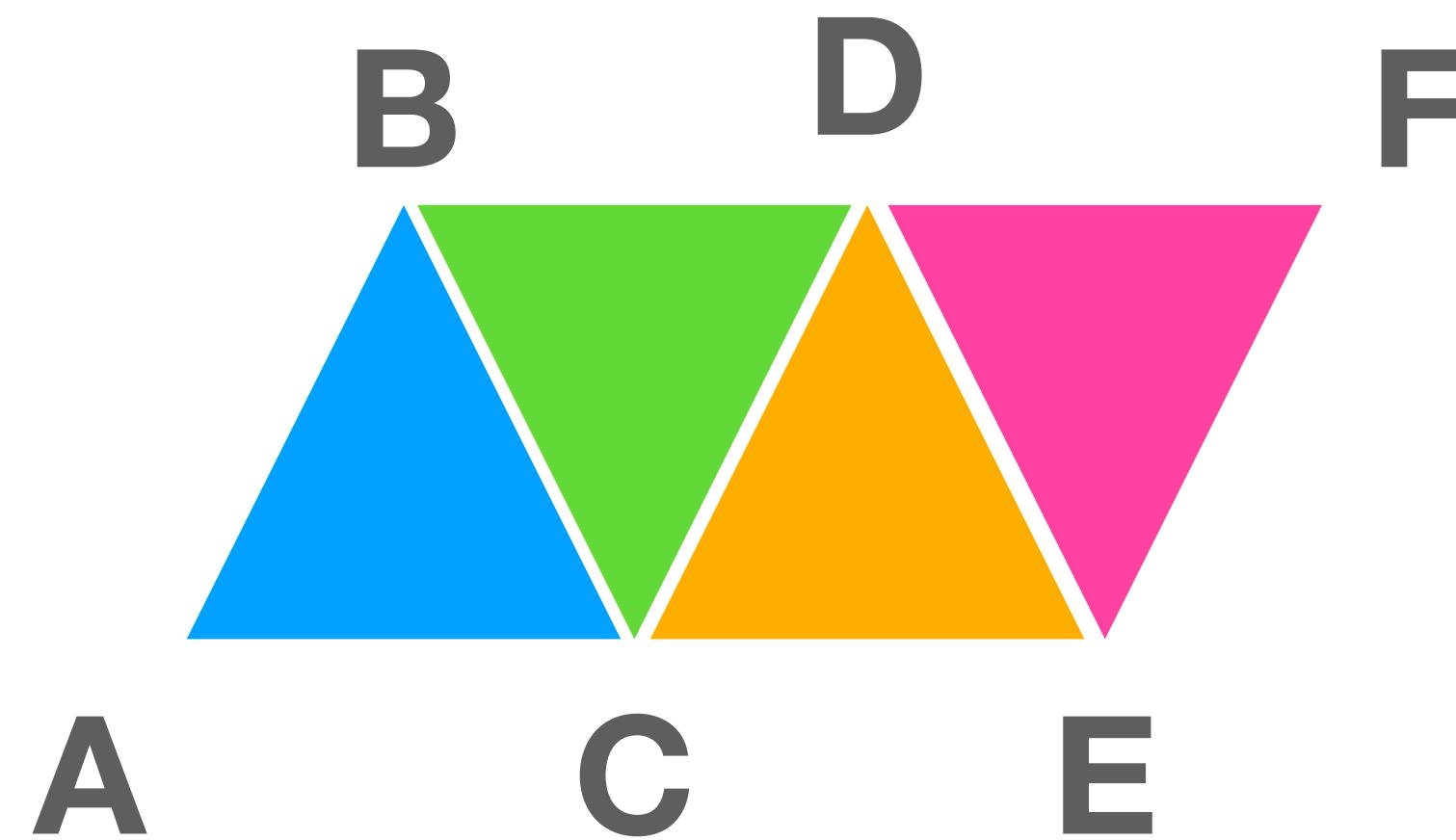


Pack your data optimally



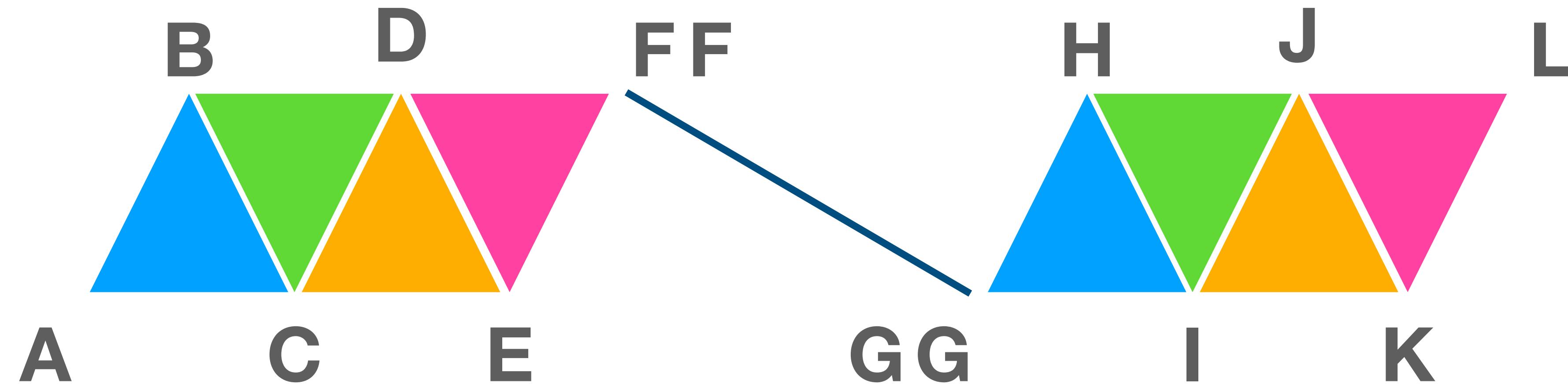
Pack your data optimally

draw calls	2
triangles	8
vertices	12



Pack your data optimally

draw calls	1
triangles	8
vertices	14

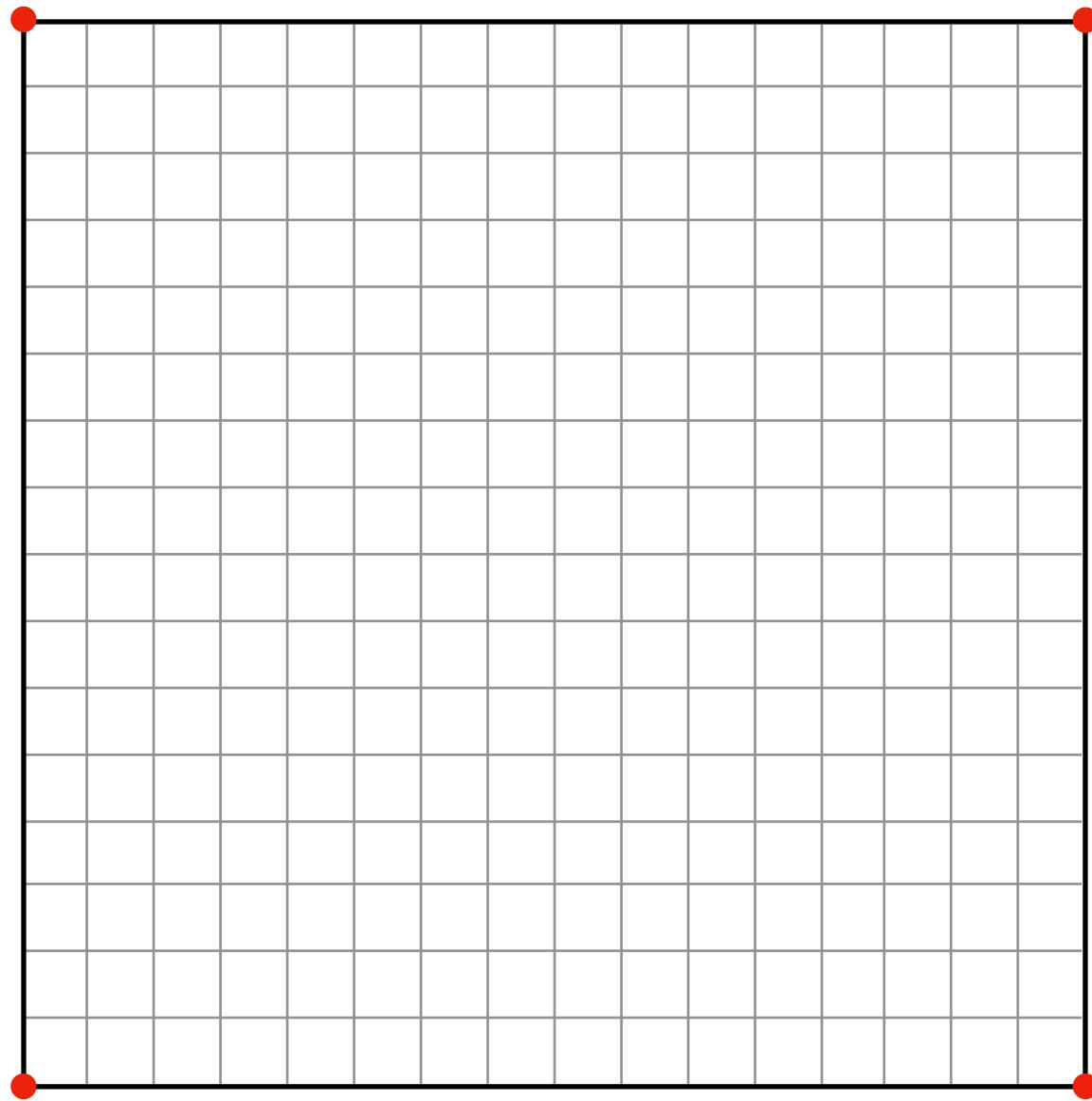


Parallel Computing

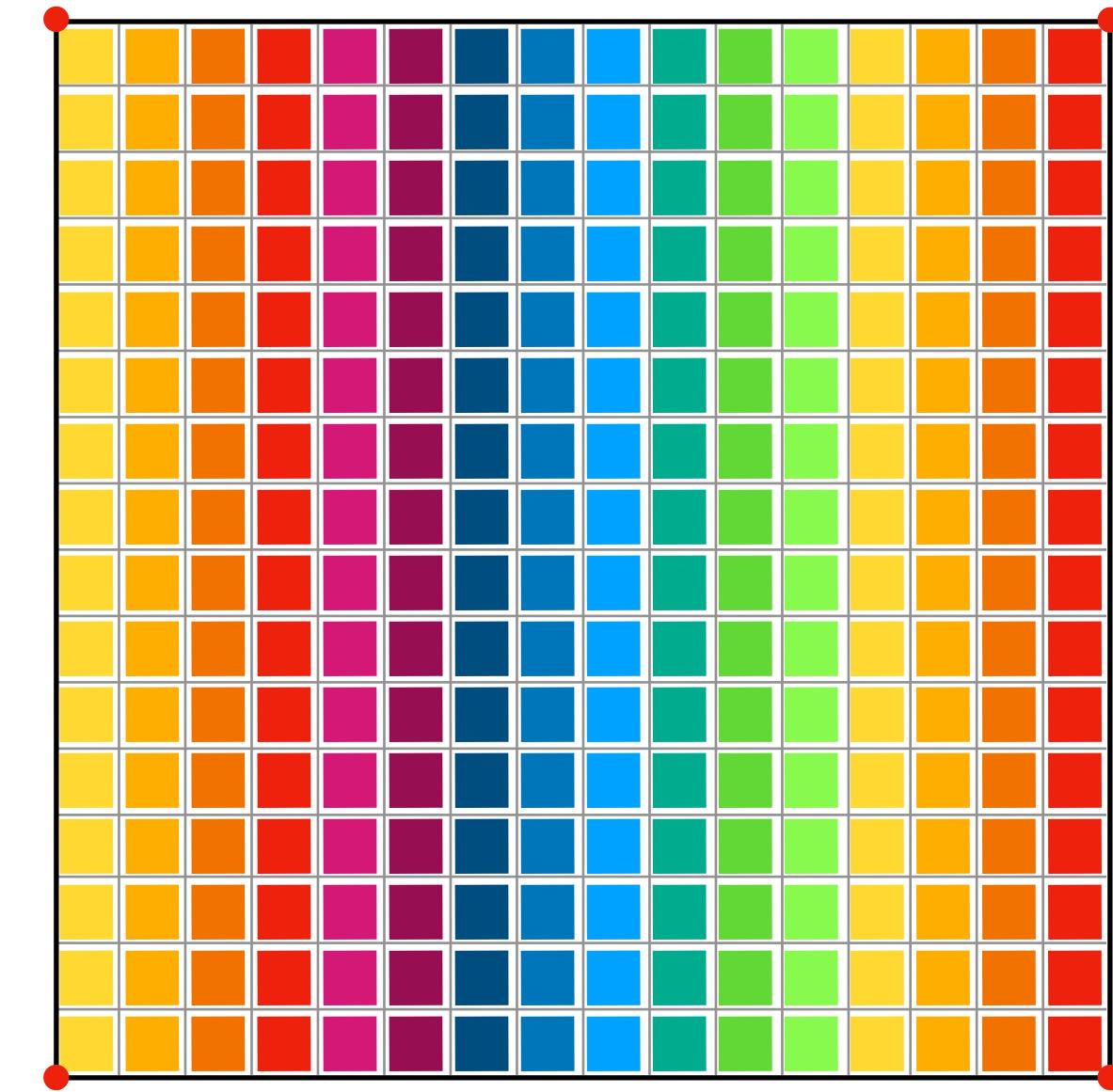
Compute - not render

Use compute for pixel to pixel processing

Vertex



Fragment



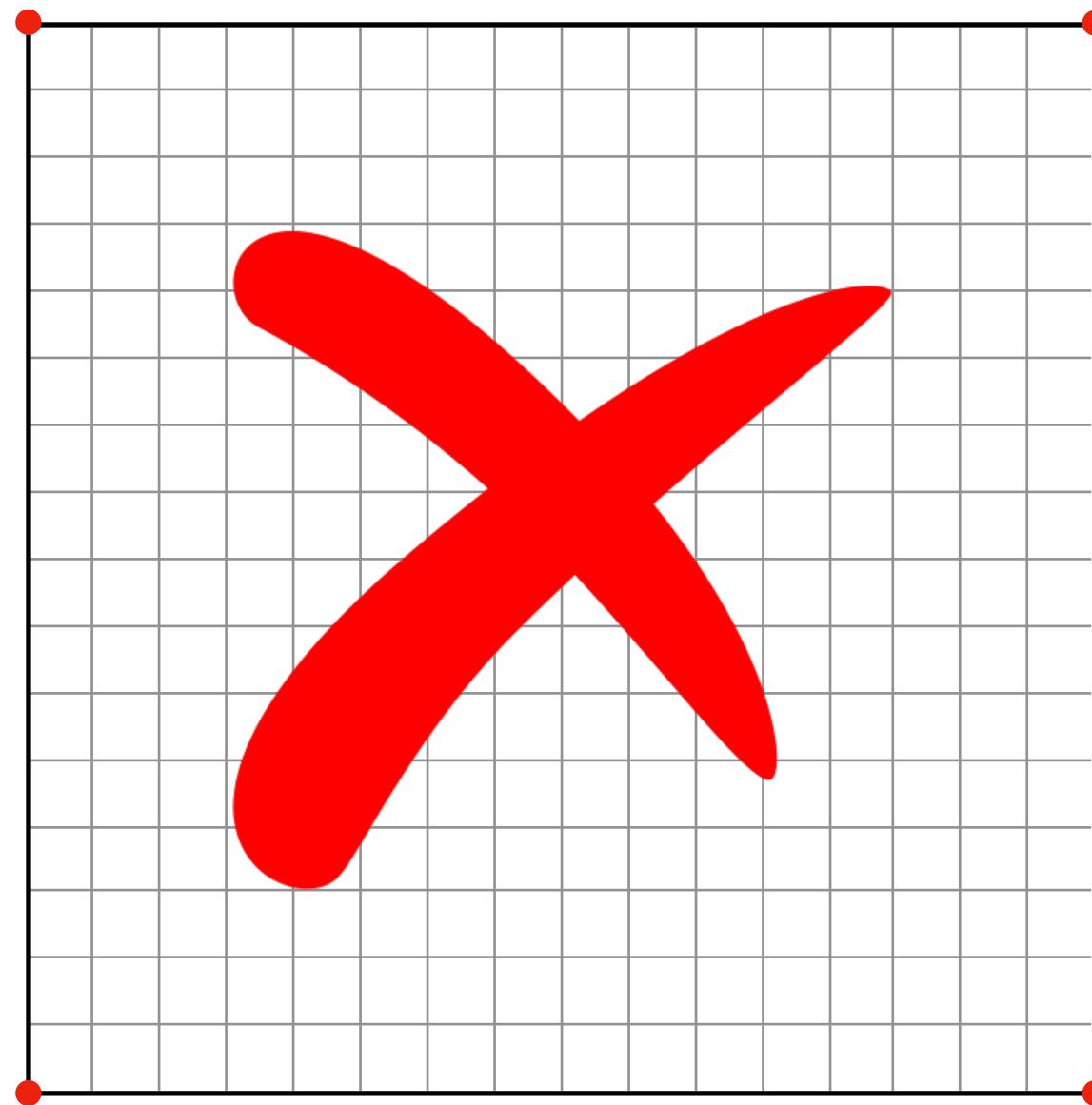
Rasterisation



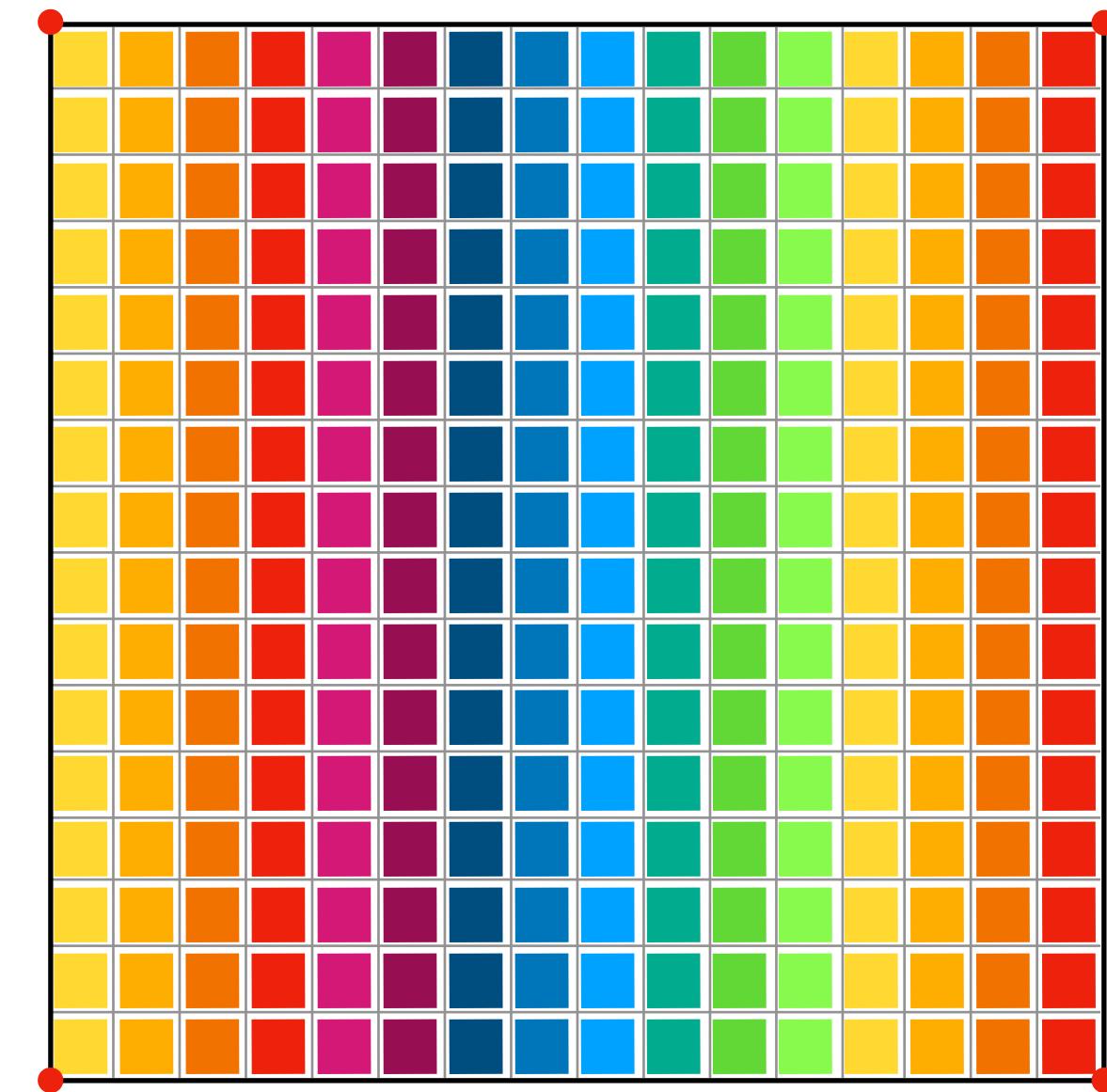
Compute - not render

Use compute for pixel to pixel processing

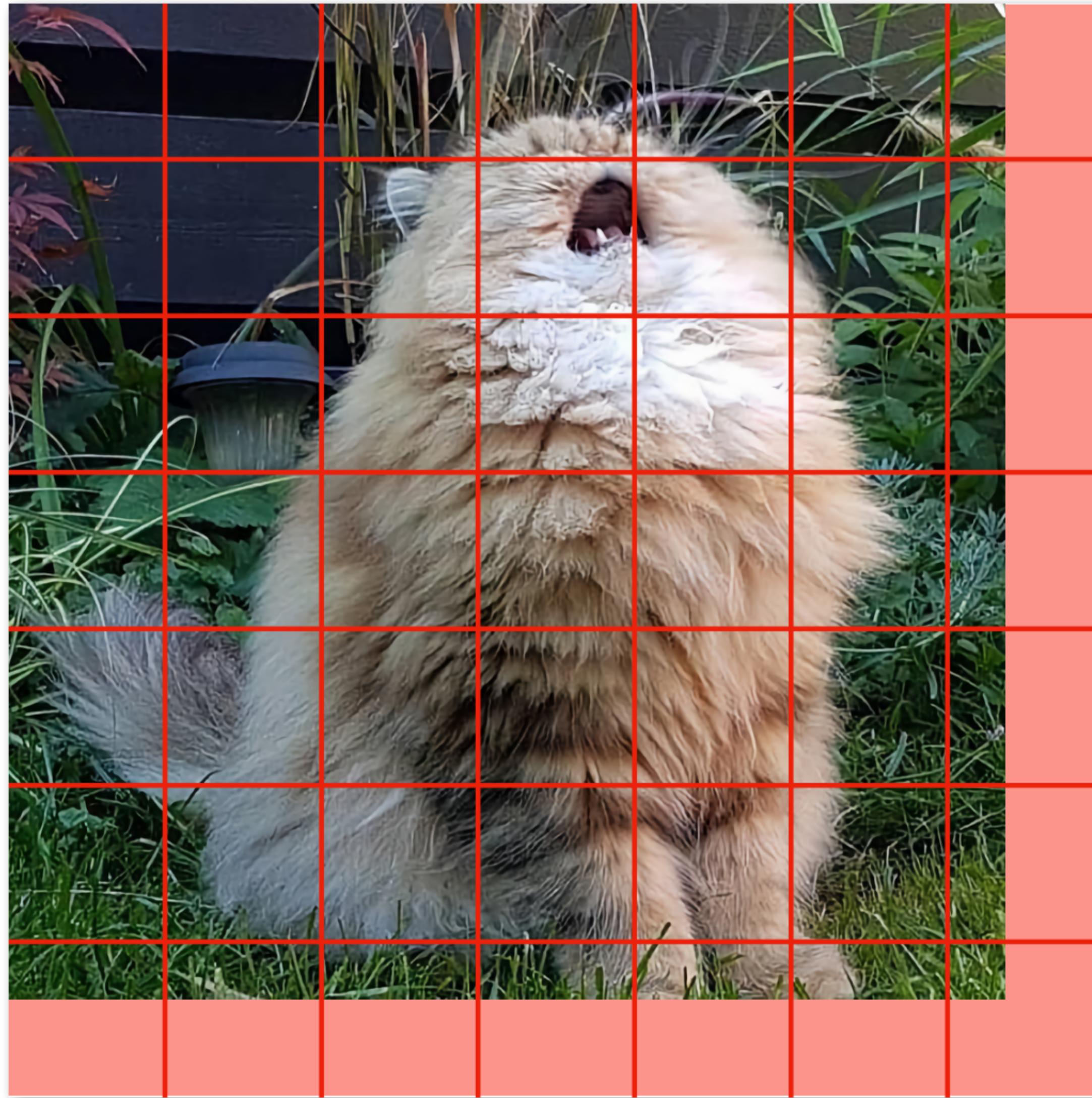
Vertex



Compute



dispatchThreadgroups

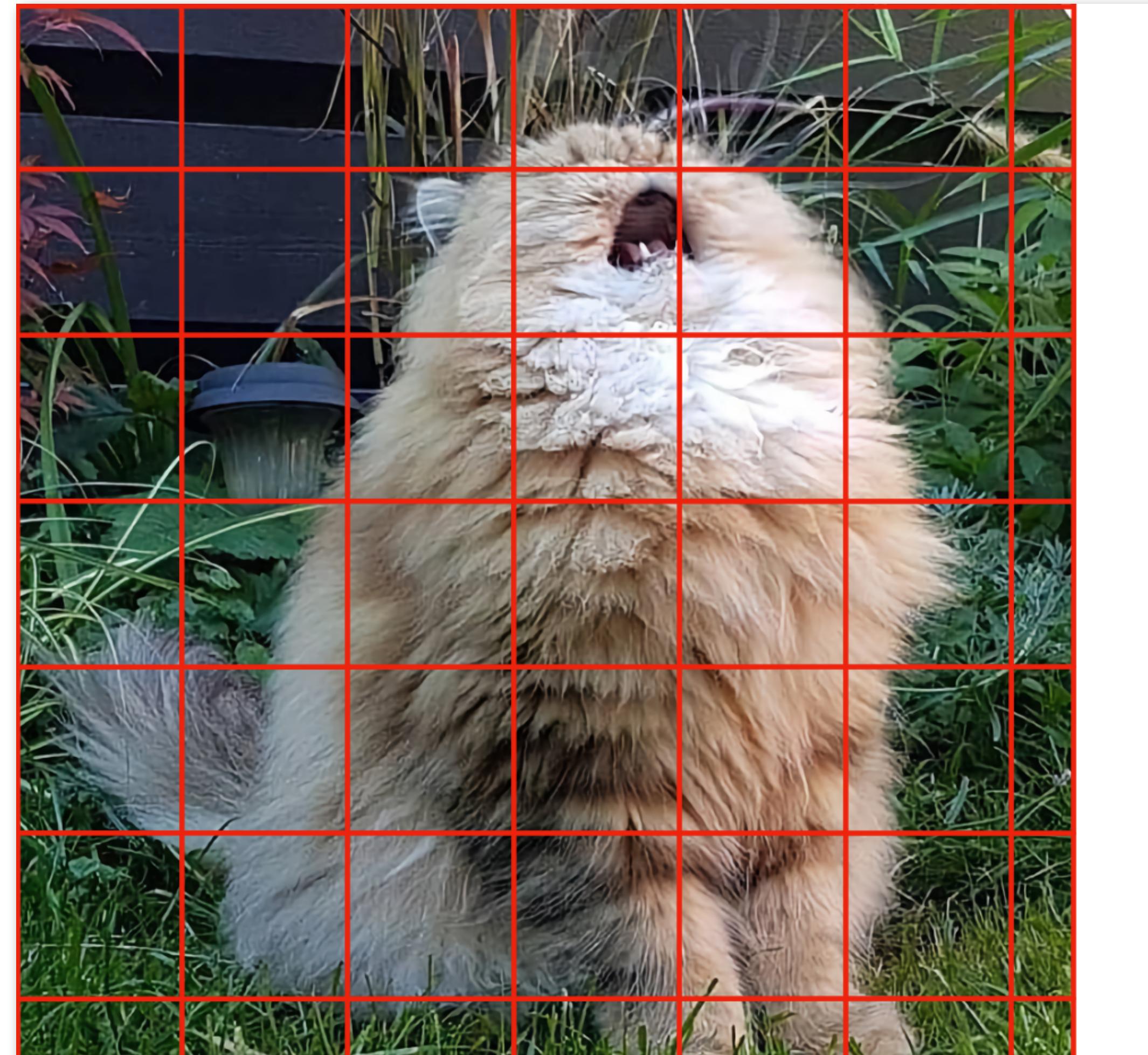


dispatchThreadgroups

```
kernel void krnFunction(texture2d<float, access::read> in [[ texture(0) ]],  
                        texture2d<float, access::write> out [[ texture(1) ]],  
                        uint2 gid [[thread_position_in_grid]])  
{  
    uint2 buf_size = uint2(in.get_width(), in.get_height());  
    if (any(gid >= buf_size)) {  
        return;  
    }  
  
    float4 src = in.read(gid);  
    out.write(src * src, gid);  
}
```

dispatchThreads

since A11 only, but...

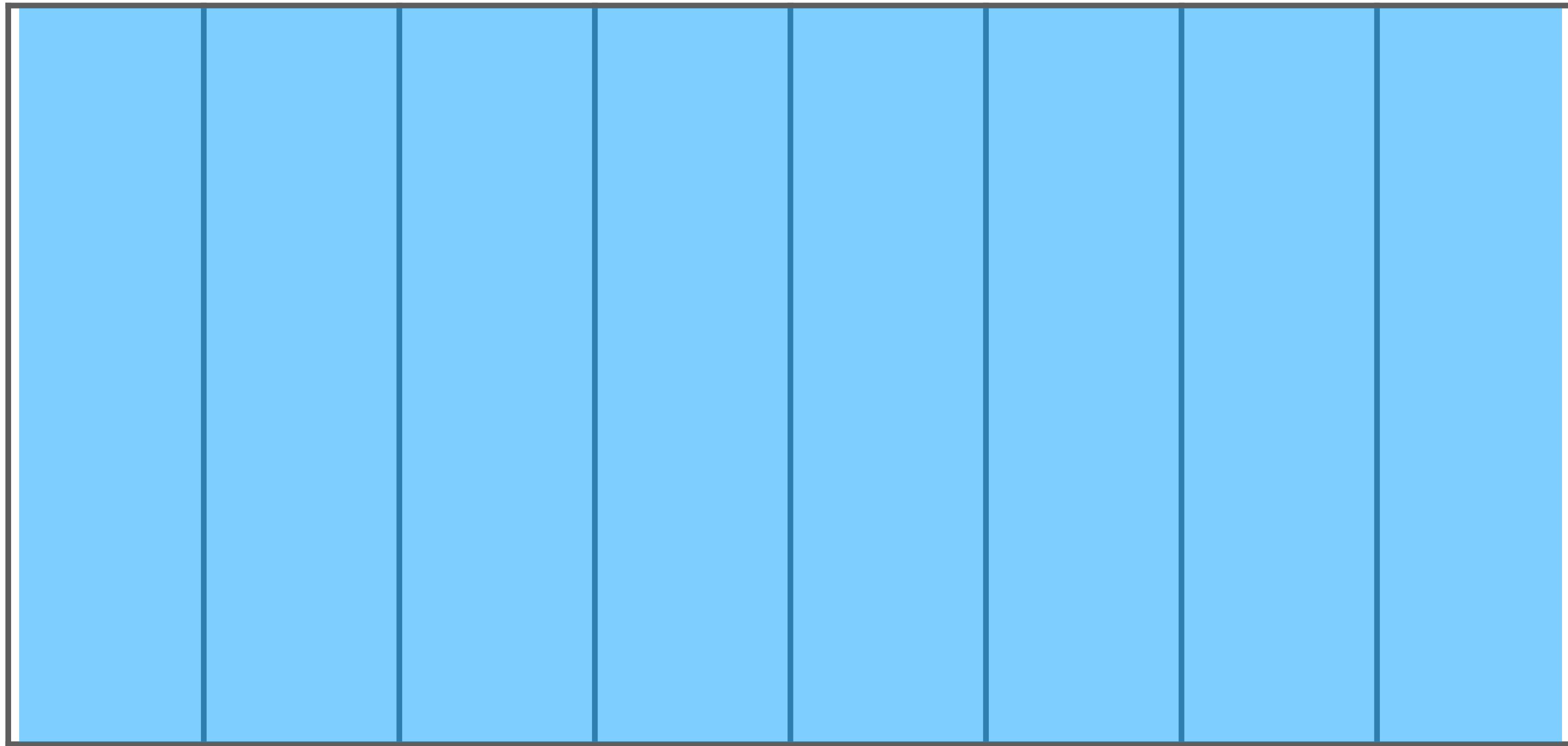


Size of thread group

```
NSUInteger w = _pipelineCompute.threadExecutionWidth;  
NSUInteger h = _pipelineCompute.maxTotalThreadsPerThreadgroup / w;
```

Size of thread group

The best 



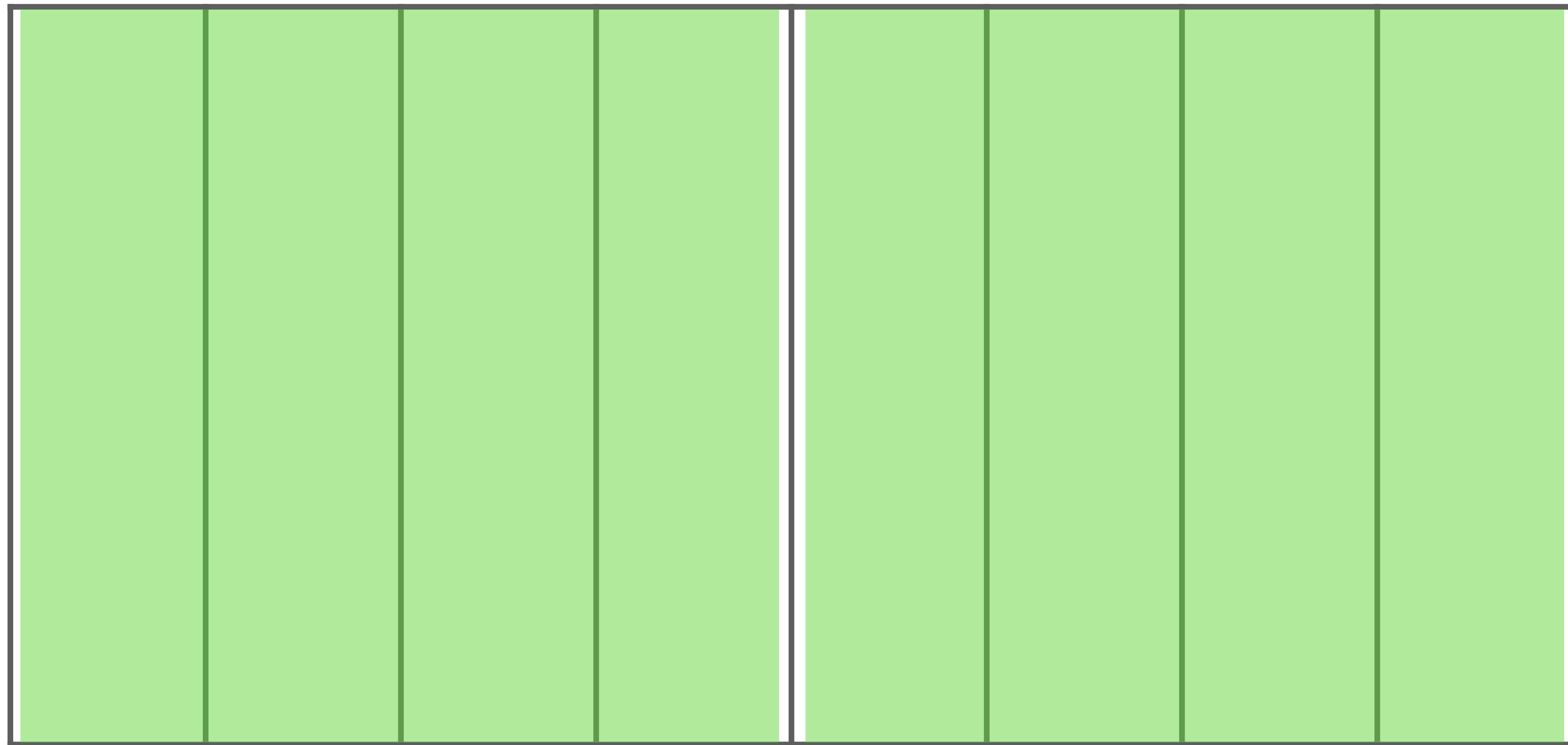
Size of thread group

The best 



Size of thread group

OK 🤷



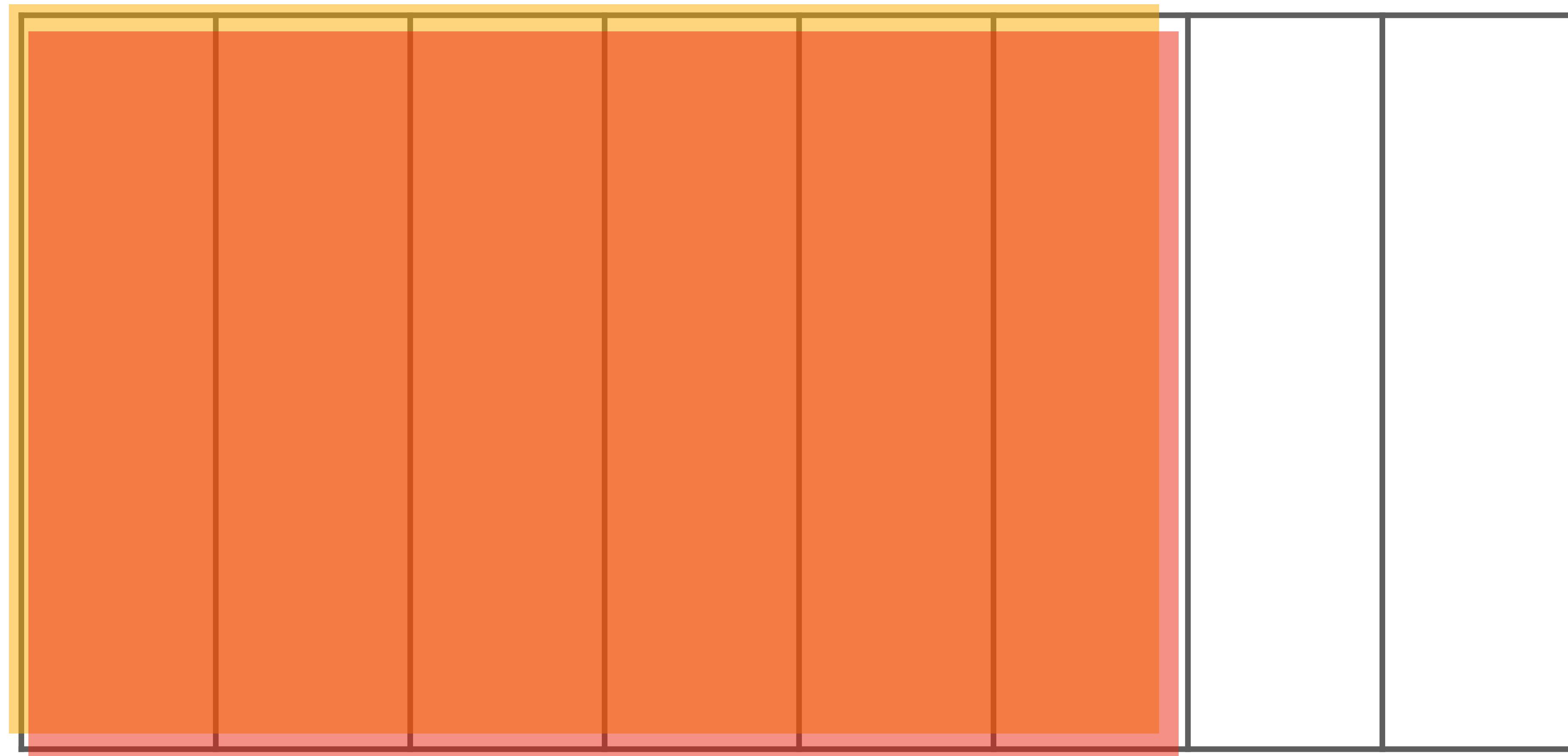
Size of thread group

OK 🤷



Size of thread group

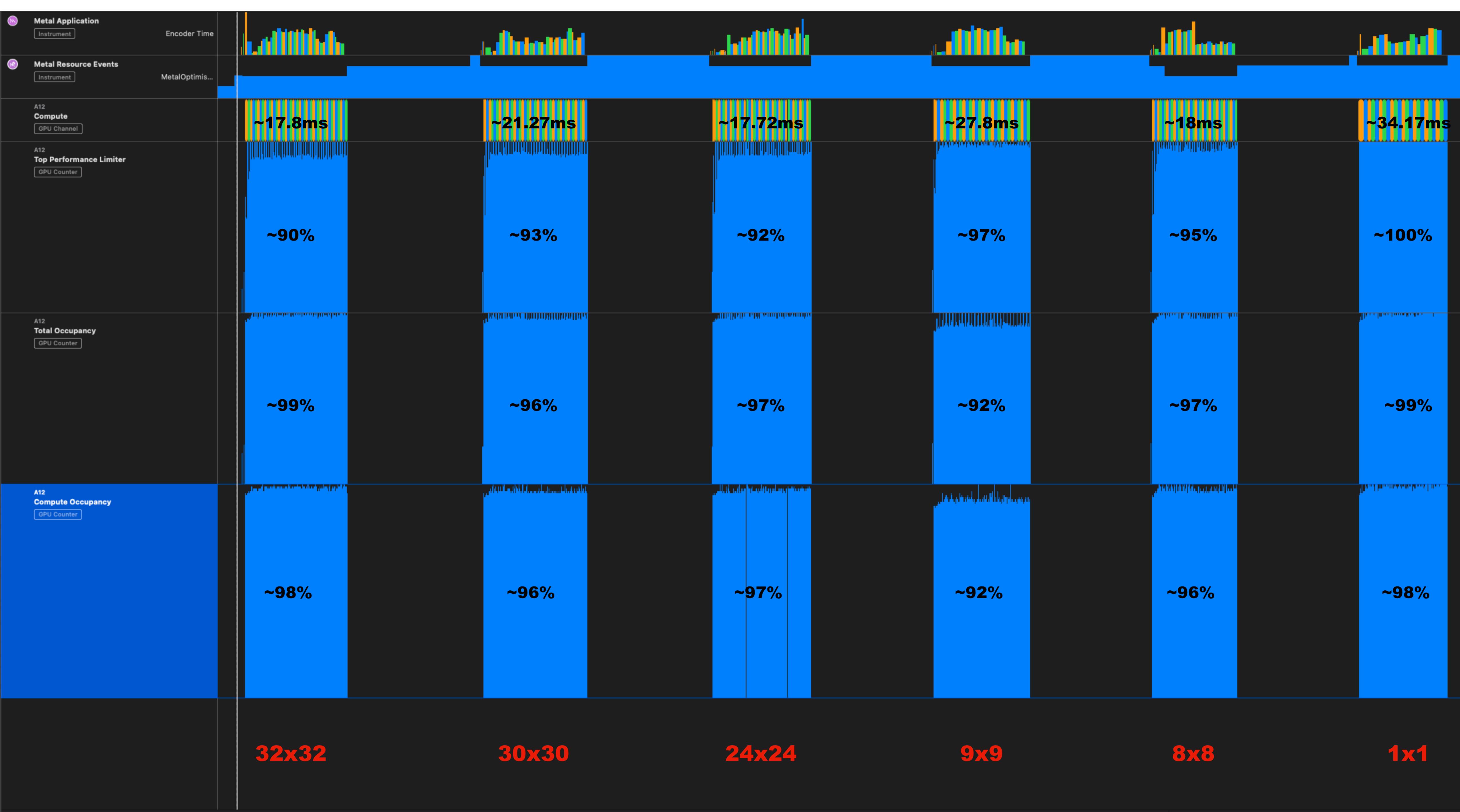
NO 



Size of thread group

NO 





A11 and TBDR

A11 and TBDR

GPU Family 4

- Asynchronous Vertex and Fragment passes
- Faster access to tile's memory
- Storing memoryless buffers in tile's memory

[WWDC + Example](#)

Tips & Tricks

Cache

everything



Swift vs C



World of C



World of Swift

Swift vs C

- SIMD
- pointers
- determinacy

Structures

padding

```
typedef struct {
    simd_float4 rect;
    bool flag;
    short index;
    simd_float3 size;
} StructureBadPadding;

//...

sizeof(StructureBadPadding)
```

```
typedef struct {
    simd_float4 rect;
    bool flag;
    short index;
    simd_float3 size;
    float f;
} StructureBadPadding;

//...

sizeof(StructureBadPadding); // 64
```

Structures

padding

```
typedef struct {
    simd_float4 rect;
    bool flag;
    // +1
    short index;
    // +12
    simd_float3 size;
    float f;
    // +12
} StructureBadPadding;
```

```
StructureBadPadding: 64
rect : 0  (16)
flag : 16 ( 1)
index: 18 ( 2)
size : 32 (16)
f     : 48 ( 4)
```

Structures

padding

```
typedef struct {
    simd_float4 rect;
    simd_float3 size;
    short index;
    bool flag;
} StructureGoodPadding;

//...

sizeof(StructureGoodPaddi
```

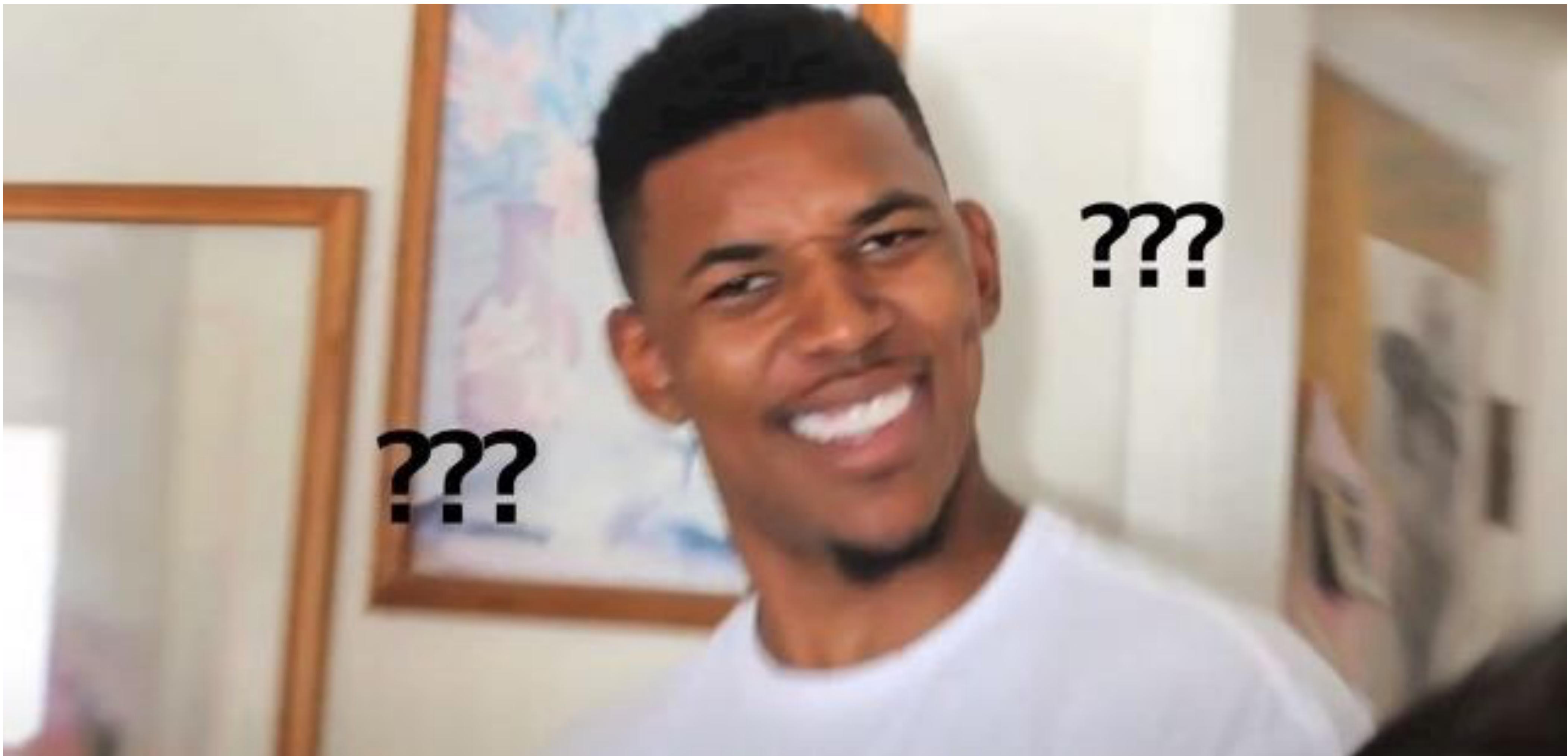
```
typedef struct {
    simd_float4 rect;
    simd_float3 size;
    float f;
    short index;
    bool flag;
} StructureGoodPadding;

//...

sizeof(StructureGoodPadding); // 48
```

Structures

padding



Structures

padding

```
typedef struct {
    simd_float4 rect;
    simd_float3 size;
    float f;
    short index;
    bool flag;
    // +9
} StructureGoodPadding;
```

```
StructureGoodPadding: 48
rect : 0  (16)
size : 16 (16)
f    : 32 ( 4)
index: 36 ( 2)
flag : 38 ( 1)
```

Constants

```
1      // int  
1.0    // double  
1.0f   // float
```

Constants

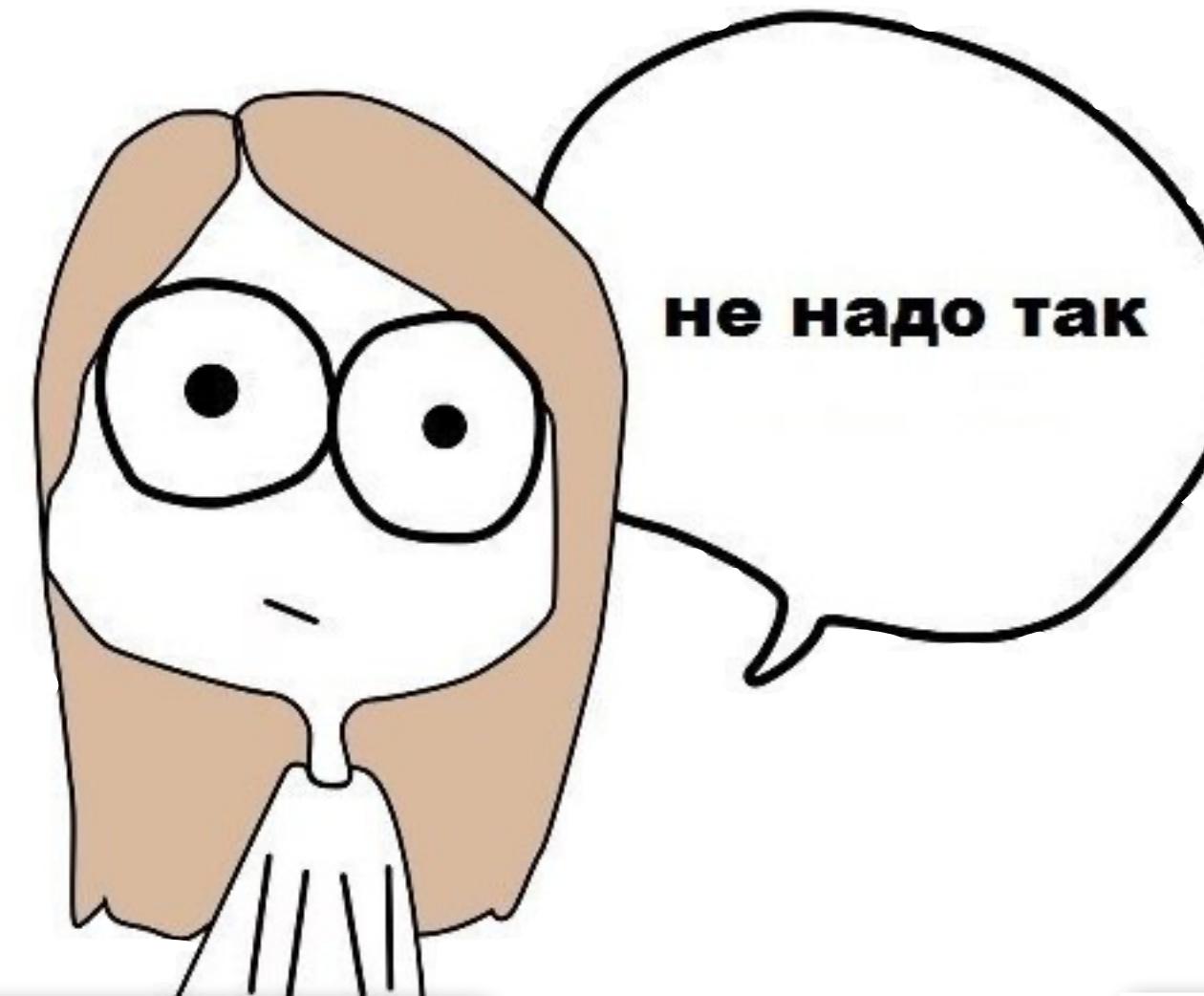
```
float foo(float a, float b) {
    return 10.0f * a * a + 5.0f * b * b;
}

foo(float, float):
    push {r4, r5, r11, lr}
    mov r4, r1
    mov r1, r0
    bl __aeabi_fmul
    mov r1, #18874368
    orr r1, r1, #1073741824
    bl __aeabi_fmul
    mov r5, r0
    mov r0, r4
    mov r1, r4
    bl __aeabi_fmul
    mov r1, #10485760
    orr r1, r1, #1073741824
    bl __aeabi_fmul
    mov r1, r5
    bl __aeabi_fadd
    pop {r4, r5, r11, lr}
    bx lr
```

```
float foo(float a, float b) {
    return 10.0 * a * a + 5.0 * b * b;
}

foo(float, float):
    push {r4, r5, r6, lr}
    mov r4, r1
    bl __aeabi_f2d
    mov r2, r0
    mov r3, r1
    bl __aeabi_dmul
    mov r3, #2359296
    mov r2, #0
    orr r3, r3, #1073741824
    bl __aeabi_dmul
    mov r5, r0
    mov r0, r4
    mov r6, r1
    bl __aeabi_f2d
    mov r2, r0
    mov r3, r1
    bl __aeabi_dmul
    mov r3, #1310720
    mov r2, #0
    orr r3, r3, #1073741824
    bl __aeabi_dmul
    mov r2, r5
    mov r3, r6
    bl __aeabi_dadd
    bl __aeabi_d2f
    pop {r4, r5, r6, lr}
    bx lr
```

Constants



```
float foo(float a, float b) {  
    return 10.0 * a * a + 5.0 * b * b;  
}
```

```
float foo(float a, float b) {  
    return 10.0f * a * a + 5.0f * b * b;  
}
```

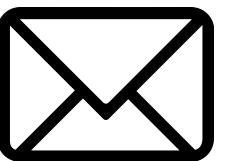


Links

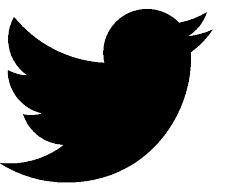
- Optimizing Performance with the Shader Profiler
- Metal System Trace
- Metal Best Practices
- Tile-Based Deferred Rendering (TBDR)
- WWDC Videos
- Demo app



Contacts



wdfeffects@gmail.com



[@GOstrobrod](#)



[@GOstrobrod](#)



<https://wdfteam.gitlab.io>