



## Harnessing Wave Intrinsic For Good (And Evil)

Alexandre Sabourin

Sharing is caring. Using wave  
intrinsic so that your threads  
aren't lonely.



 **Graphics  
Programming  
Conference**

 **Breda  
University**  
OF APPLIED SCIENCES

Hi, I'm Alex. One of the 6 Alexes at Snowed In Studios.

I'd like to thank the organizers of the conference for having me and for all of you for being here.

Today, I'll be talking about wave intrinsic. How they can be implemented and how you can use them to enable strange and powerful algorithms.

# What Will You Get Out Of This Talk?

- An intuition for how wave intrinsics are implemented
- An intuition for when you want to use wave intrinsics
- Examples of common and novel use cases
- Ideas on how you can develop your own wave-based algorithms

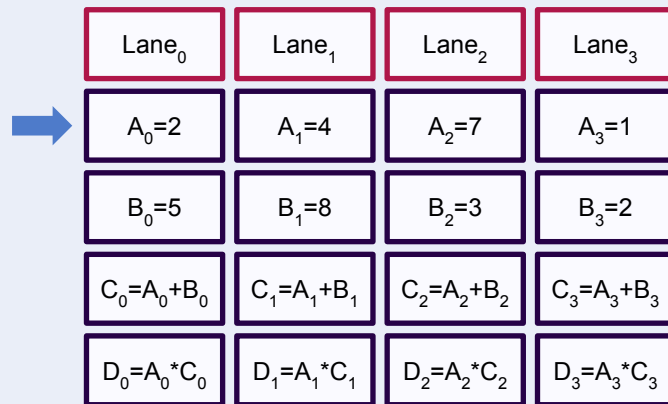
As a part of this talk, I'm hoping that you will be able to develop an intuition for how existing wave intrinsics are implemented and when you want to use them.

On top of that, we will be looking at common and novel use cases of wave intrinsics.

# Let's Meet Our GPU



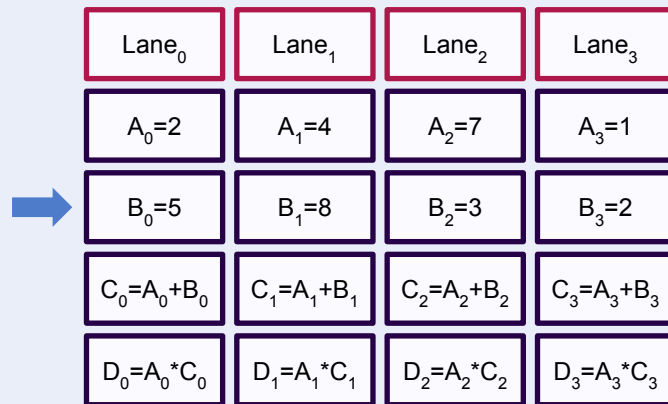
Let's meet the imaginary GPU we're going to be using for our demonstration today.



(Click through while talking) This mock GPU has 4 lanes each processing an independent stream of information.

It can do everything else a GPU can.

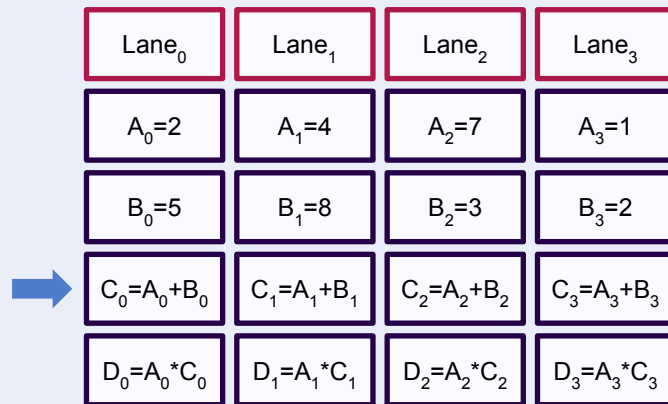
It just has fewer lanes than you would expect!



(Click through while talking) This mock GPU has 4 lanes each processing an independent stream of information.

It can do everything else a GPU can.

It just has fewer lanes than you would expect!



(Click through while talking) This mock GPU has 4 lanes each processing an independent stream of information.

It can do everything else a GPU can.

It just has fewer lanes than you would expect!

Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
A <sub>0</sub> =2	A <sub>1</sub> =4	A <sub>2</sub> =7	A <sub>3</sub> =1
B <sub>0</sub> =5	B <sub>1</sub> =8	B <sub>2</sub> =3	B <sub>3</sub> =2
C <sub>0</sub> =A <sub>0</sub> +B <sub>0</sub>	C <sub>1</sub> =A <sub>1</sub> +B <sub>1</sub>	C <sub>2</sub> =A <sub>2</sub> +B <sub>2</sub>	C <sub>3</sub> =A <sub>3</sub> +B <sub>3</sub>
→ D <sub>0</sub> =A <sub>0</sub> *C <sub>0</sub>	D <sub>1</sub> =A <sub>1</sub> *C <sub>1</sub>	D <sub>2</sub> =A <sub>2</sub> *C <sub>2</sub>	D <sub>3</sub> =A <sub>3</sub> *C <sub>3</sub>

(Click through while talking) This mock GPU has 4 lanes each processing an independent stream of information.

It can do everything else a GPU can.

It just has fewer lanes than you would expect!

# Sharing Is Caring

An Overview Of Common Intrinsic



What if we want to share across these lanes?

What if being completely independent isn't good enough?

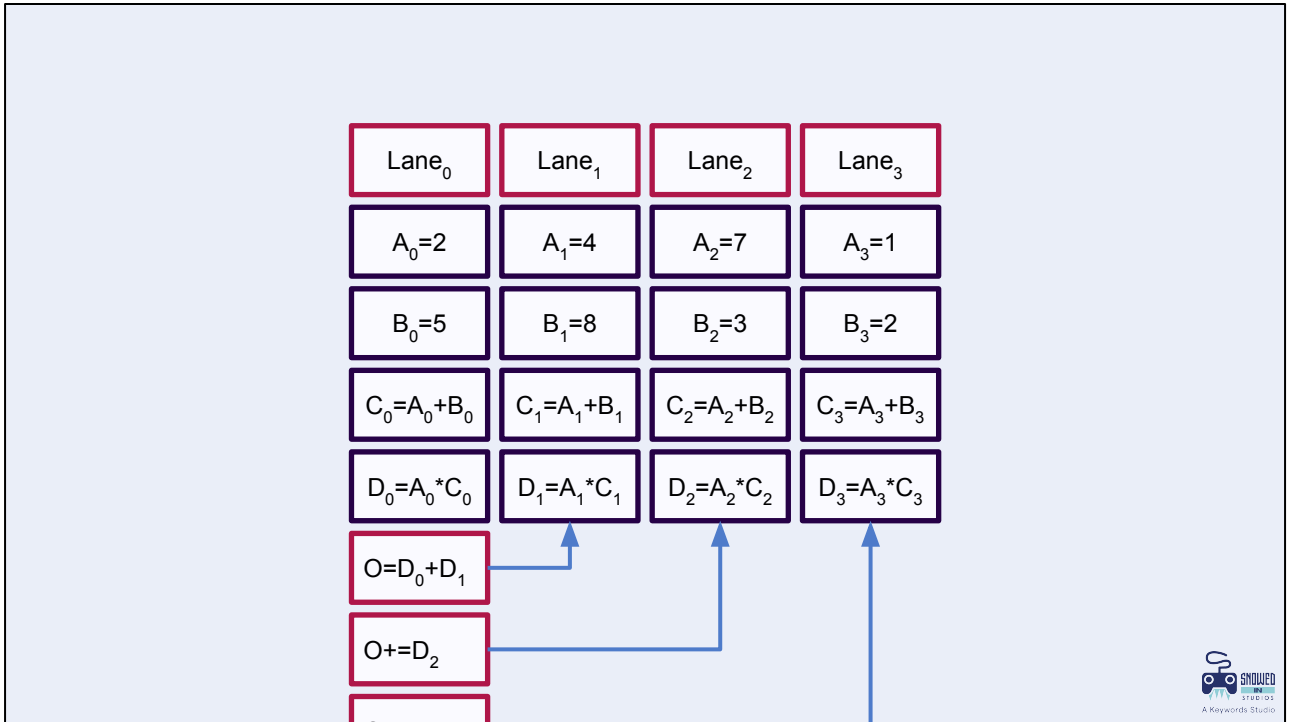


Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
A <sub>0</sub> =2	A <sub>1</sub> =4	A <sub>2</sub> =7	A <sub>3</sub> =1
B <sub>0</sub> =5	B <sub>1</sub> =8	B <sub>2</sub> =3	B <sub>3</sub> =2
C <sub>0</sub> =A <sub>0</sub> +B <sub>0</sub>	C <sub>1</sub> =A <sub>1</sub> +B <sub>1</sub>	C <sub>2</sub> =A <sub>2</sub> +B <sub>2</sub>	C <sub>3</sub> =A <sub>3</sub> +B <sub>3</sub>
D <sub>0</sub> =A <sub>0</sub> *C <sub>0</sub>	D <sub>1</sub> =A <sub>1</sub> *C <sub>1</sub>	D <sub>2</sub> =A <sub>2</sub> *C <sub>2</sub>	D <sub>3</sub> =A <sub>3</sub> *C <sub>3</sub>
Output=D <sub>0</sub> +D <sub>1</sub> +D <sub>2</sub> +D <sub>3</sub>			

Let's say that at the end of our program, we want to accumulate the values of each of our lanes in a single value.

\*click\*

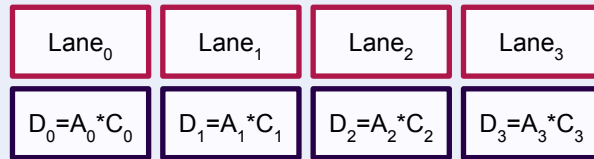
We could do that simply by iterating through each element in a single lane and adding them to our output.



Let's say that at the end of our program, we want to accumulate the values of each of our lanes in a single value.

\*click\*

We could do that simply by iterating through each element in a single lane and adding them to our output.



But there's a better way.

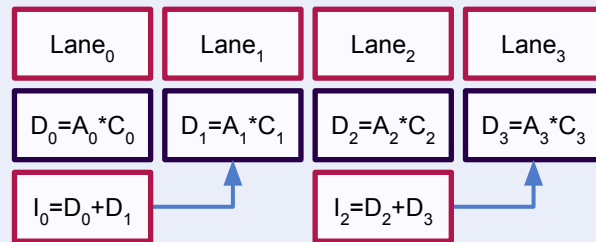
Instead, we can have some of our lanes participate in the final summation by having them accumulate the values of their neighbours.

Starting with their nearest neighbour.

\*click\*

And then doubling the distance to the next neighbour.

\*click\*



But there's a better way.

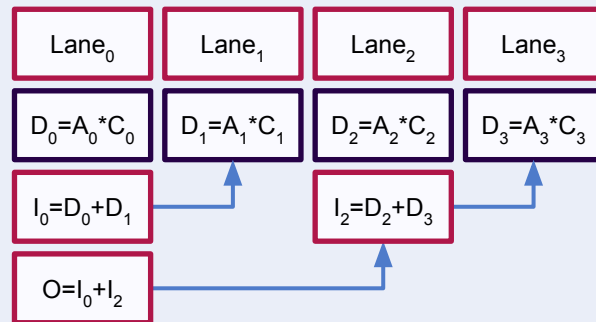
Instead, we can have some of our lanes participate in the final summation by having them accumulate the values of their neighbours.

Starting with their nearest neighbour.

\*click\*

And then doubling the distance to the next neighbour.

\*click\*



This summation ends at the second step but you can continue this process.

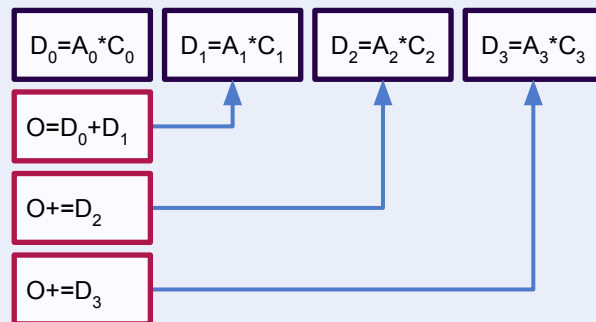
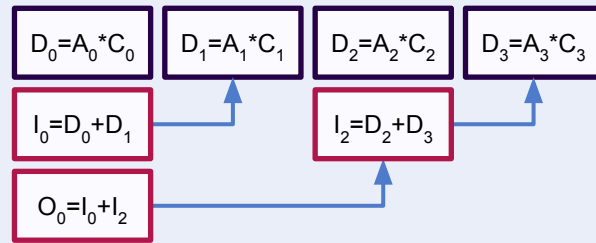
For 8 elements, the next stage would look at the neighbours at strides of 4.

For 16, an additional stage with a stride of 8.

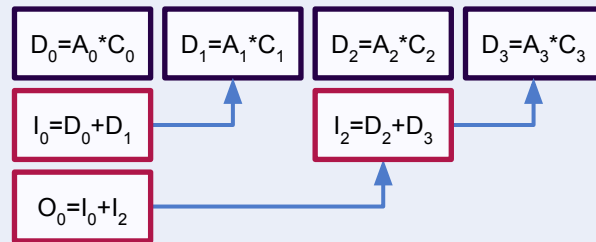
$\log_2(n)$  vs  $n$



We've turned a linear problem into a logarithmic one.



We save 1 iteration. In the case of a wave of size 32, we save 26.



WaveActiveSum

In HLSL, you may have accessed this operation using the `WaveActiveSum` intrinsic.

But why do we care about `WaveActiveSum` and its implementation?

Let's take the algorithm one step further and use a common example used in many games.

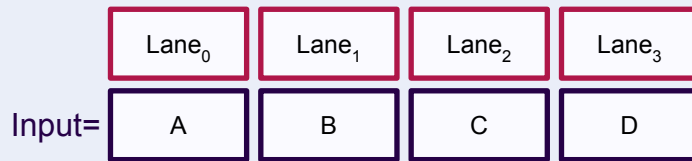


# Stream Compaction



Stream Compaction!

If you've done any form of GPU driven rendering, you almost definitely have used some form of stream compaction.



Let's look at a short example.

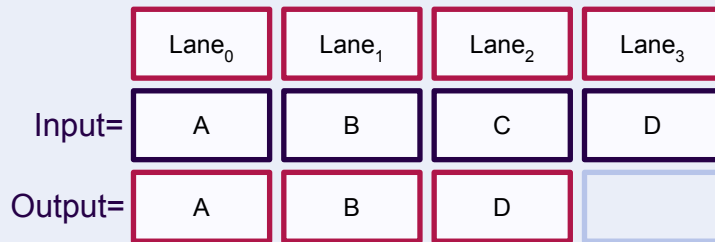
I have 4 letters.

But I want to write out only my favourite letters to an output buffer.

	Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
Input=	A	B	C	D
Output=	?	?	?	?

My favourite letters are A, B and D.

We don't like C.



We know that our output buffer should look like this.

But how do we get there?

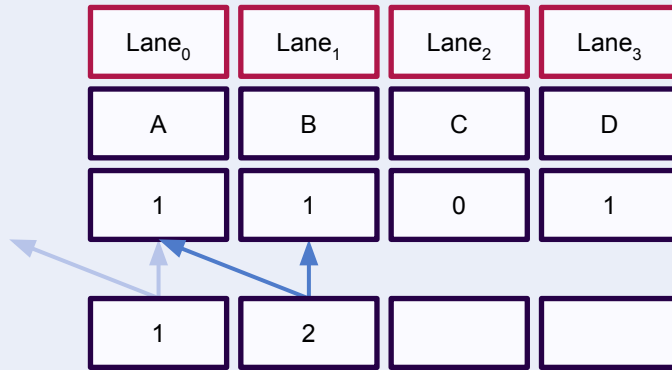
Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
A	B	C	D
1	1	0	1

Let's mark each lane with a 1 if its a favourite letter and a 0 if its not.

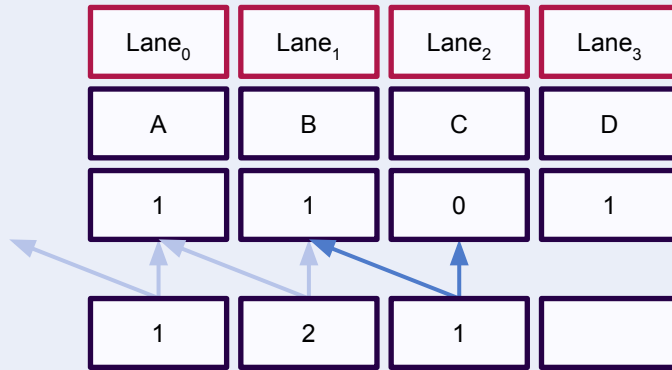
Maybe we just call an `isFavourite` function to get this result.



(Click through while talking) Starting from the same algorithm we used for WaveActiveSum, each neighbour reads and adds the value of its neighbour to itself.

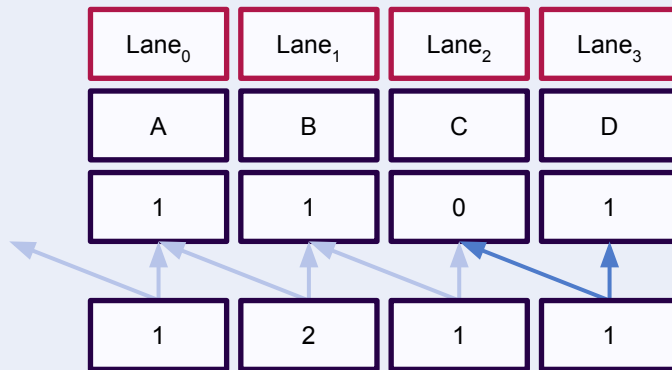


(Click through while talking) Starting from the same algorithm we used for WaveActiveSum, each neighbour reads and adds the value of its neighbour to itself.

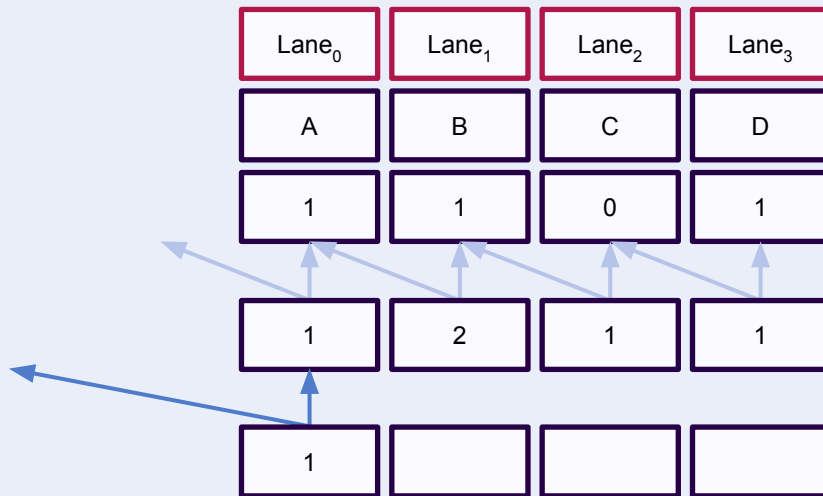


(Click through while talking) Starting from the same algorithm we used for WaveActiveSum, each neighbour reads and adds the value of its neighbour to itself.



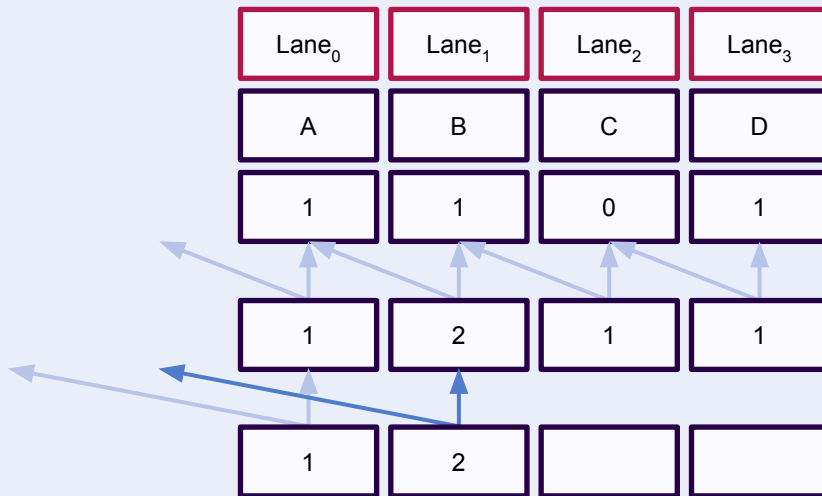


(Click through while talking) Starting from the same algorithm we used for WaveActiveSum, each neighbour reads and adds the value of its neighbour to itself.



(Click through while talking) We then accumulate the results of our neighbours at twice the distance.

You'll notice that each value is slowly accumulating the sum of all of its preceding neighbours.



(Click through while talking) We then accumulate the results of our neighbours at twice the distance.

You'll notice that each value is slowly accumulating the sum of all of its preceding neighbours.



(Click through while talking) We then accumulate the results of our neighbours at twice the distance.

You'll notice that each value is slowly accumulating the sum of all of its preceding neighbours.



(Click through while talking) We then accumulate the results of our neighbours at twice the distance.

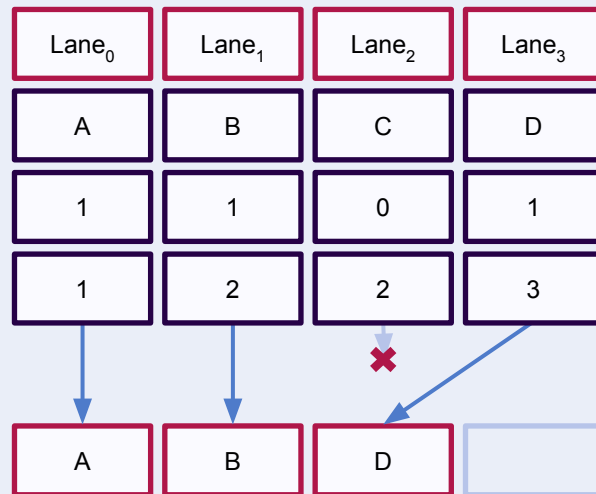
You'll notice that each value is slowly accumulating the sum of all of its preceding neighbours.

Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
A	B	C	D
1	1	0	1
1	2	2	3
?	?	?	?

Once we've accumulated the results of our neighbours, we know how many values will be written before our own.

\*click\*

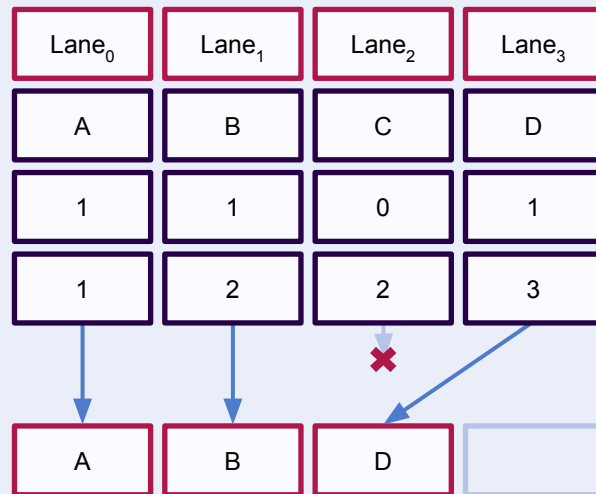
Then we subtract 1 from our sum and write to that index if it's one of my favourite letters.



Once we've accumulated the results of our neighbours, we know how many values will be written before our own.

\*click\*

Then we subtract 1 from our sum and write to that index if it's one of my favourite letters.



WavePrefixSum

To use this operation, you simply have to use `WavePrefixSum`.

Note to reader: We actually described an inclusive prefix sum which includes our current element, while WavePrefixSum is an exclusive prefix sum which excludes the current element.



**Incredibly  
Useful For  
Culling**



Stream compaction, and as a result a prefix sum, is widely used for a variety of problems. Especially for GPU culling of meshlets, lights and other primitives.

# Already Built For You!

- **WaveActiveSum->subgroupAdd**
- **WavePrefixSum->subgroupExclusiveAdd**
- **And more!**
  - <https://github.com/microsoft/directxshadercompiler/wiki/wave-intrinsics>
  - <https://www.khronos.org/blog/vulkan-subgroup-tutorial>

There are many more intrinsics for you to use and some excellent resources for discovering how to use them.

The references slide later on has a wide variety of excellent reading material.

# How Does It Compile?

(With AMD's Offline Compiler)



Now that we've had a bit of an overview. How do these wave intrinsics compile?

Knowing how these compile can provide you with valuable insight into the potential performance characteristics of your algorithm.

For the rest of this talk, disassembly will be from AMD's offline compiler.

Let's look at a simple program.

```

9   StructuredBuffer<int> Input;
10  RWStructuredBuffer<int> Output;
11
12
13  [numthreads(32,1,1)]
14  void CS()
15  {
16      int input = Input[WaveGetLaneIndex()];
17      bool selectInput = Select(input);
18      uint outputIndex = WavePrefixSum(selectInput ? 1 : 0);
19
20      if(selectInput)
21      {
22          Output[outputIndex] = input;
23      }
24  }

```

This program represents the algorithm we discussed a moment ago.

# WavePrefixSum

```
32 ;      uint outputIndex = WavePrefixSum(selectInput ? 1 : 0);
36 v_mov b32_e32 v1, v8 //
37 s_not_b64_exec, exec //
38 //
39 v_mov b32_e32 v1, 0 //
40 s_not_b64_exec, exec //
41 s_or_saveexec b64 s[4:5], -1 //
42 //
43 s_mov b32_s0, 0x6543210f //
44 v_lshlrev b64 v[2:3], v0, 1 //
45 v_permianex16 b32 v1, v1, s0, 0xedcba987 op_sel:[1,0] //
46 v_and b32_e32 v5, 0x10001, v3 //
47 v_readlane b32 s0, v1, 16 //
48 //
49 v_and b32_e32 v4, 0x10001, v2 //
50 v_and b32_e32 v3, 0xffff0000, v3 //
51 v_and b32_e32 v2, 0xffff0000, v2 //
52 v_writelane b32 v1, s0, 48 //
53 v_cmp_eq u64 e64 s0, 0, v[4:5] //
54 //
55 v_mov b32_e32 v4, 0 //
56 v_mov b32_e32 v5, 0 //
57 //
58 v_writelane b32 v1, 0, 16 //
59 v_permianex16 b32 v6, v1, 0, -1 op_sel:[1,0] //
60 v_cndmask b32_e64 v1, v6, v1, s0 //
61 v_cmp_ne u64 e64 s0, 0, v[2:3] //
62 v_add_nc u32_dpp v6, v1, v1 row_shr:1 row_mask:0xf bank_mask:0 //
63 v_mov b32_dpp v4, v1 row_shr:2 row_mask:0xf bank_mask:0xf //
64 v_mov b32_dpp v5, v1 row_shr:3 row_mask:0xf bank_mask:0xf //
65 v_add b32 v1, v6, v4, v5 //
66 v_add_nc u32_dpp v1, v1, v1 row_shr:4 row_mask:0xf bank_mask:0 //
67 v_add_nc u32_dpp v1, v1, v1 row_shr:8 row_mask:0xf bank_mask:0 //
68 //
69 v_permianex16 b32 v4, v1, -1, -1 op_sel:[1,0] //
70 v_cndmask b32_e64 v2, 0, v4, s0 //
71 //
72 v_cmp_lt u32_e64 s0, 31, v0 //
73 v_add_nc u32_e32 v1, v2, v1 //
74 v_readlane b32 s6, v1, 31 //
75 v_cndmask b32_e64 v0, 0, s6, s0 //
76 v_add_nc u32_e32 v0, v1, v0 //
77 s_mov b64_exec, s[4:5] //
78 //
79 v_mov b32_e32 v8, v0 //
```

**~40 Instructions**

If we look at the disassembly for `WavePrefixSum` you'll see that it's quite a bit of code!

Unfortunately, it's a little bit more complicated than "read your neighbour's value"...

# A Better Program



I'll admit...

I lied to you.

We don't need to use `WavePrefixSum` for this problem at all.

Since our values are 0 and 1, we can actually use the much leaner `WavePrefixCountBits`

```

9
10 StructuredBuffer<int> Input;
11 RWStructuredBuffer<int> Output;
12
13 [numthreads(32,1,1)]
14 void CS()
15 {
16     int input = Input[WaveGetLaneIndex()];
17     bool selectInput = Select(input);
18     uint outputIndex = WavePrefixCountBits(selectInput);
19
20     if(selectInput)
21     {
22         Output[outputIndex] = input;
23     }
24 }

```

`WavePrefixCountBits` will provide us the same results.

But on RDNA, compiles to a much leaner two instructions.

\*click\*

```

42 |
43 | ;      uint outputIndex = WavePrefixCountBits(selectInput);
44 |     v_mbcnt_lo_u32_b32 v1, vcc_lo, 0 //
45 |     v_mbcnt_hi_u32_b32 v1, vcc_hi, v1 //
46 |

```

**2 Instructions**

`WavePrefixCountBits` will provide us the same results.

But on RDNA, compiles to a much leaner two instructions.

\*click\*



```

36 ;      uint outputIndex = WavePrefixSum(selectInput ? 1 : 0);
37 v_mov b32_e32 v1, v8 //
38 s_not b64_exec, exec //
39 v_mov b32_e32 v1, 0 //
40 s_not b64_exec, exec //
41 s_of_saveexec_b64 s[4:5], -1 //
42 s_mov b32 s0, 0x6543210f //
43 v_lshlrev b64 v[2:3], v0, 1 //
44 v_permlanexl6 b32 v1, v1, s0, 0xedcba987 op_sel:[1,0] //
45 v_and b32_e32 v5, 0x10001, v3 //
46 v_readlane b32 s0, v1, 16 //
47 v_and b32_e32 v4, 0x10001, v2 //
48 v_and b32_e32 v3, 0xffff0000, v3 //
49 v_and b32_e32 v2, 0xffff0000, v2 //
50 v_writelane b32 v1, s0, 48 //
51 v_cmp eq_u64 e64 s0, 0, v[4:5] //
52 v_mov b32_e32 v4, 0 //
53 v_mov b32_e32 v5, 0 //
54 v_writelane b32 v1, 0, 16 //
55 v_permlanexl6 b32 v6, v1, 0, -1 op_sel:[1,0] //
56 v_cndmask b32_e64 v1, v6, v1, s0 //
57 v_cmp ne_u64 e64 s0, 0, v[2:3] //
58 v_add nc_u32_dpp v6, v1, v1 row_shr:1 row_mask:0xf bank_mask:0 //
59 v_mov b32_dpp v4, v1 row_shr:2 row_mask:0xf bank_mask:0xf //
60 v_mov b32_dpp v5, v1 row_shr:3 row_mask:0xf bank_mask:0xf //
61 v_add3 u32 v1, v6, v4, v5 //
62 v_add nc_u32_dpp v1, v1, v1 row_shr:4 row_mask:0xf bank_mask:0 //
63 v_add nc_u32_dpp v1, v1, v1 row_shr:8 row_mask:0xf bank_mask:0 //
64 v_permlanexl6 b32 v4, v1, -1, -1 op_sel:[1,0] //
65 v_cndmask b32_e64 v2, 0, v4, s0 //
66 v_cmp lt_u32 e64 s0, 31, v0 //
67 v_add nc_u32_e32 v1, v2, v1 //
68 v_readlane b32 s6, v1, 31 //
69 v_cndmask b32_e64 v0, 0, s6, s0 //
70 v_add nc_u32_e32 v0, v1, v0 //
71 s_mov b64_exec, s[4:5] //
72 v_mov b32_e32 v8, v0 //

```

~40 Instructions

VS

```

74 ;      uint outputIndex = WavePrefixCountBits(selectInput);
75 v_mbcnt_lo u32 b32 v1, vcc_lo, 0 //
76 v_mbcnt_hi u32 b32 v1, vcc_hi, v1 //

```

2 Instructions

I'll let you decide which one is likely to be faster...

# Wave-wide Interpolation

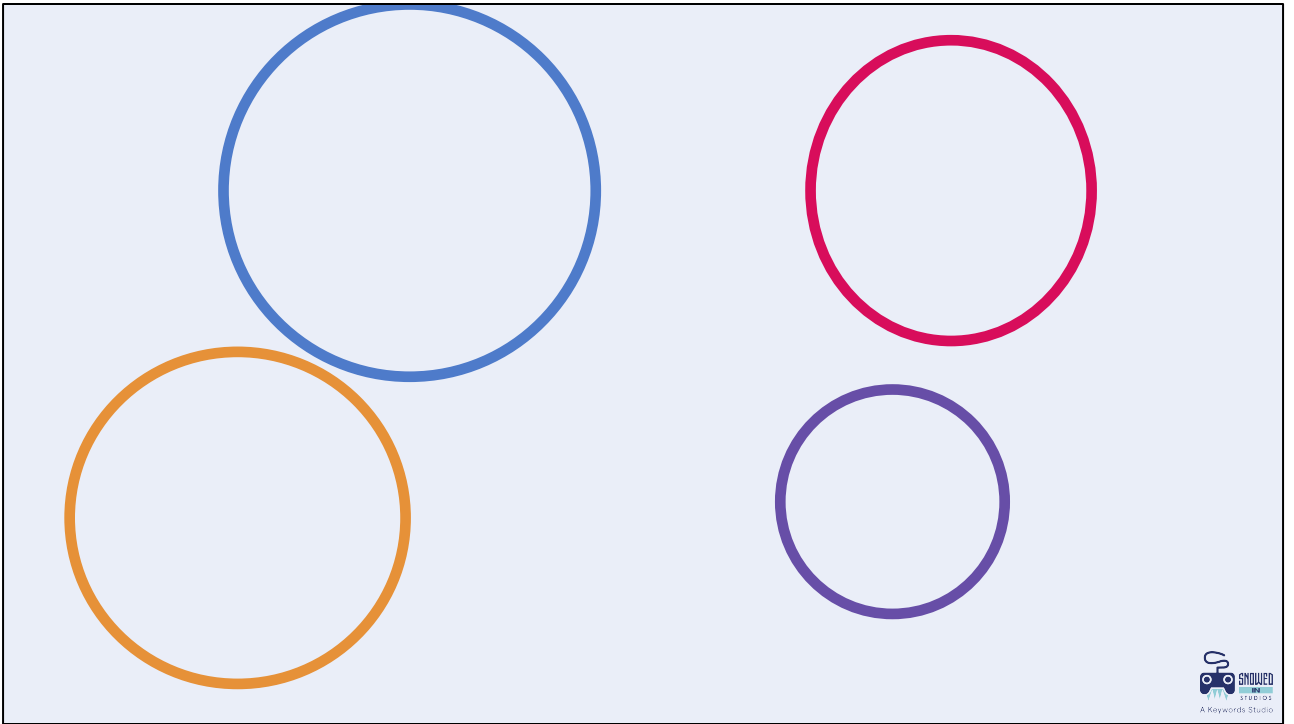


Now!

We've touched on some common use cases for wave intrinsics.

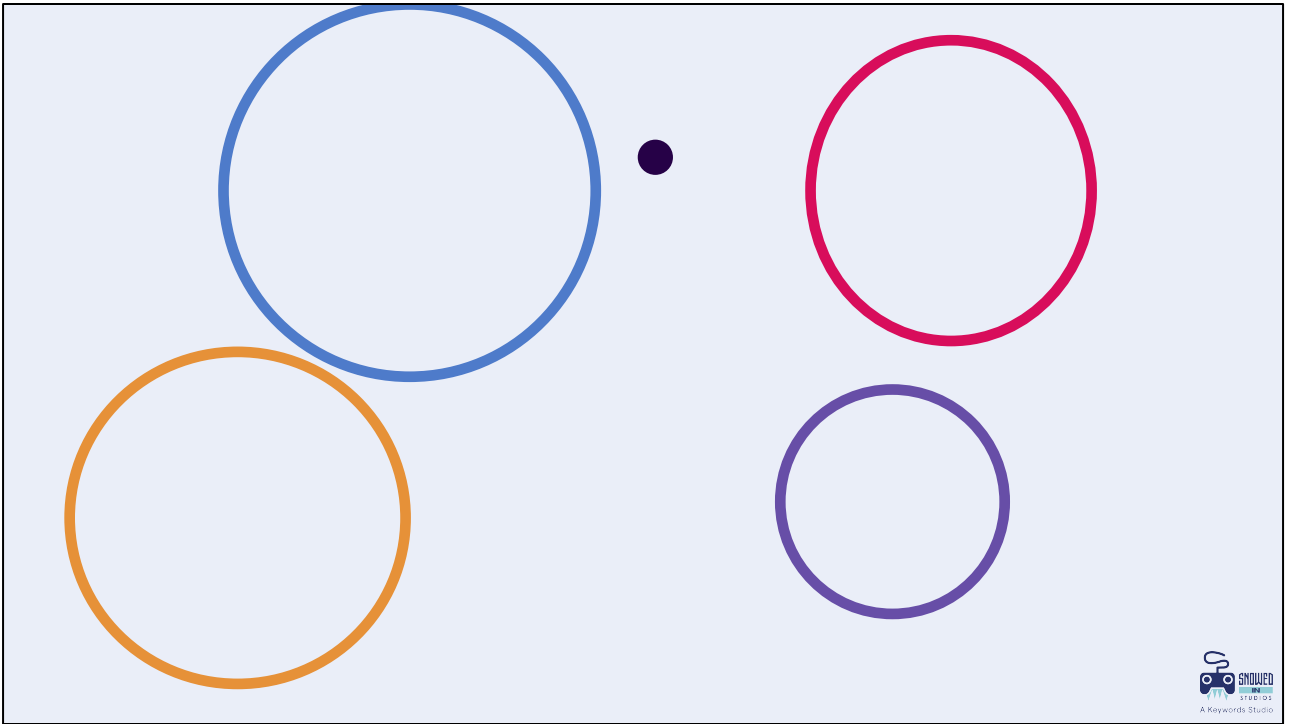
Let's look at an algorithm we used on a project currently in production.

Let's talk about wave-wide interpolation.



Imagine you have 4 spheres.

Each sphere has a radius and emits some sort of color.

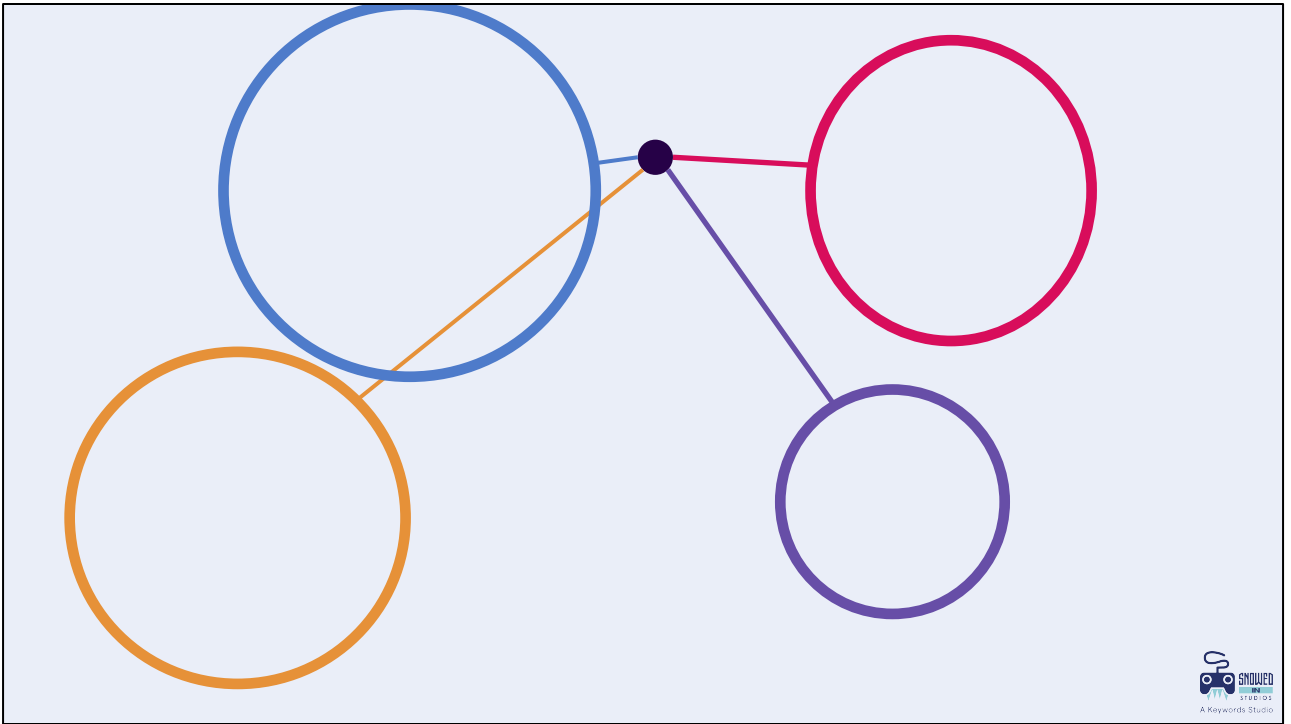


Now, let's imagine that we have a point that wants to interpolate between the colors of each of these spheres.

\*click\*

For some reason, we don't want a weighted average. We want a series of recursive linear interpolations based on some priority scheme.

This formulation of the problem is simple. We could do it on the CPU if we wanted to.

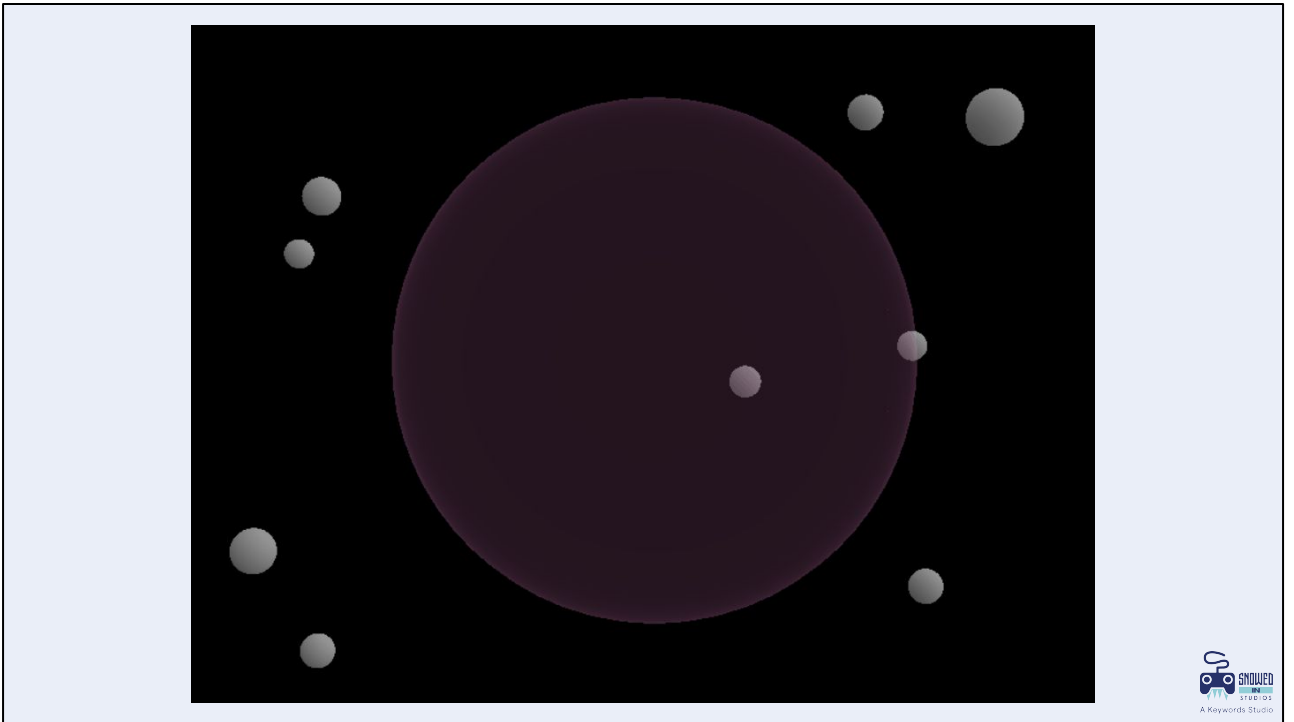


Now, let's imagine that we have a point that wants to interpolate between the colors of each of these spheres.

\*click\*

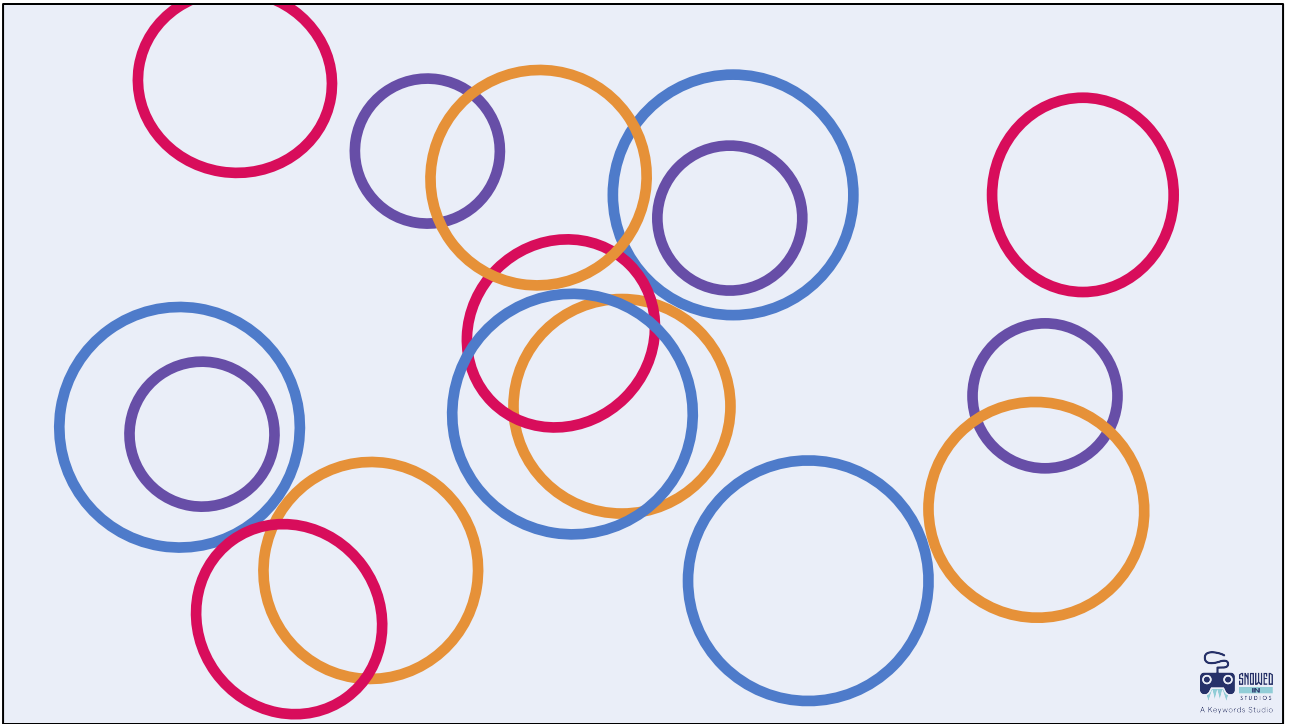
For some reason, we don't want a weighted average. We want a series of recursive linear interpolations based on some priority scheme.

This formulation of the problem is simple. We could do it on the CPU if we wanted to.



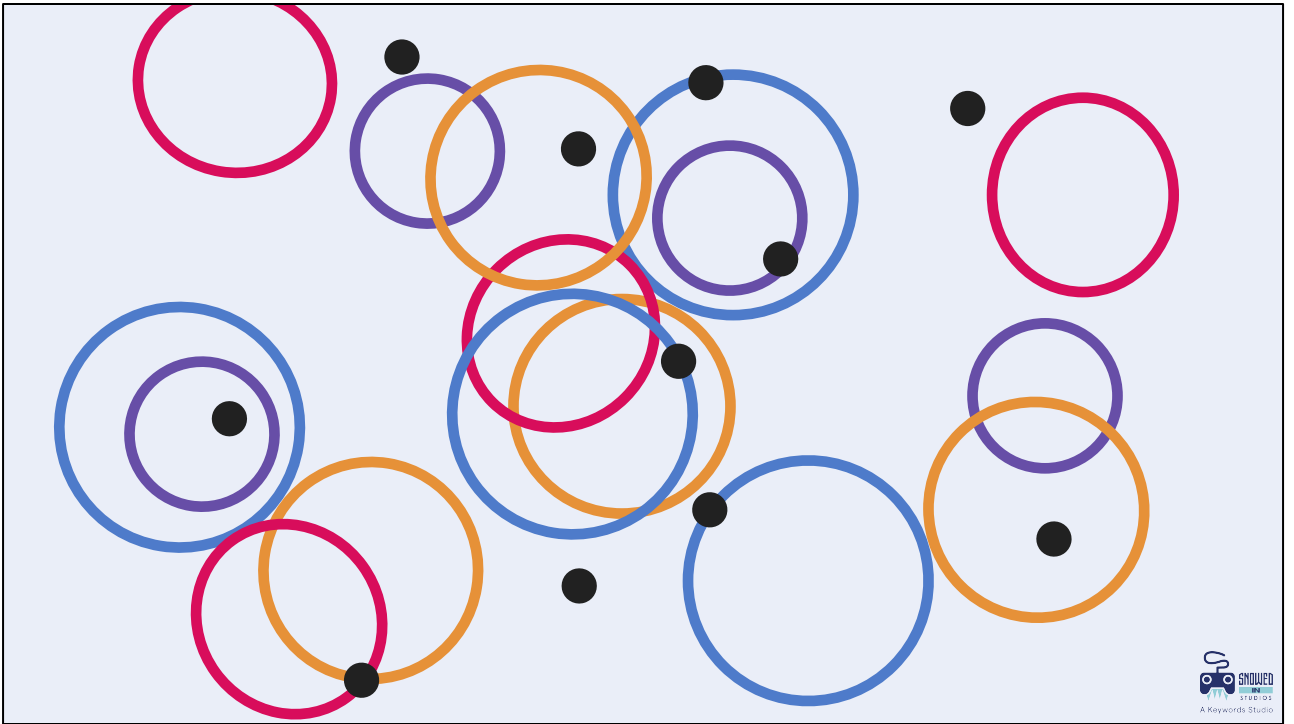
We want something that looks like this.

(Pause)



Instead, what if we have many more spheres.

Let's say 128 of them.



And let's say we have 128 points that each want to interpolate the colors of these 128 spheres.



# Points	# Spheres	Strategy	# Threads	Iterations / Thread	Approval
1 Million	4	Point Per Thread	1 Million	4	

An important detail in our process is for us to consider the amount of parallelism present in our problem.

This depends on the numbers of points and spheres.

If we have 1 million points and 4 spheres, then no problem!

\*click\*

We can just have one thread per point calculate its result with the 4 spheres.

# Points	# Spheres	Strategy	# Threads	Iterations / Thread	Approval
1 Million	4	Point Per Thread	1 Million	4	👍

An important detail in our process is for us to consider the amount of parallelism present in our problem.

This depends on the numbers of points and spheres.

If we have 1 million points and 4 spheres, then no problem!

\*click\*

We can just have one thread per point calculate its result with the 4 spheres.

# Points	# Spheres	Strategy	# Threads	Iterations / Thread	Approval
1 Million	4	Point Per Thread	1 Million	4	👍
128	128	Point Per Thread	128	128	?

But what if we only have 128 points and 128 spheres?

If we use the same strategy as above, we end up with 128 threads looping through 128 spheres.

## Example RDNA GPU

- 26 Workgroup Processors (WGP)
- 4 SIMD32 per WGP
- 20 waves per SIMD32
- $26 * 4 * 20 * 32 = 66,560$  threads for 100% saturation

As you may have guessed, this isn't a very effective use of our GPU.

If we look at an example RDNA GPU with 26 Workgroup Processors, we can schedule up to 66 thousand threads.

128 threads is a mere 0.1% of our GPU's potential capacity.

# Points	# Spheres	Strategy	# Threads	Iterations / Thread	Approval
1 Million	4	Point Per Thread	1 Million	4	👍
128	128	Point Per Thread	128	128	?

\*click\*

# Points	# Spheres	Strategy	# Threads	Iterations / Thread	Approval
1 Million	4	Point Per Thread	1 Million	4	👍
128	128	Point Per Thread	128	128	👎

\*click\*

# Points	# Spheres	Strategy	# Threads	Iterations / Thread	Approval
1 Million	4	Point Per Thread	1 Million	4	👍
128	128	Point Per Thread	128	128	👎
128	128	Point-Sphere Pair Per Thread?			

What if instead, we spawn a thread per point-sphere pair?

We would have each thread compute a single piece of our equation for each point and we would somehow combine the results.

And as a result, we be able to spawn \*click\* 16 thousand threads.

# Points	# Spheres	Strategy	# Threads	Iterations / Thread	Approval
1 Million	4	Point Per Thread	1 Million	4	👍
128	128	Point Per Thread	128	128	👎
128	128	Point-Sphere Pair Per Thread?	$128 \times 128 = 16384$	1	👍

What if instead, we spawn a thread per point-sphere pair?

We would have each thread compute a single piece of our equation for each point and we would somehow combine the results.

And as a result, we be able to spawn \*click\* 16 thousand threads.



How?



\*click\*

# Version #1


```
for(int p = 0; p < pointCount; p++)
{
    float3 pointPos = PointPositions[p];

    float3 color = float3(0.0f, 0.0f, 0.0f);
    for(uint i = 0; i < sphereCount; i++)
    {
        float influence = calcInfluence(spherePositions[i], pointPos);
        color = lerp(color, sphereColors[i], influence);
    }

    PointColors[p] = color;
}
```

If we were to write our problem as a serial program, it might look something like this.

# Version #1

```
for(int p = 0; p < pointCount; p++)  Per-Thread
{
    float3 pointPos = PointPositions[p];


    float3 color = float3(0.0f, 0.0f, 0.0f);
    for(uint i = 0; i < sphereCount; i++)
    {
        float influence = calcInfluence(spherePositions[i], pointPos);
        color = lerp(color, sphereColors[i], influence);
    }


    PointColors[p] = color;
}
```

In our base version, we simply have one thread per-point.

This works well since we don't have any dependencies between any loop iteration which means each iteration can run completely independently.

## Version #2

```
for(int p = 0; p < pointCount; p++)  Per-Thread Group
{
    float3 pointPos = PointPositions[p];

    float3 color = float3(0.0f, 0.0f, 0.0f);
    for(uint i = 0; i < sphereCount; i++)  Per-Thread
    {
        float influence = calcInfluence(spherePositions[i], pointPos);
        color = lerp(color, sphereColors[i], influence);
    }

    PointColors[p] = color;
}
```

Instead, what if we spawn a thread group for each point and spawn a thread per sphere?

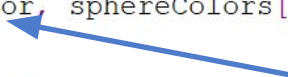
We're unrolling our loop where each individual iteration is its own thread.

## Version #2

```
for(int p = 0; p < pointCount; p++)
{
    float3 pointPos = PointPositions[p];

    float3 color = float3(0.0f, 0.0f, 0.0f);
    for(uint i = 0; i < sphereCount; i++)
    {
        float influence = calcInfluence(spherePositions[i], pointPos);
        color = lerp(color, sphereColors[i], influence);
    }

    PointColors[p] = color;
}
```



Cross-Thread Dependency

However, our lerp has a dependency between each loop iteration (aka thread).

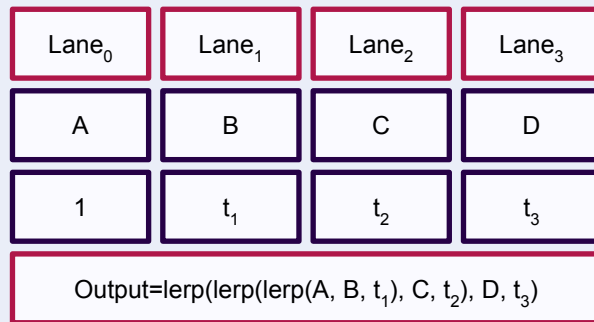
How do we address that?

Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
A	B	C	D
1	t <sub>1</sub>	t <sub>2</sub>	t <sub>3</sub>

If we look at our small 4 lane GPU again.

This set of lanes is looking to compute the result of the interpolation of A, B, C and D using 4 interpolation terms.

For convenience and so the algebra doesn't get too crazy, we've set Lane 0's interpolation term to 1.



We want to combine the results of our lanes such that we evaluate the expression above.

$\text{lerp}(A, B, t_0)$	$A*(1-t_0)+B*t_0$
$\text{lerp}(\text{lerp}(A, B, t_1), C, t_2)$	$A*(1-t_1)*(1-t_2)+B*t_1*(1-t_2)+C*t_2$
$\text{lerp}(\text{lerp}(\text{lerp}(A, B, t_0), C, t_1), D, t_2)$	$A*(1-t_1)*(1-t_2)*(1-t_3)+B*t_1*(1-t_2)*(1-t_3)+C*t_2*(1-t_3)+D*t_3$

If we expand each iteration of a recursive lerp, you might notice a pattern.



Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
A	B	C	D
1	t <sub>1</sub>	t <sub>2</sub>	t <sub>3</sub>

$$A*(1-t_1)*(1-t_2)*(1-t_3)+B*t_1*(1-t_2)*(1-t_3)+C*t_2*(1-t_3)+D*t_3$$

Each term is simply a combination of its interpolation factor and its next neighbour's 1 minus interpolation term.

\*click\*

\*click\*

\*click\*

Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
A	B	C	D
1	t <sub>1</sub>	t <sub>2</sub>	t <sub>3</sub>

$$A*(1-t_1)*(1-t_2)*(1-t_3)+B*t_1*(1-t_2)*(1-t_3)+C*t_2*(1-t_3)+D*t_3$$

Each term is simply a combination of its interpolation factor and its next neighbour's 1 minus interpolation term.

\*click\*

\*click\*

\*click\*

Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
A	B	C	D
1	t <sub>1</sub>	t <sub>2</sub>	t <sub>3</sub>

$$A*(1-t_1)*(1-t_2)*(1-t_3)+B*t_1*(1-t_2)*(1-t_3)+C*t_2*(1-t_3)+D*t_3$$

Each term is simply a combination of its interpolation factor and its next neighbour's 1 minus interpolation term.

\*click\*

\*click\*

\*click\*

Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
A	B	C	D
1	t <sub>1</sub>	t <sub>2</sub>	t <sub>3</sub>



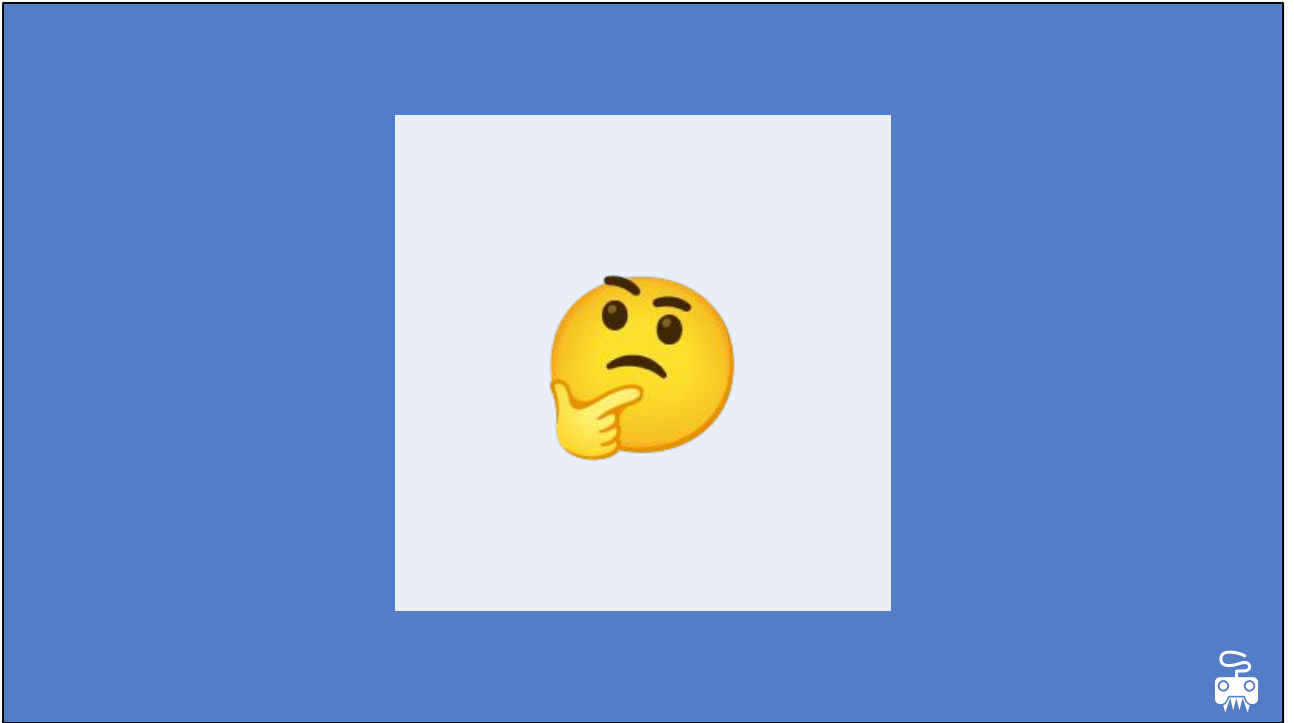
$$A*(1-t_1)*(1-t_2)*(1-t_3)+B*t_1*(1-t_2)*(1-t_3)+C*t_2*(1-t_3)+D*t_3$$

Each term is simply a combination of its interpolation factor and its next neighbour's 1 minus interpolation term.

\*click\*

\*click\*

\*click\*



Which looks a lot like a prefix operation in the reverse order.

Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
A	B	C	D
1	t <sub>1</sub>	t <sub>2</sub>	t <sub>3</sub>

$$A*(1-t_1)*(1-t_2)*(1-t_3)+B*t_1*(1-t_2)*(1-t_3)+C*t_2*(1-t_3)+D*t_3$$

So perhaps we can simply reverse the order of our lane's values.

\*click\*

Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
D	C	B	A
t <sub>3</sub>	t <sub>2</sub>	t <sub>1</sub>	1

$$D \cdot t_3 + C \cdot t_2 \cdot (1 - t_3) + B \cdot t_1 \cdot (1 - t_2) \cdot (1 - t_3) + A \cdot (1 - t_1) \cdot (1 - t_2) \cdot (1 - t_3)$$

So perhaps we can simply reverse the order of our lane's values.

\*click\*

Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
D	C	B	A
t <sub>3</sub>	t <sub>2</sub>	t <sub>1</sub>	1

$$D \cdot t_3 + C \cdot t_2 \cdot (1 - t_3) + B \cdot t_1 \cdot (1 - t_2) \cdot (1 - t_3) + A \cdot (1 - t_1) \cdot (1 - t_2) \cdot (1 - t_3)$$

Prefix Product + Wave Sum

Then we can easily apply a prefix product.



Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
D	C	B	A
t <sub>3</sub>	t <sub>2</sub>	t <sub>1</sub>	1

$$D \cdot t_3 + C \cdot t_2 \cdot (1 - t_3) + B \cdot t_1 \cdot (1 - t_2) \cdot (1 - t_3) + A \cdot (1 - t_1) \cdot (1 - t_2) \cdot (1 - t_3)$$

Prefix Product + Wave Sum

(Prefix sum but with multiplication instead of addition)

\*click\*

Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
D	C	B	A
t <sub>3</sub>	t <sub>2</sub>	t <sub>1</sub>	1

$$D \cdot t_3 + C \cdot t_2 \cdot (1 - t_3) + B \cdot t_1 \cdot (1 - t_2) \cdot (1 - t_3) + A \cdot (1 - t_1) \cdot (1 - t_2) \cdot (1 - t_3)$$

Prefix Product + Wave Sum

And once we've applied our prefix product, we simply want to combine our results using a wave sum!

$$D*t_3 + C*t_2*(1-t_3) + B*t_1*(1-t_2)*(1-t_3) + A*(1-t_1)*(1-t_2)*(1-t_3)$$

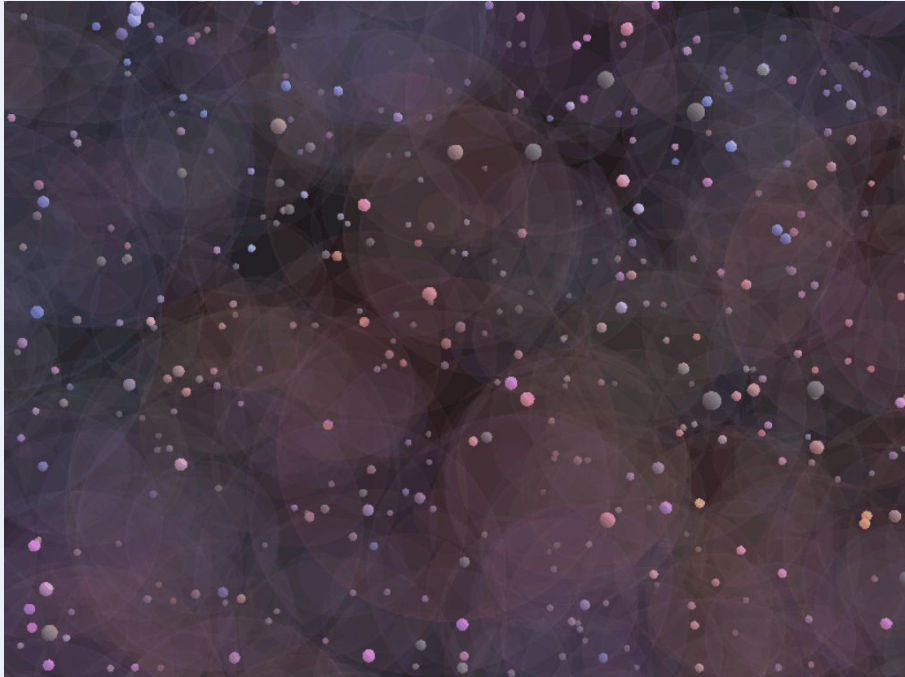
```
float WaveActiveLerp(float value, float t)
{
    float prefixProduct = WavePrefixProduct(1.0f - t);
    float laneValue = value * t * prefixProduct;
    float interpolation = WaveActiveSum(laneValue);

    return interpolation;
}
```

(Pause)

Through some simple algebra, we were able to identify a pattern that we could use to turn what might be a problem ill-suited for a GPU to one that can effectively saturate our GPU's cores.

From 128 threads to 16 thousand threads through the use of wave intrinsics.



Here we've simulated a thousand spheres influencing a thousand points at twice the speed of the naive implementation.

1024 Spheres x 1024 Points			
	Naive	Parallelized	Savings
Intel(R) Iris(R) Xe Graphics	0.6ms	0.3ms	0.3ms
NVIDIA GeForce RTX 3070 Ti	0.4ms	0.1ms	0.3ms

As you can see, the savings can be drastic between 2x and 4x in savings on various GPUs.

# Warning: Code Generation



A word of warning before you go forth.

# Code Generation - Dynamic WaveReadLaneAt

```
44 ;      uint outputIndex = WaveReadLaneAt(selectInput, WaveGetLaneIndex() + 1);
45     v_cndmask_b32_e64 v3, 0, 1, vcc_lo                                // 00000000000048:
46 000000000000000050 <L0>:
47     v_readfirstlane_b32 s6, v2                                        // 00000000000050:
48     v_cmp_eq_u32_e64 s4, s6, v2                                     // 00000000000054:
49     s_and_saveexec_b64 s[4:5], s[4:5]                               // 0000000000005C:
50     v_readlane_b32 s6, v3, s6                                       // 00000000000060:
51     v_mov_b32_e32 v1, s6                                             // 00000000000068:
52     v_xor_b64 exec, exec, s[4:5]                                     // 0000000000006C:
53     s_cbranch_execnz L0                                             // 00000000000070:
54     s_mov_b64 exec, s[2:3]                                          // 00000000000074:
55 : C:\\Users\\Alex\\Desktop\\Development\\Iccap\\Iccap\\bin\\Debug\\net8.0-wir
```

**Waterfall loop! :(**

Be mindful of your code generation. Using WaveReadLaneAt to do a shuffle boils down to a waterfall loop on AMD's offline compiler.

\*click\*

# Code Generation - Spir-V Shuffle Up?

```
40 ;      uint outputIndex = OpGroupNonUniformShuffleUp(vk::SubgroupScope, (uint)selectInput, 1);
41     v_cndmask_b32_e64 v2, 0, 1, vcc_lo                                     // 00000000000048: D5010002 01A902
42 0000000000000050 <L0>:
43     v_readfirstlane_b32 s6, v3                                           // 00000000000050: 7E0C0503
44     v_cmp_eq_u32_e64 s4, s6, v3                                          // 00000000000054: D4C20004 000206
45     s_and_saveexec_b64 s[4:5], s[4:5]                                    // 0000000000005C: BE842404
46     v_readlane_b32 s6, v2, s6                                            // 00000000000060: D7600006 00000E
47     v_mov_b32_e32 v1, s6                                                 // 00000000000068: 7E020206
48     s_xor_b64 exec, exec, s[4:5]                                         // 0000000000006C: 89FE047E
49     s_cbranch_execnz L0                                                  // 00000000000070: BF89FFF7
50     s_mov_b64 exec, s[2:3]                                               // 00000000000074: BEFE0402
```

Waterfall loop! :(

And even if you use the more constrained shuffle up spir-v instruction.

It still boils down to a waterfall loop.



# What's A Waterfall Loop?



But what's a waterfall loop?

```
for (;;)
{
    uint scalarValue = WaveReadLaneFirst(laneValue);
    [branch]
    if (scalarValue == laneValue)
    {
        // Do something.
        break;
    }
}
```

A waterfall loop iterates through each unique value within our wave.

If we have a wave with unique values A, B and D we will iterate 3 times.

Each lane that matches our first active lane will enter the branch.

Then they will mark themselves inactive.

Then each lane that matches the new first active lane will enter the branch.

Until every unique value was processed.

# Waterfall Loop

- Iterates for each unique value in a wave
- Worst case: wave lane count number of iterations
- Best case: single iteration
- Useful if you want to reduce vector register pressure

```
for (;;)
{
    uint scalarValue = WaveReadLaneFirst(laneValue);
    [branch]
    if (scalarValue == laneValue)
    {
        // Do something.
        break;
    }
}
```

In the best case, we iterate only once.

In the worst case when each lane in a wave has a different value, we iterate once for each lane.

# Conclusions



\*click\*



A lens you can use to navigate your problem spaces, is to think of your GPU as a super powered loop unrolling machine.

With that in mind, you can imagine that there's a gradient spanning a fully serial program

\*click\*

to a fully unrolled program where each iteration of the loop is a thread.

\*click\*

With a vast number of options in between.

Fully Serial



A lens you can use to navigate your problem spaces, is to think of your GPU as a super powered loop unrolling machine.

With that in mind, you can imagine that there's a gradient spanning a fully serial program

\*click\*

to a fully unrolled program where each iteration of the loop is a thread.

\*click\*

With a vast number of options in between.



A lens you can use to navigate your problem spaces, is to think of your GPU as a super powered loop unrolling machine.

With that in mind, you can imagine that there's a gradient spanning a fully serial program

\*click\*

to a fully unrolled program where each iteration of the loop is a thread.

\*click\*

With a vast number of options in between.

# How To Develop Your Own?

- **Take inspiration from CPU SIMD and integer bit tricks**
  - Hacker's Delight - Henry S. Warren, Jr
  - [Bit Twiddling Hacks](#) - Sean Eron Anderson
  - [SIMD at Insomniac Games](#) - Andreas Fredriksson
- **Look at our friends in GPGPU**
- **Play around with the algebra**
  - Maybe there's a prefix or suffix operation hiding in there!
- **Explore the literature**
  - Goes back to the 70s!
    - [A Parallel Algorithm for the Efficient Solution of a General Class of Recurrence Equations](#) - Peter M. Kogge, Harold S. Stone

There's a wealth of inspiration that you can gain from many areas to develop your own algorithms.

Here are some spaces you can explore to get some ideas!



# Where Can You Find Me?

Mastodon



<https://mastodon.gamedev.place/@AlexSneezeKing>

Blog



[https://github.com/AlexSabourinDev/cranberry\\_blog](https://github.com/AlexSabourinDev/cranberry_blog)

If you want to chat, you can find me on Mastodon.

And if you're interested in more topics like this, take a look at my blog!

# References

- [1] [SIMD at Insomniac Games](#) - Andreas Fredriksson
- [2] [Wave Programming in D3D12 and Vulkan](#) - David Lively, Holger Gruen
- [3] [Stream compaction using wave intrinsics - Interplay of Light](#) - Kostas Anagnostou
- [4] [Compute shader wave intrinsics tricks](#) - Marko Sreckovic
- [5] [GPU Programming Primitives for Computer Graphics](#) - Daniel Meister, Atsushi Yoshimura, Chih-Chen Kao
- [6] [Optimizing for the Radeon RDNA architecture](#) - Lou Kramer
- [7] [Optimizing The Graphics Pipeline with Compute](#) - Graham Wihlidal
- [8] [Stream Reduction Operations for GPGPU Applications](#) - Daniel Horn
- [9] [A Parallel Algorithm for the Efficient Solution of a General Class of Recurrence Equations](#) - Peter M. Kogge, Harold S. Stone
- [10] [Bit Twiddling Hacks](#) - Sean Eron Anderson
- [11] [The Power Of Parallel Prefix](#) - Clyde P. Kruskal, Larry Rudolph, Marc Snir
- [12] [SIMD and SWAR Techniques](#)
- [13] [Prefix Sums and Their Applications](#) - Guy E. Blelloch
- [14] [AMD GCN Assembly: Cross-Lane Operations](#) - Ben Sander
- [15] [Reading Between The Threads: Shader Intrinsics](#)
- [16] [Adaptive Exposure from Luminance Histograms](#) - Alex Tardif
- [17] <https://bruop.github.io/exposure/> - Bruno Opsenica

I recommend you take the time to explore some of these references. There's great insights in all of these.

# Thanks!



Thank you all for your time!

If you have feedback for me on my presentation, I would love to hear from you later.

# Questions?



Any questions?

# Bonus Slides: Wave Histogram



Let's explore at a problem at a different scale.

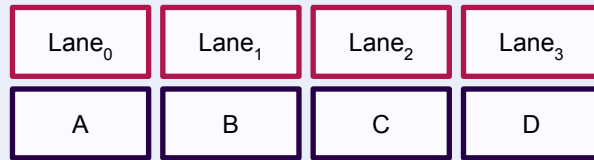
No more spheres. Instead, I want boxes.

Perhaps you're working on histogram-based adaptive exposure or perhaps you just like counting things.

We're going to investigate tool we can use to optimize atomic adds to a histogram.

```
[numthreads(1024, 1, 1)]  
void CS(uint3 dispatchId : SV_DispatchThreadId)  
{  
    uint bucket = getHistogramBucket(Input[dispatchId.x]);  
    InterlockedAdd(OutputHistogram[bucket], 1);  
}
```

Here we have a simple program that adds to a histogram.



Histogram= 

0	0	0	0
---	---	---	---

Once again, we want to look at our trusty 4 lane GPU.

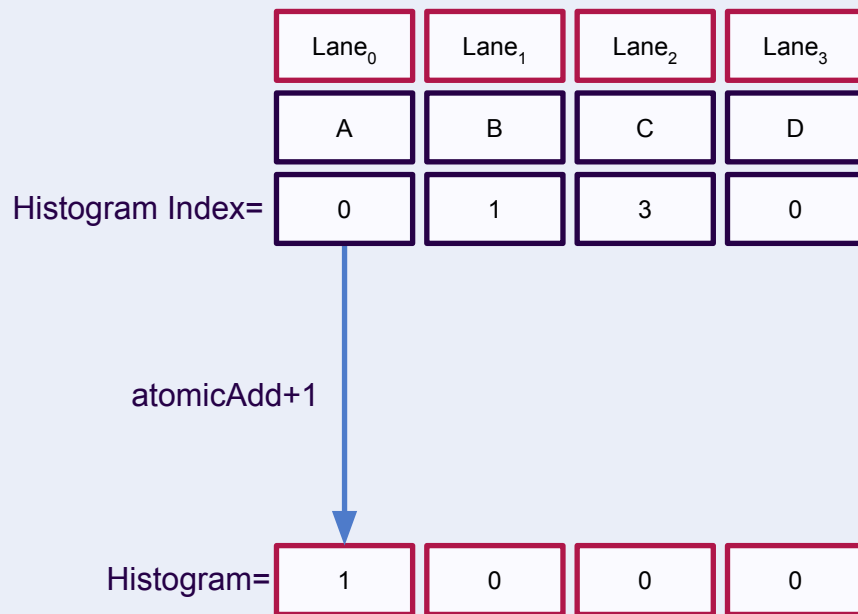
	Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
	A	B	C	D
Histogram Index=	0	1	3	0

Histogram=	0	0	0	0
------------	---	---	---	---

This program will assign an index to each value and each lane will add to their respective buckets using an atomic increment.

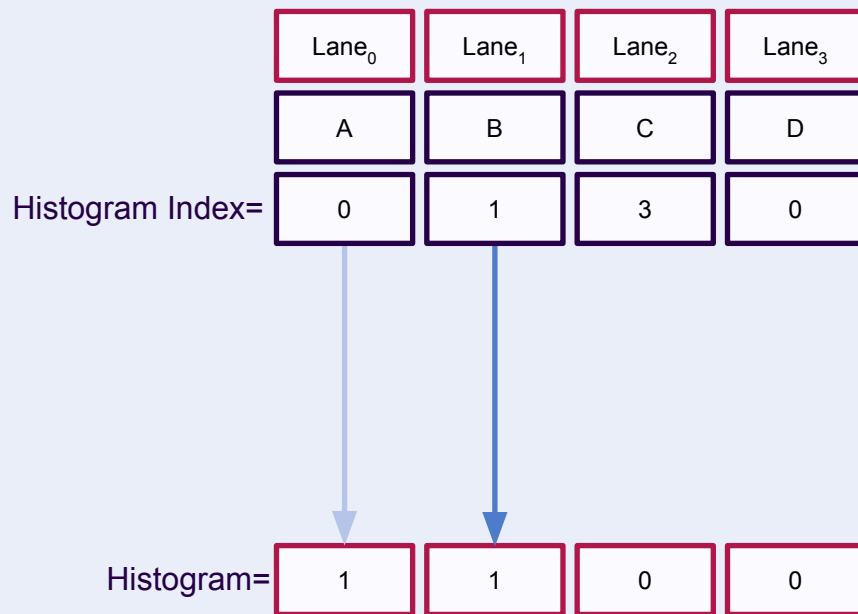
\*click\*  
\*click\*  
\*click\*  
\*click\*





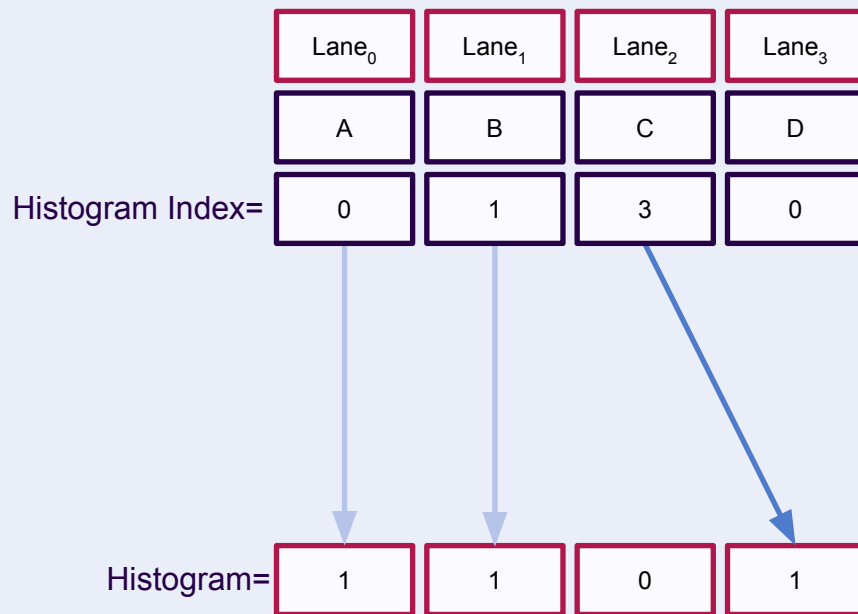
This program will assign an index to each value and each lane will add to their respective buckets using an atomic increment.

\*click\*  
\*click\*  
\*click\*  
\*click\*



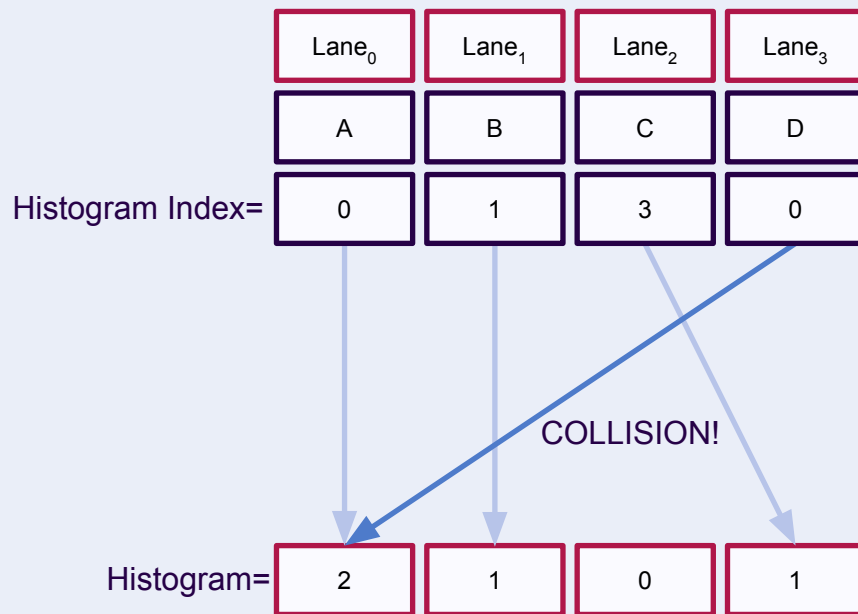
This program will assign an index to each value and each lane will add to their respective buckets using an atomic increment.

\*click\*  
\*click\*  
\*click\*  
\*click\*



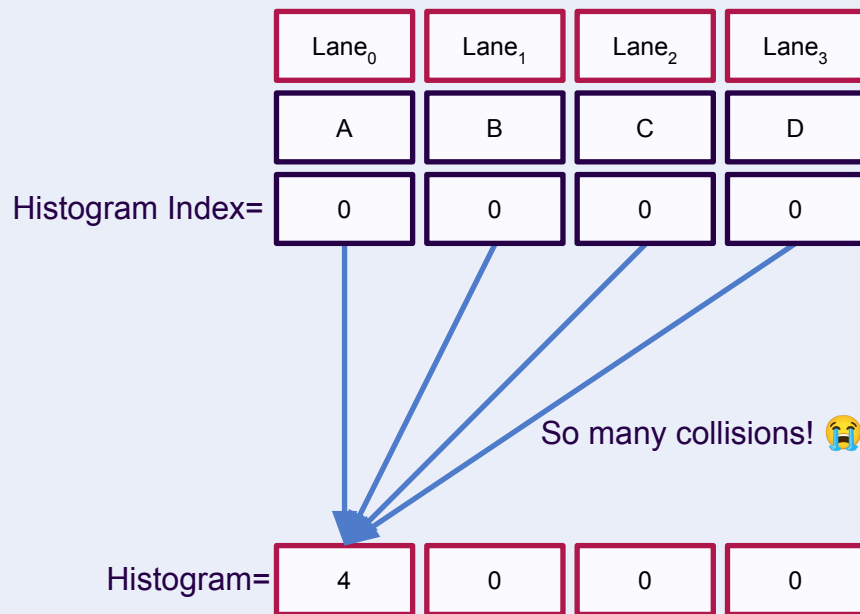
This program will assign an index to each value and each lane will add to their respective buckets using an atomic increment.

\*click\*  
\*click\*  
\*click\*  
\*click\*



Sometimes we're going to have a collision.

This can cause a small, probably unnoticeable performance penalty.



Unless everyone writes to the same location...

# 256 Bucket Histogram Benchmark

	Global Memory	+Groupshared
<b>Intel(R) Iris(R) Xe Graphics - 16 Million Elements, 256 Buckets</b>		
Minimal Collisions	30ms	1.5ms
Maximum Collisions	30ms	3.5ms
<b>NVIDIA GeForce RTX 3070 Ti - 16 Million Elements, 256 Buckets</b>		
Minimal Collisions	2ms	0.3ms
Maximum Collisions	2ms	2ms

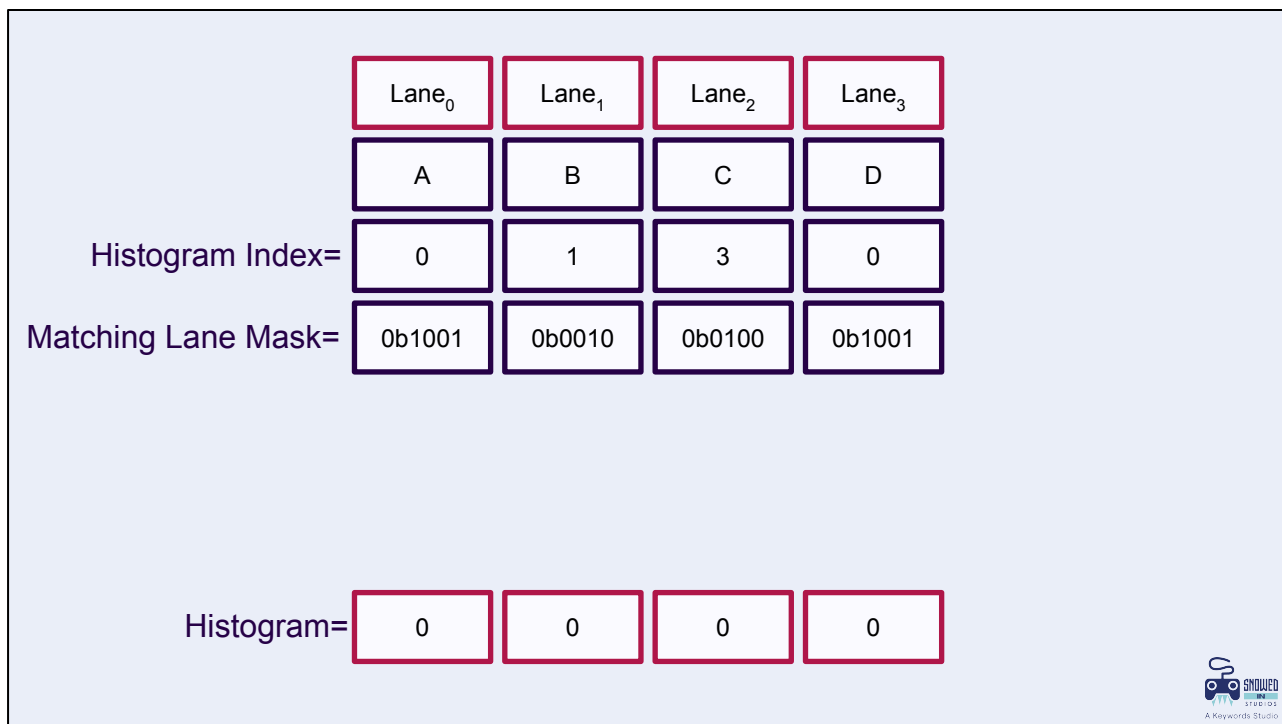
Writing directly to global memory is not bound by our atomic throughput.

But when we use groupshared memory as an intermediate data store, we start seeing some performance penalties.

	Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
	A	B	C	D
Histogram Index=	0	1	3	0

Histogram=	0	0	0	0
------------	---	---	---	---

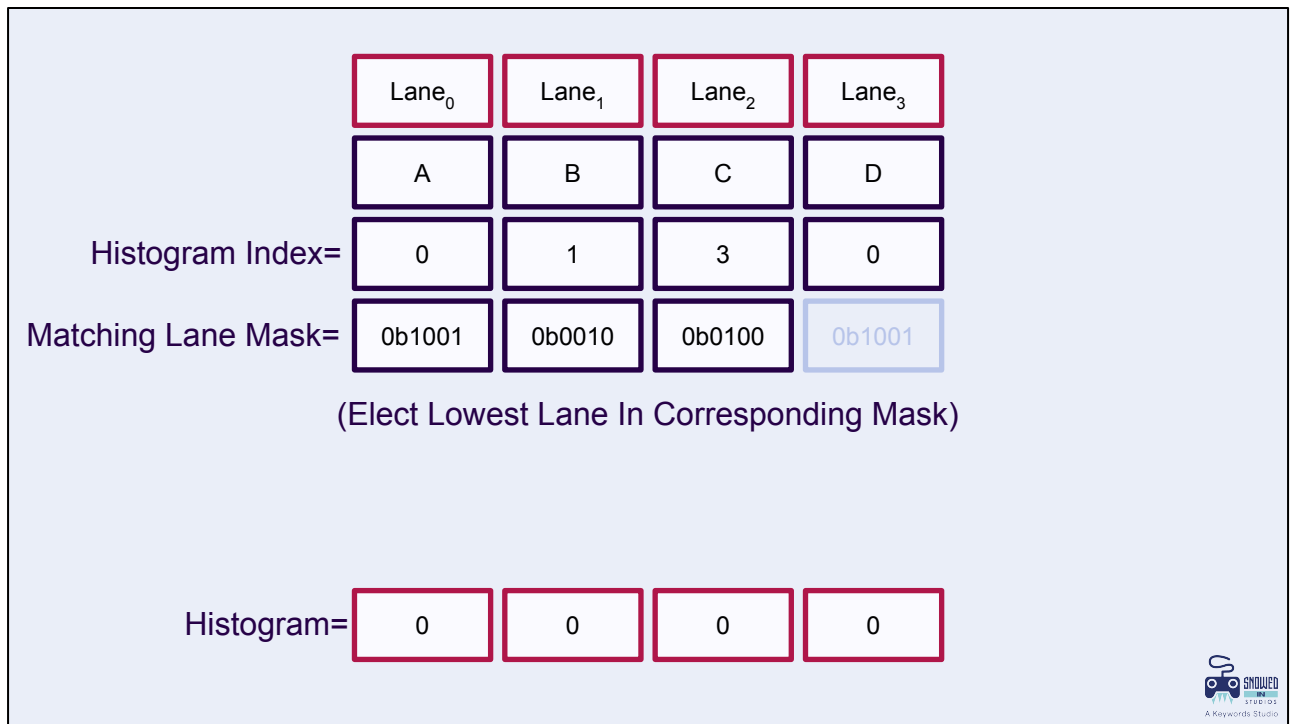
What we want, is some sort of “mask” that would tell us which lanes have the same value as our current lane.



Once we have that mask, we want to elect the lowest lane in each mask.

In this case, note that lane 0 and lane 3 have the same mask since they want to write to the same histogram bucket.



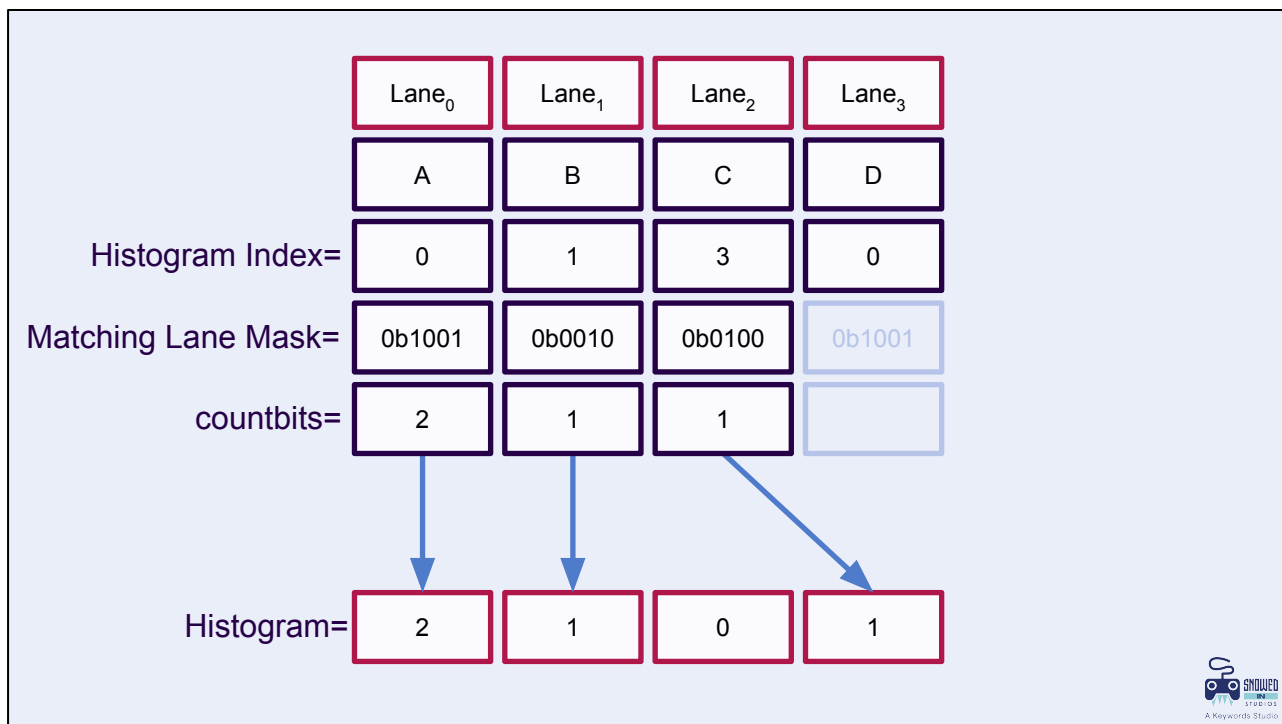


When we have a conflict, we selected Lane 0 as the lowest lane whose bit is set in its respective mask.

	Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
	A	B	C	D
Histogram Index=	0	1	3	0
Matching Lane Mask=	0b1001	0b0010	0b0100	0b1001
countbits=	2	1	1	
Histogram=	0	0	0	0

Then, it's simply a matter of counting the bits in our mask and writing to the corresponding location in our histogram.

\*click\*



Then, it's simply a matter of counting the bits in our mask and writing to the corresponding location in our histogram.

\*click\*

# How Do We Get The Lane Mask?



That's great and all, but how do we get that mask?

# How Do We Get The Lane Mask?

- We have a fixed histogram size (in this case 4 buckets)

In this case, we have a fixed histogram size with 4 buckets.

# How Do We Get The Lane Mask?

- We have a fixed histogram size (in this case 4 buckets)
- We could simply compare each lane to all of its neighbours
  - Wave lane count amount of shuffles 😞

We could compare each lane with all of its neighbours.

But then we would need an equivalent number of shuffles to the number of lanes on our GPU.

In this case, our “GPU” has 4 lanes. But in other cases we might see wave sizes of 32 or 64.

# How Do We Get The Lane Mask?

- We have a fixed histogram size (in this case 4 buckets)
- We could simply compare each lane to all of its neighbours
  - Wave lane count amount of shuffles 😞
- What if we compare each lane bit by bit compare?
  - Let's look at an example

What if we compare each lane bit by bit?

	Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
	A	B	C	D
Histogram Index=	0	1	3	0

Lane Mask=	0b1111	0b1111	0b1111	0b1111
------------	--------	--------	--------	--------

Let's start by assigning each lane a unique "Lane Mask" value and set it to 1s.



	Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
	A	B	C	D
Histogram Index=	0	1	3	0
	0b00	0b01	0b11	0b00

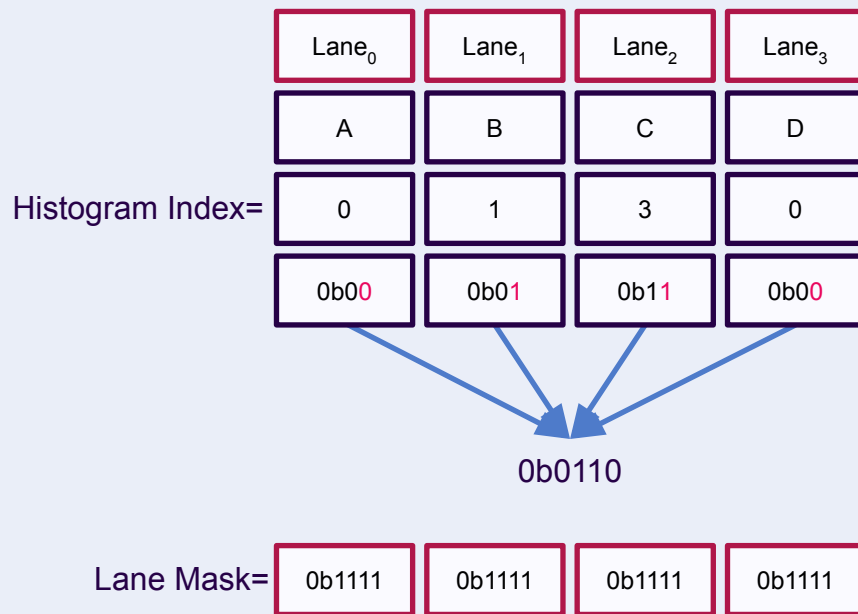
Lane Mask=	0b1111	0b1111	0b1111	0b1111
------------	--------	--------	--------	--------

Let's convert our histogram index to binary.

	Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
	A	B	C	D
Histogram Index=	0	1	3	0
	0b00	0b01	0b11	0b00

Lane Mask=	0b1111	0b1111	0b1111	0b1111
------------	--------	--------	--------	--------

We want to walk through the bits of that index one at a time.

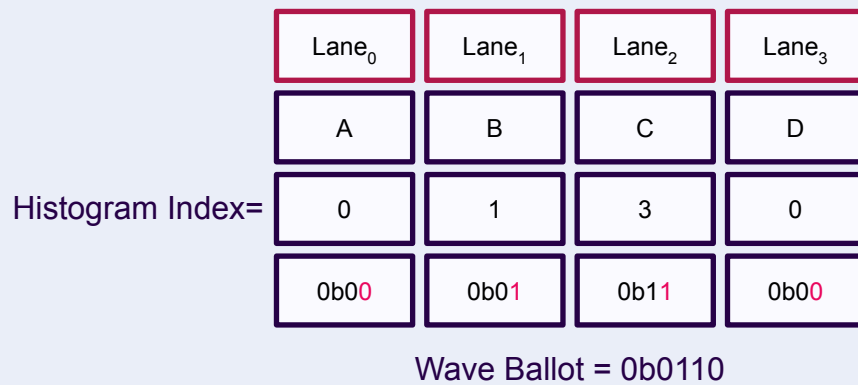


And we want to collect all the bits into a single value.

(Take some extra time to describe how this works and why we want it)

This can be achieved with a call to `WaveActiveBallot`.

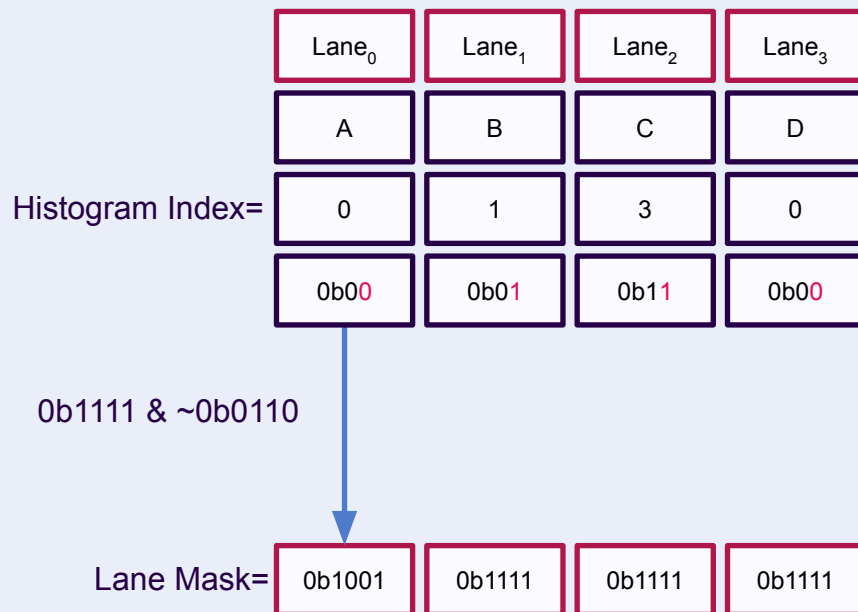
WaveActiveBallot simply allows us to collect all of the bits in our lanes into a single integer value visible to all lanes at once.



Lane Mask=

0b1111	0b1111	0b1111	0b1111
--------	--------	--------	--------

Once we have this ballot, we have a mask that tells us which lane's first bits are 1 and which lane's bits are 0.

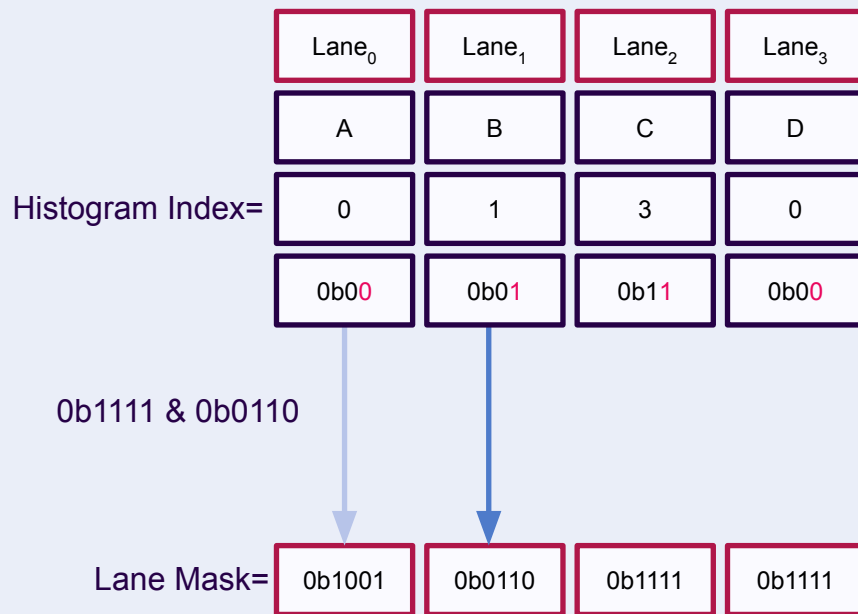


Since we're trying to find out which lanes match our own.

We use this mask as a progressive equality check.

Since our current bit is zero, then we want to invert the ballot mask.

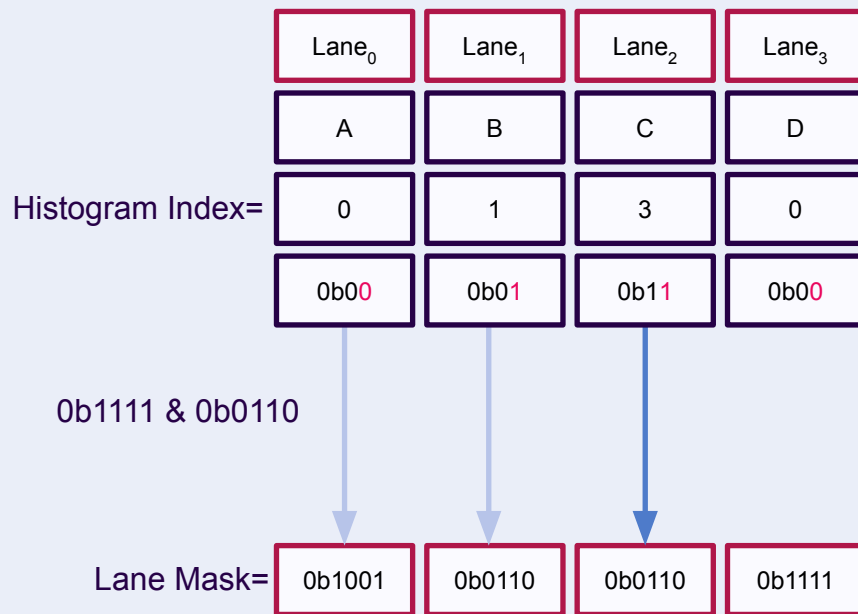
Each lane whose first bit was zero should get a "true" and each lane whose first bit was one should get a "false".



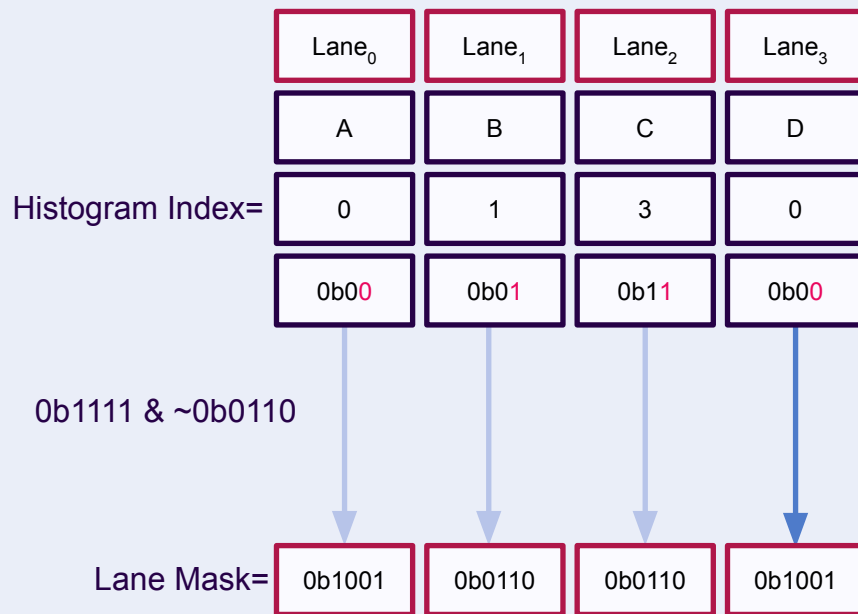
Then we continue this process.

In this case, our first bit is one.

As a result, we keep our ballot mask as it is.



\*click\*



(Pause so people can look at the current masks)

\*click\*



	Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
	A	B	C	D
Histogram Index=	0	1	3	0
	0b00	0b01	0b11	0b00

Wave Ballot = 0b0100

Lane Mask=	0b1001	0b0110	0b0110	0b1001
------------	--------	--------	--------	--------

You'll notice that we're progressively disqualifying neighbouring lanes from matching our current one.

We repeat this process up to the number of bits that we want to use in our equality test.

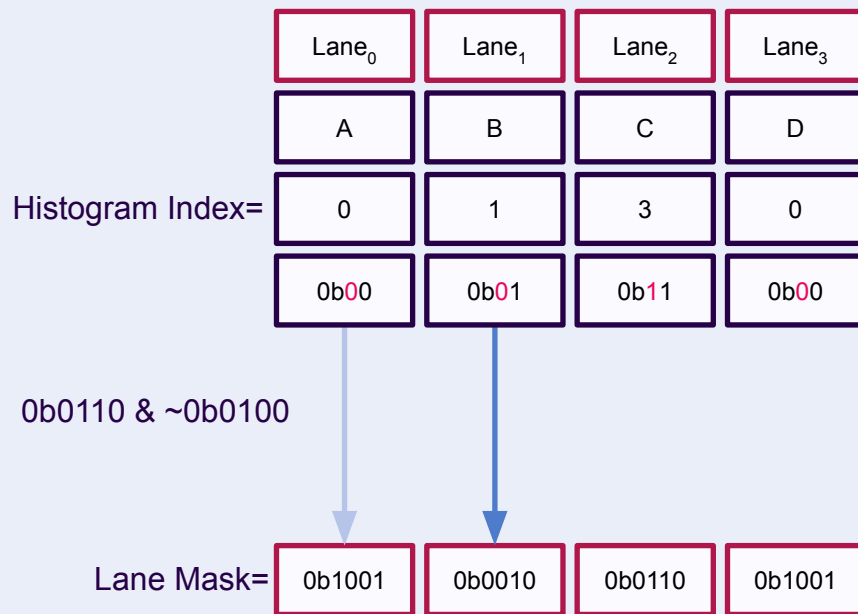
	Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
	A	B	C	D
Histogram Index=	0	1	3	0
	0b00	0b01	0b11	0b00

0b1001 & ~0b0100

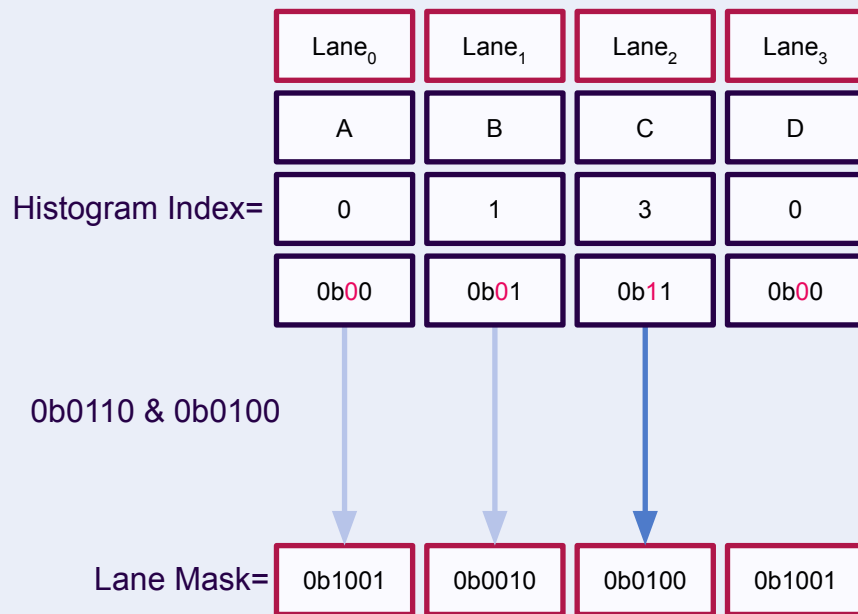
Lane Mask=

0b1001	0b0110	0b0110	0b1001
--------	--------	--------	--------

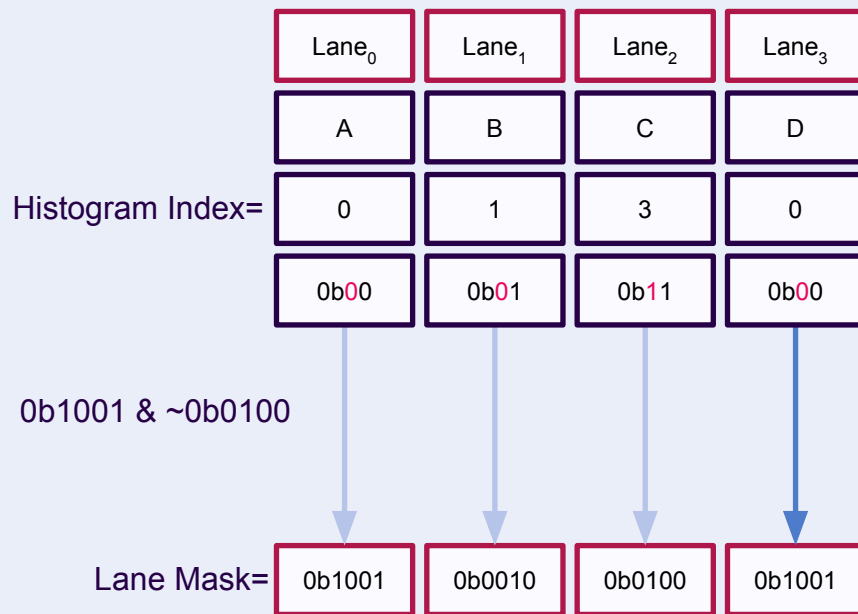
\*click\*



\*click\*



\*click\*

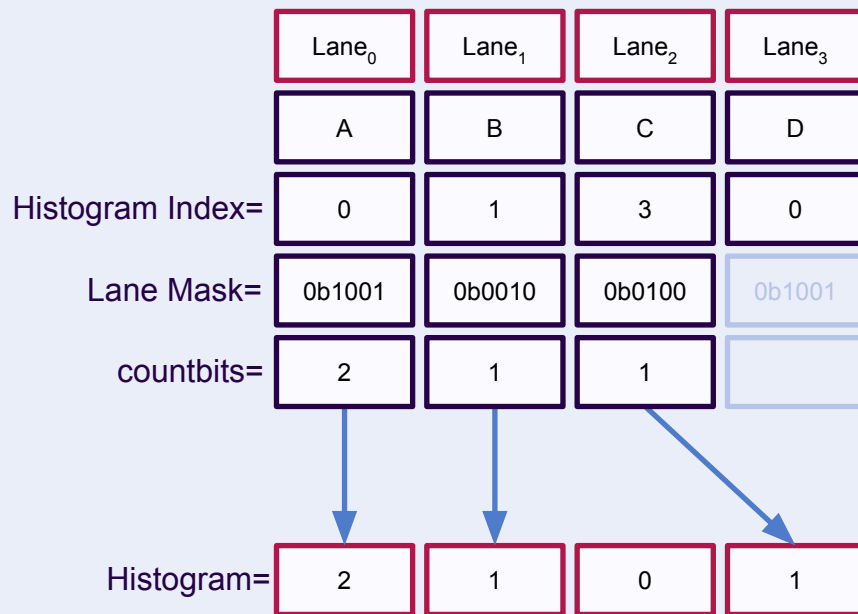


(Pause for a moment so people can look at the lane mask)

\*click\*

	Lane <sub>0</sub>	Lane <sub>1</sub>	Lane <sub>2</sub>	Lane <sub>3</sub>
	A	B	C	D
Histogram Index=	0	1	3	0
Lane Mask=	0b1001	0b0010	0b0100	0b1001

With two iterations, each lane now has a mask that represents the lane's whose value match our own.



Now we simply go through the routine that we described above and write to our histogram.

# How Do We Get The Lane Mask?

- We have a fixed histogram size (in this case 4 buckets)
- We could simply compare each lane to all of its neighbours
  - Wave lane count amount of shuffles 😞
- What if we compare each lane bit by bit compare?
  - countbits(bucketCount) number of iterations
    - Supporting 8 bits takes 8 iterations 😊

For a histogram of 256 values, you simply need to iterate 8 times.



# How Do We Get The Lane Mask?

- We have a fixed histogram size (in this case 4 buckets)
- We could simply compare each lane to all of its neighbours
  - Wave lane count amount of shuffles 😞
- What if we compare each lane bit by bit compare?
  - countbits(bucketCount) number of iterations
    - Supporting 8 bits takes 8 iterations 😊
    - Supporting 32 bits takes 32 iterations 😞

But for supporting 32 bits you need 32 iterations...

# How Do We Get The Lane Mask?

- We have a fixed histogram size (in this case 4 buckets)
- We could simply compare each lane to all of its neighbours
  - Wave lane count amount of shuffles 😞
- What if we compare each lane bit by bit compare?
  - countbits(bucketCount) number of iterations
    - Supporting 8 bits takes 8 iterations 😊
    - Supporting 32 bits takes 32 iterations 😞
    - Depends on a fast WaveActiveBallot

And this depends on a fast ballot method.

If requesting a ballot is secretly a waterfall loop, or even simply a reduction the cost of this approach soars.

```
43 | ;         uint4 ballot = WaveActiveBallot(laneBit);  
44 |     v_cmp_eq_u32_e64 s0, 0, v3  
45 |     v_cmp_ne_u32_e32 vcc_lo, 0, v3  
46 | . C:\Users\Alex\Desktop\Development\Iccan\Iccan\bin\
```

Which thankfully, in the case of RDNA and AMD's offline compiler is a nice and tidy 2 instructions.

```

uint4 WaveActiveMatch(uint value, uint bitCount)
{
    uint4 laneMask = uint4(0xFFFFFFFF, 0xFFFFFFFF, 0xFFFFFFFF, 0xFFFFFFFF);
    for (uint i = 0; i < bitCount; i++)
    {
        uint laneBit = (value >> i) & 1;
        uint4 ballot = WaveActiveBallot(laneBit);

        // Make sure to flip our ballot if our bit is 0.
        laneMask.x &= laneBit ? ballot.x : ~ballot.x;
        laneMask.y &= laneBit ? ballot.y : ~ballot.y;
        laneMask.z &= laneBit ? ballot.z : ~ballot.z;
        laneMask.w &= laneBit ? ballot.w : ~ballot.w;
    }

    return laneMask;
}

```

Here is the algorithm that we've just described to capture the lane mask.

As well as the election process to determine which lanes will write to the histogram.

\*click\*

\*click\*

```
uint4 ballot = WaveActiveMatch(bucket, 8);  
uint matchCount = WaveBallotBitCount(ballot);  
uint electedLane = WaveBallotLSB(ballot);  
  
[branch]  
if (WaveGetLaneIndex() == electedLane)  
{  
    InterlockedAdd(OutputHistogram[bucket], matchCount);  
}
```

Here is the algorithm that we've just described to capture the lane mask.

As well as the election process to determine which lanes will write to the histogram.

\*click\*

\*click\*

# 256 Bucket Histogram Benchmark

	Global Memory	+Groupshared	+Wave Match
<b>Intel(R) Iris(R) Xe Graphics - 16 Million Elements, 256 Buckets</b>			
Minimal Collisions	30ms	1.5ms	2ms
Maximum Collisions	30ms	3.5ms	2ms
<b>NVIDIA GeForce RTX 3070 Ti - 16 Million Elements, 256 Buckets</b>			
Minimal Collisions	2ms	0.3ms	0.4ms
Maximum Collisions	2ms	2ms	0.4ms

If we look at some performance numbers.

You'll note that regardless of the number of collisions, you get the same performance characteristics.

2ms on my intel integrated GPU and 0.4ms on my RTX 3070 Ti.

Using wave intrinsics you can be confident that you won't be surprised by sudden spikes in your histogram calculation.

# Bonus Slides: Things I Dream About



# Things I Dream About

- **More expressive intrinsics**
  - HLSL Shuffle Up/Shuffle Down/Xor Shuffle
    - Already available with subgroups operations!

Especially in HLSL, it would be nice to be able to access more expressive intrinsics.

Thankfully, for my primary use case I can use inline Spir-V to access these intrinsics.



# Things I Dream About

- **More expressive intrinsics**
  - HLSL Shuffle Up/Shuffle Down/Xor Shuffle
    - Already available with subgroups operations!
- **Don't waterfall loop my shuffles (please)**

Additionally, it can be difficult to develop more complex intrinsics when many intrinsics degenerate to a waterfall loop in a variety of cases.

Understandably, I can see that hardware limitations might lead to needing to use a waterfall loop.

But this is why I called them dreams.

# Things I Dream About

- **More expressive intrinsics**
  - HLSL Shuffle Up/Shuffle Down/Xor Shuffle
    - Already available with subgroups operations!
- **Don't waterfall loop my shuffles (please)**
- **More shader compilation tools**
  - RGA Good!
    - No source correlation for SPIR-V Debug Info
    - Currently using a modified version of LLPC to get source code correlation
  - Intel disassembly on Windows?



Finally, more shader compilation tools would be incredible.

RGA is a great tool, but the lack of source code correlation with graphics Spir-V makes it harder to correlate the disassembly to the source code that generated it.

You may have noticed that I have source correlations in the disassembly I've presented. I use a modified version of DXC and LLPC to allow for debug instructions to make it all the way to the binary.

On top of that, more tools for different vendors would be incredible.

Perhaps something like that exists! Let me know if that's the case.