



Arm[®] GPU Best Practices

Revision 3.3

Developer Guide

Non-Confidential

Copyright © 2017, 2019–2024 Arm Limited (or its affiliates).
All rights reserved.

Issue 09

101897_0303_09_en



Arm® GPU Best Practices Developer Guide

This document is Non-Confidential.

Copyright © 2017, 2019–2024 Arm Limited (or its affiliates). All rights reserved.

This document is protected by copyright and other intellectual property rights.

Arm only permits use of this document if you have reviewed and accepted [Arm's Proprietary Notice](#) found at the end of this document.

This document (101897_0303_09_en) was issued on 2024-10-10. There might be a later issue at <https://developer.arm.com/documentation/101897>

The product revision is 3.3.

See also: [Proprietary notice](#) | [Product and document information](#) | [Useful resources](#)

Start reading

If you prefer, you can skip to [the start of the content](#).

Intended audience

This guide is for experienced software engineers who want to optimize the performance of the graphics in their application.

Inclusive language commitment

Arm values inclusive communities. Arm recognizes that we and our industry have used language that can be offensive. Arm strives to lead the industry and create change.

We believe that this document contains no offensive language. To report offensive language in this document, email terms@arm.com.

Feedback

Arm welcomes feedback on this product and its documentation. To provide feedback on the product, create a ticket on <https://support.developer.arm.com>.

To provide feedback on the document, fill the following survey: <https://developer.arm.com/documentation-feedback-survey>.

Contents

1. Overview.....	7
1.1 Before you begin.....	7
1.2 Arm GPU datasheet and performance counters.....	7
1.3 The graphics rendering pipeline.....	8
2. Optimization basics.....	10
2.1 Optimization process.....	10
2.2 Basic optimization checklist.....	14
2.3 Memory bandwidth.....	15
2.4 Converting from desktop to mobile.....	16
3. Optimizing application logic.....	17
3.1 Basic application optimizations.....	17
3.2 Draw call batching best practices.....	18
3.3 Draw call culling best practices.....	19
3.4 Optimizing the draw call render order.....	20
3.5 Avoid using depth prepasses.....	23
3.6 OpenGL ES GPU pipelining.....	23
3.7 OpenGL ES Separate Shader Objects.....	25
3.8 Vulkan GPU pipelining.....	25
3.9 Vulkan pipeline synchronization.....	27
3.10 Pipelined resource updates.....	30
3.11 Optimize attachment grouping.....	32
3.12 Queries.....	32
4. CPU overheads.....	34
4.1 Compiling shaders in OpenGL ES.....	34
4.2 Pipeline creation in Vulkan.....	34
4.3 Allocating memory in Vulkan.....	36
4.4 OpenGL ES CPU memory mapping.....	36
4.5 Vulkan CPU memory-mapping.....	37
4.6 Command pools for Vulkan.....	39
4.7 Optimizing command buffers for Vulkan.....	40

4.8 Secondary command buffers.....	41
4.9 Optimizing descriptor sets and layouts for Vulkan.....	42
5. Vertex shading.....	44
5.1 Basic vertex shader optimizations.....	44
5.2 Index draw calls.....	45
5.3 Index buffer encoding.....	46
5.4 Index sparsity.....	47
5.5 Attribute precision.....	48
5.6 Attribute layout.....	49
5.7 Varying precision.....	50
5.8 Triangle density.....	51
5.9 Instanced vertex buffers.....	51
6. Tessellation, geometry shading, and tiling.....	53
6.1 Tessellation.....	53
6.2 Geometry shading.....	54
6.3 Tiling and effective triangulation.....	55
7. Fragment shading.....	56
7.1 Basic fragment shader optimizations.....	56
7.2 Efficient render passes with OpenGL ES.....	57
7.3 Efficient render passes with Vulkan.....	58
7.4 Multisampling for OpenGL ES.....	59
7.5 Multisampling for Vulkan.....	60
7.6 Multipass rendering.....	61
7.7 HDR rendering.....	65
7.8 Stencil updates.....	66
7.9 Blending.....	66
7.10 Transaction elimination.....	67
7.11 Variable rate shading.....	69
8. Buffers and textures.....	72
8.1 Buffer update for OpenGL ES.....	72
8.2 Robust buffer access.....	73
8.3 Staging buffers.....	74
8.4 Texture sampling performance.....	75

8.5 Anisotropic sampling performance.....	77
8.6 Texture and sampler descriptors.....	79
8.7 sRGB textures.....	81
8.8 AFBC textures.....	82
8.9 AFBC textures for Vulkan.....	82
8.10 AFRC.....	85
8.11 Runtime compression.....	86
9. Compute shading.....	88
9.1 Image processing.....	88
9.2 Workgroup sizes.....	90
9.3 Shared memory.....	91
10. Shader code.....	93
10.1 Minimize precision.....	93
10.2 Check precision.....	95
10.3 Vectorized arithmetic code.....	95
10.4 Vectorize memory access.....	96
10.5 Manual source code optimization.....	97
10.6 Generating SPIR-V.....	98
10.7 Instruction caches.....	99
10.8 Uniforms.....	100
10.9 Uniform subexpressions.....	101
10.10 Uniform control-flow.....	101
10.11 Branches.....	102
10.12 Discards.....	103
10.13 Atomics.....	104
11. Ray tracing.....	105
11.1 Acceleration structures.....	105
11.2 Efficient ray tracing.....	107
11.3 Ray query.....	110
11.4 Ray tracing pipeline.....	113
12. System integration.....	116
12.1 Using EGL buffer preservation in OpenGL ES.....	116
12.2 Android blob cache size in OpenGL ES.....	117

12.3 Optimizing the swapchain surface count for Vulkan.....	118
12.4 Optimizing the swap chain surface rotation for Vulkan.....	119
12.5 Optimizing swapchain semaphores for Vulkan.....	120
12.6 Window buffer alignment.....	121
12.7 Vulkan private data.....	121
12.8 Vulkan extensions to avoid.....	122
Proprietary notice.....	123
Product and document information.....	125
Product status.....	125
Revision history.....	125
Conventions.....	128
Useful resources.....	131

1. Overview

Our guide provides recommendations on how to optimize your applications when developing for Arm Immortalis and Mali GPUs. The recommendations are intended for established developers who want to begin working with Arm GPUs.

1.1 Before you begin

Real-world applications can be complicated and there are always exceptions to this generalized advice. Therefore, we recommend that you regularly make measurements of any applied optimizations to check that they are performing as intended on your target devices.



This guide is aimed at established developers. We therefore assume that you have experience as a graphics developer, and that you have a general familiarity with the underlying APIs.

1.2 Arm GPU datasheet and performance counters

Arm GPUs have evolved as technology has improved. Over time, new standards have been supported, features added, and performance and efficiency improved.



The following recommendations are correct for Midgard, Bifrost, Valhall and 5th Gen, which includes the Mali™-G725 and Immortalis™-G725. The recommendations are also correct for DDK drivers up to the r51 release of OpenGL ES and also the r51 release of Vulkan.

Arm GPU datasheet

For an overview of the different Arm GPU capabilities, see [Arm GPU Datasheet](#).

Arm GPU performance counters

Arm GPUs implement a comprehensive range of performance counters that enable you to closely monitor GPU activity as your application runs. To help you to identify the cause of bottlenecks or inefficient workloads, the Arm Streamline profiler can capture and present Arm CPU and GPU performance counter data as a series of easy-to-read charts.

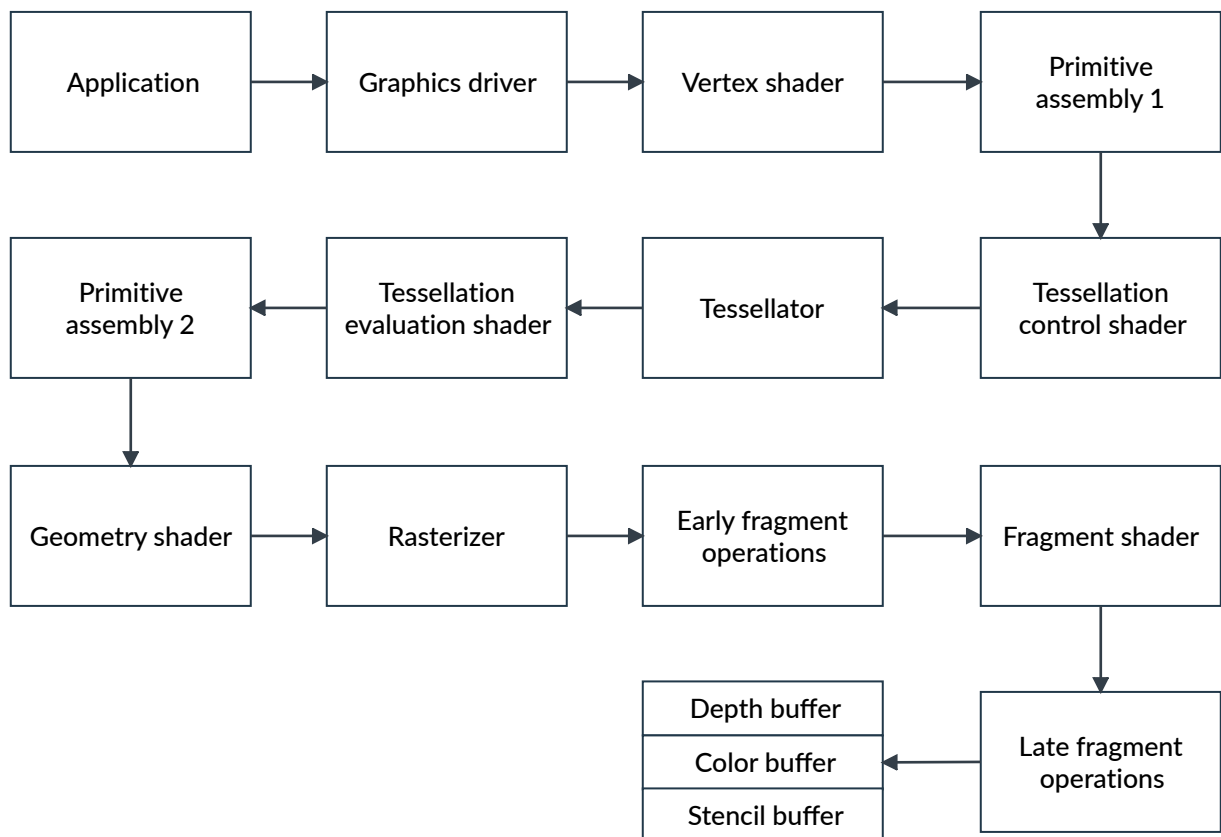
For an overview of the performance counters, see [Arm GPU Performance Counters](#).

1.3 The graphics rendering pipeline

Graphics processing can be represented as a pipeline of processing stages that includes the application, the graphics driver, and the various hardware stages inside the GPU. Most stages follow a strict architectural pipeline, with outputs from one stage becoming the inputs into the next stage.

The following figure shows the graphics pipeline beginning at the application, and ending at the depth, color, and stencil buffers:

Figure 1-1: The stages of the graphics pipeline



Compute shaders are the exception to this strict pipeline because shader results are written back to system memory. Any stage of the pipeline that can consume a buffer or texture resource can use the compute shader outputs.

Guide structure

This guide has been structured to follow the graphics rendering pipeline. Each topic gives recommendations that can be applied to workloads that are running in that pipeline stage.

There are some generic topics, such as shader program authoring advice, that applies to multiple pipeline stages.

We provide explanations for our recommendations, along with technical points that can be considered during development. Where possible, we document the likely impact of ignoring the advice and also which techniques to use when debugging.

We provide relatively specific advice in this guide. For example, we tell you not to use `discard` in fragment shaders. There are use-cases where you must use a feature that we would otherwise tell you not to use, because the feature is required for some algorithms. Therefore, make your implementation decisions as per your specific use-case. However, always consider that there could be an underlying performance risk with certain techniques.

2. Optimization basics

This chapter covers general information on how to optimize for mobile GPUs.

2.1 Optimization process

It is important to have a good process to work out where needs optimizing and closer inspection in your code and shaders. While in theory it is important to follow best practices everywhere, there are some aspects that are more important than others.

Prerequisites

You must understand the following concepts:

- Application code
- Vertex shaders
- Fragment shaders

Set a computation budget

First, try to set a computation budget to measure your performance against. There are multiple measures you can set a budget for. For example:

- Frame time
- Triangle count
- Application processor cycles
- Vertex processing cycles
- Fragment processing cycles
- Memory bandwidth

While frame time is the key budget, measuring the other components helps you to locate where other bottlenecks are. In turn, highlighting where extra performance gains can be found.

Example: Fragment shader budget

The calculation to work out the fragment shader budget is:

- Multiply the number of GPU shader cores by the GPU clock speed. This gives the maximum theoretical number of cycles per-second. Multiply this by 0.8 to give a more realistic number of available fragment processing cycles per-second. This is result A.
- Multiply the frame height by the frame width. This gives the number of pixels per frame. Now multiply this by the required frame rate. This gives the number of pixels required per-second. To take account of average overdraw, multiply this number by 2.5. This gives the number of fragments required per second. This is result B.
- Divide the value of result A by the value of result B. The result is the average number of cycles a fragment shader can take.

However, you do not have to make all of your fragment shaders take this number of cycles. For example, you can use longer, more complex, shaders on objects closer to the camera as well as shorter, and less complex, shaders on more distant objects.

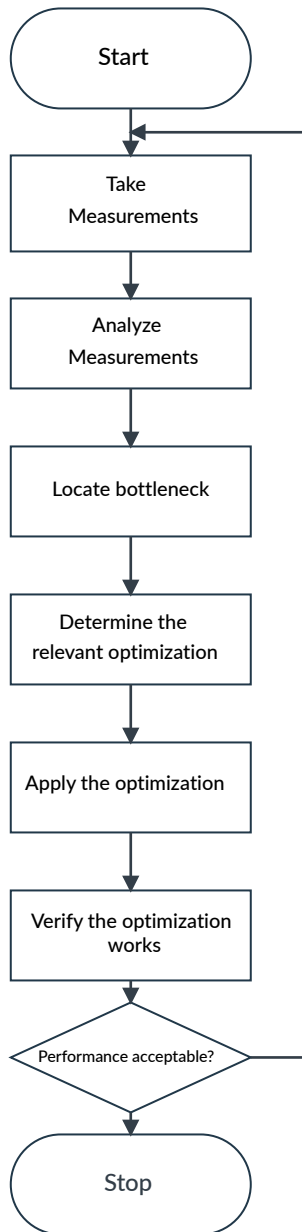
Do not assume the number of fragment processing cycles equals the number of fragment processing instructions. The processors in Arm GPUs can do many operations per cycle. The Mali Offline Compiler can be used to determine the number of cycles a shader requires.

Basic optimization workflow

Try repeating the following process until your application fulfills its computation budget:

- Take measurements.
- Locate the bottleneck.
- Determine and apply an optimization.
- Verify the optimization improves performance.

Measurements must be taken with real hardware, with the Arm GPU you want to test the performance on. Verifying performance gains from an optimization is important, as there can be multiple bottlenecks, and some optimizations are context-dependent.

Figure 2-1: Optimization process

Consider the following areas when determining where an optimization is needed:

- Application code.
- Misuse of API.
- Use of blocking API calls.
- Vertex processing.
- Triangle setup.

- Fragment processing.
- Memory bandwidth.

Optimization pitfalls to avoid

There are several common pitfalls to be aware of when optimizing:

- Do not use frames per second as a measurement. Use frame times, because improvements will be linear and easier to interpret. Frames per second gives a skewed inverse measurement.
- Do not assume that because your GPU bottleneck is fixed that your CPU is running fast enough, or too slow. Bottlenecks move between processors, so be sure to measure all of them to monitor performance.
- Do not keep optimizing fragment shaders when the application is either vertex shader or application-bound, or vice versa. Work on the bottleneck. Indeed, you can improve performance by moving work to the less heavily loaded parts of the pipeline.
- Make sure VSYNC is off when taking measurements.

Negative impacts of incorrect optimization process

An inefficient optimization process means:

- You cannot find the correct performance bottlenecks that need fixing.
- You cannot reach the target performance that you want.

Debugging tools to use in your optimization process

Arm Mobile Studio is available with various tools to help you in your optimization process:

- Streamline provides the performance profile needed for deep-dive analysis of your application's CPU usage, GPU usage, and memory bandwidth. Its Performance Advisor feature gives you easy-to-read reports with actionable advice about how to optimize your application.
- Graphics Analyzer lets you debug OpenGL ES and Vulkan API calls in your application. It can visualize the GPU rendering and overdraw, as well as the API state and data resources used, for each draw call.
- Mali Offline Compiler performs static analysis of your shader programs, for any current Arm GPU, giving you resource usage and an approximate performance cost.
- Frame Advisor gives analysis of a problem frame, showing the API calls, rendering and comprehensive geometry metrics to show what might be slowing your application.

For more about Arm Mobile Studio, along with useful tutorials and training videos, please visit:

Web Link: [Arm Mobile Studio](#)

2.2 Basic optimization checklist

This chapter contains a set of standard things to check to ensure your application is getting the optimization basics correct.

Prerequisites

You must understand the following concepts:

- Application code

Checklist

There are many basic optimization approaches that you can try. These include:

- Ensure target display settings are correct: Check that you have requested a compatible resolution and color depth for the screen, and a drawing surface format that works with the framebuffer. Also check that what you have received from the API is what you requested.
- Build for the correct architecture: Work out which you are targeting and accept the limitations of older chips if necessary.
- Use the latest tools: Improvements are always being added. Remember to rebuild after updating your tools or targets.
- Use Neon or other hardware facilities that are available: Be aware of the abilities of the hardware you are targeting and be sure to take full advantage of their features.
- Remove debugging information: Optimize on release builds, with a minimum of `printf` calls and no more than one `glGetError()` call per-frame.
- Use VSYNC - but *not* while doing performance measurements. While enabling VSYNC reduces power consumption in your final build, it ruins your optimization measurements.
- Use mipmapping for textures in 3D scenes: Mipmapping is a rare optimization that improves both quality *and* performance. See [Texture sampling performance](#) for more information.
- Use texture compression, ideally ASTC: See [Texture sampling performance](#) for more information.
- Use vertex buffers to save bandwidth: For more details on the best use of vertex buffers, see [Instanced vertex buffers](#).
- Ensure you are not CPU-bound. See [Optimizing application logic](#) and [CPU overheads](#) for more information.

Avoid optimization pitfalls

There are also some things to avoid for optimal running performance:

- Avoid calls that stall the graphics pipeline. `glReadPixels()`, `glCopyTexImage()`, and `glTexSubImage()` cause the entire image to be rendered if there are queued texture references. See [OpenGL ES GPU pipelining](#) for more information.
- Do not compile shaders within a frame. It is more efficient to compile them on application start, or ship pre-compiled shaders. If necessary, compile them asynchronously on a background thread. See [Compiling Shaders in OpenGL ES](#) for more information.

Negative impacts of non-optimized applications

Running non-optimized code results in increased power use and a less pleasant experience for users.

Related information

The [Arm GPU Training Series](#) on the Arm Software YouTube Channel.

2.3 Memory bandwidth

Memory bandwidth requires a lot of power, so it is very restricted in mobile devices compared to desktop systems.

Prerequisites

You must understand the following concepts:

- Memory
- Caching

Reduce memory bandwidth

Memory bandwidth can easily become a bottleneck limiting the performance of your application. It is a shared resource, so using too much can affect the performance of the system in unpredictable ways.

Accessing data in caches reduces power usage and improves performance. If you must read from memory a lot, use techniques such as mipmapping, and texture compression to ensure your data is cache-friendly. Other methods to reduce memory bandwidth use include:

- Activate back face culling.
- Use view frustum culling.
- Ensure textures are not too large.
- Use a texture resolution that fits the object on screen.
- Use low bit depth textures where possible.
- Use lower resolution textures if the texture does not contain sharp detail.
- Use fewer textures. Pre-bake effects in, utilize unused channels, and use texture atlases.
- Use texture compression.
- Only use trilinear filtering on specific objects.
- Only use anisotropic filtering on specific objects.
- Utilize Level of Detail (LOD) for meshes and mipmapping for textures.

Avoid overdraw

Overdraw causes excess memory bandwidth use, so it should be limited as much as possible.

Negative impacts of excess bandwidth use

Memory bandwidth impacts everything, which means that issues can appear anywhere. If it seems that a processor is the bottleneck, but optimizations make no difference, it may in fact be bandwidth issues.

Excess memory bandwidth use causes excess power use, so battery life becomes impacted.

Debugging memory bandwidth issues

To track memory bandwidth use, use [Arm Streamline](#) to profile your application.

2.4 Converting from desktop to mobile

When converting an application from desktop to mobile, there are several standard considerations to improve how the application runs.

Prerequisites

You must understand the following concepts:

- Tile-based graphics processing

Use appropriate graphics assets

When porting to mobile, the whole [Basic optimization checklist](#) applies along with reducing [Memory bandwidth](#), but also ensure you use graphic assets that are appropriate for your platform.

It is also recommended on mobile to:

- Draw non-transparent objects in front-to-back order, before any transparent objects.
- Use texture compression.

Avoid inappropriate graphics assets

Large assets slow your application.

- Avoid high numbers of triangles. See [Triangle density](#).
- Avoid long shaders.
- Avoid high bit depth and high-resolution textures.
- Avoid overly conservative Vulkan dependencies. These can block the pipeline on the tile-based GPUs used in mobile. See [Vulkan pipeline synchronization](#) for how to handle different shader stages dependencies.

Negative impacts of not considering mobile

Running desktop content on mobile leads to unnecessarily fast use of power, and a less smooth experience for users.

3. Optimizing application logic

It is important to optimize application logic to minimize the CPU load that the graphics driver generates.

The CPU is used to process the application logic, drive the graphics stack, and run the graphics driver. The graphics driver is the first potential source of performance issues that your application can run into.

3.1 Basic application optimizations

There are several standard practices to ensure the work being done on the CPU does not slow down your application.

Prerequisites

You must understand the following concepts:

- Application code
- Caching

Common optimizations

When programming there are many ways of implementing solutions. It is worth spending time thinking about whether there is a simpler way or an acceptable approximation. For instance, for 3D collision detection it can be useful to use hierarchical bounding boxes to test against.

Other important considerations include:

- Align data: This makes it more cacheable and faster to copy between the application and the graphics driver.
- Optimize loops: Move computations out of the loop where possible. Simplify the code and minimize the blocks of data worked on. Where possible avoid branches and function calls, especially in inner loops. Some level of unrolling can be useful.
- Use vector instructions: SIMD technologies like Neon allow the processor to do multiple calculations simultaneously.
- Use fast data structures: Keep data together rather than using pointers, so it can be copied from memory efficiently. Optimize the data elements to fit in a cache.
- Use multi-processing: Working on different threads can allow work to be done in parallel, meaning a quicker execution time.

These tips, while often handy in themselves, can also mean the compiler is able to make further optimizations.

API optimizations

The main API optimization is to minimize draw calls, combining what you can. To aid optimization, try combining textures in a texture atlas and multiple texture atlases together. Text textures can be combined in a font atlas. For more on optimizing draw calls see [Draw call batching best practices](#), [Draw call culling best practices](#) and [Optimizing the draw call render order](#).

Try minimizing API state changes. You can reduce state changes by:

- Grouping objects using the same texture, allowing shared texture binding calls.
- Grouping objects using the same state together, allowing shared state setting calls.
- Using texture atlases, allowing more grouping and reduced numbers of texturing binding calls.
- Keeping track of the current state, allowing only making necessary state changes.
- For OpenGL ES, removing redundant `glEnable()` and `glDisable()` calls if no state change occurred.

Avoid the graphics pipeline stalling

Do not make API calls that cause the image to be rendered unnecessarily.

Negative consequences of inefficient application code

Inefficient application code results in a slower program.

3.2 Draw call batching best practices

Committing draw calls to the command stream is an expensive operation in terms of graphics driver overheads on the application processor. Draw call runtime costs are higher on OpenGL ES than on Vulkan.

Prerequisites

You must understand the following concepts:

- Draw calls.
- Instancing.

Batch rendering and command buffers to reduce draw calls

Draw calls that contain few vertices and fragments are quicker to process on the GPU, compared to the CPU time that is required to dispatch the workload. Therefore, the performance of the application is CPU-limited because the CPU is unable to keep the GPU busy.

These issues encourage draw call batching. Draw call batching merges rendering for multiple objects that use the same render state into a single draw call. This reduces the total number of draw calls required to render each frame, which lowers the CPU processing cost and energy consumption.

How to optimize draw call batching

Try using the following optimization techniques:

- Batch objects to reduce the draw call count.
- Use batches, even if not CPU-limited, to reduce system power consumption.
- Use instanced draw calls when drawing multiple copies of a mesh. Instancing allows the application to test and cull individual instances that are, for example, outside the view frustum.
- For OpenGL ES, aim for fewer than 500 draw calls per frame.
- For Vulkan, aim for fewer than 1000 draw calls per frame.



Because CPU performance can vary widely between chipsets, these draw call count recommendations are approximate guidelines for current Arm GPU hardware.

Behaviors to avoid when draw call batching

Arm recommends that you:

- Do not make batches so large that frustum culling, and sorting for front-to-back rendering order, become inefficient.
- Do not request many small draw calls for rendering, for example single points or quads, without batching.

Negative impacts of unoptimized draw call batches

The different types of impact you can see are:

- Higher application CPU load, because of a high draw call count.
- A reduction in overall performance if the application is CPU bound.

Debugging when draw call batching

Try the following debugging tips:

1. Profiling the application CPU load.
2. Counting the number of draw calls per frame.

3.3 Draw call culling best practices

The fastest draw calls that an application can process are the ones that are discarded before they reach the graphics API.

Prerequisites

You must understand the following concepts:

- Culling.

- Primitives.
- Vertex shading.

Minimizing the number of draw calls

Application culling of entire meshes is always more efficient than per-primitive culling in the GPU, because it can exploit scene knowledge. GPU culling can only be performed after the primitive clip-space coordinates are known. Therefore, the vertex shading must be executed even if the primitive is eventually culled.

How to optimize culling draw calls

Try using the following optimization techniques:

- Cull objects that are out of frustum on the CPU. For example, by using bounding box frustum checks.
- Cull objects that are known to be occluded on the CPU. For example, by using portal culling.
- Experiment to find a balance between batch size and culling efficiency.

Outcome to avoid when draw call culling

Do not send every object to the graphics API.

Negative impacts on the CPU and GPU

Unoptimized draw call culling can lead to the following problems:

- Higher application CPU load, because of unnecessary draw calls.
- High GPU vertex shading and memory bandwidth, because of redundant vertex shading.

Debugging draw call culling problems

Try the following debugging tips:

- Profile the application CPU load.
- Use GPU performance counters to verify tiler primitive culling rates. Expect around 50% of the triangles to be culled because they are back-facing inside the frustum. Higher culling rates can indicate that the application needs an improved draw call culling approach.

3.4 Optimizing the draw call render order

When rendering draw calls, the GPU can reject occluded fragments efficiently using the early ZS test. An efficient draw call order can maximize the benefit of the early ZS test by removing as many occluded fragments as possible.

Prerequisites

You must understand the following concepts:

- Command buffers.
- Draw calls.

- Early and late ZS testing.
- Render targets.
- Geometry buffer (G-buffer).

Increasing culling rates using early ZS testing and Forward Pixel Kill

To get the highest fragment culling rate from the early ZS unit, first render all opaque meshes in a front-to-back render order. To ensure blending works correctly, render all transparent meshes in a back-to-front render order, over the top of the opaque geometry.

All Arm GPUs since the Mali-T620 GPU include the Forward Pixel Kill (FPK) optimization. FPK provides automatic hidden surface removal of fragments that are occluded, which early ZS testing does not kill.

The removal of occluded fragments occurs due to the use of a back-to-front render order for opaque geometry. However, do not rely on the FPK optimization alone. An early ZS test is always more energy-efficient, consistent, and works on older Arm GPUs that do not include hidden surface removal.

Although it is best to utilize early ZS rather than FPK, it can be useful to know what stops FPK working. Situations where FPK commonly fails include draw calls using:

- Alpha blended transparency
- Shader programmatic framebuffer access.
- Late ZS testing.
- Small triangles.

Fragment Prepass

From Mali-G725 and Immortalis-G925 onwards, there is a new feature called Fragment Prepass. Fragment Prepass reliably removes overdrawn fragments, removing redundant shading which improves performance and energy efficiency. Fragment Prepass' culling efficiency is not sensitive to application draw order, enabling applications to disable front-to-back sorting of objects in software. This reduces the CPU cost of issuing draw calls to the GPU.

The prepass handles many common usage patterns. This includes arbitrary depth and stencil testing, write-only side-effects, and modifying coverage with alpha-to-coverage or shader code use of `discard`. Attachment read-back, including for depth and stencil, is handled well, meaning in OpenGL ES framebuffer object use without clearing initial attachments, or in Vulkan initial attachments using `loadOp=LOAD`.

However, you should still follow a number of best practices to get the most out of Fragment Prepass.

The prepass terminates at the first incompatible draw in a tile. Incompatible draw calls are:

- Transparent draws that write ZS
- Draws that write ZS-only and are after a transparent draw.

- Draws with side-effects where the value read from memory is used for color computation.

Also, incompatible are draws that have shader execution impacted by their coverage mask, for example, using centroid mapping, using coverage in shader code, or reading from a color attachment for cross-lane operations in shader code. Draws should write to a superset of the previous draws' render targets.

Avoid mid-frame depth clears, and avoid late-Z where possible, and failing that make late-Z calculations small. Try to keep the G-buffer at 128 bits per pixel (bpp) or less, and if not keep it within 256bpp.

Many general GPU best practices still apply and are even more important with the prepass:

- Simplify your geometry meshes as much as possible.
- Render transparents after opaques.
- Minimize the computations needed to compute vertex position.
- Minimize the computations needed to compute fragment depth and/or coverage, if it is programmatically determined.
- Minimize the number of draws that would terminate the prepass, and if unavoidable draw them as late as possible.
- Minimize use of modifiable coverage (discard, alpha-to-coverage).

How to optimize draw call render ordering

Try using the following optimization techniques:

- Render opaque objects in a front-to-back order.
- Render opaque objects with blending disabled.

Behaviors to avoid when optimizing draw call render ordering

Arm recommends that you:

- Do not use `discard` in the fragment shader, as it forces a late ZS test.
- Do not use alpha-to-coverage, as it forces a late ZS test.
- Do not write to fragment depth in the fragment shader, as it forces a late ZS test.

The negative impact of an unoptimized draw call render order

Using a suboptimal draw call render order and late ZS testing incurs a higher fragment shading load. This occurs because of visually occluded fragments that were not killed before fragment shading.

Debugging draw call render order problems

Try the following debugging techniques:

- Render your scene without your transparent elements. Now use GPU performance counters to check the number of fragments that is being rendered per output pixel. If the number of fragments is higher than 1, you have opaque fragment overdraw that the early ZS test can remove.

- Use the GPU performance counters to check the number of fragments that require late ZS testing. Also check the number of fragments that late ZS testing kills.

3.5 Avoid using depth prepasses

A depth prepass is a common technique in PC and console games development. It is used in scenarios where there might be significant overdraw and expensive fragment shaders. It is also used where you cannot reliably get front-to-back sorted opaque geometry.

Prerequisites

You must understand the following concepts:

- Early ZS testing.
- Render passes.

An optimization that reduces performance

The aim of a depth prepass is to quickly set the depth for all geometry, without incurring fragment shading cost. The color shading pass, which follows the prepass, only shades the fragments where the depth exactly matches. This gives the ideal of one fragment shader invocation per pixel, minimizing the amount of redundant fragment processing that occurs. However, to gain this behavior, depth prepasses must double the draw call count and the processed vertex count.

Arm GPUs already include optimizations, such as Forward Pixel Kill (FPK) hidden surface removal, to reduce redundant fragment processing automatically. Therefore, the performance cost of the additional draw calls, vertex shading, and memory bandwidth needed to implement a depth prepass usually outweigh the benefits of reduced overdraw.

Behaviors to avoid when draw call batching

Do not use depth prepass algorithms to remove any fragment overdraw.

Negative impacts of using depth prepasses

The impact on performance when using depth prepasses are:

- The CPU incurs a higher load due to duplicated draw calls.
- There is a higher vertex shading and memory bandwidth cost due to duplicated geometry.

3.6 OpenGL ES GPU pipelining

OpenGL ES exposes a synchronous rendering model to users of the API, despite using asynchronous execution on the GPU. Users must be aware of this to avoid stalling the GPU rendering pipeline.

Prerequisites

You must understand the following concepts:

- The differences between synchronous and asynchronous execution.
- Pipeline draining

Keeping the GPU busy, fences, and queries

To get the best performance, the application must use the synchronous OpenGL ES API to keep the GPU busy during asynchronous execution of the workloads. Therefore, avoid using operations that cause the driver to drain the GPU pipeline and starve the GPU of work.

How to optimize OpenGL ES GPU pipelining on Arm GPUs

Try using the following optimization techniques:

- Do not use API calls that cause the driver to block and wait for the GPU.
- Pipeline use of fences and query objects, only waiting on them when you know the result will already be available.
- Pipeline resource updates, avoiding modifications to textures and buffers that are referenced by in-flight draws.
- Use `GL_MAP_UNSYNCHRONIZED` to enable the use of `glMapBufferRange()` on a buffer that is referenced by an in-flight draw.
- Use asynchronous `glReadPixels()` calls to read data into a pixel buffer object.

Behaviors to avoid when optimizing OpenGL ES GPU pipelining on Arm GPUs

Arm recommends that you:

- Avoid using the following synchronous OpenGL ES operations:
 - `glFinish()`.
 - `glReadPixels()`.
 - `glCopyTexImage()`, while the target texture is still referenced by an in-flight draw call.
 - `glTexSubImage()`, while the target texture is still referenced by an in-flight draw call.
 - `glMapBufferRange()`, without using `GL_MAP_UNSYNCHRONIZED`, while the target buffer is still referenced by an in-flight draw call.
- Avoid using `glMapBufferRange()` with either `GL_MAP_INVALIDATE_RANGE`, or `GL_MAP_INVALIDATE_BUFFER`. Both of these flags can trigger the creation of a resource ghost on some Arm GPU driver versions.
- Avoid using `glFlush()` to split render passes because the driver automatically flushes when needed.

Negative impacts of inefficient OpenGL ES GPU pipelining

The different types of impact you can see are:

- If the pipeline is drained, the GPU becomes partially idle during the resulting bubble, causing a loss of performance.
- Depending on the interaction with the system Dynamic Voltage and Frequency Scaling power management logic, there might be some performance instability.

Debugging OpenGL ES issues for Arm GPUs

CPU and GPU activity can be monitored using system profilers. For example, the [Arm Streamline Performance Analyzer](#). Pipeline drains are visible as periods of busy time oscillating between the CPU and GPU, without the CPU or the GPU being fully utilized.

3.7 OpenGL ES Separate Shader Objects

OpenGL ES allows the use of Separate Shader Objects (SSOs), which allow a program pipeline to be built at draw time rather than relying on full program linkage.

Prerequisites

You must understand the concept of Shader Objects.

Only use SSOs if necessary

Using SSOs prevents several link-time optimizations from being implemented, so the resulting shaders can execute more slowly on the GPU. However, in some situations they can avoid a substantial increase in the number of shader programs, which can help reduce CPU overhead.

If you do use SSOs, optimize your shader programs to manually implement the link-time optimizations. For example, remove redundant outputs from a pipeline stage if later pipeline stages do not consume them. For this use case, the developer takes on the added optimization responsibility that the compiler and driver cannot do.

Behaviors to avoid

Arm recommends that you do not use SSOs.

How to optimize SSOs

If SSOs are needed, remove redundant computation, and ensure shader stages only pass on the values that are needed by later pipeline stages.

Negative impacts of SSOs

Some shader optimizations are not available, and rendering performance may be impacted.

3.8 Vulkan GPU pipelining

Arm GPUs can run either compute or vertex work while the fragment processing from another render pass is running. To ensure high performance, applications must not unnecessarily create bubbles in this pipeline.

Prerequisites

You must understand the following concepts:

- Using command buffers.

- The different shader stages.

Pipeline bubbles



Note

For Arm GPUs, it is important to overlap the vertex or compute work from one render pass with the fragment work from earlier render passes.

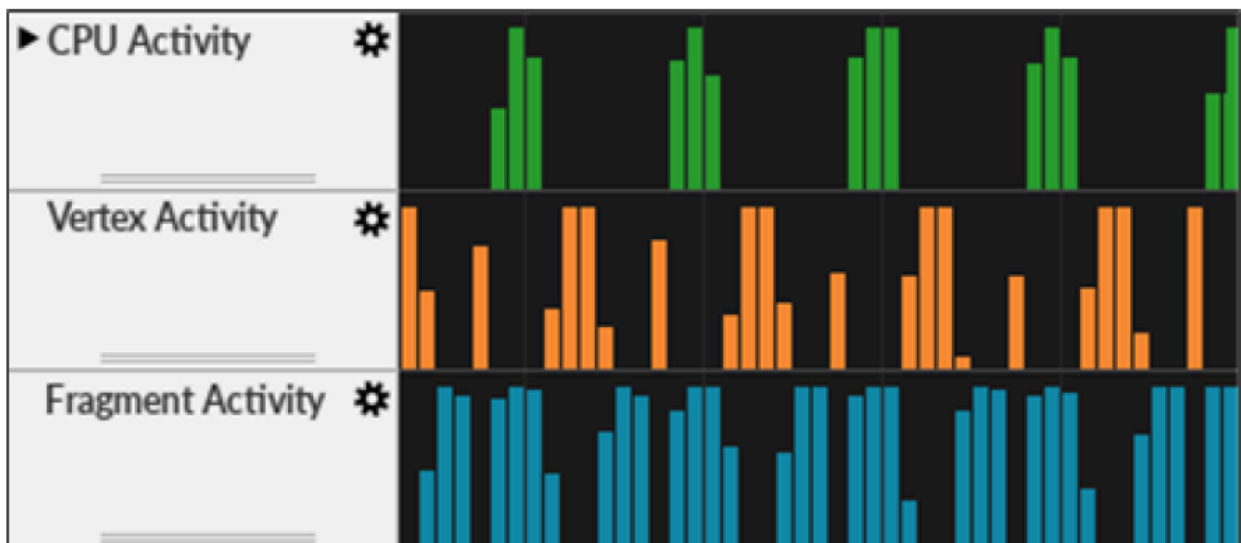
When using Vulkan, the following reasons can lead to pipeline bubbles in applications:

1. Command buffers not submitted often enough: Not submitting command buffers often enough reduces the amount of work in the GPU processing queue. Restricting the possible scheduling opportunities.
2. Data dependency: For example, consider two render passes N and M. Render pass M occurs at a later stage in the pipeline. Data dependency results when N is consumed earlier in the pipeline by M. Data dependency causes a delay during which enough work must be done to hide the latency in the result generation.

In the following image, which is taken from our Streamline profiler, you can see that we are not hitting the targeted 60FPS. Instead, the CPU and the vertex and fragment activity in the GPU all have idle time.

This pattern of work and idle time oscillating between the processors is a good indicator of the presence of API calls, or data dependencies that are limiting GPU scheduling opportunities.

Figure 3-1: Pipeline bubble example



How to prevent pipeline bubbles

To prevent the occurrence of pipeline bubbles:

- Submit command buffers for processing frequently, for example, for each major render pass in a frame.
- If you have a case that causes a bubble, experiment with bubble filling techniques. For example, by inserting independent workloads between the two render passes.
- Consider generating dependent data in an earlier pipeline stage than the stage that consumes it. For example, the compute stage is suited to generate input data for the vertex shading stage. The fragment stage is poorly suited because it occurs later than the vertex shading stage in the pipeline.
- Consider processing dependent data later in the pipeline. For example, fragment shading consuming output from other fragment shading works better than compute shading consuming fragment shading.
- Use fences to asynchronously read back data to the CPU. Do not block synchronously and cause the pipeline to drain.

Actions to avoid while optimizing the GPU pipeline

Arm recommends that you:

- Do not unnecessarily wait for GPU data anywhere in the pipeline.
- Do not wait until the end of the frame to submit all of the render passes.
- Do not create any backwards data dependencies in the pipeline without sufficient intervening work to hide the result generation latency.
- Do not use either `vkQueueWaitIdle()` or `vkDeviceWaitIdle()`.

Debugging your application

The [Arm Streamline](#) system profiler visualizes the Arm CPU and GPU activity on both GPU queues. The system profiler quickly shows bubbles in scheduling that occur in two ways:

- Locally to the GPU queues. This type of bubble is indicative of a stage dependency issue.
- Globally across both the CPU and GPU. This type of bubble is indicative of a blocking CPU call that is being used.

3.9 Vulkan pipeline synchronization

Arm GPUs expose two hardware processing slots. Each slot implements a subset of the rendering pipeline stages and runs each slot in parallel with the other slot.

Prerequisites

You must understand the following concepts:

- Command synchronization barriers.
- Shader stages.

Using parallel processing with Vulkan

The GPU allows the maximum amount of parallel processing across the two hardware slots.

**Note**

Arm tile-based GPUs differ from desktop intermediate-mode renderers as tile-based GPUs have two independently scheduled hardware slots present for different types of workload. Tune your pipeline to work well on a tile-based GPU when porting content from a desktop GPU.

Make sure to overlap vertex or compute processing with fragment processing to get maximum efficiency from your GPU.

The following lists show the mappings of Vulkan stages to the Arm GPU processing slots:

Vertex or compute hardware slot

- `VK_PIPELINE_STAGE_DRAW_INDIRECT_BIT`
- `VK_PIPELINE_STAGE_VERTEX_*_BIT`
- `VK_PIPELINE_STAGE_TESSELLATION_*_BIT`
- `VK_PIPELINE_STAGE_GEOMETRY_SHADER_BIT`
- `VK_PIPELINE_STAGE_COMPUTE_SHADER_BIT`
- `VK_PIPELINE_STAGE_TRANSFER_BIT`

Fragment hardware slot

- `VK_PIPELINE_STAGE_EARLY_FRAGMENT_TESTS_BIT`
- `VK_PIPELINE_STAGE_FRAGMENT_SHADER_BIT`
- `VK_PIPELINE_STAGE_LATE_FRAGMENT_TESTS_BIT`
- `VK_PIPELINE_STAGE_COLOR_ATTACHMENT_OUTPUT_BIT`
- `VK_PIPELINE_STAGE_TRANSFER_BIT`

Vulkan places the application in control of how dependencies between commands are expressed. An application must make sure that a pipeline stage for one command has produced its results before a later command can consume the results.

The API includes multiple primitives that are available for command synchronization, for example:

- Subpass dependencies, pipeline barriers, and events, which are used to express fine-grained synchronization in a single queue.
- Semaphores, which are used to express heavier weight dependencies across queues.

All fine-grained dependency tools allow the application to specify a restricted scope for their synchronization. The `srcStage` mask indicates which pipeline stages must be waited for. The `dstStage` mask indicates which pipeline stages must wait for synchronization before processing starts.

To get best parallel processing across the two Arm GPU hardware processing slots begin by minimizing the scope for synchronization. Set `srcStage` as early as possible in the pipeline, and set `dstStage` as late as possible.

Semaphores allow control when dependent commands run using `pWaitDstStages`. However, semaphores assume that the `srcStage` is the worst case `VK_PIPELINE_STAGE_BOTTOM_OF_PIPE_BIT`. Therefore, only use semaphores when no fine-grained alternative is available.

Synchronization at low levels

Two kinds of synchronization are needed at the low level:

- Synchronization in a single hardware processing slot.
- Synchronization across the two hardware processing slots.

As fragment shading always comes after vertex or compute in the Arm GPU rendering pipeline, synchronization from a `srcStage` running in the vertex or compute processing slot to a `dstStage` running in the fragment processing slot is low cost.

Synchronization from a `srcStage` in the fragment hardware slot to a `dstStage` in the vertex or compute hardware slot is expensive. The synchronization creates a pipeline bubble unless the extra latency for `srcStage` result generation is accounted for. For example, by having sufficient non-dependent work to fill the bubble.

The `TRANSFER` stage is an overloaded term in the Vulkan pipeline. The driver can implement the transfer operations in either hardware processing slot. This means that the direction of a dependency through the pipeline, either backwards or forwards, is not obvious.

Transfers from buffer-to-buffer are implemented in the vertex or compute processing slot. Other transfers can be implemented in either processing slot and determining which hardware processing slot is used depends on the state of the data resource that is being written at the time.

The processing slot that is used materially impacts the rendering workload of an application pipeline. Therefore, always review the performance of transfer operations.

How to optimize your Vulkan pipeline synchronizations

Try using the following optimization techniques:

- Set up `srcStageMask` as early as possible in the pipeline.
- Set up `dstStageMask` as late as possible in the pipeline.
- Check whether dependencies point forwards, so source is vertex or compute and destination is fragment, or backwards, so source is fragment and destination is vertex or compute, through the pipeline. Minimize the use of backwards-facing dependencies.
- If backwards-facing dependency is needed, then add sufficient latency between the generation and the consumption of the resource that is used to hide the added scheduling bubble.
- Use `srcStageMask = ALL_GRAPHICS_BIT` and `dstStageMask = FRAGMENT_SHADER_BIT` to synchronize render passes with each other.
- Zero-copy algorithms are the most efficient, so minimize the use of `TRANSFER` copy operations. Always review how `TRANSFER` copies are impacting the hardware pipelining.
- Only use intra-queue barriers when they are needed and put as much work as possible between barriers.

Vulkan pipelining techniques to avoid

Arm recommends that you:

- Do not leave the vertex / compute processing queue empty while doing fragment processing or vice versa.
- Do not use the following `srcStageMask` to `dstStageMask` synchronization pairings because they completely drain the pipeline:
 - `BOTTOM_OF_PIPE_BIT` to `TOP_OF_PIPE_BIT`
 - `ALL_GRAPHICS_BIT` to `ALL_GRAPHICS_BIT`
 - `ALL_COMMANDS_BIT` to `ALL_COMMANDS_BIT`
- Be careful not to introduce false dependencies if merging pipeline barriers.
- Do not use more stages in masks for barriers than is necessary for synchronization. Conservative masks may cause unexpected delays.
- Do not use more barriers than needed. Each barrier comes with some processing cost, so fewer barriers are better if merging them does not cause synchronization delays.
- Do not use `vkEvent` to signal completion, and then wait for that event right away. Use `vkCmdPipelineBarrier()` instead. `vkCmdPipelineBarrier2()` can also be used, allowing several `vkCmdPipelineBarrier()` calls to be rolled into one.
- Do not use `vkSemaphore` for dependency management in a single queue.

The negative impact of inefficient Vulkan pipelining

The wrong pipeline barrier can either starve the GPU of work with too much synchronization, or cause rendering corruption with too little synchronization.

Debugging your Vulkan pipeline to identify bubbles

The [Arm Streamline](#) system profiler tool quickly shows bubbles in scheduling either locally to the GPU hardware, or globally across both the CPU and GPU. GPU local bubbles are indicative of a stage dependency issue. Global bubbles suggest that a blocking CPU call is being used.

Enabling the [Vulkan Validation Layers](#) while developing can highlight some inefficiencies in synchronization as well as many other errors.

3.10 Pipelined resource updates

OpenGL ES must make sure that resources are in the correct state when the draw call is executed on the GPU. Doing so correctly prevents the modification of resources that are still referenced by in flight draw call.

Prerequisites

You must understand the following concepts:

- Resource ghosting.
- N-buffering.

Referencing resources

**Note**

Depending on how the driver handles the conflict, attempting to modify a resource that is still referenced can cause either pipeline bubbles, or increased CPU load.

OpenGL ES presents a synchronous rendering model to the application developer, even though the underlying execution is asynchronous. Rendering must reflect the state of the data resources at the point that the draw call was made. If an application modifies a resource while a pending draw call is still referencing it, then the driver must take some action to ensure correctness.

The Arm GPU driver avoids blocking and waiting for the resource reference count to hit zero, because doing so drains the pipeline and leads to poor performance. To reflect the new state, the driver creates a new version of the resource. The old, or ghost, version of the resource is kept until the pending draw calls have completed and its reference count drops to zero.

This process is expensive because it requires memory allocation for the new resource, and the cleaning up of the ghost resource when it is no longer needed. If the update is not a total replacement, then a data copy from the old resource into the new one is also required.

How to optimize your pipelined resources

Try using the following optimization techniques:

- To prevent modifying resources that are referenced by the queued draw calls, use N-buffered resources, and pipeline your dynamic resource updates.
- Use `GL_MAP_UNSYNCHRONIZED` to allow the use of `glMapBufferRange()` to patch an unreferenced region of a buffer that is still referenced by in-flight draw calls.

Approaches to avoid when updating pipelined resources

Arm recommends that you:

- Do not modify resources that are still being referenced by in-flight draw calls.
- Do not use `glMapBufferRange()` with either `GL_MAP_INVALIDATE_RANGE` or `GL_MAP_INVALIDATE_BUFFER` on some older Arm GPU driver versions. These flags trigger the creation of an unnecessary resource ghost.

Negative impacts of unoptimized pipelined resource updates

Suboptimal resource update pipelining can result in the following negative impacts:

- Increased CPU load due to memory allocation overhead, and the copies needed to build new versions of resources.
- Unstable memory footprint, due to the constant allocation and freeing of the ghost resources.

Debugging resource updates

The [Arm Streamline](#) system profiler visualizes the Arm CPU and GPU activity. Failing to pipeline resource updates normally show as elevated CPU activity.

3.11 Optimize attachment grouping

With Vulkan subpasses, attachment grouping is generally explicit. However, for OpenGL ES and some other Vulkan uses, it is good to consider what attachments different groups of draw calls can read and write to.

Prerequisites

You must understand the following concepts:

- Draw calls
- Attachments

Parallel execution

For the Mali-G710 and later GPUs there is increased parallelization of draw calls, but consideration of which attachments are read and written is required to make best use of that parallelization.

Parallelization is considered at attachment level, not channel level. This means that multiple draw calls with reads of an attachment can be done in parallel, if there are no pending writes. It also means that it can be worth splitting an attachment if there are separate reads and writes to different channels, to allow parallelization of the reads.

Avoid overlapping writes

Try to avoid writing all draw calls to the same attachment, as that requires all calls to be executed serially.

Negative consequences of non-optimal attachment grouping

If draw call attachments are not grouped optimally, more calls must be executed in series, rather than in parallel, slowing down the execution.

3.12 Queries

OpenGL ES and Vulkan both have query objects that are used for such things as testing occlusion. Using them correctly improves performance.

Prerequisites

You must understand the following concepts:

- Query objects.

How to optimize your queries

Try using the following optimization techniques:

- Pipeline the use of query objects, only reading the query result when you know that it is already available.

- For occlusion queries, prefer use of any-samples queries, rather than the precise sample count mode. Use `GL_ANY_SAMPLES_PASSED` for OpenGL ES, or `VK_QUERY_CONTROL_PRECISE_BIT = false` for Vulkan, unless you must know the amount of occlusion.

The negative impacts of inefficient queries on Arm GPUs

The different types of impact you see are:

- If you wait too early for a query result, the pipeline will drain waiting for the result, which wastes cycles.
- If you use a precise count occlusion query unnecessarily, the GPU does additional work, which also wastes power and cycles.

4. CPU overheads

There are various ways to reduce CPU overheads to increase efficiency and reduce software processing costs.

4.1 Compiling shaders in OpenGL ES

Both shader compilation and program linkage are expensive operations. Generating new programs while trying to render at high frame rates can cause dropped frames. Therefore, avoid shader compilation and program linkage in the interactive part of your application.

Prerequisites

You must understand the following concepts:

- Compiling shaders.
- Linking programs.

How to optimize compiled shaders

Arm recommends that you:

- Compile all shaders, and link all programs, when starting an activity or when loading a level of a game.
- Store program binaries back to disk, avoiding re-compilation on subsequent runs of the application.

Behaviors to avoid when compiling shaders

Arm recommends that you:

- Do not attempt to compile shaders, or link programs, during interactive gameplay.
- Do not rely on the Android blob cache for program caching, because it is also often too small to contain all shader programs for a complex application.

The negative impact of not compiling shaders correctly

The cost of trying to compile, and link, shaders during the interactive portions of the application is a high CPU load and dropped frames.

4.2 Pipeline creation in Vulkan

Vulkan pipelines have similar compile requirements to OpenGL ES shaders. Vulkan also provides no automatic caching that is equivalent to the Android blob cache that is available for OpenGL ES

programs. You are responsible for providing persistent storage of compiled programs for use across program invocations.

Prerequisites

You must understand the following concepts:

- Vulkan pipelines.

How to optimize pipeline creation

To optimize pipeline creation, use the following optimization techniques:

- Create pipelines when starting an activity or loading a game level.
- Use a pipeline cache to speed up pipeline creation.
- Serialize the pipeline cache to disk and then reload it the next time the application is used, providing end users with a faster load time.



Arm GPUs ignore the following flags: `VK_PIPELINE_CREATE_DERIVATIVE_BIT`, `VK_PIPELINE_CREATE_ALLOW_DERIVATIVES_BIT`, and also `VK_PIPELINE_CREATE_DISABLE_OPTIMIZATION_BIT`.

Pipeline caches

As well as standard pipeline cache use, loading a whole pipeline from a previous run, pipeline caches can be useful for partial loading.

A pipeline cache can help where there still needs a partial recompilation. For example, combining a previously compiled vertex shader with a new fragment shader. This includes changes that requires a recompilation from the API, but ultimately does not require a recompilation, like updating an unused descriptor set.

Where pipeline caches risk becoming too large it is possible to request from the API a smaller size than the full cache. In this case Arm's Vulkan implementation will prioritize returning the entries with the highest amount of hits.

Negative impacts on the application

Incorrect pipeline creation can have the following impacts on your application:

- Attempting to create pipelines during the interactive portions of your application increases the CPU load and can cause skipped frames.
- Failure to serialize, and then reload, a pipeline cache increases load times on later application runs.

Example

Example code of this best practice is available in the Vulkan samples repository on GitHub, see the [Pipeline cache tutorial](#).

4.3 Allocating memory in Vulkan

For frequent allocations, do not use the `vkAllocateMemory()` allocator. All allocations by `vkAllocateMemory()` are expensive kernel calls.

Prerequisites

You must understand the following concepts:

- Memory allocation.

How to optimize memory allocation

When allocating memory, use your own allocator to manage block allocations.

Dedicated allocations

For potential performance improvements, use the extension `VK_KHR_dedicated_allocation` to determine when a dedicated allocation may be beneficial for your use-case.

Something to avoid when allocating memory

Do not use `vkAllocateMemory()` as a general-purpose allocator.

Negative impacts of suboptimal memory allocation

Using `vkAllocateMemory()` as a general-purpose allocator increases the load on the application CPU load.

Debugging memory allocation issues

Monitor the frequency and allocation size of all calls to `vkAllocateMemory()` at runtime.

4.4 OpenGL ES CPU memory mapping

By using `glMapBufferRange()`, OpenGL ES provides direct access to buffer objects that are mapped into the application address space. It is efficient to stream writes, but reads are expensive for mapped buffers because they are uncached on the CPU.

Prerequisites

You must understand the following concepts:

- Buffers.
- Memory mapping.

How to optimize CPU memory mapping

To benefit from store merging in the CPU write buffer, make write-only buffer updates to sequential addresses.

Outcome to avoid when using CPU memory mapping

Arm recommends that you do not read values from mapped buffers.

Negative impacts of not using CPU memory mapping correctly

Reading from uncached buffers shows up as increased CPU load in the application functions that read them.

4.5 Vulkan CPU memory-mapping

Vulkan gives applications support for multiple buffer memory types.

Prerequisites

You must understand the following concepts:

- CPU and GPU memory mapping.
- Cached versus uncached memory.

Memory Types

On Midgard architecture GPUs, the Arm GPU driver exposes the following memory types:

1. `DEVICE_LOCAL_BIT | HOST_VISIBLE_BIT | HOST_COHERENT_BIT`
2. `DEVICE_LOCAL_BIT | HOST_VISIBLE_BIT | HOST_CACHED_BIT`
3. `DEVICE_LOCAL_BIT | LAZILY_ALLOCATED_BIT`

On Bifrost and later architecture GPUs, the Arm GPU driver exposes the following memory types:

1. `DEVICE_LOCAL_BIT | HOST_VISIBLE_BIT | HOST_COHERENT_BIT`
2. `DEVICE_LOCAL_BIT | HOST_VISIBLE_BIT | HOST_CACHED_BIT`
3. `DEVICE_LOCAL_BIT | HOST_VISIBLE_BIT | HOST_COHERENT_BIT | HOST_CACHED_BIT`
4. `DEVICE_LOCAL_BIT | LAZILY_ALLOCATED_BIT`

The four memory types and their purposes

The purposes that these memory types are useful for are:

Not cached, coherent

The `HOST_VISIBLE`, `HOST_COHERENT_BIT` memory type:

- Is guaranteed to be supported.
- Provides uncached storage on the CPU.
- Avoids polluting the CPU caches with unnecessary data.
- Efficiently merges small writes that are then sent on to the external memory device by using the CPU write buffer hardware.
- It is the optimal memory type for CPU write-only resources.

Cached, incoherent

The `HOST_VISIBLE`, `HOST_CACHED` memory type:

- Provides cached storage on the CPU but does not guarantee that the CPU and the GPU get coherent views of the underlying memory. The CPU view of the memory is not coherent with the GPU view of the memory, therefore you must call:
 - `vkFlushMappedRanges()` when the CPU has finished writing data for the GPU.
 - `vkInvalidateMappedRanges()` when reading back the data that the GPU has written.
 - However, both calls are expensive to use, so must be used sparingly.
 - Must only be used for resources that are both mapped and read by the application software on the CPU.

Cached, coherent

The `HOST_VISIBLE`, `HOST_COHERENT`, `HOST_CACHED` memory type:

- Provides cached storage on the CPU, which is also coherent with the GPU view of the memory without needing manual synchronization.
- Supported by Arm Bifrost and later GPUs if the chipset supports the hardware coherency protocol between the CPU and GPU.
- Due to hardware coherency, it avoids the overheads of manual synchronization operations. Cached, coherent memory is preferred over the Cached, incoherent memory type when available.
- Must be used for resources that are both mapped and read by the application software on the CPU.
- Hardware coherency has a small power cost, so must not be used for resources that are write-only on the CPU. For write-only resources, bypass the CPU cache by using the Not Cached, coherent memory type.

Lazily allocated

The `LAZILY_ALLOCATED` memory type:

- It is a special memory type that is initially only backed by GPU virtual address space and not physical memory pages. Physical backing pages are allocated on demand if the memory is accessed.
- It must be used alongside transient image attachments created using `VK_IMAGE_USAGE_TRANSIENT_ATTACHMENT_BIT`. Transient images are intended for use as framebuffer attachments which only exist during a single render pass. Doing so avoids backing the attachment with physical memory. You can do this by backing the attachment image with a lazy allocation, and a `VK_ATTACHMENT_STORE_OP_DONT_CARE` storage operation. Common use cases for this include depth or stencil buffers for simple renders. Another use case and G-buffer attachments, that are used for deferred lighting, that are then consumed by a later subpass and are not written back to memory.
- Must not be used to provide storage for resources that are expected to be written back to memory.

How to optimize Vulkan CPU memory mapping

Try using the following optimization techniques:

- Use `HOST_VISIBLE` | `HOST_COHERENT` memory for immutable resources.

- Use `HOST_VISIBLE | HOST_COHERENT` memory for resources which are write-only on the CPU.
- Use `memcpy()` to write updates to `HOST_VISIBLE | HOST_COHERENT` memory, or write sequentially to get best efficiency from the CPU write-combine unit.
- Use `HOST_VISIBLE | HOST_COHERENT | HOST_CACHED` memory for resources which are read back on to the CPU. Use `HOST_VISIBLE | HOST_CACHED` if it is not available.
- Use `LAZILY_ALLOCATED` memory for transient framebuffer attachments which only exist during a single render pass.
- Only use `LAZILY_ALLOCATED` memory for `TRANSIENT_ATTACHMENT` framebuffer attachments.
- Mapping and unmapping buffers have a CPU cost. Therefore, persistently map buffers which are accessed often. Example: uniform buffers, data buffers, or dynamic vertex data buffers.
- When using the `VK_EXT_host_image_copy` extension use `HOST_CACHED` flagged memory.

Vulkan CPU memory-mapping techniques to avoid

Arm recommends that you:

- Do not read back data from uncached memory on the CPU.
- Do not store suballocator metadata, which must be read on the CPU, inside an uncached memory buffer.

Negative impacts of inefficient Vulkan CPU memory mapping

Uncached readbacks can be much slower than cached reads due to increased CPU processing costs.

Debugging Vulkan CPU memory-mapping performance problems

There are a couple of techniques you can take:

- Check that all CPU-read buffers are using cached memory.
- Design interfaces for buffers to encourage implicit flushes or invalidates as needed. It is difficult and time consuming to debug coherency failures due to missing maintenance operations without this infrastructure already in place.

4.6 Command pools for Vulkan

Command pools do not automatically recycle memory from deleted command buffers unless created with the `RESET_COMMAND_BUFFER_BIT` flag. Pools without this flag do not recycle their memory until the application resets the pool.

Prerequisites

You must understand the following concepts:

- Command pools.

How to optimize command pools

Periodically call `vkResetCommandPool()` to release the memory.

Command pool techniques to avoid

- Read [Optimizing command buffers for Vulkan](#) for an explanation why using `RESET_COMMAND_BUFFER_BIT` is not recommended. However, using `RESET_COMMAND_BUFFER_BIT` is better than not releasing memory.
- When creating command pools, Arm GPUs ignore the flag `VK_COMMAND_POOL_CREATE_TRANSIENT_BIT`. Arm GPUs also ignore the `vkTrimCommandPool()` command.

The negative impact of inefficient Vulkan command pools

An increase in memory usage until a manual command pool reset is triggered.

4.7 Optimizing command buffers for Vulkan

Command buffer usage flags affect performance.

Prerequisites

You must understand the following concepts:

- Command buffers.
- Command pools.

How to optimize command buffers

Try using the following optimization techniques:

- For best performance set the `ONE_TIME_SUBMIT_BIT` flag.
- Build per-frame command buffers instead of using simultaneous command buffers.
- If the alternative is to replay the same command sequence every time in application logic, then use `SIMULTANEOUS_USE_BIT`. It is more efficient than an application replaying commands manually, but less efficient than a one-time submit buffer.

Command buffer techniques to avoid

Arm recommends that you:

- Do not set `SIMULTANEOUS_USE_BIT`, unless required to do so.
- Do not use command pools with `RESET_COMMAND_BUFFER_BIT` set. Doing so increases the memory management overhead, as it prohibits the driver from using a single large allocator for all command buffers in a pool.

Negative impacts of inefficient Vulkan command buffers

Keep the following points in mind:

- There is a risk of an increase in CPU load if inappropriate flags are used, or if command buffer resets are too frequent.
- Avoid calling `vkResetCommandBuffer()` on a high frequency call path.

Debugging your Vulkan command buffer performance problems

There are a couple of approaches you can try, to speed up your debugging process. These are:

- Review and evaluate the performance impact every use of any command buffer flag other than `ONE_TIME_SUBMIT_BIT`.
- Evaluate every use of `vkResetCommandBuffer()` and assess if it could be replaced with `vkResetCommandPool()` instead.

Example

Example code of command buffer usage and multi-threaded recording is available in the Vulkan Samples repository on GitHub: [Command Buffer usage tutorial](#)

4.8 Secondary command buffers

Arm GPU hardware prior to the Mali-G710 series does not have native support for invoking commands in a secondary command buffer. Therefore, there is extra overhead that is incurred when using secondary command buffers.

Prerequisites

You must understand the following concepts:

- Command buffers.
- Command pools.

Using secondary command buffers

It is expected that applications must use secondary command buffers, typically to allow multi-threaded command buffer construction. However, the total number of secondary command buffer invocations must be minimized. As with primary command buffers, we recommend avoiding creating command buffers with the `SIMULTANEOUS_USE_BIT` due to increased overheads.

How to optimize secondary command buffers

Try using the following optimization techniques:

- Use secondary command buffers to allow multi-threaded render pass construction.
- Minimize the number of secondary command buffer invocations that are used per frame.

Secondary command buffer step to avoid

Do not set `SIMULTANEOUS_USE_BIT` on secondary command buffers.

The negative impact of inefficient secondary command buffers

Keep in mind that an increased CPU load is incurred.

Example

Example code of command buffer usage and multi-threaded recording is available in the Vulkan samples repository on GitHub: [Command buffer usage tutorial](#)

4.9 Optimizing descriptor sets and layouts for Vulkan

Midgard and Bifrost Arm GPUs support four simultaneous bound descriptor sets at the API level. However, they require a single physical descriptor table per draw call internally.

Prerequisites

You must understand the following concepts:

- Descriptor sets.
- Binding spaces.
- Uniform Buffer Objects (UBOs).
- Shader Storage Buffer Objects (SSBOs).

Descriptor sets and layouts

If any of the four source descriptor sets have changed, then the driver rebuilds the internal table for a draw call. The first draw call, after a descriptor changes, has a higher CPU overhead than following draw calls that reuse the same descriptor set. Larger descriptor sets cause a more expensive rebuild. With current drivers, the descriptor set pool allocations are not pooled. Do not call `vkAllocateDescriptorSets()` on a performance critical code path.



Before Valhall, descriptor set pool allocations were never pooled.

Table rebuilds are much smaller for Valhall Arm GPUs. However, our advice to not allocate on critical code paths still applies.

For ease of programming, `VkDescriptorSetLayoutBinding::stageFlags` can always be set to `VK_SHADER_STAGE_ALL` with no performance loss.

The `VkDescriptorPool` creation flag

Arm GPUs ignore the `VK_DESCRIPTOR_POOL_CREATE_FREE_DESCRIPTOR_SET_BIT` flag, but applications must still support the specification. This means that if `VK_DESCRIPTOR_POOL_CREATE_FREE_DESCRIPTOR_SET_BIT` is set, then you must still check for pool fragmentation. If the flag is not set, then you must not free individual descriptor sets, and only use `vkAllocateDescriptorSets()` and `vkResetDescriptorPool()`.

How to optimize descriptor sets and layouts

Try using the following optimization techniques:

- Pack the descriptor set binding space as much as possible.
- Instead of resetting descriptor pools and reallocating new descriptor sets, update descriptor sets that are already allocated, but no longer referenced.

- Reuse pre-allocated descriptor sets and do not update them with the same information every time.
- Use `VK_DESCRIPTOR_TYPE_UNIFORM_BUFFER_DYNAMIC` OR `VK_DESCRIPTOR_TYPE_STORAGE_BUFFER_DYNAMIC` to bind the same UBO or SSBO with different offsets. The alternative is building more descriptor sets.

Descriptor indexing

Using the extension `VK_EXT_descriptor_indexing` allows the use of bindless resources. This provides many useful functionality improvements, especially with regards to added flexibility and reduced CPU load. When using `VK_EXT_descriptor_indexing`, try to use SSBOs instead of UBOs, as SSBOs offer better performance.

Descriptor set and layout techniques to avoid

Arm recommends that you:

- Do not leave holes in the descriptor set.
- Do not leave unused entries as copying and merging still has a computational cost.
- Do not allocate descriptor sets from descriptor pools on performance critical code paths.
- Do not use `DYNAMIC_OFFSET` UBOs/SSBOs if you never plan on changing the binding offset, as there is a small, extra, cost for handling the dynamic offset. The negative impact of inefficient descriptor sets and layout unoptimized Vulkan descriptor sets and layout leads to a risk of increased CPU load for draw calls.

Debugging your descriptor sets and layout performance problems

Ways to speed up your debugging process:

- Monitor the pipeline layout for unused entries.
- Monitor for contention-related performance problems on `vkAllocateDescriptorSets()`.

Example and further information

Example code of descriptor and buffer management is available in the Vulkan Samples repository on GitHub: [Descriptor Management tutorial](#)

For more information on `VK_EXT_descriptor_indexing` [this Arm Community blog](#) provides further details.

5. Vertex shading

The efficiency of vertex processing, in terms of both buffer packing and vertex shading, is important when rendering a scene. Avoid poor mesh selection and inefficient vertex data encoding on mobile devices, because it can significantly increase DRAM bandwidth.

5.1 Basic vertex shader optimizations

Most of the basic vertex optimizations are implemented by making improvements in the input mesh data, or in the application rendering logic.

Prerequisites

You must understand the following concepts:

- Vertices and meshes

Artistic improvements

There are various artistic optimizations to make for mobile that reduce the number of vertices processed, making vertex shaders much faster.

- Reduce the number of vertices, using only the number of triangles needed to preserve the desired object silhouette.
- Ensure that vertices are unnecessarily duplicated during asset creation.
- Use textures and normal maps to simulate fine details of color and lighting.
- Use dynamic level of detail for meshes that change distance from the camera.

Culling improvements

Try processing only those triangles that can be seen. There are several culling types that can be used to reduce vertex cost:

- Use view frustum culling in the application, discarding all objects that are outside of the view frustum. This type of testing can use simple bounding volumes around objects to make these tests more efficient.
- Use portal culling in the application, discarding all objects that are inside the frustum but which are inside a region of the level that can be easily proven to be occluded from the current camera position.
- Enable back-face culling during rendering to avoid shading the occluded back faces of your meshes.

Negative consequences of not making artistic optimizations

Processing more vertices takes more time, so the scene runs slower, consuming more bandwidth and power.

Reference

For more details, see [Real-time 3D Art Best Practices - Geometry](#), and [Triangle density](#).

5.2 Index draw calls

Indexed draw calls allow for reuse of vertices and are more efficient than non-indexed draw calls.

Prerequisites

You must understand the following concepts:

- Draw calls.
- Index buffers.

How to optimize index draw calls

Try using the following optimization techniques:

- Use indexed draw calls whenever vertex reuse is possible.
- Optimize index locality for a post-transform cache.
- To minimize redundant processing and data fetch, ensure that the index buffer references the entire range of index values.
- Create tightly packed vertex ranges for lower level of detail meshes.
- Use `glDrawRangeElements()` for volatile resources.
- To implement geometry Level of Detail (LOD), create a contiguous set of vertices for each detail level.

Index draw call techniques to avoid

Arm recommends that you:

- Do not use client-side index buffers.
- Do not modify the contents of index buffers. Doing so causes the driver to rescan the index buffer to determine the active index range.
- Do not use indices that sparsely access vertices in the vertex buffer.
- Do not use indexed draws for simple geometry which has no vertex reuse, such as simple quads or point lists.
- Do not create a low detail mesh by sparsely sampling vertices from the full detail mesh.

Negative impacts of inefficient index draw calls

Unoptimized index draw calls can lead to the following problems:

- The use of client-side index buffers increases the CPU load. Client-side index buffers first allocate server-side buffer storage, then copies the data, and finally, scans the contents to determine the active index ranges.

- The use of index buffers with frequently modified contents shows up as increased load on the CPU. The increased load is due to the need to rescan the buffer to determine the active index ranges.
- Inefficient mesh indexing that is due to index sparseness, or poor spatial locality. This shows up as extra GPU vertex processing time or extra memory bandwidth that is being used. The severity of the impact depends on the complexity of the mesh and the layout of the index buffer that is in memory.

Debugging index draw calls issues

Scan through all index buffers before submission. Optimize any buffers that include any unused indices.

5.3 Index buffer encoding

The index buffer data is one of the primary data sources into the primitive assembly and tiling process in Arm GPUs. To minimize the cost of tiling, efficiently pack the index buffer.

Prerequisites

You must understand the following concepts:

- Index buffers.
- Strip formats and simple list formats.

How to optimize index buffer encodes

Try using the following optimization techniques:

- Use the lowest precision index data type possible to reduce index list size.
- To reduce index list size, use strip formats over simple list formats.
- To reduce index list size, use primitive restart instead of degenerate triangles.
- For a post-transform cache, optimize the index locally.

Index buffer encoding techniques to avoid

Arm recommends that you:

- Do not use 32-bit index values.
- Do not use index buffers with low spatial coherency because they hurt caching.

The negative impact of inefficient index buffer encoding

Inefficient index buffer encoding does not usually have a significant impact, but can be part of small improvements that add up to improve application performance. A significant negative impact is possible if draw calls are small, and vertex shading and tiling do not pipeline cleanly.

5.4 Index sparsity

For Midgard and Bifrost, very sparse index references in indexed draw calls can cause Out Of Memory (OOM) errors, leading to `VK_ERROR_DEVICE_LOST`.

Prerequisites

You must understand the following concepts:

- Index buffers.
- Draw calls.

Pack mesh indices

Arm GPUs have a memory region which is available to store the intermediate geometry output from a render pass. The region is fixed to 180MB on the Midgard and Bifrost Arm GPUs that have this issue, that is, before the Mali-G77.

For an indexed draw call, Midgard and Bifrost GPUs would allocate memory to span the highest and lowest indices referenced. In extreme cases, this means that a draw call for a single triangle with indices `[0, 1, 100000]` uses a substantial amount of memory. A second neighboring draw call for another triangle with `[0, 2, 100000]` allocates the same amount again. A single draw call covering the two triangles only allocates the memory once.

You can avoid this problem by:

- Ensuring that indices are spatially coherent. Ideally within the original meshes, or with a static layer of indirection. Alternatively, a compute shader can be used to create a coherent mesh, at the cost of performance.
- Ensure draw calls are combined when possible
- Splitting the render pass. Be warned that this has a performance cost, especially on bandwidth with further readbacks, but it is another way around the issue.

Very large geometry

It is also possible to hit the OOM or `VK_ERROR_DEVICE_LOST` error if you have huge geometry that exceeds the 180MB limit, even without indexing issues. Note that this limit has been significantly increased on recent Valhall GPUs.

Avoid sparse mesh index references

To avoid memory issues, try to:

- Avoid sparse vertex index references in a draw call.
- Not have excessively large geometry.
- Avoid encoding metadata into index buffer values.

Negative impacts of index sparsity

For Midgard and Bifrost, very sparse index buffers can cause OOM errors, leading to `VK_ERROR_DEVICE_LOST`.

Further reading

There is a useful blog on [Vulkan memory limits with Mali GPUs](#).

5.5 Attribute precision

Full FP32 high precision of vertex attributes is unnecessary for many uses of attribute data. For example, color computation. Asset pipelines must keep the data at the minimum precision that is required. Doing so reduces bandwidth and improves performance.

Prerequisites

You must understand the following concepts:

- Attribute data.
- Asset pipelines.
- FP32 and FP16.
- Vector types.

Attribute precision

OpenGL ES and Vulkan provide different attribute data types that fit every precision that is needed. Lower precision attributes are available in 8-bit, 16-bit, and packed formats such as RGB10_A2. Arm GPU hardware can convert attributes to FP16 and FP32 for free on data load. Therefore, there is no shader processing overhead from using narrower data types.

How to optimize attribute precision

Try using the following optimization techniques:

- Stable geometry positioning needs extra precision, so use FP32 for computing vertex positions.
- Use low precision for other attributes, increasing to high precision only when it is needed.
- For integer vectors, prefer `uvec` over `ivec` as it enables extra compiler optimizations.

Relaxed precision FP16 outputs

If FP16 is sufficient for your vertex shader outputs, it is recommended to use the `RelaxedPrecision` decoration to allow them to be FP16 on Arm GPUs. If the fragment shader needs FP32 as input, then this is interpolated as needed, but still saves bandwidth. Note that doing so can flag up some warnings in older validation layers, but this usage is now permitted by the Vulkan specification.

Attribute precision techniques to avoid

Arm recommends that you:

- Do not use FP32 for everything.
- Do not upload FP32 data into a buffer and then read it as a mediump attribute. Doing so wastes both memory storage and bandwidth as the extra precision is discarded.

Negative impacts of inefficient attribute precision

Higher memory bandwidth, large memory footprint, energy inefficiency, and reduced vertex shader performance.

5.6 Attribute layout

Vertices are shaded using an Index Driven Vertex Shading (IDVS) flow on Bifrost Arm GPU architectures and above.

Prerequisites

You must understand the following concepts:

- Vertices.
- Vertex buffers.

IDVS vertex shading order

The order in which vertices are shaded using the IDVS are as follows:

1. Positions.
2. Varyings of vertices of primitives which have survived culling.

Good buffer layout maximizes the benefit of this geometry pipeline.

How to optimize the attribute layout

Try using the following optimization techniques:

- Use a dedicated, tightly packed, vertex buffer for position data.
- Use a second dedicated vertex buffer for non-position data, if any exists.
- Remove unused attributes for specific uses to optimize the meshes. For example: Generating a tightly packed buffer consisting of only position-related attributes for shadow mapping.

Attribute layout techniques to avoid

Arm recommends that you:

- Do not use one buffer per attribute as this wastes bandwidth by fetching data for culled vertices when they share a cache line with visible vertices.
- Do not pad interleaved vertex buffer strides up to the power-of-two. Arm GPU hardware does not require it, and doing so increases memory bandwidth.
- Do not store any constant data in attributes, use uniforms. For best performance in more complicated cases, for example, where vertex color is constant at one Level of Detail (LoD) and varies at another, write a custom shader. This creates a significant reduction in memory bandwidth for the LoDs where data is constant.

The negative impacts of an inefficient attribute layout

Here are two potential problems that your application can experience:

- Increased memory bandwidth due to redundant data fetches.
- Increased pressure on caches which cause loss of performance.

5.7 Varying precision

Vertex shader outputs are commonly called varying outputs. They are written back to main memory before fragment shading commences.

Prerequisites

You must understand the following concepts:

- `mediump`.
- `highp`.

Vertex shader outputs

It is important to minimize the precision of varying outputs. Doing so reduces the amount of memory storage and bandwidth that is needed for the varying output data.

Typical uses of varying data that usually have sufficient precision in `mediump` include:

- Normals.
- Vertex color.
- Texture coordinates for non-tiled textures that go up to 512x512 pixels.

Typical uses of varying data that need `highp` precision are:

- World-space positions.
- Texture coordinates for large textures, or textures with high levels of UV coordinate wrapping.

How to optimize varying outputs

Use the `mediump` qualifier on varying outputs if the precision is acceptable.

Varying output precision techniques to avoid

Arm recommends that you:

- Do not use varyings with more components and precision than is needed.
- Do not use vertex shaders with outputs that are left unused in the fragment shader.

The negative impacts of inefficient varying outputs

The different types of impact you can see are:

- Increased GPU memory bandwidth.
- Reduced vertex shader and fragment shader performance.

5.8 Triangle density

A vertex requires more bandwidth and processing power to process than a fragment. Ensure that there are many pixels worth of fragment work for each primitive that is rendered. Doing so spreads the processing expense of the vertices over multiple pixels of output.

Prerequisites

You must understand the following concepts:

- Memory bandwidth and vertex processing costs.
- Mesh Levels of Detail (LOD).

How to optimize triangle density

Try using the following optimization steps:

- Use models that create at least 10-20 fragments per primitive.
- Use dynamic mesh level-of-detail, using simpler meshes when objects are further away from the camera.
- Use techniques, such as normal mapping, to bake the required complex geometry during asset creation. Turning the geometry into a simpler runtime mesh with a supplemental texture for per-pixel lighting.
- To improve final image quality, favor improved lighting effects and textures instead of increased geometry.

Vertex processing techniques to avoid

Do not generate micro triangles, as they increase bandwidth and processing power costs for little visual benefit.

The negative impacts of inefficient vertex processing

If the triangle density of a mesh is not properly optimized, then high volumes of geometry cause many problems for a tile-based renderer. For example: poor shading performance, high memory bandwidth, and high system energy consumption due to the memory traffic.

5.9 Instanced vertex buffers

Both OpenGL ES and Vulkan have support for instanced drawing. Instanced drawing uses attribute instance divisors to divide a buffer to locate the data for each instance. There are some hardware limitations which determine the optimal usage of instanced vertex buffers.

Prerequisites

You must understand the following concepts:

- Instance divisors.
- Instanced vertex buffers.

How to optimize the use of vertex buffers

Try using the following optimization techniques:

- Use a single interleaved vertex buffer for all instance data.
- Use instanced attributes to work around any uniform buffer size limitations. For example, 16KB uniform buffers.
- Use several vertices per instance that are a power-of-two.
- Prefer indexed lookups using `gl_InstanceID` into uniform buffers or shader storage buffers. Rather than per-instance attribute data.

Instanced vertex buffer techniques to avoid

Do not use more than one instanced vertex buffer per draw call.

The negative impacts of an inefficient instanced vertex buffer

If the vertex buffer is not properly optimized, then you can expect a reduced performance on all impacted instanced draw calls.

6. Tessellation, geometry shading, and tiling

This chapter covers ways on how to optimize tessellation, geometry shading, and tiling instances in your application.

6.1 Tessellation

Tessellation is a powerful brute-force technique that creates higher-quality meshes, at the cost of GPU memory bandwidth and power consumption. Therefore, tessellation can become expensive to run on tile-based GPUs, like Arm's, that write the output of the geometry processing back to system memory.

Prerequisites

You must understand the following concepts:

- Tessellation.
- Using and optimizing mesh Levels of Details (LODs).
- How to monitor GPU memory bandwidth and power consumption.

Alternatives to using tessellation

Try using the following optimization techniques instead:

- Add more static mesh LoD to character meshes.
- Use the geo-mipmap morphing technique for terrain meshes.

However, if you must use tessellation:

- Only add geometry where it benefits most. For example, on the silhouette edges of your meshes.
- Tessellation can also be used to reduce geometry complexity.
- Cull patches to avoid redundant evaluation shader invocations. In the control shader, use either `glPrimitiveBoundingBox()`, or frustum cull patches.
- To ensure that there is a sensible upper-bound on complexity, clamp your maximum tessellation factor.
- Pre-tessellate a new static mesh instead of using fixed tessellation factors.

Techniques to avoid when using tessellation

Arm recommends that you:

- Do not use tessellation until you have evaluated other options.
- Do not use tessellation together with geometry shaders in the same draw call.
- Do not use transform feedback with tessellation.

- Do not use microtriangulation. There is little perceptual quality benefit from triangle densities of fewer than ten fragments per primitive.

Negative impacts of unnecessarily implementing tessellation

Using tessellation to significantly increase triangle density leads to poor performance, high memory bandwidth, and increased system power consumption.

Debugging tessellation issues

Try the following debugging tips:

- It is easy to overlook that tessellation is generating millions of triangles because it is procedurally generated and is hidden from the application. Therefore, always check the GPU performance counters to monitor the number of triangles that are being generated and the number of ineffective microtriangles that are being generated.
- To perform a controlled exploration of the performance versus visual benefit trade-offs, apply different levels of `min()` to the tessellation factors in the control shader.

6.2 Geometry shading

Before using geometry shading, keep in mind that tile-based GPU architecture is sensitive to geometry bandwidth levels.

Prerequisites

You must understand the following concepts:

- Geometry shading.

How to optimize geometry shading

Use compute shaders instead of geometry shading, as compute shaders are more flexible and avoid the unnecessary duplication of vertices that most geometry shaders trigger.

Techniques to avoid when using geometry shading

Arm recommends that you do not use:

- Geometry shaders.
- Geometry shaders to filter primitives, such as triangles, or to pass down more attributes to the fragment stage.
- Geometry shaders to expand points to quads. Instance a quad instead.
- Transform feedback with geometry shaders.
- Primitive restart with geometry shading.

Negative impacts of implementing geometry shading

Using geometry shading leads to increased memory bandwidth. This reduces both performance and energy efficiency.

6.3 Tiling and effective triangulation

Tiling and rasterization both work on fragment patches that are larger than a single pixel. For Arm GPUs, the tiling uses bins that are usually 16x16 pixels or more. Fragment rasterization emits 2x2 pixel quads for fragment shading.

Prerequisites

You must understand the following concepts:

- Tiling and rasterization techniques.
- Fragment quads.

Minimizing number of triangles needed

Optimal performance is achieved when mesh triangulation uses the minimum number of triangles to cover the necessary pixels.

How to optimize tiling and effective triangulation

Use triangles that are almost equilateral. Using equilateral triangles maximizes the ratio of the area to edge length. Doing so reduces the number of generated partial fragment quads.

Tiling and rasterization techniques to avoid

The center point of a triangle fan has a high triangle density for low pixel coverage per triangle loaded. Therefore, avoid the use of fans or similar geometry layouts.

Negative impacts of inefficient tiling and rasterization

There is a risk of increasing the fragment shading overhead due to the triangle fetch and partial sample coverage of the 2x2 pixel fragment quads.

Debugging triangulation issues

Use the Arm Graphics Analyzer. The Graphics Analyzer contains mesh visualization tools that allow the outline of the mesh, submitted to a draw call, to be visualized in object-space. [Graphics Analyzer](#)

7. Fragment shading

This chapter provides multiple areas of guidance on how to optimize the fragment shading elements of your application on Arm GPUs.

7.1 Basic fragment shader optimizations

Fragment shading can often be the bottleneck in graphics rendering. If that is the case, simplify fragment shaders and textures.

Prerequisites

You must understand the following concepts:

- Fragment shading

Simplify the shader

Look to see if there is simpler arithmetic or acceptable approximations to achieve the required shader effects. If work can be done in the vertex shader or application, it can be more efficient or avoid a bottleneck. If bandwidth is not the problem, using textures instead of calculations is another possible solution.

Reduce the number of branches. Branches are not expensive in themselves, but can lead to a shader being too long and sections needing to be read from cache or RAM.

Reduce texture bandwidth. Textures can use a large amount of memory bandwidth. When too much bandwidth is used the fragment shaders cannot get sufficient data, and stalls. Do not use bigger textures than needed.

Avoid overdraw

Overdraw occurs when the GPU draws over the same pixel multiple times, wasting compute and memory bandwidth. There are multiple things to do to avoid overdraw:

- Enable depth testing.
- Enable back face culling.
- Sort objects by distance to camera. Draw non-transparent objects in front-to-back order, and then draw transparent objects from back-to-front.
- Optimize use of transparency. Transparent objects are more expensive, so minimize their use and draw them last, if still visible.

Negative consequences of overdraw and complex shaders

Using complex shaders slows down performance, especially if they are long enough to have to be read from RAM. If there is too much processing to be done, fragment shading is the bottleneck.

Overdraw causes an unnecessary slowdown of your application, as compute and memory bandwidth resource is wasted drawing pixels that are overwritten before being seen.

Debugging fragment processing

To find whether there is overdraw, or whether fragment shaders are your bottleneck, use Arm Mobile Studio to see what is happening on your Arm GPU.

Web Link: [Arm Mobile Studio](#)

7.2 Efficient render passes with OpenGL ES

Tile-based rendering operates on the concept of render passes. Each render pass has an explicit start and end and produces an output in memory only at the end of the pass.

Prerequisites

You must understand the following concepts:

- OpenGL ES rendering APIs.
- Render passes.
- Tile-based rendering.

Render pass handling

At the start of the pass, the tile memory is initialized inside the GPU. At the end of the pass, the required outputs are written back to system memory. The intermediate framebuffer working state lives entirely inside the tile memory.

Efficient render passes

The driver infers the OpenGL ES render passes based on framebuffer binding calls, as they are not explicit in the API. A render pass for the framebuffer starts when it is bound as the `GL_DRAW_FRAMEBUFFER` target. A render pass normally ends when the draw framebuffer binding changes to another framebuffer.

How to optimize render passes

Try using the following optimization techniques:

- When starting a render pass, clear or invalidate every attachment. This does not apply if the content of a render target is used as the starting point for rendering.
- To clear the tile memory quickly, clear the entire content of the attachment, ensuring that the color, depth, or stencil writes are not masked.
- At the end of the render pass, invalidate any attachments that are not needed outside of the pass, before changing the framebuffer binding to the next FBO.
- For rendering to a subregion of framebuffer, use a scissor box to restrict the area of clearing and rendering required.

Render pass techniques to avoid

Arm recommends that you:

- Do not switch back and render to the same FBO multiple times in a frame. Complete each of your render passes in a single `glBindFramebuffer()` call before moving on to the next. On Valhall, it is especially important to minimize FBO switching as the driver does a flush in `glBindFramebuffer()`.
- Do not split a render pass by using either `glFlush()` or `glFinish()`.
- Do not create a packed depth-stencil texture, `D24_S8` or `D32F_S8`, and only attach one of the two components as an attachment.

Negative impacts of inefficient render passes

Incorrect handling of render passes causes worse fragment shading performance and increased memory bandwidth. Therefore, avoid lowering fragment shading performance and increasing memory bandwidth. At the start of rendering, read non-cleared attachments into tile memory, and then write out noninvalidated attachments at the end of rendering.

Debugging render pass issues

Review your API usage of framebuffer binding, clears, draws, and invalidates.

7.3 Efficient render passes with Vulkan

Tile-based rendering operates on the concept of render passes. Each render pass has an explicit start and end and produces an output in memory only at the end of the pass.

Prerequisites

You must understand the following concepts:

- Vulkan rendering APIs.
- Render passes.
- Tile-based rendering.

Render pass handling

At the start of the pass, the tile memory is initialized inside the GPU. At the end of the pass, the required outputs are written back to system memory. The intermediate framebuffer working state lives entirely inside the tile memory.

Efficient render passes

Unlike with OpenGL ES, Vulkan render passes are explicit in the API. There are defined `loadOp` and `storeOp` operations. `loadOp` defines how GPUs initialize the tile memory at the start of the pass. `storeOp` defines what is written back at the end of a pass.

Vulkan introduces lazily allocated memory, meaning that transient attachments existent during a single render pass do not need physical storage.

How to optimize render passes

Try using the following optimization techniques:

- Clear or invalidate each attachment at the start of a render pass using `loadOp = LOAD_OP_CLEAR` or `loadOp = LOAD_OP_DONT_CARE`.
- Set up any attachment that is only live during a single render pass as a `TRANSIENT_ATTACHMENT` that is backed by `LAZILY_ALLOCATED` memory.
- Ensure that the contents are invalidated at the end of the render pass using `storeOp = STORE_OP_DONT_CARE`.
- Ensure that the contents of read-only attachments that you did not modify during the pass, but still need to keep, are not written at the end of the render pass using `storeOp = STORE_OP_NONE`.

Render pass techniques to avoid

Arm recommends that you:

- Do not clear an attachment inside a render pass using `vkCmdClearAttachments()`. This is not free, unlike a clear or invalidate `loadOp` operation.
- Do not write a constant color using a shader program to manually clear a render pass.
- Do not use `loadOp = LOAD_OP_LOAD` unless your algorithm relies on the initial framebuffer state.
- Do not set `loadOp` or `storeOp` for attachments that are not needed in the render pass to avoid generating a needless round trip through the tile-memory for that attachment.
- Do not use `vkCmdBlitImage()` as a way of upscaling a low-resolution game frame to native resolution. Especially if you render the UI or HUD directly on top of the frame with `loadOp = LOAD_OP_LOAD`, as this is an unnecessary round trip to memory.

Negative impacts of inefficient render passes

Incorrect handling of render passes causes worse fragment shading performance and increased memory bandwidth.

Debugging render pass issues

Review the API usage of render pass creation and any use of `vkCmdClearColorImage()`, `vkCmdClearDepthStencilImage()`, and `vkCmdClearAttachments()`.

7.4 Multisampling for OpenGL ES

For most multisampling, all data for the additional samples are kept in the tile memory, which is inside the GPU. This data is resolved to a single pixel color as part of the tile write-back. This is efficient because the bandwidth for the additional samples never enters the external main memory.

Prerequisites

You must understand the following concepts:

- Anti-aliasing.
- OpenGL ES APIs.
- [The `EXT_multisampled_render_to_texture` extension](#)

Optimal multisampling performance

To get optimal render-to-texture multisampling performance, use the `EXT_multisampled_render_to_texture` extension. This extension can render multisampled data directly into a single-sampled image in memory, without needing a second pass.

How to optimize the use of MSAA

Try using the following optimization techniques:

- Use 4x MSAA because it is not expensive and provides good image quality improvements.
- Use the `EXT_multisampled_render_to_texture` extension render-to-texture multisampling.

MSAA techniques to avoid

Arm recommends that you:

- Do not use `glBlitFramebuffer()` to implement a multisample resolve.
- Do not use more than 4x MSAA without checking performance.

Negative impacts of implementing MSAA incorrectly

Failing to resolve multisampling inline results in higher memory bandwidth and reduced performance. For example, manually writing and resolving a 4xMSAA 1080p surface at 60 FPS requires 3.9GB/s of memory bandwidth. This is compared to 500MB/s when using the extension.

Debugging MSAA issues more effectively

Review any use of `glBlitFramebuffer()`.

7.5 Multisampling for Vulkan

For most multisampling, all data for the additional samples are kept in the tile memory, which is inside the GPU. This data is resolved to a value of a single pixel color as part of the tile write-back. This is efficient because the bandwidth for the additional samples never enters the external main memory.

Prerequisites

You must understand the following concepts:

- Anti-aliasing.
- Vulkan APIs.

Optimal multisampling performance for Vulkan

Multisampling is fully integrated with Vulkan render passes. Allowing the multisample resolve to be explicitly specified at the end of the subpass using the end of pass `resolveOp`.

Be aware that subpasses with different Multisample Anti-Aliasing (MSAA) levels cannot be fused. For example, in a scenario with 2 subpasses: `subpass 0` renders 4xMSAA to transient attachment; `subpass 1` reads `subpass 0` output and resolves to 1xMSAA output buffer which is written out. For

Arm GPUs, performance can be improved by making the second subpass also use 4xMSAA and then using write out resolve to 1xMSAA.

How to optimize the use of MSAA with Vulkan

Try using the following optimization techniques:

- Use 4x MSAA as it is not expensive and provides good image quality improvements.
- Use `loadOp = LOAD_OP_CLEAR` OR `loadOp = LOAD_OP_DONT_CARE` for the multisampled image.
- Use `pResolveAttachments` in a subpass to automatically resolve a multisampled color buffer into a single-sampled color buffer.
- Use `storeOp = STORE_OP_DONT_CARE` for the multisampled image.
- Use `LAZILY_ALLOCATED` memory to back the allocated multisampled images. No physical backing storage is required as they do not need to be stored in the main memory.

Vulkan MSAA steps to avoid

Arm recommends that you:

- Do not use `vkCmdResolveImage()`. Bandwidth and performance are negatively impacted.
- Do not use `storeOp = STORE_OP_STORE` for multisampled image attachments.
- Do not use `storeOp = LOAD_OP_LOAD` for multisampled image attachments.
- Do not use more than 4x MSAA without checking performance.
- Do not have subpasses with different MSAA levels unnecessarily.

The negative impact of implementing MSAA with Vulkan incorrectly

Failing to resolve multisampling inline results in higher memory bandwidth and reduced performance. For example, manually writing and resolving a 4x MSAA 1080p surface at 60 FPS requires 3.9GB/s of memory bandwidth. This is compared to 500MB/s when using an inline resolve.

Vulkan MSAA example

Example code of MSAA is available in the Vulkan Samples repository on github: [Vulkan MSAA example](#)

7.6 Multipass rendering

Multipass rendering is an important feature of Vulkan which enables applications to exploit the full power of tile-based architectures using the standard API.

Prerequisites

You must understand the following concepts:

- Anti-aliasing.
- Vulkan APIs.

- Using late ZS testing.
- Render passes and subpasses.

Enabling powerful algorithms

Arm GPUs can take color attachments and depth attachments from one subpass, and use them as input attachments in a later subpass without going through the main memory. This enables powerful algorithms, for example, deferred shading or programmable blending, to be used efficiently. However, you must have the correct set up.

Per-pixel storage requirements

Most Arm GPUs are designed for rendering 16x16 pixel tiles, with older Arm GPUs having 128 bits per pixel of tile buffer color storage. From Mali-G72, this count increased to up to 256 bits per pixel, and recent GPUs have up to 1024 bits per pixel. G-buffers, which require more color storage, can be used at the expense of requiring smaller tiles during fragment shading, which can reduce performance.

For example, a sensible G-buffer layout that fits neatly into a 128-bit budget could be:

- Light: B10G11R11_UFLOAT
- Albedo: RGBA8_UNORM
- Normal: RGB10A2_UNORM
- PBR material parameters/misc: RGBA8_UNORM

Image layouts

Multipass rendering is one of the few cases where image layout matters because it impacts driver-enabled optimizations. Initial layouts should be “safe” layouts, see [Transaction elimination](#) for more information.

Initial layouts

The following list is a sample multipass layout which hits all the good paths:

- Light: COLOR_ATTACHMENT_OPTIMAL
- Albedo: COLOR_ATTACHMENT_OPTIMAL
- Normal: COLOR_ATTACHMENT_OPTIMAL
- PBR: COLOR_ATTACHMENT_OPTIMAL
- Depth: DEPTH_STENCIL_ATTACHMENT_OPTIMAL

G-buffer pass (subpass #0) output attachments

- Light: COLOR_ATTACHMENT_OPTIMAL
- Albedo: COLOR_ATTACHMENT_OPTIMAL
- Normal: COLOR_ATTACHMENT_OPTIMAL
- PBR: COLOR_ATTACHMENT_OPTIMAL
- Depth: DEPTH_STENCIL_ATTACHMENT_OPTIMAL

To boost performance, the eventual output must occupy the first render target slot in the hardware, in this case, the light attachment. To do this, light attachment must be `attachment #0` in the `VkRenderPass`.

To enable the render to write out emissive parameters from the opaque material, light is included as an output from the G-buffer in this example. There is no extra bandwidth to write out an extra render target because the subpasses are merged.

Unlike a desktop GPU, there is no need to invent schemes to forward emissive light contributions through the other G-buffer attachments.

Lighting pass (subpass #1) input attachments

- Albedo: `SHADER_READ_ONLY_OPTIMAL`
- Normal: `SHADER_READ_ONLY_OPTIMAL`
- PBR: `SHADER_READ_ONLY_OPTIMAL`
- Depth: `DEPTH_STENCIL_READ_ONLY`

From the point that any pass starts to read from the tile buffer, optimize multipass performance by marking every depth or stencil attachment as read-only. `DEPTH_STENCIL_READ_ONLY` is designed for this read-only depth or stencil testing. You can use `DEPTH_STENCIL_READ_ONLY` concurrently as an input attachment to the shader program for programmatic access to depth values.

Lighting pass (subpass #1) output attachments

- Light: `COLOR_ATTACHMENT_OPTIMAL` - Lighting that is computed during `subpass #1`, is blended on top of the pre-computed emissive data from `subpass #0`. If needed, the application also blends transparent objects after the lighting passes have completed.

Subpass dependencies

Dependencies between the subpasses use `VkSubpassDependency` which sets the `DEPENDENCY_BY_REGION_BIT` flag. This dependency tells the driver that each subpass depends on the previous subpasses at that pixel coordinate.

For the example above, the subpass dependency setup would look like:

```
VkSubpassDependency subpassDependency = {};
subpassDependency.srcSubpass = 0;
subpassDependency.dstSubpass = 1;
subpassDependency.srcStageMask = VK_PIPELINE_STAGE_COLOR_ATTACHMENT_OUTPUT_BIT |
VK_PIPELINE_STAGE_EARLY_FRAGMENT_TESTS_BIT |
VK_PIPELINE_STAGE_LATE_FRAGMENT_TESTS_BIT;
subpassDependency.dstStageMask = VK_PIPELINE_STAGE_FRAGMENT_SHADER_BIT |
VK_PIPELINE_STAGE_COLOR_ATTACHMENT_OUTPUT_BIT |
VK_PIPELINE_STAGE_EARLY_FRAGMENT_TESTS_BIT |
VK_PIPELINE_STAGE_LATE_FRAGMENT_TESTS_BIT;
subpassDependency.srcAccessMask = VK_ACCESS_COLOR_ATTACHMENT_WRITE_BIT |
VK_ACCESS_DEPTH_STENCIL_ATTACHMENT_WRITE_BIT;
subpassDependency.dstAccessMask = VK_ACCESS_INPUT_ATTACHMENT_READ_BIT |
VK_ACCESS_COLOR_ATTACHMENT_READ_BIT |
VK_ACCESS_COLOR_ATTACHMENT_WRITE_BIT |
VK_ACCESS_DEPTH_STENCIL_ATTACHMENT_READ_BIT |
VK_ACCESS_DEPTH_STENCIL_ATTACHMENT_WRITE_BIT;
```

```
subpassDependency.dependencyFlags = VK_DEPENDENCY_BY_REGION_BIT;
```

Subpass merge considerations

The driver merges subpasses if the following conditions are met:

- Merging can save a write-out or read-back. Two unrelated subpasses which do not share any data do not benefit from multipass and are not merged.
- The number of unique `VkAttachments` used for input and color attachments in all considered subpasses is less than nine. However, keep in mind that depth or stencil does not count towards this limit.
- The depth or stencil attachment does not change between subpasses.
- Multisample counts are the same for all attachments.

How to optimize the use of multipass rendering with Vulkan

Try using the following optimization techniques:

- Use multipass.
- Use a 128-bit G-buffer budget for color.
- Use by-region dependencies between subpasses.
- Use `DEPTH_STENCIL_READ_ONLY` image layout for depth after the G-buffer pass is done.
- Use `LAZILY_ALLOCATED` memory to back images for every attachment except for the light buffer, which is the only texture that is written out to memory.
- Follow the basic render pass best practices, with `LOAD_OP_CLEAR` or `LOAD_OP_DONT_CARE` for attachment loads, and `STORE_OP_DONT_CARE` for transient stores.

Multipass rendering techniques to avoid

You must not store G-buffer data to memory, only to write the final color output.

Dynamic Rendering incompatibility

From Vulkan 1.3, Arm GPUs supports the `VK_KHR_dynamic_rendering` extension. The `VK_KHR_dynamic_rendering` extension simplifies renderpass building, but it also automatically disables subpass fusion. This created a trade-off between ease of use and optimal performance, that has been solved in driver version r50 with the `VK_KHR_dynamic_rendering_local_read` extension. This extension enables similar functionality to subpass fusion, and shaders created for a classic renderpass can be ported unmodified to dynamic rendering with tile buffer access preserved.

The negative impact of implementing multipass rendering incorrectly

Not using multipass correctly forces the driver to use multiple physical passes, sending intermediate image data back to the main memory between passes. In turn, losing all of the benefits of multipass rendering.

How to debug multipass rendering issues more effectively

To debug potential issues you encounter, try the following techniques:

- To determine if passes are being merged, refer to the GPU performance counters for information about the number of physical tiles that are rendered.
- The GPU performance counters also provide information about the number of fragment threads using late ZS testing. A high value in the late ZS test can be indicative of your application not using `DEPTH_STENCIL_READ_ONLY` correctly.

Reference

To learn about the tile bits per pixel of particular Arm GPUs, see [Arm GPU Datasheet](#).

Example

Example code of render subpasses is available in the Vulkan Samples repository on GitHub, see [Render Subpasses Tutorial](#) for more information.

7.7 HDR rendering

For mobile content, the relevant formats for rendering HDR images are `RGB10_A2`, for unorm data, and `B10R11G11` and `RGBA16F`, for floating-point data.

Prerequisites

You must understand the following concepts:

- HDR rendering.

How to optimize the use of HDR rendering

Try using the following optimization techniques:

- Use `RGB10_A2 UNORM` formats for rendering where small increases in dynamic range are required.
- Use `B10G11R11` for floating-point rendering as it is only 32bpp, compared to 64bpp with full fp16 float.

HDR rendering techniques to avoid

Arm recommends that you:

- Do not use `RGBA16F` unless:
 - `B10G11R11` is not providing suitable image quality.
 - Alpha is a requirement in the framebuffer.

Negative impacts of implementing HDR rendering incorrectly

The different types of impact you can see are:

- Increased bandwidth usage and reduced application performance.
- Fitting into 128bpp is difficult for multipass rendering to achieve efficiently.

7.8 Stencil updates

Many stencil masking algorithms toggle the stencil between a few values when creating and using a mask. When designing a stencil mask algorithm that uses multiple draw calls, the aim is to minimize the number of stencil buffer updates that occur.

Prerequisites

You must understand the following concepts:

- Stencil masking algorithms.

How to minimize stencil buffer updates

Try using the following optimization techniques:

- If the values are the same, then use `KEEP` rather than `REPLACE`.
- Some algorithms use pairs of draws: one to create the stencil mask, and one to color the unmasked fragments. In this case, use the second draw to reset the stencil value so it is ready for the next pair. Doing so avoids the need for a separate clear operation.

HDR rendering techniques to avoid

Do not waste performance by writing a new stencil value unnecessarily.

The negative impact of implementing stencil updates incorrectly

A fragment that writes a stencil value cannot be rapidly discarded. In turn, introducing an extra fragment processing cost that affects the performance of your application.

7.9 Blending

Blending is efficient on Arm GPUs because the `dstColor` is available on-chip, inside the tile buffer. However, it is important to remember that blending is more efficient for some formats than others.

Prerequisites

You must have a foundational knowledge of the following key graphics development concepts:

- Tile buffers.
- Blending.
- [Forward Pixel Kill \(FPK\)](#)
- Early ZS testing.

Optimized float blend modes

From Bifrost onwards float blending is enabled, but in Valhall GPUs there is accelerated hardware blending for FP16 and R11G11B10 formats. Simple blends of those formats are accelerated, but advanced blends are not. Advanced blends examples include those with logical operations or min/max.

If it is suspected that there is a costly blend, be sure to measure whether it is actually the bottleneck. If blending is the bottleneck, reduce overdraw where possible, before inspecting the blend steps.

How to optimize blending

Here are several optimizations that you can try to optimize your application on Arm GPUs when using blending:

- Prefer blending on unorm formats, rather than floating-point values.
- Always disable blending and alpha-to-coverage if an object is opaque.
- Monitor the number of blended layers that are being generated on a per-pixel basis. Even if the shaders are simple, high layer counts quickly consume cycles due to the number of fragments that must be processed.
- Consider splitting large UI elements into opaque and transparent portions. The opaque and transparent portions can then be drawn separately, allowing either early ZS, or FPK, to remove the overdraw beneath the opaque parts.

Avoiding suboptimal blending

Arm recommends that you:

- Do not use blending on floating-point framebuffers.
- Do not use blending on multisampled floating-point framebuffers.
- Do not generalize the user interface rendering code so that blending is always enabled.
- Do not just set alpha to 1.0 in the fragment shader to disable blending.

The negative impact of inefficient blending

Blending has a significant impact on performance since blending disables many of the important optimizations that remove fragment overdraw, such as early ZS testing and FPK. The negative impact is especially noticeable for user interfaces and 2D games that use multiple layers of sprites.

Debugging OpenGL ES issues for Arm GPUs

When the time comes to debugging any blending issues, use the Graphics Analyzer tool to step through the construction of a frame. You must also monitor which draws calls are being blended and the amount of overdraw that the blends create.

7.10 Transaction elimination

Transaction elimination is an Arm GPU technology that is used to avoid framebuffer write bandwidth for static regions of the framebuffer. Transaction elimination is beneficial for games that contain many static opaque overlays.

Prerequisites

You must understand the following concepts:

- Vulkan APIs.

- Image layouts.

Avoiding framebuffer write bandwidth in Vulkan

Transaction elimination is used for an image if the following conditions are met:

- The sample count is 1.
- The mipmap level is 1.
- The image uses `COLOR_ATTACHMENT_BIT`.
- The image does not use `TRANSIENT_ATTACHMENT_BIT`.
- The effective tile size is 16x16 pixels. Pixel data storage determines the effective tile size.

For GPUs before Mali-G51 there was the additional condition that only a single color attachment is being used.

The difference between safe and unsafe image layouts

Transaction elimination is a rare case where the driver uses the image layout. Whenever the image transitions from a safe to an unsafe image layout, the driver invalidates the transaction elimination signature buffer.

Safe image layouts are defined as image layouts that are either for read-only access or for color attachment writes.

These layouts are:

- `VK_IMAGE_LAYOUT_COLOR_ATTACHMENT_OPTIMAL`.
- `VK_IMAGE_LAYOUT_SHADER_READ_ONLY_OPTIMAL`.
- `VK_IMAGE_LAYOUT_TRANSFER_SRC_OPTIMAL`.
- `VK_IMAGE_LAYOUT_PRESENT_SRC_KHR`.
- `VK_IMAGE_LAYOUT_ATTACHMENT_OPTIMAL_KHR`.
- `VK_IMAGE_LAYOUT_READ_ONLY_OPTIMAL_KHR`.

All other layouts, including `VK_IMAGE_LAYOUT_UNDEFINED`, are considered unsafe as they allow writes to an image that is outside of the write path of the tile. For instance, `VK_IMAGE_LAYOUT_RENDERING_LOCAL_READ_KHR` is unsafe because storage image operations are not in the tile write path.

If the color attachment reference layout is different from the final layout, then the signature buffer invalidation can happen as part of `VkImageMemoryBarrier`, `VkCmdPipelineBarrier()`, `VkCmdWaitEvents()`, or `VkRenderPass`.

How to optimize transaction elimination on Arm GPUs

To optimize transaction elimination, try using the following optimization techniques:

- For color attachments, use the `COLOR_ATTACHMENT_OPTIMAL` image layout.
- To avoid unnecessary signature invalidation, use the safe image layouts for color attachments.

- To skip unneeded render target writes, use `storeOp=DONT_CARE` rather than `VK_IMAGE_LAYOUT_UNDEFINED`.

Vulkan input attachments requiring `VK_IMAGE_LAYOUT_GENERAL`

Some extensions, including `VK_EXT_rasterization_order_attachment_access`, require the input or color attachments to be in the `VK_IMAGE_LAYOUT_GENERAL` image layout for certain operations. This would normally stop transaction elimination. However, you can use the subpass layout functionality with the latest Arm GPU drivers to keep transaction elimination. To keep transaction elimination:

- Enter the render pass with safe layouts.
- Use safe layouts everywhere possible - avoid `VK_IMAGE_LAYOUT_GENERAL` where possible.
- Make sure `VK_IMAGE_LAYOUT_GENERAL` is the only unsafe layout used.
- Leave the render pass with safe layouts.



Note

`VK_KHR_dynamic_rendering` does not offer a similar implicit image transition path, so there is no equivalent of this method for dynamic rendering.

Transaction elimination techniques to avoid

Do not transition color attachments from safe to unsafe, unless the algorithm requires it.

The negative impact of not using transaction elimination

The loss of transaction elimination increases external memory bandwidth for scenes with static regions across frames. In turn, reducing the performance on systems that have limited memory bandwidth.

Debugging transaction elimination

To determine if transaction elimination is triggered, make the GPU performance counters count the number of tile writers that have been killed.

7.11 Variable rate shading

Variable Rate Shading (VRS) adds the ability to control the shading rate of fragment shaders. VRS improves power consumption and also gives a performance increase when your application is fragment bound.

You can specify the shading rate in the following ways:

- Drawcall level
- Primitive level
- Using screen-space attachments for the whole renderpass.

You can also use any combination of these three methods.

Prerequisites

You must understand the following concepts:

- Fragment and compute shading
- Attachments, especially screen-space attachments
- Drawcalls

Explore whether attachment or drawcall VRS works best

Specifying VRS per drawcall, or per screen-space attachment, is use-case dependent. Use the method that is appropriate for your usecase.

Specifying VRS per drawcall works well for the following cases:

- Fast-moving objects
- Out-of-focus objects
- With smoke and similar particle effects.

VRS screen-space attachments can involve combining a number of different factors. The following examples are when you would specify VRS screen-space attachments:

- Luminance variance
- Depth discontinuities
- Motion vectors
- Many other features



You can combine drawcall-level and attachment-level VRS without issue.

Use compute or fragment shading to generate screen-space attachments for your usecase

When generating screen-space attachments per frame, we generally advise to not use compute shaders. However, VRS screen-space attachments cannot use Arm Frame Buffer Compression (AFBC) because `VK_IMAGE_USAGE_FRAGMENT_SHADING_RATE_ATTACHMENT_BIT_KHR` disables AFBC. This means there is no loss in using storage images for the attachments when written by compute shaders. However, you must make sure that synchronization is handled correctly, so that compute jobs overlap with the fragment shading of other passes.

Use subgroup operations to perform reductions

Generating VRS screen-space attachments often involves shaders that perform reductions. For better performance, use subgroup operations instead of shared memory.

Profile to find best shading rates and tile sizes

8x8 tile size for VRS screen-space attachments often gives better performance in its renderpass. An 8x8 tile size gives better performance because it is more granular than bigger tile sizes, which is especially important in lower resolutions. However, generating the VRS screen-space attachment for an 8x8 tile size may be more costly, in terms of compute, compared to bigger tile sizes. For these reasons, start with an 8x8 tile size but then also experiment with bigger tile sizes.

An 8x8 tile size for VRS screen-space attachments often gives the best performance because it is more granular than bigger tile sizes. However, if the screen-space attachment is generated every frame, then the higher cost of generating smaller tile sizes can be more significant. Additionally, lower resolutions require a smaller tile size.

Similarly, shading rates greater than 2x2 do not always give the substantial improvement required to justify the quality loss, although GPUs after the Immortalis-G720 have improved performance at higher rates. To see what performance gains are achieved, experiment on the hardware you are targeting.

Avoid reprojecting render targets

VRS screen-space attachments are often generated with data from the previous frame. In this case, reprojecting the image from the previous frame is likely to be unnecessary overhead. If the lack of reprojection causes artifacts, you can try limiting the shading rate to 2x2.

Negative consequences of incorrect VRS usage

You must profile and monitor the quality of your output. If you use the VRS effectively, it allows significant performance gains without any noticeable loss of quality. However, if you use VRS incorrectly, you can lose both quality and performance.

8. Buffers and textures

This chapter contains information on how to optimize the performance of your in-game textures on an Arm GPU.

8.1 Buffer update for OpenGL ES

OpenGL ES implements a synchronous programming model with each API call behaving as if the rendering triggered by earlier API calls has completed. In reality, this is an illusion as rendering is handled asynchronously, often completing milliseconds after the triggering API call was made.

Prerequisites

You must understand the following concepts:

- Buffers.
- Pipeline draining.
- Resource ghosting.
- OpenGL ES APIs.

Buffer update

To maintain the illusion, the driver must track resources that are referenced by pending render commands. The driver locks them to prevent modification until those rendering commands have been completed.

If the application attempts to modify a locked resource, then the driver must take some evasive action. Either draining the pipeline until the lock is released, or creating a new ghost copy of the resource to contain the modifications. Both choices incur an overhead that the application can avoid.

How to optimize buffer updates

Try using the following optimization techniques:

- Perform modification of buffers using `glMapBufferRange()` and `MAP_UNSYNCHRONIZED` to explicitly bypass the synchronous rendering and resource lock semantics.
- Design buffer use to cycle through N buffers. Where the value of N is high enough to ensure that any resource locks have been released before a buffer is used again.
- Prefer complete buffer replacement using `glBufferData()` over partial replacement of buffers that are still referenced by pending commands.

Buffer update techniques to avoid

Do not use `glBufferSubData()` to modify buffers that are still referenced by pending commands. Practically this means not calling it more than once every frame or two for any given buffer. A second call can stall the GPU until any updated objects are through the rendering pipeline.

The negative impacts of inefficient buffer updates

Keep the following points in mind:

- Pipeline draining stalls the application until the resource lock is released. While this does not increase CPU load, GPU processing efficiency is reduced.
- CPU loads increase when resource ghosting requires a new allocated buffer and any parts that are not overwritten are to be copied from the original.

8.2 Robust buffer access

Vulkan devices support bound checking to the GPU memory accesses by using `robustBufferAccess`.

Prerequisites

You must understand the following concepts:

- Bounds checking.
- Vulkan APIs.
- Uniform buffers.

Adding bounds checking

Bounds checking ensures that accesses cannot fall outside of the buffer, preventing hard failure modes such as a GPU segmentation fault. However, do note that out-of-bounds access behavior, with bounds checking, is defined in the implementation. In turn, creating a render that cannot be predicted.



Enabling bounds checking causes loss in performance for accesses to uniform buffers and shader storage buffers.

How to optimize robust buffer access on Vulkan

Try using the following optimization techniques:

- Use `robustBufferAccess` as a debugging tool during development.
- Disable `robustBufferAccess` in release builds. Only leave it enabled if the application use-case requires the additional level of reliability due to the use of unverified user-supplied draw parameters.
- Push constants have no advantage over uniform buffers on Arm GPUs.
- Avoid dynamic indexing where possible, whether using push constants, or uniform buffer objects.

Robust buffer access techniques to avoid

Arm recommends that you do not enable:

- `robustBufferAccess` unless it is needed.
- `robustBufferAccess` without reviewing the performance impact that it can have on your application.

The negative impact of using robust buffer access

Even though `robustBufferAccess` has advantages, the use of `robustBufferAccess` causes measurable performance loss to both uniform buffers and shader storage buffers.

Debugging robust buffer access

Try the following debugging tips:

- To verify performance input, compare two test runs. One with `robustBufferAccess` enabled, and one without.
- The `robustBufferAccess` feature is a useful debug tool during development. If the application has problems with crashes or `DEVICE_LOST` errors being returned, then enable the robust access feature and see if the problem stops. If the problem does stop, there is either a draw call, or compute dispatch, that is making an out-of-bounds access.

8.3 Staging buffers

Most buffer use does not need staging buffers and you should avoid them when possible. However, tiled-optimal content needs staging buffers, and there may be other staging buffer requirements.

Prerequisites

You must understand the following concepts:

- Vulkan buffers
- Vulkan images
- Memory

VK_EXT_host_image_copy

The `VK_EXT_host_image_copy` extension is available from driver r50. The extension uses the CPU to copy direct between host memory and an image, and image to image copy. The `VK_EXT_host_image_copy` extension saves GPU memory and processor use by avoiding a staging buffer between them.

Only certain image layouts are supported, and you can retrieve these from the API.

The currently supported layouts are the following:

Copy source:

- `VK_IMAGE_LAYOUT_GENERAL`
- `VK_IMAGE_LAYOUT_TRANSFER_DST_OPTIMAL`

- `VK_IMAGE_LAYOUT_UNDEFINED`
- `VK_IMAGE_LAYOUT_PREINITIALIZED`

Copy destination:

- `VK_IMAGE_LAYOUT_GENERAL`
- `VK_IMAGE_LAYOUT_TRANSFER_DST_OPTIMAL`

Because only a limited number of layouts are supported, you may need a device-side layout transition with appropriate barriers.

Copying to or from host memory is most efficient with linear tiling. Image to image copy is fastest with the same tiling used for source and destination. If you are using optimal tiling for source and destination, we recommend to have the extents aligned with block size, or at least its subdivisions. Block size will be 16x16 elements for uncompressed formats, or 4x4 for compressed ones like ASTC; elements are 1 texel single format blocks for uncompressed and N texels for compressed.

Memory backing the images should be created with the `VK_MEMORY_PROPERTY_HOST_CACHED_BIT` flag.

The `VK_EXT_host_image_copy` extension is not compatible with AFBC or AFRC.

How to optimize staging buffers

Avoid using staging buffers where possible. Use staging buffers mainly for tiled-optimal content.

The negative impact of using staging buffers incorrectly

Unnecessary staging buffer use will result in wasted GPU memory, processing, bandwidth and power use.

8.4 Texture sampling performance

Depending on texture format and filtering mode, the Arm GPU texture unit can spend variable amounts of cycles sampling a texture. The texture unit is designed to give full-speed performance for both nearest sample and bilinear filtered, `LINEAR_MIPMAP_NEAREST`, texel sampling.

Prerequisites

You must understand the following concepts:

- Texture sampling.
- Texture formats.

Cases where extra texture sampling cycles are required

Ignoring data cache effects, the cases that require extra cycles are:

- Trilinear, `LINEAR_MIPMAP_LINEAR`, filtering has a 2x cost.
- 3D formats have a 2x cost.

- FP32 formats have a 2x cost.
- Depth formats have a 2x cost for Midgard GPUs, but only a 1x cost on Bifrost GPUs. For Valhall GPUs, Depth is generally a 1x cost. But, if there is a reference or comparison depth, then it is a 2x cost, for example Shadow Maps.
- Cubemap formats have a 1x cost per cube face that is accessed.
- For older Midgard, or first-generation Bifrost GPUs, YUV formats have an Nx cost, where N is the number of texture planes that are required. For the second-generation Bifrost and Valhall GPUs, meaning Mali-G51 onwards, YUV has a 1x cost, irrespective of the plane count.

For example, a trilinear filtered RGBA8 3D texture access takes four times longer to filter than a bilinear 2D RGBA texture access.

How to optimize texture sampling performance

Try using the following optimization techniques:

- Use the lowest texture resolution that you can.
- Use the narrowest precision that still retains a level of texture quality that is acceptable to you.
- For static resources, use offline texture compression, such as Ericsson texture Compression (ETC), Ericsson texture Compression 2 (ETC2), or ideally Adaptive Scaleable Texture Compression (ASTC).
- To improve both texture caching and image quality, always use mipmaps for textures that are in 3D scenes.
- Be careful to use trilinear filtering selectively if your content is texture-rate limited. Trilinear filtering gives the most benefits to textures with fine detail, such as text rendering.
- Use `texelFetch()` or `texture()` instead of the slower `imageLoad()` wherever possible.
- Use `mediump` samplers. `highp` samplers can be half the speed due to their wider data path requirements.
- If you need higher dynamic range, then consider using packed 32-bit formats. Such 32-bit formats include `RGB10_A2` or `RGB9_E5`, as an alternative to FP16 or FP32 textures.
- To lower the intermediate precision of the ASTC texture filtering, which can further improve performance and energy efficiency.
 - [Use the `EXT_texture_compression_astc_decode_mode` extension for OpenGL ES](#)
 - [Or the `VK_EXT_astc_decode_mode` extension for Vulkan](#)

Texture sampling techniques to avoid

Arm recommends that you:

- Do not use wider data types unless essential.
- Do not use trilinear filtering for everything. Use it on a case-by-case basis only.
- Do not use `textureGrad()` unless absolutely necessary. It is much slower than `texture()` and `textureLod()`. Replace with trilinear or bilinear filtering where possible. `textureGrad()` can be improved by ensuring the derivatives specified are the same for all threads in an aligned 2x2 pixel-quad.

Negative impacts of unoptimized texture sampling

The different types of impact you can see are:

- Applications can experience problems due to either texture cache pressure, or general external memory bandwidth issues when content is doing any of the following:
 - Loading too much texture data.
 - Is using sparse sampling due to missing mipmaps.
 - Is using wide data types.
- Content that makes effective use of texture compression and mipmapping are typically limited by filtering performance rather than external memory bandwidth. Such content only sees a measurable impact if the texture unit is the critical path unit for the shaders. However, if a shader is arithmetically expensive, then the texture filtering cost can be hidden beneath.

Debugging texture sampling

Try the following debugging tips:

- The GPU performance counters can show you the utilization of the texture unit to determine if your application is texture filtering limited. Also, external bandwidth counters can monitor traffic to the external system memory.
- Try disabling trilinear filtering to see if it improves the performance.
- Try clamping the texture resolution to see if the performance improves.
- Try narrowing the texture formats that you use to see if performance improves.

One final note

The Mali Offline Compiler can help you identify if your important shaders are texture-rate limited through the statistics about functional unit usage that it presents. However, the texture cycle counts from the tool assumes one cycle per texel performance. The compiler does not have visibility of the precise sampler or format that you plan on using. You must manually de-rate the texture performance and base the performance on your own texture and sampler usage.

Further reading

[Arm GPU Datasheet](#)

8.5 Anisotropic sampling performance

When determining the sample pattern to use, anisotropic filtering enables the texture sampling unit to account for the orientation of the triangle. Doing so improves image quality, in particular for primitives that are viewed at a steep angle regarding the view plan. However, anisotropic filtering comes at the cost of needing extra samples for a texture operation.

Prerequisites

You must understand the following concepts:

- Texture sampling.

- Bilinear filtering.
- Trilinear filtering.

Anisotropic Filtering (AF)

The following figure shows a cube that has been textured as a wooden crate. The image on the left uses traditional trilinear filtering, and the image on the right uses a bilinear filter with 2x AF. You can see the improved fidelity of the right-hand face of the cube when AF is used.

Figure 8-1: Trilinear vs Bilinear AF



From a performance perspective, the worst-case cost of AF is one sample per level of maximum anisotropy, which the application controls. Samples can be bilinear and use `LINEAR_MIPMAP_NEAREST`, or trilinear and use `LINEAR_MIPMAP_LINEAR`. Therefore, a 2x bilinear AF makes up to two bilinear samples. One significant advantage of anisotropic filtering is that actual number of samples that are made can be dynamically reduced. The reduction is based on the orientation of the primitive under the current sample position.

How to optimize anisotropic filtering

Try using the following optimization techniques:

- Start using a max anisotropy of two, and then assess if it provides sufficient quality. Higher numbers of samples can improve quality, but they also give diminishing returns which are often not worth the performance cost.
- Consider using 2x bilinear AF in preference to isotropic trilinear filtering. 2x bilinear is faster, and has better image quality in regions of high anisotropy. Note, that by switching to bilinear filtering, you can see some visible seams at the decision point between mipmap levels.
- Only use AF and trilinear filtering for objects that benefit from it the most.

Anisotropic filtering techniques to avoid

Arm recommends that you:

- Do not use higher levels of max anisotropy without reviewing performance. 8x bilinear AF costs eight times more GPU computational power than a simple bilinear filter.
- Do not use trilinear AF without reviewing performance. 8x trilinear AF costs 16 times more than a simple bilinear filter.

Negative impacts of using incorrect anisotropic filtering

The different types of impact you can see are:

- Using 2x bilinear AF, instead of trilinear filtering, increases image quality and can also improve performance.
- Using high levels of max anisotropy can improve image quality, but at the cost of performance.

Debugging anisotropic filtering performance

Arm recommends that to debug a performance problem with texture filtering, first try disabling AF completely and measuring any improvement. Then, incrementally increase the amount of max anisotropy that is allowed, and assess whether the quality is worth the additional cost to performance.

8.6 Texture and sampler descriptors

Arm GPUs cache texture and sampler descriptors in a control structure cache that can store a variable number of descriptors, depending on their content.

Prerequisites

You must understand the following concepts:

- Texture caching.
- Sampler descriptors.

Directions for maximizing the cache capacity of descriptor entries

Arm recommends that you use the following descriptor settings. Doing so ensures that you reach the maximum cache capacity in terms of descriptor entries, and therefore, the best performance:

For OpenGL ES only

- Set `GL_TEXTURE_WRAP_(S|T|R)` to identical values



The OpenGL ES driver can specialize the sampler state that is based on the current texture, unlike with Vulkan. So there is no need to set `GL_TEXTURE_WRAP_R` for 2D textures.

-
- Do not use `GL_CLAMP_TO_BORDER`.
 - Set `GL_TEXTURE_MIN_LOD` to the default of `-1000.0`
 - Set `GL_TEXTURE_MAX_LOD` to the default of `+1000.0`

- Set `GL_TEXTURE_BASE_LEVEL` to the default of 0
- Set `GL_TEXTURE_SWIZZLE_R` to the default of `GL_RED`.
- Set `GL_TEXTURE_SWIZZLE_G` to the default of `GL_GREEN`.
- Set `GL_TEXTURE_SWIZZLE_B` to the default of `GL_BLUE`.
- Set `GL_TEXTURE_SWIZZLE_A` to the default of `GL_ALPHA`.
- If the `EXT_texture_filter_anisotropic` filtering extension is available, then set `GL_TEXTURE_MAX_ANISOTROPY_EXT` to 1.0.

For Vulkan only

For Vulkan, when populating the `vkSamplerCreateInfo` structure:

- Set sampler `addressMode(U|V|W)` so they are all the same.



`addressModeW` must be set to be the same as `U` and `V`, even when sampling a 2D texture.

- Set sampler `mipLodBias` to 0.0
- Set sampler `minLod` to 0.0
- Set sampler `maxLod` to `VK_LOD_CLAMP_NONE`
- Set sampler `anisotropyEnable` to `VK_FALSE`
- Set sampler `maxAnisotropy` to 1.0
- Set sampler `borderColor` to `VK_BORDER_COLOR_FLOAT_TRANSPARENT_BLACK`
- Set sampler `unnormalizedCoordinates` to `VK_FALSE`.

When populating the `vkImageViewCreateInfo` structure:

- Set every field that is in a view component to either `VK_COMPONENT_SWIZZLE_IDENTITY` or to the explicit per-channel identity-mapping equivalent.
- Set view `subresourceRange.baseMipLevel` to 0.

Think whether it makes sense for you to have separate or combined texture and sampler descriptors. Either works for Arm GPUs, so use separate where it reduces the number of samplers.

Using the Nearest and Linear filters in OpenGL ES and Vulkan

For emulating `GL_NEAREST` and `GL_LINEAR` sampling for mipmapped textures, the requirements for maximizing descriptor storage conflicts with the Vulkan recommended specification approach.

The Vulkan specification states that there are no Vulkan filter modes that directly correspond to OpenGL minification filters of `GL_LINEAR` or `GL_NEAREST`. However, required filters can be emulated using `VK_SAMPLER_MIPMAP_MODE_NEAREST`, `minLod = 0`, `maxLod = 0.25`, and using `minFilter = VK_FILTER_LINEAR` or `minFilter = VK_FILTER_NEAREST`, respectively.

To emulate these two texture filtering modes for a texture with multiple mipmaps levels, while also being compatible with the requirements for compact samplers, the recommended application behavior is to create a unique `VkImageView` instance. The `VkImageView` instance references only the level 0 mipmap and uses a `VkSampler` with `pCreateInfo.maxLod` setting to `VK_LOD_CLAMP_NONE` in accordance with the compact sampler restrictions.

Direct access to textures through `imageLoad()` and `imageStore()` in shader programs, or the equivalent in SPIR-V, are not impacted by this issue. However, remember to always prefer `texelFetch()` OVER `imageLoad()` where possible.

Behaviors to avoid when optimizing texture and sampler descriptors

Do not set `maxLod` to the maximum mipmap level in the texture chain. Instead, use `VK_LOD_CLAMP_NONE`. Otherwise, you can experience a reduced texture filtering throughput.

The negative impact of unoptimized texture and sampler descriptors

Expect reduced texture filtering throughput.

8.7 sRGB textures

sRGB textures are natively supported in Arm GPU hardware. Both sampling from and rendering or blending to an sRGB surface comes at no performance cost. sRGB textures have a better perceptual color resolution than non-gamma corrected formats at the same bit depth. Therefore, Arm encourages the use of sRGB textures.

Prerequisites

You must understand the following concepts:

- sRGB textures.
- Color formats.

How to optimize sRGB texture performance

Try using the following optimization techniques:

- Use sRGB textures for improved color quality.
- Use sRGB framebuffers for improved color quality.
- Remember that Adaptive Scalable Texture Compression (ASTC) supports sRGB compression modes for offline compressed textures.

Something to avoid when using sRGB textures

Do not use 16-bit linear formats to gain perceptual color resolution when 8-bit sRGB can suffice.

The negative impact of using an unoptimized texture format

The different types of impact you can see are:

- Not using sRGB textures, where appropriate, can reduce the quality of the image that is being rendered.

- Using wider float formats in place of sRGB textures increases bandwidth and reduces performance.

8.8 AFBC textures

Arm FrameBuffer Compression (AFBC) is a lossless image compression format, which is supported from the Mali-T760 GPU onwards. AFBC can be used for compressing framebuffer outputs from the GPU.

Prerequisites

You must understand the following concepts:

- Texture compression techniques.
- Framebuffer attachments.

Using AFBC

When enabled, AFBC is automatic and functionally transparent to the application. However, be aware of some areas where AFBC cannot be used, and others that can require the driver to insert runtime decompression from AFBC back to an uncompressed pixel format.

How to optimize the use of AFBC textures

Try using the following optimization techniques:

- Use the `texture()` and `texelFetch()` functions in shaders to access textures and images that the GPU previously rendered as framebuffer attachments.
- When packing data into color channels, to get the best compression rates, store the most volatile bits in the least significant bits of the channel.

Something to avoid when using AFBC textures

Do not use `imageLoad()` or `imageStore()` to read or write into a texture or image that the GPU has rendered as a framebuffer attachment. Doing so triggers decompression.

The negative impact of not using AFBC textures correctly

Arm recommends that you keep in mind that the incorrect use of AFBC can trigger decompression. Decompression increases memory bandwidth usage and decreases performance.

8.9 AFBC textures for Vulkan

Arm FrameBuffer Compression (AFBC) is a lossless image compression format. AFBC can be used for compressing framebuffer outputs from the GPU. Partial Vulkan support for AFBC is available from Mali-G71 onwards, and full support from Mali-G31, Mali-G51, and Mali-G76.

Prerequisites

You must understand the following concepts:

- Texture compression techniques.
- Framebuffer attachments.

Using AFBC

When enabled, AFBC is automatic and functionally transparent to the application. However, be aware of some areas where AFBC cannot be used. And of others that can require the driver to insert runtime decompression from AFBC back to an uncompressed pixel format.

For AFBC usage, `VkImage` requires:

- `VkSampleCountFlagBits` must be `VK_SAMPLE_COUNT_1_BIT`.
- `VkImageType` must be `VK_IMAGE_TYPE_2D`.
- `VkImageTiling` must be `VK_IMAGE_TILING_OPTIMAL`.

As a rule, only enable image flags that are needed. Setting the `VkImageUsageFlags` `VK_IMAGE_USAGE_STORAGE_BIT` or setting the `VkImageCreateFlags` `VK_IMAGE_CREATE_ALIAS_BIT` stops AFBC being used. Similarly, setting `VkImageUsageFlags` `VK_IMAGE_USAGE_TRANSIENT_ATTACHMENT_BIT` or `VkImageCreateFlags` `VK_IMAGE_CREATE_MUTABLE_FORMAT_BIT` stops AFBC use on GPUs before Valhall. Other `VkImage` flags may also stop AFBC, depending on the device.

To query the compression properties of an image, add a `VkImageCompressionPropertiesEXT` structure to the pNext chain of the `VkSubresourceLayout2EXT` structure in a call to `VkGetImageSubresourceLayout2KHR` or `VkGetImageSubresourceLayout2EXT`. `VkImageCompressionPropertiesEXT::imageCompressionFlags` will contain `VK_IMAGE_COMPRESSION_DEFAULT_EXT` if AFBC is used, `VK_IMAGE_COMPRESSION_DISABLED_EXT` if there was no compression, and `VK_IMAGE_COMPRESSION_FIXED_RATE_EXPLICIT_EXT` if Arm Fixed Rate Compression (AFRC) is used.

From the Valhall generation of Arm GPUs, `VkImage` can have any `VkFormat` 32 bits or smaller. From the Mali-G710 onwards, all `VK_FORMAT_R16G16B16A16_*` `VkFormats` are supported, including float. For Arm GPUs released before Valhall, only a subset was supported.

The supported subset in Bifrost GPUs is as follows:

- `VK_FORMAT_R4G4B4A4_UNORM_PACK16`
- `VK_FORMAT_B4G4R4A4_UNORM_PACK16`
- `VK_FORMAT_R5G6B5_UNORM_PACK16`
- `VK_FORMAT_R5G5B5A1_UNORM_PACK16`
- `VK_FORMAT_B5G5R5A1_UNORM_PACK16`
- `VK_FORMAT_A1R5G5B5_UNORM_PACK16`
- `VK_FORMAT_B8G8R8_UNORM`
- `VK_FORMAT_B8G8R8A8_UNORM`
- `VK_FORMAT_B8G8R8A8_SRGB`

- `VK_FORMAT_A8B8G8R8_UNORM`
- `VK_FORMAT_A8B8G8R8_SRGB`
- `VK_FORMAT_A8R8G8B8_SRGB`
- `VK_FORMAT_B10G10R10A2_UNORM`
- `VK_FORMAT_R4G4B4A4_UNORM`
- `VK_FORMAT_R5G6B5_UNORM`
- `VK_FORMAT_R5G5B5A1_UNORM`
- `VK_FORMAT_R8_UNORM`
- `VK_FORMAT_R8G8_UNORM`
- `VK_FORMAT_R8G8B8_UNORM`
- `VK_FORMAT_R8G8B8A8_UNORM`
- `VK_FORMAT_R8G8B8A8_SRGB`
- `VK_FORMAT_A8R8G8B8_UNORM`
- `VK_FORMAT_R10G10B10A2_UNORM`
- `VK_FORMAT_D24_UNORM_S8_UINT`
- `VK_FORMAT_D16_UNORM`
- `VK_FORMAT_D32_SFLOAT`

How to optimize the use of AFBC textures

Try using the following optimization techniques:

- Only enable `vkImage` flags that are needed.
- Only use 2D images and optimal tiling, along with `VK_SAMPLE_COUNT_1_BIT`.
- Only use supported `VkFormatS`.
- When packing data into color channels, to get the best compression rates, store the most volatile bits in the least significant bits of the channel.

Something to avoid when using AFBC textures

Minimize use of storage images as they cannot use AFBC compression.

The negative impact of not using image flags correctly

Incorrect use of Vulkan image flags can cause AFBC to not be used or textures to be decompressed. Decompression increases memory bandwidth usage and decreases performance.

Example

Example code for AFBC is available in the Vulkan samples repository on GitHub: [AFBC Tutorial](#)

8.10 AFRC

Arm Fixed Rate Compression (AFRC) is a lossy image compression format. You can enable AFRC in Vulkan and OpenGL ES with extensions.

Prerequisites

You must understand the following concepts:

- Texture compression techniques.
- Framebuffer attachments.

Using AFRC

AFRC is available in Arm Immortalis™ GPUs and in Arm Mali GPUs from the Mali-G510, Mali-G310 and Mali-G715 onwards, but it is not enabled by default.

You can enable AFRC for different GPU tasks:

- To enable AFRC on display images, use the `VK_EXT_image_compression_control_swapchain` extension in Vulkan, or the OpenGL ES equivalent `EGL_EXT_surface_compression`.
- To enable AFRC for render to texture usages, use `VK_EXT_image_compression_control` in Vulkan, or the OpenGL ES equivalent `GL_EXT_texture_storage_compression`.
- To share AFRC images between GPU and ISP, use the `VK_EXT_image_compression_control` extension or the OpenGL ES equivalent `GL_EXT_EGL_image_storage_compression`.

Using AFRC gives smaller and predictable memory footprints compared to Arm FrameBuffer Compression (AFBC). However, AFRC is a lossy compression format, so you must use caution when using AFRC for render to texture, for example, because losses are magnified with each compression pass. The main, and safest, use of AFRC is for display images.

Table 8-1: Examples of what compression is achievable

Format	Bitrates	Vulkan setting	Comp Ratio
RGBA8	16 bpp 12 bpp 8 bpp	VK_IMAGE_COMPRESSION_FIXED_RATE_4BPC_BIT_EXT	2:1
		VK_IMAGE_COMPRESSION_FIXED_RATE_3BPC_BIT_EXT	3:1
		VK_IMAGE_COMPRESSION_FIXED_RATE_2BPC_BIT_EXT	4:1
RGB8	15 bpp 12 bpp 6 bpp	VK_IMAGE_COMPRESSION_FIXED_RATE_5BPC_BIT_EXT	3:2
		VK_IMAGE_COMPRESSION_FIXED_RATE_4BPC_BIT_EXT	2:1
		VK_IMAGE_COMPRESSION_FIXED_RATE_2BPC_BIT_EXT	4:1
RG8	8 bpp 6 bpp 4 bpp	VK_IMAGE_COMPRESSION_FIXED_RATE_4BPC_BIT_EXT	2:1
		VK_IMAGE_COMPRESSION_FIXED_RATE_3BPC_BIT_EXT	3:1
		VK_IMAGE_COMPRESSION_FIXED_RATE_2BPC_BIT_EXT	4:1
R8	4 bpp 3 bpp 2 bpp	VK_IMAGE_COMPRESSION_FIXED_RATE_4BPC_BIT_EXT	2:1
		VK_IMAGE_COMPRESSION_FIXED_RATE_3BPC_BIT_EXT	3:1
		VK_IMAGE_COMPRESSION_FIXED_RATE_2BPC_BIT_EXT	4:1

AFRC support limitations

AFRC has the following limitations:

- Only works with 2D images
- Does not support mutable views.
- You cannot use AFRC for storage images
- It has no MSAA support.
- Only works on a subset of Vulkan formats, so you should always query supported compression properties with `VkImageCompressionPropertiesEXT` before enabling. In OpenGL ES, you can check supported compression properties using the `SURFACE_COMPRESSION_EXT` and `EGL_SURFACE_COMPRESSION_EXT` attributes.

AFRC has all the restrictions of [AFBC textures](#) and [AFBC textures for Vulkan](#). For example, to access textures and images, you should use the `texture()` and `texelFetch()` functions in shaders. Do not use `imageLoad()` or `imageStore()`.

With EGL window surfaces in OpenGL ES, `EGL_KHR_partial_update` can cause visible corruption if used with AFRC and so these commands should not be used together.

Best practices when using AFRC textures

When turning off AFRC, do not accidentally turn off AFBC in the process. To turn off AFRC in Vulkan, use the flag `VK_IMAGE_COMPRESSION_DEFAULT_EXT`. You must not use the `VK_IMAGE_COMPRESSION_DISABLED_EXT` flag, because this deactivates AFBC.

The negative impact of using AFRC incorrectly

Overuse of AFRC on textures can lead to a loss of quality. However, if you do not use AFRC, you can miss out on potential bandwidth and performance savings.

Example

Example code for AFRC is available in the Vulkan samples repository on GitHub, see [Image Compression Control Tutorial](#) for more information. There is also an [Arm Learning Path on AFRC](#).

8.11 Runtime compression

Where possible, use offline texture compression over runtime compression. If the offline approach is not possible, for example, when texture data is generated at runtime, you can use Arm compression technologies. A final option is implementing compression manually in a compute shader.

Prerequisites

You must understand the following concepts:

- Texture compression formats:
 - Adaptive Scalable Texture Compression (ASTC)
 - Ericsson Texture Compression (ETC)
- Arm compression technologies:

- Arm Fixed Rate Compression (AFRC)
- Arm FrameBuffer Compression (AFBC)

Use offline compression if possible

Any static resources should be compressed offline. Only use runtime compression where offline compression is not an option. For offline texture compression, we recommend ASTC, but ETC and ETC2 are other available options. See [Texture sampling performance](#).

Consider AFBC

If you are unable to use the offline approach, AFBC is lossless and can help reduce bandwidth, but does not reduce the overall memory footprint. The driver automatically enables AFBC compression when possible. See [AFBC textures for Vulkan](#) for more information about AFBC.

Consider AFRC

If a lossy compression is acceptable, AFRC is a good solution for this case. Because AFRC is fixed rate, you get the benefit of reduced memory footprint as well as reduced bandwidth. See [AFRC](#) for more information on using AFRC.

If the lossiness is an issue, or AFRC is not available, consider using AFBC instead.

Shader based compression

If AFBC or AFRC do not suit your usecase, you can implement an ASTC encoder in shader code.

ASTC is a complex format. We recommend that you pick a subset to implement to make the encoder less complicated. Alternatively, you can try ETC2, which is less complex.

Implementing a shader is not a trivial task. The additional compute load may impact the performance of your application. Consider how efficiently it can be done and what trade-offs are necessary.

9. Compute shading

This chapter covers how to optimize workgroup sizes, how to correctly use shared memory on an Arm GPU, and optimized ways to process images.

9.1 Image processing

One common use case for compute shaders is for image post-processing effects. Remember however, that fragment shaders have access to many fixed-function features in the hardware. Such features can speed up performance, reduce power, and even reduce bandwidth. Ideally use fragment shaders for image processing where possible.

Prerequisites

You must understand the following concepts:

- Fragment shading.
- Compute shading.
- Pipeline bubbles.
- Thread quads.

Image processing advantages

Here are some advantages to using fragment shading for image processing:

- Texture coordinates are interpolated using fixed function hardware when using varying interpolation. In turn, freeing up shader cycles for more useful workloads.
- Write out to memory can be done using the tile-writeback hardware in parallel to shader code.
- There is no need to range check `imageStore()` coordinates. Doing so can be a problem when you are using workgroups that do not subdivide a frame completely.
- It is possible to do framebuffer compression and transaction elimination.

Here are some advantages to using compute shading for image processing:

- It can be possible to exploit shared data sets between neighboring pixels. Doing so avoids extra passes for some algorithms.
- It is easier to work with larger working sets per thread, avoiding extra passes for some algorithms.
- For complicated algorithms such as Fast Fourier Transforms (FFTs), which require multiple fragment render passes, it is often possible to merge into a single compute dispatch.

How to optimize the use of image processing

Try using the following optimization techniques:

- Use fragment shaders for simple image processing.

- For more complicated scenarios try using, and monitoring the performance of, compute shaders.
- Use `texture()` or `texelFetch()` instead of `imageLoad()` for reading read-only texture data. `texture()` works with Arm FrameBuffer Compression (AFBC) textures that have been rendered by previous fragment passes. Using `texture()` also load balances the GPU pipelines better because `texture()` operations use the texture unit and both `imageLoad()` and `imageStore()` use the load or store unit. The load/store units are often already being used in compute shaders for generic memory accesses.
- To improve `imageStore()`, memory coherence across threads in a quad is important. Therefore, make sure to group four adjacent lanes to use addresses in the same 64 byte cache line. For compute shaders, threads are packed linearly across X then Y then Z. Correctly combining the thread packing with the texture storage pattern allows you to layout your compute shader thread group in the same shape as the `imageStore()` accesses. This more efficient memory use provides a significant performance improvement. This is also true for `imageLoad()` if it cannot be avoided.

For more detail on the last point, consider that for the storage pattern `VK_IMAGE_TILING_LINEAR` the texture is laid out linearly. For a quad of threads, the y and z coordinates must stay the same and only the x coordinate can vary, in groups of 0-3, 4-7, and so on. For `VK_IMAGE_TILING_OPTIMAL` the texture is block-interleaved, which means the y co-ordinate can vary as well, and there are options for your quad layout.

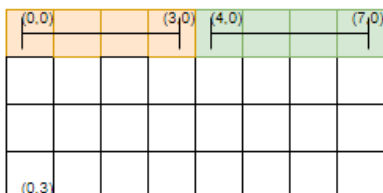
If there are 8 bytes per texel then y can only vary in groups of 2 (0-1, 2-3, and so on); but if there are 4 or fewer bytes per texel, then y can vary in groups of 0-3, 4-7, like the x co-ordinate. See the below figure for possible thread quad to texel layouts for the different storage patterns. Working with this knowledge, you can then either layout your workgroup in the same shape as your accesses must be, or layout linearly and manually swizzle your `gl_GlobalInvocationID` into coordinates for `imageStore()`.

Note that in the following figure, the colored areas represent memory blocks that your quad must be selected within.

Figure 9-1: Texel to quad mapping

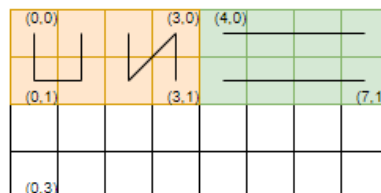
Potential texel mapping of Quads with:

`VK_IMAGE_TILING_LINEAR`



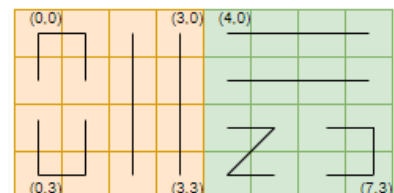
`VK_IMAGE_TILING_OPTIMAL`

8 bytes/texel



`VK_IMAGE_TILING_OPTIMAL`

1,2 or 4 bytes/texel



Things to avoid when optimizing your image processing implementation

Arm recommends that you:

- Do not use `imageLoad()` in compute unless you must use coherent read and writes within the dispatch.
- Avoid `textureGrad()` unless it is essential. It is much slower than `texture()` and `textureLod()`. Instead, replace with trilinear or bilinear filtering where possible. Note: `textureGrad()` can be improved by ensuring the derivatives specified are the same for all threads in an aligned 2x2 pixel-quad.
- Do not use compute to process images generated by fragment shading. Doing so creates a backwards dependency that can cause a bubble. If fragment shader outputs are consumed by fragment shaders of later render passes, then render passes go through the pipeline more cleanly.

The negative impact of not using the correct image processing method

Compute shaders can be slower and less energy-efficient than fragment shaders for simple post-processing workloads. Examples of such workloads include: downscaling, upscaling, and blurs.

9.2 Workgroup sizes

Arm Midgard, Bifrost, and Valhall GPUs have a fixed number of registers available in each shader core. These GPUs can split those registers across a variable number of threads depending on the register usage requirements of the shader program.

Prerequisites

You must understand the following concepts:

- Shader core resource scheduling.
- Workgroups.
- Stack memory.

Using workgroups

The GPU hardware can split up, and then merge, workgroups during shader core resource scheduling. If barriers or shared memory are used, then GPUs cannot do this with workgroups. In such a case, all work items in the workgroup must be executed concurrently in the shader core.

Large workgroup sizes restrict the number of registers that are available to each work item in this scenario. In turn, forcing shader programs to use stack memory if insufficient registers are available.

Arm GPUs currently use subgroup sizes of 16 and smaller, so implementations which assume larger warp-sizes may require modification to function correctly. Be especially careful with optimizations which attempt to avoid barrier calls.

How to optimize the use of workgroup sizes

Try using the following optimization techniques:

- Use 64 as a baseline workgroup size.
- Use a multiple of 4 as a workgroup size.

- Try smaller workgroup sizes before larger ones, especially if using barriers or shared memory.
- When working with images or textures, use a square execution dimension, for example 8x8, to exploit optimal 2D cache locality.
- If a workgroup has per-workgroup work to be done, consider splitting the work into two passes. Doing so avoids barriers and kernels that contain portions where most threads are idle.
- Compute shader performance is not always intuitive, so keep measuring the performance levels.

Things to avoid when optimizing workgroup sizes

Arm recommends that you:

- Do not use more than 64 threads per workgroup.
- Do not assume that barriers with small workgroups are free from performance costs.

Negative impacts of not using workgroup sizes correctly

The different types of impact you can see are:

- Be careful with large workgroups. If a high percentage of work items are waiting on a barrier, then the shader core can be starved of work.
- Shaders that spill to the stack incur a higher load and store unit utilization, along with a higher cost to external memory bandwidth.

9.3 Shared memory

Arm GPUs do not implement dedicated on-chip shared memory for compute shaders. The shared memory that is available to use is system RAM that is backed up by the load-store cache.

Prerequisites

You must understand the following concepts:

- Cache memory.
- Memory allocation.
- Reduction shaders.

How to optimize the use of shared memory on Arm GPUs

Try using the following optimization techniques:

- Use shared memory to share significant computation between threads in a workgroup.
- Keep your shared memory as small as possible, as it reduces the chance of thrashing the data cache.
- To reduce the size of the shared memory that is needed, reduce the precision and data widths.
- You need barriers to synchronize access to shared data. Shader code that has been ported from desktop development can sometimes omit barriers due to GPU-specific assumptions on warp width. Such an approach is not safe to use on mobile GPUs.

- It can be computationally cheaper splitting an algorithm over multiple shaders when compared to inserting barriers.
- For barriers, smaller workgroups are less expensive.

Things to avoid when optimizing shared memory use

Arm recommends that you:

- Do not copy data from global memory to shared memory on Arm GPUs. Doing so pollutes the caches.
- Do not use shared memory to implement code. For example:

```
if (localInvocationID == 0) {
    common_setup();
}
barrier();
// Per-thread workload here
barrier();
if (localInvocationID == 0) {
    result_reduction();
}
```

Splitting the example problem into three shaders would be an improvement. The setup and reduction shaders would need fewer threads.

The negative impact of using shared memory incorrectly

Depending on the algorithmic design of the compute shader that you use, the negative impact of using shared memory incorrectly is specific to your application.

10. Shader code

Writing optimal shader code through correct precision, vectorizing, uniforms and other techniques is key to optimizing the graphics performance of your application.

10.1 Minimize precision

Arm GPUs have full support for reduced precision in the shader core register file and arithmetic units. Also, reducing precision on inputs and outputs saves data bandwidth. Using 16-bit precision is normally sufficient for computer graphics, especially for fragment shading when computing an output color.

Prerequisites

You must understand the following concepts:

- Different precision types, including `lowp` and `mediump`, FP16, and FP32.
- Maintaining correct values through calculations involving reduced precision variables and temporaries.
- HLSL and GLSL.

Marking variables and temporaries

Both ESSL and Vulkan GLSL support marking variables and temporaries to use reduced precision level with `mediump`. There is no benefit to using `lowp` for Arm GPUs as it is functionally identical to `mediump`.

Forcing 16-bit floating point

There are two aspects to 16-bit float support: storage and arithmetic. 16-bit storage support reduces bandwidth and memory requirements. 16-bit arithmetic support improves throughput.

In Vulkan, enforced 16-bit support is available through the extension `VK_KHR_shader_float16_int8` for arithmetic and the extension `VK_KHR_16bit_storage`, which is core to Vulkan 1.1, for storage. There are examples of their use at the end of the chapter.

Where graphics drivers can ignore `mediump`, with explicit FP16 through the extensions drivers producing known performance. But where both FP32 and FP16 variant shaders must be provided, `mediump` gives convenience if different results on different devices are acceptable. Explicit FP16 is especially useful for specialized compute shaders with FP16 kernels.

Which precision to use

FP16 is not usually accurate enough to use for every vertex output in your application. Use cases which work well with FP16 precision:

- Normals / Tangent / Bi-tangent
- Vertex colors

- Any auxiliary data which is centered around 0 and does not need exceptionally high precision.

Use cases which can work well with FP16 precision:

- Local world position. Precision is significantly improved with `delta_pos = f16vec3(world_position - camera_position);`. This gives the property that the closer you get to the camera, the better precision you get. For mobile, precision is probably acceptable.
- Texture coordinates with smaller texture resolutions and constrained UV range. If UVs can be kept between `[-1, 1]`, there is reasonable resolution in FP16.

Use cases unlikely to work:

- Global world position.
- UI texture coordinates.

Beware of overflowing the precision, and getting *Infinity* or *NaN*. In ESSL 1 propagation of *Infinity* and *NaN* through calculations is usually suppressed, but with ESSL 3 and Vulkan values are likely to propagate through all affected calculations. `VK_KHR_shader_float_controls` can be used to change default *Infinity* and *NaN* propagation behaviour.

How to optimize the use of minimized precision

Try using the following optimization techniques:

- Use `mediump` when the resulting precision is acceptable.
- Use `mediump` for inputs, outputs, variables, and samplers where possible.
- Extra precision for `mediump` values can be gained by centering around zero and making use of the floating-point sign bit. For example, with angles try to use a range of `-PI` to `+PI`, rather than `0` to `2PI`.
- Use `mediump` where you must make shader variants for FP16 and FP32, for example graphics fragment shaders. Use explicit FP16 for cases like specially optimized compute shaders.
- If using explicit FP16, make sure to vectorize the code by using `f16vec2` or `f16vec4`. Modern GPU architectures use packed `f16x2` instructions to improve arithmetic performance. Scalar `float16_t` does not gain the same benefit.

HLSL

HLSL compiles to SPIR-V with the right tool, or pipeline of tools. Whether relaxed, or explicit, precision is required can depend on the tool and what other platforms are targeted. However, `minFloat16` works well for relaxed precision and `float16` works well for explicit precision.

Things to avoid when optimizing your use of minimized precision

Be aware of the following pitfalls when using minimized precision:

- Do not test the correctness of `mediump` precision on desktop GPUs. Desktop GPUs usually ignore `mediump` and process it as `highp`. There is no difference in function or performance, so the test is worthless.
- Do not cast between FP16 and FP32 too much. Most GPUs must spend cycles when converting between FP16 and FP32.

The negative impact of using minimized precision incorrectly

If you choose to use full FP32 precision, performance and power efficiency can be negatively impacted. `highp` can use up to twice the resources and halve the speed.

How to debug precision-related performance issues

Try forcing `mediump` or FP16 for everything except for the contributors to `gl_Position`, and then compare the performance difference afterwards.

Vulkan Examples

Example code for 16-bit arithmetic is available in the Vulkan Samples repository on GitHub: [16-Bit Arithmetic tutorial](#)

Example code for 16-bit storage is available in the Vulkan Samples repository on GitHub: [16-Bit Storage input/output tutorial](#)

10.2 Check precision

In extensions and elsewhere, make sure that attribute precisions are as expected.

Prerequisites

You must understand the concept of precision types such as `mediump` and `highp`.

Predefined attributes

You can get overflows if you use an attribute as though it was a different precision to what it is defined as. For example, in the extension `EXT_shader_framebuffer_fetch`, `gl_LastFragData` is `mediump`, as is `gl_LastFragColorARM` in the extension `GL_ARM_shader_framebuffer_fetch`. Both of these variables have caused overflows when used incorrectly by developers.

How to optimize

Check the definition of attributes, pre-defined or otherwise, before use.

Negative impact

Overflows from wrong precision errors can cause random crashes and other types of corruption.

10.3 Vectorized arithmetic code

The Arm Midgard GPU architecture implements Single Instruction Multiple Data (SIMD) maths units, exposing vector instructions to each thread of execution. The Arm Bifrost and Valhall GPU architectures switch to scalar warp-based arithmetic instructions, but still implement vector access to memory.

Prerequisites

You must understand the following concepts:

- Vector and scalar arithmetic instructions.
- Vector processing units.

How to optimize the use of vectorized arithmetic code on Arm GPUs

Try using the following optimization techniques:

- Write vector arithmetic code in your shaders. While doing so is less critical since the Arm Bifrost GPUs and upwards, there are still large numbers of devices that are using Arm Midgard GPU architecture.
- Write compute shaders so that work items contain enough work to fill the vector processing units.

Something to avoid when optimizing your use of vectorized arithmetic code

Do not write scalar code and hope that the compiler optimizes it. While the compiler can, it is more reliably vectorized if the input code starts out in vector form.

10.4 Vectorize memory access

The Arm GPU shader core load-store data cache has a wide data path capable of returning multiple values in a single clock cycle. It is important to make vector data accesses to get the highest access bandwidth from the data caches. Shader programs expose direct access to the underlying memory.

Prerequisites

You must understand the following concepts:

- The load-store data cache.
- Vector data accesses.
- Thread quads.

How to optimize the use of vectorized memory accesses on Arm GPUs

Try using the following optimizations:

- For memory accesses with a single thread, use vector data types.
- Bifrost GPUs can run four neighboring threads in lock-step, which is known as a quad. You must access overlapping or sequential memory ranges across the four threads to allow load merging. Mali-G52 and Mali-G76 GPUs can do eight threads in a warp, and Valhall GPUs can do 16.

Things to avoid when optimizing your use of vectorized memory accesses

Arm recommends that you:

- Do not use scalar loads if vector loads are possible.
- Do not access divergent addresses across a thread quad where possible.

The negative impact of using vectorized memory accesses incorrectly

Many types of common compute programs perform relatively light arithmetic on large data sets. Getting memory access correct for them can have a significant impact on performance.

10.5 Manual source code optimization

There is no guarantee that the shader compiler can safely perform code transforms. Rendering errors can result from a floating-point infinity or Not-a-Number (NaN) that would not have occurred in the original program order.

Prerequisites

You must understand the following concepts:

- Code transforms.
- Matrix mathematics.
- Floating-point limitations.

Reducing the number of computations needed

Where possible, refactor your source code to reduce the number of computations that are needed, rather than relying on the compiler to apply the optimizations.

How to optimize your source code

Try using the following optimizations:

- Based on your knowledge of values, refactor your source code to optimize the code as much as possible by hand.
- Graphics are not bit exact. Therefore, you must be willing to approximate when it helps refactoring. For example, simplify $(A * 0.5) + (B * 0.45)$ by using $(A + B) * 0.5$ instead, saving a multiply.
- Use the built-in function library where possible. Often, a hardware implementation that is faster or, lower power than the equivalent hand-written shader code backs up the function library.

Beware timestamp effects

Timestamps are useful to see how long different parts of code are taking, but:

- Timestamps have an overhead; be mindful of usage and check GPU performance impact if using a lot.
- Timestamps can impact work scheduling and introduce serialization. Check how the overall frame time is affected when adding them in, and do not leave them in released code.

Things to avoid when optimizing your use of your source code

- Do not reinvent the built-in function library in your custom shader code.
- Do not do clears or fills of buffers with shader code where `vkCmdFillBuffer` can be used.

The negative impact of not using minimized precision correctly

Less efficient shader programs cause reduced application performance.

How to debug source-code related issues

Use the Mali Offline Compiler to measure the impact of your shader code changes, including analysis of shortest and longest path through the programs.

Web link: [Mali Offline Compiler](#)

10.6 Generating SPIR-V

With Vulkan, whether using HLSL or GLSL, there are a number of SPIR-V generators - or pipelines to generate SPIR-V. Be aware that they will create different code, and be prepared to analyze it.

Prerequisites

- SPIR-V
- Microsoft DXC and FXC
- glslang
- Mali Offline Compiler

Analyzing SPIR-V

To check SPIR-V, there is Khronos' spirv-dis tool, to disassemble the binary into assembly language text. Analysis of compiled code can be done with the Mali Offline Compiler. However, if you wish to try other compilers, then they will either need to be run separately, or you can use other third-party tools, such as Shader Playground.

Something to avoid in compiled SPIR-V

Often the fastest SPIR-V involves the fewest instructions. However, be aware of the `NoContraction` decoration, which disables optimizations in an entire shader. Avoid this if at all possible.

As an example, currently the DXC HLSL compiler uses the `NoContraction` decoration if fused multiply-adds are used (`fma()` or `mad()`). Normal math instructions of `A*B+C` therefore end up unexpectedly faster, as they do not generate `NoContraction`.

The negative impact of not checking generated SPIR-V

If you do not look at the SPIR-V generated by your shader compiler, then you can end up with less efficient shaders than you expected.

Additional reading

Further information on SPIR-V can be found at the following website: [Official Khronos SPIR overview](#).

10.7 Instruction caches

The shader core instruction cache is a performance-impacting area that is often overlooked. Due to the number of threads running concurrently, it is a critically important part to be aware of.

Prerequisites

You must understand the following concepts:

- Instruction caches.
- Early ZS testing.

How to optimize the use of instruction caches

Try using the following optimizations:

- Use shorter shaders with many threads over long shaders with few threads. A shorter program is more likely to be hot in the cache.
- Use shaders that do not have control-flow divergence. Divergence can reduce temporal locality and increase cache pressure.

Things to avoid when optimizing your use of instruction caches

Arm recommends that you:

- Do not unroll loops too aggressively, although some unrolling can help.
- Do not generate duplicate shader programs or pipeline binaries from identical source code.
- Beware of fragment shading with many visible layers in a tile. The shaders for all layers that are not killed by early ZS or Forward Pixel Kill (FPK), must be loaded and executed, increasing cache pressure.

Web link: [Killing Pixels blog post](#).

How to debug instruction cache-related performance issues

Try the following debugging techniques:

- Use the Mali Offline Compiler to statically determine the sizes of the programs being generated for any given Arm GPU.

Web link: [Mali Offline Compiler](#)

- The Arm Mobile Studio tool suite can be used to step through draw calls and visualize how many transparent layers are building up in your render passes.

Web link: [Arm Mobile Studio](#)

10.8 Uniforms

Arm GPUs can promote data from API-set uniforms and uniform buffers into shader core registers. The data is then loaded on a per-draw basis, instead of on every shader thread. In turn, removing many load operations from the shader programs.

Prerequisites

You must understand the following concepts:

- Uniforms.
- Uniform Buffer Objects (UBOs).

Using uniforms

Not all uniforms can be promoted into registers. Uniforms that are dynamically accessed cannot always be promoted to register-mapped uniforms, unless the compiler is able to make the uniforms constant expressions. For example, expression array indices that are made constant by loop unrolling a fixed iteration for-loop.

How to optimize the use of uniforms on Arm GPUs

Try using the following optimizations:

- Keep your uniform data small. 128 bytes is a good general rule for how much data can be promoted to registers in any given shader.
- Avoid uniform vectors or matrices that are padded with constant elements that are used in computation. For example, elements that are always zero or one.
- Promote uniforms to compile-time constants with `#defines` for OpenGL ES, specialization constants for Vulkan, or literals in the shader source if they are static.
- Prefer uniforms set by `glUniform()` on OpenGL ES, rather than uniforms loaded from buffers.
- Prefer column major format for matrices that are stored in UBOs with OpenGL ES.
- Vulkan push constants are equal to UBOs for Arm GPUs.

Things to avoid when optimizing your use of uniforms

Arm recommends that you:

- Do not dynamically index into uniform arrays.
- Do not over use instancing. Instanced uniforms that are indexed using `gl_InstanceID` count as being dynamically indexed and cannot use register mapped uniforms.

The negative impact of using uniforms incorrectly

Register mapped uniforms cost little to use computationally. Any spilling to buffers in memory increases the load-store cache accesses to the per-thread uniform fetches.

How to debug uniform-related performance issues on Arm GPUs

The Mali Offline Compiler provides statistics about the number of uniform registers that are being used. And also about the number of load and store instructions that are being generated.

Web link: [Mali Offline Compiler](#)

10.9 Uniform subexpressions

One common source of inefficiency is the presence of uniform subexpressions in the shader source. Uniform subexpressions are pieces of code that only depend on the value of literals or other uniforms. Therefore, the results are always the same.

Prerequisites

You must understand the following concepts:

- Uniforms.

How to optimize the use of uniform subexpressions on Arm GPUs

Minimize the number of uniform-on-uniform or uniform-on-literal computations. Compute the result of the uniform subexpression on the CPU and then upload that as your uniform.

The negative impact of using uniform subexpressions incorrectly

The Mali GPU drivers can optimize the cost of most uniform subexpressions so that they are only computed a single time per draw. While it can seem like the benefit is not as large as it appears, the optimization pass still incurs a small cost for every draw call. The draw call cost can be avoided by removing the redundancy.

How to debug uniform sub-expression-related performance issues on Arm GPUs

Use the Mali Offline Compiler to measure the impact of your shader code changes, including analysis of shortest and longest path through the programs.

Web link: [Mali Offline Compiler](#)

10.10 Uniform control-flow

One common source of inefficiency is the presence of conditional control-flow, such as if blocks and for loops, which are parameterized by uniform expressions.

Prerequisites

You must understand the following concepts:

- `#defines` in OpenGL ES.
- Specialization constants in Vulkan.

How to optimize the use of a uniform control-flow

Use `#defines` at compile time in OpenGL ES, and specialization constants in Vulkan for all control flow. Doing so allows the compilation to completely remove unused code blocks and statically unroll loops.

Things to avoid when optimizing your use of a uniform control-flow

Do not use uniform values that parametrize control-flows. Instead, specialize shaders for each control path that is needed.

The negative impact of not using control-flow correctly

You can expect a reduced performance in your application due to less efficient shader programs.

How to debug uniform-control-flow-related performance issues

Use the Mali Offline Compiler to measure the impact of your shader code changes. Include an analysis of the shortest and longest path through the programs.

Web link: [Mali Offline Compiler](#)

10.11 Branches

Branching can be expensive on a GPU. Branching either restricts how the compiler can pack groups of instructions in a thread. Another example is when there is a divergence across multiple threads, which introduces cross-thread scheduling restrictions.

Prerequisites

You must understand the following concepts:

- Branching.
- Thread scheduling.

How to optimize the use of branches

Try using the following optimizations:

- Minimize the use of complex branches in shader programs.
- Minimize the amount of control-flow divergence in spatially adjacent shader threads.
- Use `min()`, `max()`, `clamp()`, and `mix()` functions to avoid small branches.
- Check the benefits of branching over computation. For example, skipping lights that are above a threshold distance from the camera. Often, it is faster just doing the computation.

Things to avoid when optimizing your use of branches

Do not implement multiple expensive data paths that are selected from using a `mix()`. Branching is usually the best solution for minimizing the overall cost in this particular scenario.

The negative impact of not using branches correctly

You can expect to experience a reduced performance in your application due to less efficient shader programs.

How to debug branch-related performance issues

Use the Mali Offline Compiler to measure the impact of your shader code changes. Include an analysis of shortest and longest path through the programs.

Web link: [Mali Offline Compiler](#)

10.12 Discards

Using a discard in a fragment shader, or when using alpha-to-coverage, are commonly used techniques. For example, alpha-testing complex shapes for foliage and trees.

Prerequisites

You must understand the following concepts:

- Fragment shaders.
- Alpha-to-coverage.
- Late ZS updates.
- Depth and stencil writing.

Using discards

These techniques force fragments to use late ZS updates. It is not known beforehand if the fragments survive the fragment operations stage of the pipeline until after shading. You must run the shader to determine the discard state of each sample. Doing so can cause redundant shading or pipeline starvation due to pixel dependencies.

How to optimize the use of discards and alpha-to-coverage

Try using the following optimizations:

- Minimize your use of shader discard and alpha-to-coverage.
- Computing lights in deferred lighting commonly uses fragment discard to cull fragments that are too far from the light source. To minimize execution bubbles, Arm recommends that you do not do depth & stencil writes for these lighting passes.
- Render alpha-tested geometry front-to-back with depth-testing enabled. Doing so causes as many fragments as possible to fail during early ZS testing. In turn, minimizing the number of late ZS updates that are needed.

The negative impact of not using discards and alpha-to-coverage correctly

Keep the following in mind:

- Extra fragment shading costs can cause a performance loss or bandwidth increase.
- Pipeline starvation waiting for pixel dependencies to resolve can also cause a loss in performance.

10.13 Atomics

Atomic operations are common to many compute algorithms and some fragment algorithms. With some slight modifications, atomic operations allow many algorithms to be implemented on highly parallel GPUs that would otherwise be serial.

Prerequisites

You must understand the following concepts:

- Atomics.
- Contention.
- L1 and L2 caches.

Atomics and contention

The key performance problem with atomics is contention. Atomic operations from different shader cores. Hitting the same cache line requires data coherency snooping through L2 cache, which is computationally expensive.

Optimized compute applications that use atomics must aim to spread out the contention by keeping the atomic operations local to a single shader core. Atomics are efficient when a shader core controls the necessary cache line in its L1.

How to optimize atomics

Try using the following optimization techniques:

- Consider how to avoid contention when using atomics in algorithm design.
- Consider spacing atomics 64 bytes apart to avoid multiple atomics contending on the same cache line.
- Consider whether it is possible to amortize the contention by accumulating into a shared memory atomic. Then, have one thread push the global atomic operation at the end of the workgroup.
- For OpenCL, consider the use of the `cl_arm_get_core_id` extension to allow explicit management of per-shader-core atomic variables.

Things to avoid when optimizing atomics

If better solutions that use multiple passes are available, then do not use atomics.

The negative impact of not using atomics correctly

Heavy contention on a single atomic cache entry significantly reduces overall throughput. It also impacts how well problems scale up when running on a GPU implementation with more shader cores.

How to debug atomics related performance issues

The GPU performance counters include counters for monitoring the frequency of L1 cache snoops from the other shader cores in the system.

11. Ray tracing

This chapter covers the best practices that we recommend you use when adding ray tracing and ray queries to your graphics application, particularly on platforms powered by Arm GPUs. Ray tracing requires support for the Vulkan extension `VK_KHR_acceleration_structure` and either `VK_KHR_ray_query` or `VK_KHR_ray_tracing_pipeline`. All Immortalis™ GPUs support ray tracing.

The ray tracing pipeline extension introduces a new set of shader types that work differently to the shaders in the traditional shader pipeline.

The ray query extension adds support for tracing rays in all shader types, and can be used in the traditional pipeline of vertex, compute, and fragment shaders. Although the two extensions are orthogonal, for the Arm GPU hardware that supports ray tracing, ray queries within the traditional pipeline is the preferred method as this currently has greater performance acceleration.

Reference

For more about Vulkan ray tracing extensions, see Khronos' introductory blog on [ray tracing in Vulkan](#) and the [Vulkan Guide Ray Tracing topic](#).

11.1 Acceleration structures

Acceleration Structures (AS) are important to make ray tracing efficient. To quickly find ray intersections, AS are used to represent geometry in a spatially sorted way.

Prerequisites

You must understand the following concepts:

- Ray tracing
- Bounding volumes
- Bottom-Level Acceleration Structures (BLAS)
- Top-Level Acceleration Structures (TLAS)

Use triangle geometry and minimize triangle count

We recommend using acceleration structures with triangles rather than Axis-Aligned Bounding Boxes (AABBs). In `VkAccelerationStructureBuildGeometryInfoKHR` objects, use `VkAccelerationStructureGeometryKHR` structures of `VK_GEOMETRY_TYPE_TRIANGLES_KHR` rather than `VK_GEOMETRY_TYPE_AABBs_KHR`. Ray intersections with both triangles and AABBs are hardware-accelerated on Immortalis™ GPUs, but intersections with AABBs require more interaction with user shaders, which is slower.

To reduce intersection tests and accelerate the traversals, try to use a minimal number of triangles for the structure.

Use the same acceleration structure in all threads

There is a significant performance improvement if all the threads in a warp can traverse their rays together because the GPU can calculate them more efficiently.

Warps can only traverse together if they are traversing the same acceleration structure. This is only a problem if you have multiple TLAS. If you have multiple TLAS, avoid dynamically selecting a TLAS in the shader.

Use build flags

The GPU can improve performance if it knows the expected usage of acceleration structures, which you can define with `VkBuildAccelerationStructureFlagBitsKHR`.

For example, when setting the `VkAccelerationStructureBuildGeometryInfoKHR` flags, use the following:

- `VK_BUILD_ACCELERATION_STRUCTURE_PREFER_FAST_TRACE_BIT_KHR` for static BLAS. This is usually in combination with `VK_BUILD_ACCELERATION_STRUCTURE_ALLOW_UPDATE_BIT_KHR` not being set.
- `VK_BUILD_ACCELERATION_STRUCTURE_PREFER_FAST_BUILD_BIT_KHR` for TLAS, and BLAS that require frequent updates, for example, animations. In this case also use `VK_BUILD_ACCELERATION_STRUCTURE_ALLOW_UPDATE_BIT_KHR` to re-fit the acceleration structure rather than building it from scratch every frame.

Use geometry flags

To achieve performance improvements, set the `VkGeometryFlagsKHR` flags.

When possible, set the `VK_GEOMETRY_OPAQUE_BIT_KHR` flag in `VkAccelerationStructureGeometryKHR` because it allows the GPU to immediately confirm ray intersections and save ray query calls or any-hit shader invocations.

Group acceleration structures when rebuilding

To achieve better parallelization:

- Group several AS builds into a single command
- Use deferred operations in host builds.

You must consider the amount of scratch memory required. Try to group AS so that several can be built together without requiring high memory usage. For example, you may have a single large scratch buffer and reserve chunks with an offset for each AS build. The buffer would only be reset, with an AS barrier, when it runs out. When tuning scratch buffer size, you must test the performance on target Arm GPU devices.

Group static geometry into a single BLAS

In general, you should group static geometry into a single BLAS to test against. However, if there are large holes in the static geometry, resulting in a largely empty bounding box, you can split the geometry into regions. In this instance, place the geometry into sensible groupings and profile their performance.

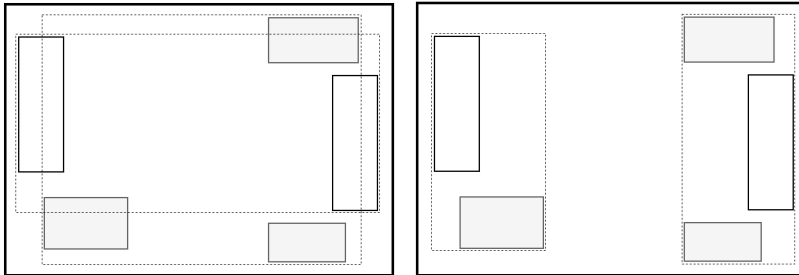
Separately, skinned geometry is also best grouped into a single BLAS with `VK_BUILD_ACCELERATION_STRUCTURE_PREFER_FAST_BUILD_BIT_KHR` and `VK_BUILD_ACCELERATION_STRUCTURE_ALLOW_UPDATE_BIT_KHR` set in `VkAccelerationStructureBuildGeometryInfoKHR::flags`.

Avoid overlapping geometries

Grouping geometries can result in faster ray traversals. Conversely, if a BLAS has disjoint geometries that have significant space between them, it often makes sense to split the BLAS. Doing so ensures that rays do not intersect with the empty region of a BLAS.

For example, the following images are two almost completely overlapping BLAS, that should be joined together, and then split into 2 regions, right and left.

Figure 11-1: Overlapping and Non-overlapping BLAS division



Negative consequences of non-optimal acceleration structures

AS that use the following take longer to test against and slow your rendering down:

- Non-triangle geometry
- More triangles than needed
- Inappropriate build flags

11.2 Efficient ray tracing

There are several practical tips to get good ray tracing performance on mobile devices. Here are some aims to consider.

Prerequisites

You must understand the following concepts:

- Ray tracing
- Acceleration Structures (AS)
- Ray query
- Ray tracing pipeline (RTP)
- Level of Detail (LOD)

Minimize count and size of dynamic geometries

Rebuilding and updating AS for animated objects has a significant cost. Try to minimize the number of dynamic models, and their primitive count.

Dynamic geometries often need their acceleration structure to be updated every frame, and can use `VK_BUILD_ACCELERATION_STRUCTURE_MODE_UPDATE_KHR` if they do so, as this is faster. However, they will still need an occasional full rebuild with `VK_BUILD_ACCELERATION_STRUCTURE_PREFER_FAST_BUILD_BIT_KHR` to ensure that the quality of the AS does not drop beyond an acceptable level.

Minimize rays per pixel

Every ray cast has a performance cost. Therefore, try to minimize the number of rays used and only cast rays that are essential.

To achieve this aim, consider the following techniques:

- Render at a lower resolution and then upscale appropriately.
- Checkerboard rendering allows you to only trace half the rays each frame.
- Consider using temporal super-sampling techniques to distribute rays across multiple frames.
- Hybrid rendering methods. For example, try combining Screen-Space Reflections (SSR) or even reflection probes with ray traced reflections, and shadow-maps with ray traced shadows.

Optimize shadow rays

Try to avoid casting unnecessary rays. For shadows, there are many cases where you can know the result without casting a ray. Such examples include:

- Avoid tracing a ray if the surface normal is facing away from the light and is therefore automatically in shadow.
- Detect if a given pixel is in range of a light, or on the far plane, so you can avoid casting the ray.
- Use a simple shadow map to detect shadow edges, where shadow rays are more valuable.
- Use `gl_RayFlags`. See [Ray query](#) or [Ray tracing pipeline](#) for further information.

Minimize divergence between rays

GPUs work best when work can be parallelized. As such, it is best to avoid shader divergence across neighboring threads to maximize shader throughput and exploit caching strategies. This is particularly important in ray tracing since spatially adjacent shader threads can result in different ray traversals and intersection logic. Therefore, try to maximize ray coherency. As an example, with warp-swizzling, a noise function could be written to diverge less per warp.

Aiming for better ray coherency means:

- Hard shadows are faster than soft shadows.
- Mirror reflections are faster than glossy/rough reflections.
- Stochastic sampling gives lower performance.

Performance-quality trade-offs may be required: for instance reduce roughness of reflections to reduce divergence and improve performance, consider ray-binning techniques, and combining ray tracing with rasterization-based techniques where possible.

However, 'divergence' through early discard of rays is good, especially if it results in whole warps being dismissed early.

Minimize ray length

When casting rays, it is best practice to keep the parameters t_{Min} as large as possible, and t_{Max} as small as possible. Make sure to also keep the parameters coherent across a warp while doing this. Keeping these values optimized can sometimes reduce traversal times and avoid precision issues.

Minimize ray payload

Additional ray payloads add to both bandwidth and memory use.

Use simple shaders

Ray tracing uses more GPU resources, so it is very important to simplify shaders to avoid register spilling.

You may also want to avoid combining ray tracing with other features that use a lot of bandwidth.

Avoid non-opaque traversals and geometry

Traversal hints may be provided with flags at the shader level (compile time) and/or at AS build time (runtime). The table below compares some options. In this case, shader flags take precedence over AS flags:

Shader flag (traceRayEXT/ rayQueryInitializeEXT)	Instance flag (VkAccelerationStructureInstanceKHR)	Geometry flag (VkAccelerationStructureGeometryKHR)
gl_RayFlagsOpaqueEXT	VK_GEOMETRY_INSTANCE_FORCE_OPAQUE_BIT_KHR	VK_GEOMETRY_OPAQUE_BIT_KHR
gl_RayFlagsNoOpaqueEXT	VK_GEOMETRY_INSTANCE_FORCE_NO_OPAQUE_BIT_KHR	0 (absence of opaque flag) or VK_GEOMETRY_NO_DUPLICATE_ANY_HIT_INVOCATION_BIT_KHR

Transparency requires more complex traversals and shader logic, which are both expensive. Only use transparent geometry where needed.

Avoid using `VK_GEOMETRY_NO_DUPLICATE_ANY_HIT_INVOCATION_BIT_KHR`, only enable it only if there is a strong need, for example, with partially transparent materials like glass or smoke particles. For alpha-test transparency, flags may be set to 0. This allows for fast traversals since the GPU does not have to track hits, but may result in alpha-test code executing more than once.

Consider materials and texturing

After a ray hit, generally the next step is to resolve the material. Different rays within a warp may hit different objects or materials.

Fragment shaders are run in groups of 2x2 fragments. This means, for instance, a full screen triangle is preferred over a full screen quad to minimise the number of helper invocations. Also,

fragment shaders will want an explicit texture LOD to stop rendering artifacts where all 4 rays in a 2x2 quad hit different materials and UV interpolation will fail.

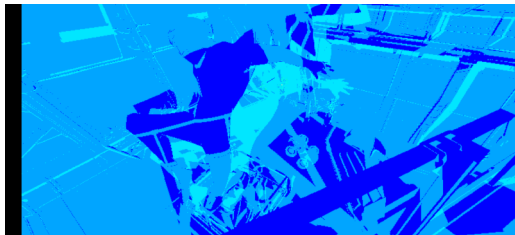
Negative consequences of too many rays cast

Ultimately, ray casting is expensive, so casting too many rays at once results in your frame taking too long to render. Make sure that every ray cast is required and provides useful information to the final pixel color.

Debugging techniques to use with ray tracing

Consider implementing a heatmap system to visualize the amount of hit-candidates for ray query, or the amount of hit-shader invocations for RTP.

Figure 11-2: Heatmap



Use [Arm Mobile Studio suite](#) for profiling.

Extensive counter instrumentation is available in Streamline, and shader best practice static analysis is available when using Mali Offline Compiler.

Figure 11-3: Streamline



11.3 Ray query

Ray query is the preferred Arm GPU ray tracing solution because it is more performant.

Prerequisites

You must understand the following concepts:

- Ray casting

Use a single call to `rayQueryProceed` for each `rayQueryInitialize`

For each `rayQueryEXT` object, ensure that there is exactly one call to `rayQueryProceedEXT` following a call to `rayQueryInitializeEXT`. For example:

```
rayQueryEXT ray_query;
for (int i = 0; i < iter; i++) {
    rayQueryInitializeEXT(ray_query, ...);
    rayQueryProceedEXT(ray_query);
}
```

Use `gl_RayFlags`

`gl_RayFlags` flags help minimize the work done during ray traversal. For use-cases such as shadows, where you only need to determine whether there is a hit rather than find the closest one, use `gl_RayFlagsTerminateOnFirstHitEXT` and `gl_RayFlagsOpaqueEXT`. The fastest way to get the closest hit is to use only `gl_RayFlagsCullNoOpaqueEXT` Or `gl_RayFlagsOpaqueEXT` as appropriate for the scene. `gl_RayFlagsCullNoOpaqueEXT` is fastest, but does not include any hits on non-opaque objects.

If you do not use Axis-Aligned Bounding Boxes (AABBs) in your Acceleration Structures, you can use `gl_RayFlagsSkipAABBEXT` to improve performance. AABBs are often not used, see [Acceleration structures](#) for more information. There is a significant performance uplift if both `gl_RayFlagsOpaqueEXT` and `gl_RayFlagsSkipAABBEXT` flags can be used. Using `gl_RayFlagsCullNoOpaqueEXT` instead of `gl_RayFlagsOpaqueEXT` provides an even larger performance improvement.

Using the flags `gl_RayFlagsOpaqueEXT` Or `gl_RayFlagsCullNoOpaqueEXT` with `gl_RayFlagsSkipAABBEXT` guarantees that `rayQueryProceedEXT` completes the traversal and return false. Therefore, using its output value is unnecessary.

```
rayQueryEXT rq;
rayQueryInitializeEXT(rq, accStruct, gl_RayFlagsTerminateOnFirstHitEXT |
    gl_RayFlagsCullNoOpaqueEXT | gl_RayFlagsSkipAABBEXT, cullMask, origin, tMin,
    direction, tMax);
```



Note

You should use Static values for ray flags if possible. Do not dynamically modify flags or pass them in from Uniform Buffer Objects (UBOs). This ensures all optimizations can be applied.

For hybrid renderers `gl_RayFlagsCullFrontFacingTrianglesEXT` can be useful to avoid self-intersection on rays spawned from the GBuffer. Culling may also be specified with AS flags, which take precedence over shader flags:

Shader flag (<code>traceRayEXT/rayQueryInitializeEXT</code>)	Instance flag (<code>VkAccelerationStructureInstanceKHR</code>)
<code>gl_RayFlagsCullBackFacingTrianglesEXT</code>	<code>VK_GEOMETRY_INSTANCE_TRIANGLE_FRONT_COUNTERCLOCKWISE_BIT_KHR</code>
<code>gl_RayFlagsCullFrontFacingTrianglesEXT</code>	<code>VK_GEOMETRY_INSTANCE_TRIANGLE_FACING_CULL_DISABLE_BIT_KHR</code>



Triangle face culling implies more work for GPU. You must only use triangle face culling if there is an algorithmic need. Alternatively, you can use a bias to ray origin to avoid self-interactions by setting t_{\min} to a small value.

Shading

Arm GPUs allow you to use ray queries from most shader types:

- Fragment shader-based queries allow the pipeline to take advantage of ray query workloads running without breaking the render pass, avoiding additional barriers. It also benefits from optimizations such as AFBC and AFRC.
- Compute shader-based pipelines can benefit from the flexibility of compute shaders and cross-thread communication, if required.
- Vertex shader-based queries lack the benefits of the above and are discouraged.
- Ray query can also be used within ray tracing pipeline shaders, but this is strongly discouraged.

Use Mali Offline Compiler to test if shaders are performant

The Arm Mobile Studio Mali Offline Compiler can analyze whether ray tracing code can be hardware accelerated with Immortalis™ GPUs or is on a slower path. This is an important and easy check.

Avoid `rayQueryProceed` conditionals

Avoid using either the `rayQueryInitializeEXT` or `rayQueryProceedEXT` call within a conditional or recursion. However, `rayQueryProceedEXT` will be needed as the condition in a `while` loop unless geometry is entirely opaque. Ensure `rayQueryProceedEXT` calls have the `rayQuery` object initialized with the flags as previously described, and if possible use single calls.

Minimizing the number of return trips between the traversal logic and the calling shader improves performance. Use `VK_GEOMETRY_OPAQUE_BIT_KHR` / `VK_GEOMETRY_INSTANCE_FORCE_OPAQUE_BIT_KHR` to minimize return trips and get better performance.

Ray query loop example: Alpha-tested shadows

```
rayQueryInitializeEXT(ray_query, top_level_as,
    gl_RayFlagsTerminateOnFirstHitEXT, ...);
// Use gl_RayFlagsTerminateOnFirstHitEXT to indicate that traversal can stop after
// first confirmed hit was found.
// Without this flag GPU will search for the closest hit, resulting into extra work,
// that is not needed for shadows.
while (rayQueryProceedEXT(ray_query))
{
    const uint instance_id =
        rayQueryGetIntersectionInstanceCustomIndexEXT(ray_query, false);
    const uint primitive_id = rayQueryGetIntersectionPrimitiveIndexEXT(ray_query,
        false);
    const vec2 barycentrics = rayQueryGetIntersectionBarycentricsEXT(ray_query,
        false);
    if (alpha_test(instance_id, primitive_id, barycentrics))
    {
        rayQueryConfirmIntersectionEXT(ray_query);
    }
}
```


Negative consequences of complex shaders and `rayQueryProceed` loops

The compiler can optimize for best performance when the shaders meet the above conditions. Not following the advice results in the compiler being forced into a slow path and a loss of performance.

Example

Example code for ray query is available in the Khronos Vulkan Samples on GitHub, see [Basic ray queries](#).

11.4 Ray tracing pipeline

The ray tracing pipeline defines a dedicated set of shaders including ray generation, closest-hit, any-hit and miss shaders.

Prerequisites

You must understand the following concepts:

- Ray casting

Use `gl_RayFlags`

These flags help minimize the work done during ray traversal. For use-cases such as shadows, where you only need to determine whether there is a hit rather than find the closest one, use `gl_RayFlagsTerminateOnFirstHitEXT` and `gl_RayFlagsOpaqueEXT`. The fastest way to get the closest hit is to use just `gl_RayFlagsOpaqueEXT`.

If Axis-Aligned Bounding Boxes (AABBs) are not used in your Acceleration Structures, you can use `gl_RayFlagsSkipAABBEXT` to improve performance. AABBs are often not used, as covered in [Acceleration structures](#). There is a significant performance uplift if both `gl_RayFlagsOpaqueEXT` and `gl_RayFlagsSkipAABBEXT` flags can be used.

```
traceRayEXT(accStruct, gl_RayFlagsTerminateOnFirstHitEXT | gl_RayFlagsOpaqueEXT |  
gl_RayFlagsSkipAABBEXT, cullMask, 0, 0, 0, origin, tMin, direction, tMax, 0);
```



Note

Static values should be used for ray flags if possible. Do not dynamically modify flags or pass them in from Uniform Buffer Objects (UBOs). This ensures all optimizations can be applied.

For hybrid renderers, using `gl_RayFlagsCullFrontFacingTrianglesEXT` can be useful to avoid self-intersection on rays spawned from the GBuffer. Culling may also be specified with AS flags, which take precedence over shader flags:

Shader flag (traceRayEXT/rayQueryInitializeEXT)	Instance flag (VkAccelerationStructureInstanceKHR)
gl_RayFlagsCullBackFacingTrianglesEXT	VK_GEOMETRY_INSTANCE_TRIANGLE_FRONT_COUNTERCLOCKWISE_BIT_KHR
gl_RayFlagsCullFrontFacingTrianglesEXT	VK_GEOMETRY_INSTANCE_TRIANGLE_FACING_CULL_DISABLE_BIT_KHR

However, be aware that triangle face culling implies more work for GPU, and should only be used if there is an algorithmic need. Alternatively, you can use a bias to ray origin to avoid self-intersections, by setting t_{\min} to a small value.

Avoid pipeline libraries

We do not recommend pipeline libraries for ray tracing pipeline. You must have a monolithically created pipeline. A monolithically created pipeline does not have pipeline libraries, so that the full state is known at creation and compiler optimizations can be applied. To achieve the best performance, you must minimise the number of shaders, state, and payloads.

Do not combine ray tracing and ray query

When using the ray tracing pipeline, do not use ray query from any shader stage because this results in a significant performance penalty. If multiple rays are required, they should all be done at the ray generation stage and, if needed, feed back information about other rays in their payload.

Minimize the use of any-hit and closest-hit shader calls

Additional shader calls can be computationally expensive and any-hit shaders can be repeatedly called. Setting a maximum number of, or eliminating, any-hit calls can give a large improvement in performance, but you must monitor the effect on quality. This trade-off requires consideration when non-opaque geometry is used. In such instances, hybrid techniques may be worth considering.

You must minimise the number of different closest-hit shaders, to reduce the chance of divergent execution.

Generally, minimise the entries in the Shader Binding Table (SBT).

Consider using a single bindless descriptor set

If materials suit it, a single bindless descriptor set is a good choice, rather than using unique descriptors for each SBT entry. A bindless descriptor set will work well with combining resources, allowing indexing into texture and buffer arrays.

Use the efficient ray tracing advice

You should read and follow the advice in the [Efficient ray tracing](#) section, especially minimizing ray payload.

Negative consequences of combining ray tracing pipeline and ray query

Using ray query within the ray tracing pipeline is not hardware accelerated in early Immortalis™ GPUs and is therefore less performant.

Example

Example code for ray tracing pipeline is available in the Khronos Vulkan Samples on GitHub:

- [Ray tracing basic](#)
- [Ray tracing extended](#)

12. System integration

This chapter covers how to optimize your Vulkan or OpenGL ES system integration.

12.1 Using EGL buffer preservation in OpenGL ES

Creating window surfaces that are configured with `EGL_BUFFER_PRESERVED` allows applications to logically update only part of the screen. With `EGL_BUFFER_PRESERVED`, applications always start rendering with the framebuffer state from the previous frame, rather than rendering from scratch.

Prerequisites

You must understand the following concepts:

Framebuffers:

- [How to use the Partial update extension](#)
- [How to use the Swap Buffers with Damage extension](#)

Preserving the EGL buffer

If your application only wants to update a subregion of the screen, then using the EGL buffer can appear like an efficiency boost. However, in real systems, the use of double, or triple buffering means that rendering starts with a full screen blit from the window surface for *frame n*, into the window surface for *frame n+1*.

How to optimize the EGL buffer

Try using the following optimization techniques:

- To minimize client rendering, use the *Partial Update* extension, instead of `EGL_BUFFER_PRESERVED`. The *Partial Update* extension allows incremental updates and true partial rendering of a frame.
- To minimize server-side composition overheads, use the *Swap Buffers with Damage* extension instead of `eglSwapBuffers()`.
- If the previous two extensions are not available, review the cost of using `EGL_BUFFER_PRESERVED` against the cost of just re-rendering a whole frame. When using simple UI content that uses compressed texture inputs and procedural fills, it can be more efficient to re-render the entire frame from the source material.
- If you know that an entire frame is being overdrawn when using `EGL_BUFFER_PRESERVED`, then insert a full-screen `glClear()` at the start of the render pass. Doing so removes the readback of the previous frame into the GPU tile memory.

Outcomes to avoid when using the EGL buffer

Arm recommends that you:

- Do not use `EGL_BUFFER_PRESERVED` surfaces without considering if it makes sense for the content involved. Always try using `EGL_BUFFER_DESTROYED` and then measure the benefits.

- Do not use either `EGL_KHR_partial_update` Or `EGL_KHR_swap_buffers_with_damage` for applications that always re-render the whole frame. There is an extra cost to the software in doing so.
- Do not use `EGL_KHR_partial_update` with [AFRC](#) on an EGL window surface. If update and AFRC regions are misaligned there will be visible corruption.

The negative impact of not using the EGL buffer correctly

Unnecessary readbacks can reduce performance and increase memory bandwidth.

12.2 Android blob cache size in OpenGL ES

The Arm GPU OpenGL ES drivers use the `EGL_ANDROID_blob_cache` extension to persistently cache the results of shader compilation and linkage.

Prerequisites

You must understand the following concepts:

- [BlobCache extension](#)

Reconsider the size of the Android blob cache

For many applications, shaders are only compiled and linked when the application is first used. Subsequent application runs use binaries from the cache, benefiting from a faster startup and level load times.

The Android blob cache defaults to a relatively small size, 64KB per application. 64KB per application is insufficient for many modern games and applications that can use hundreds of shaders.

To increase the number of applications benefiting from program caching, we recommend that system integrators significantly increase the maximum size of the blob cache for their platforms.

How to optimize the size of the Android blob cache

Try increasing the size of the Android blob cache for each application. For example, up to 512KB or even 1MB.

The negative impact of not setting the Android blob cache size correctly

As shader programs must be compiled and linked at runtime, any games and applications that exceed the size of the blob cache begin more slowly.

12.3 Optimizing the swapchain surface count for Vulkan

The application has control over the window surface swapchain when using Vulkan. In particular, the application can decide how many surfaces to use.

Prerequisites

You must understand the following concepts:

- Vulkan windows surface swapchain.

The swapchain and vsync

Most mobile platforms use the vsync signal of the display to prevent screen tearing on buffer swap. If the GPU is rendering more slowly than the vsync period, then a swapchain that contains only two buffers is prone to stalling the GPU.

How to optimize the swapchain surface count

Try using the following optimization techniques:

- If your application always runs faster than vsync, then use two surfaces in the swapchain. Doing so reduces your memory consumption.
- If your application sometimes runs more slowly than vsync, then use three surfaces in the swapchain. Doing so gives you the best performance for your application.

Outcomes to avoid when using the window surface swapchain

If your application runs more slowly than vsync, then do not use two surfaces in the swapchain.

The negative impact of not using the window surface swapchain correctly

Using double buffering and vsync locks rendering to an integer fraction of the vsync rate. If a frame renders slower than vsync, then the performance of the application is reduced. For example, on a device that uses a 60FPS panel refresh rate, an application that is otherwise capable of running at 50FPS, drops down to 30FPS.

How to debug swapchain related performance issues

System profilers can show when the GPU is going idle. If the GPU is going idle, and the frame period is a multiple of the vsync period, then it can indicate rendering blocking and that the GPU is waiting for the vsync signal to release a buffer.

Example

Example code for N-buffering is available in the Vulkan Samples repository on GitHub: [Swapchain images tutorial](#)

12.4 Optimizing the swap chain surface rotation for Vulkan

For Vulkan, the application has control over the window surface orientation. The application is responsible for handling the differences between the logical and the physical orientation of the window when a mobile device is rotated.

Prerequisites

You must understand the following concepts:

- How swap chain surface rotation works.

Swap chain surface rotation

It is more efficient for the presentation subsystem if the application renders into a window surface whose orientation matches the physical orientation of the display panel.

How to optimize swap chain surface rotation

Try using the following optimization techniques:

- To avoid presentation engine transformation passes, ensure that swap chain `preTransform` value matches the `currentTransform` value that is returned by `vkGetPhysicalDeviceSurfaceCapabilitiesKHR()`.
- If a swap chain image acquisition returns `VK_SUBOPTIMAL_KHR` OR `VK_ERROR_OUT_OF_DATE_KHR`, then recreate the swap chain. When doing so, consider any updated surface properties, including potential orientation updates reported using `currentTransform`.

Outcomes to avoid when using swap chain surface rotation

Do not assume that supported presentation engines transforms, other than `currentTransform`, are free. Many presentation engines can handle rotation or mirroring. However, it can come with extra processing cost.

The negative impact of not using swap chain surface rotation correctly

Non-native orientation can require extra transformation passes in the presentation engine. Therefore, some systems must use the GPU as part of the presentation engine to handle cases that the display controller cannot handle natively.

How to debug swap chain surface rotation performance issues

System profilers, such as the kernel-integrated version of Streamline, can track the use of the Arm GPU to specific processes. Monitoring such tracking allows for the attribution of any extra GPU workload to the compositor process. However, the previous step assumes that the GPU is being used by the compositor to apply the presentation transformation.

Example

Example code for surface rotation is available in the Vulkan Samples repository on GitHub: [Surface Rotation Tutorial](#)

12.5 Optimizing swapchain semaphores for Vulkan

In Vulkan, semaphores allow synchronization and safe access of swapchain data.

Prerequisites

You must understand the following concepts:

- Semaphores
- Swap chains
- Windows System Integration (WSI)
- Pipeline bubbles

Swapchain semaphores

The following is a typical example of what a Vulkan frame looks like:

1. Create a `vkSemaphore`, #1, for the start of frame acquire.
2. Create a separate `vkSemaphore`, #2, for the end of frame release.
3. Calling `vkAcquireNextImage()` gives you the swapchain index #N and then associates with semaphore #1.
4. Wait for all fences that are associated with swapchain index #N.
5. Build the command buffers.
6. Submit the command buffers rendering to the window surface to `vkQueue`. Tell the command buffers to wait for semaphore #1 before rendering can begin, and to then signal semaphore #2 when the command buffers have completed.
7. Call `vkQueuePresent()`, configuring it to wait for semaphore #2.

The critical part of the previous example occurs when setting up the wait for any command buffers on semaphore #1. We must also specify which pipeline stages must wait for the WSI semaphore.

Along with `pWaitSemaphores[i]`, there is also `pWaitDstStageMask[i]`. The `pWaitDstStageMask[i]` mask specifies which pipeline stages must wait for the WSI semaphore. You must wait for the final color buffer using `VK_PIPELINE_STAGE_COLOR_ATTACHMENT_OUTPUT_BIT`.

When rendering, you must transition the WSI image from either an `UNDEFINED`, or `PRESENT_SRC_KHR` layout to a different layout. The layout transition must wait for `COLOR_ATTACHMENT_OUTPUT_BIT`, creating a dependency chain to the semaphore.

How to optimize swapchain semaphores in WSI

When the application is waiting for a semaphore from WSI, use `pWaitDstStageMask = VK_PIPELINE_STAGE_COLOR_ATTACHMENT_OUTPUT_BIT`.

The negative impact of not using swapchain semaphores correctly

Large pipeline bubbles are created when the vertex, or compute, stalls.

12.6 Window buffer alignment

If rows are insufficiently aligned, then linear format framebuffers can suffer from dramatically reduced write performance. To ensure an efficient write performance, align all memory system allocated surfaces that are imported into Arm GPU drivers.

Prerequisites

You must understand the following concepts:

- Framebuffer formats.
- Row alignments.

How to optimize the window buffer alignment

Try using the following optimization techniques:

- Align rows in linear format framebuffers to the smallest possible alignment of either a multiple of 16 pixels, or 64 bytes. For example, for RGB565 framebuffers you can use a 32-byte alignment. You can use a 64-byte alignment for RGBA8 and RGBA fp16 framebuffers.
- When an alpha channel is not required, then use power-of-two formats as they have better memory alignment properties. Use a format such as RGBX8 with a dummy channel, instead of using RGB8.

Outcomes to avoid when using the window buffer alignment

Arm recommends that you:

- Do not use any row alignment other than either a multiple of 16 pixels, or 64 bytes.
- Do not use framebuffer formats that are not power-of-two. For example, do not use true 24-bit RGB8.

The negative impact of not using the window buffer alignment correctly

AXI Bursts must be aligned to the Burst size. Unaligned rows must be broken into multiple smaller AXI Bursts to meet this requirement. Doing so makes less efficient use of the memory bus and, often causes the GPU to stall. The stalling occurs because the GPU exhausts the pool of available transaction IDs.

12.7 Vulkan private data

The Vulkan extension `VK_EXT_private_data` adds the ability to associate application data with a Vulkan object. Try to only use private data slots reserved at device creation time.

Prerequisites

You must understand the following concepts:

- Vulkan extensions

Pre-reserve slots

If private data is required to be stored with a Vulkan object through the Vulkan private data extension, ensure that all the necessary slots are reserved at device creation time. Pre-reserved slots are *not* recycled and are accessed efficiently within the Vulkan object's memory area.

If you do not reserve enough slots, you can affect slot reservations of other layers. Conversely, you can be affected by other poorly implemented layers - even if you reserve the slots correctly.

Avoid unreserved slots

If more slots are needed than were reserved at device creation time, it is better to implement a Vulkan object data storage system on the app side. The app has a better overview of lifetime and thread access patterns and can produce better performance.

Also avoid using more pre-reserved slots than is necessary. In general, only one should be required. Each pre-reserved slot increases the size of all created Vulkan objects by eight bytes.

Negative consequences of using unreserved private data slots

Unreserved slots have extra overhead to track them, and reduces performance below what an efficient app-side implementation can achieve. Ultimately, they are a convenience feature for when performance is not the primary concern.

12.8 Vulkan extensions to avoid

There are many Vulkan extensions. However, as Vulkan is a cross-platform API, not every extension is suited to mobile and tile-based GPUs. Therefore, some extensions are non-optimal and must therefore be avoided.

Prerequisites

You must understand the following concepts:

- Vulkan extensions

Non-optimal extensions

Extensions to be avoided include:

- `VK_EXT_transform_feedback`: An extension which supports old OpenGL ES and DirectX applications running on top of Vulkan.

How to avoid using non-optimal extensions

To avoid using these extensions, try the following alternative:

- Use compute shaders rather than the transform feedback buffers of `VK_EXT_transform_feedback`.

The negative impact of using non-optimal extensions

The impact differs by extension, but performance is significantly affected.

Proprietary Notice

This document is protected by copyright and other related rights and the use or implementation of the information contained in this document may be protected by one or more patents or pending patent applications. No part of this document may be reproduced in any form by any means without the express prior written permission of Arm Limited ("Arm"). No license, express or implied, by estoppel or otherwise to any intellectual property rights is granted by this document unless specifically stated.

Your access to the information in this document is conditional upon your acceptance that you will not use or permit others to use the information for the purposes of determining whether the subject matter of this document infringes any third party patents.

The content of this document is informational only. Any solutions presented herein are subject to changing conditions, information, scope, and data. This document was produced using reasonable efforts based on information available as of the date of issue of this document. The scope of information in this document may exceed that which Arm is required to provide, and such additional information is merely intended to further assist the recipient and does not represent Arm's view of the scope of its obligations. You acknowledge and agree that you possess the necessary expertise in system security and functional safety and that you shall be solely responsible for compliance with all legal, regulatory, safety and security related requirements concerning your products, notwithstanding any information or support that may be provided by Arm herein. In addition, you are responsible for any applications which are used in conjunction with any Arm technology described in this document, and to minimize risks, adequate design and operating safeguards should be provided for by you.

This document may include technical inaccuracies or typographical errors. THIS DOCUMENT IS PROVIDED "AS IS". ARM PROVIDES NO REPRESENTATIONS AND NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, INCLUDING, WITHOUT LIMITATION, THE IMPLIED WARRANTIES OF MERCHANTABILITY, SATISFACTORY QUALITY, NON-INFRINGEMENT OR FITNESS FOR A PARTICULAR PURPOSE WITH RESPECT TO THE DOCUMENT. For the avoidance of doubt, Arm makes no representation with respect to, and has undertaken no analysis to identify or understand the scope and content of, any patents, copyrights, trade secrets, trademarks, or other rights.

TO THE EXTENT NOT PROHIBITED BY LAW, IN NO EVENT WILL ARM BE LIABLE FOR ANY DAMAGES, INCLUDING WITHOUT LIMITATION ANY DIRECT, INDIRECT, SPECIAL, INCIDENTAL, PUNITIVE, OR CONSEQUENTIAL DAMAGES, HOWEVER CAUSED AND REGARDLESS OF THE THEORY OF LIABILITY, ARISING OUT OF ANY USE OF THIS DOCUMENT, EVEN IF ARM HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

Reference by Arm to any third party's products or services within this document is not an express or implied approval or endorsement of the use thereof.

This document consists solely of commercial items. You shall be responsible for ensuring that any permitted use, duplication, or disclosure of this document complies fully with any relevant

export laws and regulations to assure that this document or any portion thereof is not exported, directly or indirectly, in violation of such export laws. Use of the word “partner” in reference to Arm’s customers is not intended to create or refer to any partnership relationship with any other company. Arm may make changes to this document at any time and without notice.

This document may be translated into other languages for convenience, and you agree that if there is any conflict between the English version of this document and any translation, the terms of the English version of this document shall prevail.

The validity, construction and performance of this notice shall be governed by English Law.

The Arm corporate logo and words marked with ® or ™ are registered trademarks or trademarks of Arm Limited (or its affiliates) in the US and/or elsewhere. Please follow Arm’s trademark usage guidelines at <https://www.arm.com/company/policies/trademarks>. All rights reserved. Other brands and names mentioned in this document may be the trademarks of their respective owners.

Arm Limited. Company 02557590 registered in England.

110 Fulbourn Road, Cambridge, England CB1 9NJ.

PRE-1121-V1.0

Product and document information

Read the information in these sections to understand the release status of the product and documentation, and the conventions used in Arm documents.

Product status

All products and services provided by Arm require deliverables to be prepared and made available at different levels of completeness. The information in this document indicates the appropriate level of completeness for the associated deliverables.

Product completeness status

The information in this document is Final, that is for a developed product.

Revision history

These sections can help you understand how the document has changed over time.

Document release information

The Document history table gives the issue number and the released date for each released issue of this document.

Document history

Issue	Date	Confidentiality	Change
0303-09	10 October 2024	Non-Confidential	First release of version 3.3
0302-08	30 January 2024	Non-Confidential	First release of version 3.2
0301-07	30 January 2023	Non-Confidential	First release of version 3.1
0300-06	1 August 2022	Non-Confidential	First release of version 3.0
0202-05	14 May 2021	Non-Confidential	First release of version 2.2
0201-04	26 May 2020	Non-Confidential	First release of version 2.1
0200-03	30 October 2019	Non-Confidential	First release of version 2.0
0101-02	30 June 2019	Non-Confidential	First release of version 1.1
0100-01	27 October 2017	Non-Confidential	Initial release

Change history

The Change history tables describe the technical changes between released issues of this document in reverse order. Issue numbers match the revision history in [Document release information](#) on page 125.

Table 2: Version 3.3

Change	Location
Added Staging buffers chapter	Staging buffers
Updated Ray tracing pipeline chapter	Ray tracing pipeline
Updated Ray query chapter	Ray query
Updated Arm Fixed Rate Compression	AFRC
Added Fragment Prepass advice	Optimizing the draw call render order
Added Pipeline cache advice	Pipeline creation in Vulkan
Added Scratch buffer advice	Acceleration structures
Updated Runtime compression	Runtime compression
Updated Dynamic Rendering	Multipass rendering
Updated Safe image layouts	Transaction elimination
Updated Variable Rate Shading	Variable rate shading

Table 3: Version 3.2

Change	Location
Added Runtime compression chapter	Runtime compression
Renamed AFBC textures chapter	AFBC textures
Changed All Vulkan example code links	Multiple Chapters
Changed Extensive Raytracing advice updates	All Raytracing chapters
Updated Manual source code optimization	Manual source code optimization
Updated Minimize precision	Minimize precision
Updated Optimization process	Multiple Chapters
Updated Optimize attachment grouping	Multiple Chapters
Updated Vulkan pipeline synchronization	Vulkan pipeline synchronization
Updated Workgroup sizes	Workgroup sizes

Table 4: Version 3.1

Change	Location
Utgard advice removed from this guide	–
Added additional tips	Efficient render passes with Vulkan
Info on using static values for ray flags	Ray query
Added section on relaxed precision FP16 outputs	Attribute precision
Added a section on HLSL	Minimize precision
Added <i>Arm Fixed Rate Compression</i> (AFRC) section	AFRC
Added section about dedicated allocations	Allocating memory in Vulkan
Added minimize ray length section	Efficient ray tracing

Change	Location
Added information on <code>VK_EXT_descriptor_indexing</code>	Optimizing descriptor sets and layouts for Vulkan
Added new ray tracing pipeline section	Ray tracing pipeline
Added new section on generating SPIR-V	Generating SPIR-V
Added Dynamic Rendering incompatibility section	Multipass rendering
Added new <i>Variable Rate Shading</i> (VRS) section	Variable rate shading

Table 5: Version 3.0

Change	
Changed guide name from Mali to Arm GPU Best Practices	-
Added Ray tracing chapter and sub-chapters	Ray tracing
Integrated Arm Mali GPU OpenGL ES Application Optimization Guide, added new sub-chapters from integrated guide	Optimization basics Basic application optimizations Basic vertex shader optimizations Basic fragment shader optimizations
Added section on Sparse Index Buffers	Index sparsity
Added section on Private Data	Vulkan private data
Added section on Optimizing Attachment Grouping	Optimize attachment grouping
Re-added Vectorize Memory Access section	Vectorize memory access
Updated developer web links	Multiple chapters
Updated Vulkan AFBC textures	AFBC textures for Vulkan
Updated Attribute layout	Attribute layout
Updated Attribute precision	Attribute precision
Updated Image processing	Image processing
Updated Arm GPU datasheet and performance counters	Arm GPU datasheet and performance counters
Updated Multipass rendering	Multipass rendering
Updated Texture sampling performance	Texture sampling performance
Updated Transaction elimination	Transaction elimination
Updated Uniforms	Uniforms

Table 6: Version 2.2

Change	Location
Updated Minimize precision	Minimize precision
AFBC chapter split into Vulkan and GLES chapters	Multisampling for OpenGL ES
Updated Multisampling for Vulkan	Multisampling for Vulkan
Updated Image processing	Image processing
Updated Texture and sampler descriptors	Texture and sampler descriptors
Updated Texture sampling performance	Texture sampling performance
Updated Vulkan pipeline synchronization	Vulkan pipeline synchronization
Updated Buffer update for OpenGL ES	Buffer update for OpenGL ES
Updated Draw call batching best practices	Draw call batching best practices
Added Vulkan extensions to avoid chapter	Vulkan extensions to avoid
Updated Multipass rendering	Multipass rendering

Change	Location
Updated Blending	Blending
Updated Optimizing the draw call render order	Optimizing the draw call render order

Table 7: Version 2.1

Change	Location
Added Queries to the Optimizing application logic chapter	Queries
Updated Pipeline creation in Vulkan	Pipeline creation in Vulkan
Updated Command pools for Vulkan	Command pools for Vulkan
Added a Git Hub example to Optimizing command buffers	Optimizing command buffers for Vulkan
Added a Git Hub example to Secondary command buffers	Secondary command buffers
Added a Git Hub example to Optimizing descriptor sets	Optimizing descriptor sets and layouts for Vulkan
Added a Git Hub example to Multipass rendering	Multipass rendering
Added AFBC Vulkan support details and GitHub example	AFBC textures
Added a Git Hub example to Optimizing the swapchain surface count for Vulkan	Optimizing the swapchain surface count for Vulkan

Table 8: Version 2.0

Change	Location
No changes	-

Table 9: Version 1.1

Change	Location
Added information about command pools	Command pools for Vulkan
Added information about secondary command buffers	Buffer update for OpenGL ES
Added information about buffer update	Buffer update for OpenGL ES

Table 10: Version 1.0

Change	Location
Initial Release	-

Conventions

The following subsections describe conventions used in Arm documents.

Glossary

The Arm Glossary is a list of terms used in Arm documentation, together with definitions for those terms. The Arm Glossary does not contain terms that are industry standard unless the Arm meaning differs from the generally accepted meaning.

See the Arm Glossary for more information: developer.arm.com/glossary.

Typographic conventions

Arm documentation uses typographical conventions to convey specific meaning.

Convention	Use
<i>italic</i>	Citations.
bold	Interface elements, such as menu names. Terms in descriptive lists, where appropriate.
monospace	Text that you can enter at the keyboard, such as commands, file and program names, and source code.
monospace <u>underline</u>	A permitted abbreviation for a command or option. You can enter the underlined text instead of the full command or option name.
<and>	Encloses replaceable terms for assembler syntax where they appear in code or code fragments. For example: <div>MRC p15, 0, <Rd>, <CRn>, <CRm>, <Opcode_2></div>
SMALL CAPITALS	Terms that have specific technical meanings as defined in the <i>Arm® Glossary</i> . For example, IMPLEMENTATION DEFINED , IMPLEMENTATION SPECIFIC , UNKNOWN , and UNPREDICTABLE .



We recommend the following. If you do not follow these recommendations your system might not work.



Your system requires the following. If you do not follow these requirements your system will not work.



You are at risk of causing permanent damage to your system or your equipment, or harming yourself.



This information is important and needs your attention.



A useful tip that might make it easier, better or faster to perform a task.



A reminder of something important that relates to the information you are reading.

Useful resources

This document contains information that is specific to this product. See the following resources for other useful information.

Access to Arm documents depends on their confidentiality:

- Non-Confidential documents are available at developer.arm.com/documentation. Each document link in the following tables goes to the online version of the document.
- Confidential documents are available to licensees only through the product package.

Arm product resources	Document ID	Confidentiality
<i>Arm GPU Datasheet</i>	–	Non-Confidential
<i>Arm GPU Performance Counters</i>	–	Non-Confidential