

ARM® Guide for Unity Developers

Version 3.3

Optimizing Mobile Gaming Graphics

ARM®

ARM® Guide for Unity Developers

Optimizing Mobile Gaming Graphics

Copyright © 2014–2017 ARM Limited or its affiliates. All rights reserved.

Release Information

Document History

Issue	Date	Confidentiality	Change
0100-00	05 September 2014	Non-Confidential	First release of version 1.0
0200-00	23 June 2015	Non-Confidential	First release of version 2.0
0201-00	28 July 2015	Non-Confidential	First release of version 2.1
0300-00	18 September 2015	Non-Confidential	First release of version 3.0
0300-01	05 November 2015	Non-Confidential	Second release of version 3.0
0301-00	07 April 2016	Non-Confidential	First release of version 3.1
0301-01	25 April 2016	Non-Confidential	Second release of version 3.1
0302-00	31 May 2016	Non-Confidential	First release of version 3.2
0303-00	17 March 2017	Non-Confidential	First release of version 3.3
0303-01	01 June 2017	Non-Confidential	Second release of version 3.3

Non-Confidential Proprietary Notice

This document is protected by copyright and other related rights and the practice 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. **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 implementations infringe any third party patents.

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, third party patents, copyrights, trade secrets, or other rights.

This document may include technical inaccuracies or typographical errors.

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.

This document consists solely of commercial items. You shall be responsible for ensuring that any 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.

If any of the provisions contained in these terms conflict with any of the provisions of any signed written agreement covering this document with ARM, then the signed written agreement prevails over and supersedes the conflicting provisions of these terms. 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 the Agreement shall prevail.

Words and logos marked with ® or ™ are registered trademarks or trademarks of ARM Limited or its affiliates in the EU and/or elsewhere. All rights reserved. Other brands and names mentioned in this document may be the trademarks of their respective owners. Please follow ARM's trademark usage guidelines at <http://www.arm.com/about/trademark-usage-guidelines.php>

Copyright © 2014–2017, ARM Limited or its affiliates. All rights reserved.

ARM Limited. Company 02557590 registered in England.

110 Fulbourn Road, Cambridge, England CB1 9NJ.

LES-PRE-20349

Confidentiality Status

This document is Non-Confidential. The right to use, copy and disclose this document may be subject to license restrictions in accordance with the terms of the agreement entered into by ARM and the party that ARM delivered this document to.

Unrestricted Access is an ARM internal classification.

Product Status

The information in this document is Final, that is for a developed product.

Web Address

<http://www.arm.com>

Contents

ARM® Guide for Unity Developers Optimizing Mobile Gaming Graphics

Preface

<i>About this book</i>	7
<i>Feedback</i>	9

Chapter 1

Introduction

1.1 <i>About Unity</i>	1-11
1.2 <i>About ARM® Mali™ GPUs</i>	1-12
1.3 <i>About optimization</i>	1-13
1.4 <i>About the Ice Cave demo</i>	1-14

Chapter 2

Optimizing applications

2.1 <i>The optimization process</i>	2-16
2.2 <i>Unity quality settings</i>	2-17

Chapter 3

Profiling your application

3.1 <i>About profiling</i>	3-21
3.2 <i>Profiling with the Unity profiler</i>	3-22
3.3 <i>The Unity Frame Debugger</i>	3-24

Chapter 4

Optimization lists

4.1 <i>Application processor optimizations</i>	4-27
4.2 <i>GPU optimizations</i>	4-33

4.3	<i>Asset optimizations</i>	4-53
4.4	<i>Optimizing with the Mali™ Offline Shader Compiler</i>	4-55
Chapter 5	<i>Global illumination in Unity with Enlighten</i>	
5.1	<i>About Enlighten</i>	5-61
5.2	<i>The structure of Enlighten</i>	5-62
5.3	<i>Setting up a scene with Enlighten</i>	5-71
5.4	<i>Example: Enlighten setup of the Ice Cave</i>	5-73
5.5	<i>Using Enlighten in custom shaders</i>	5-80
Chapter 6	<i>Advanced graphics techniques</i>	
6.1	<i>Custom shaders</i>	6-85
6.2	<i>Implementing reflections with a local cubemap</i>	6-98
6.3	<i>Combining reflections</i>	6-114
6.4	<i>Dynamic soft shadows based on local cubemaps</i>	6-120
6.5	<i>Refraction based on local cubemaps</i>	6-128
6.6	<i>Specular effects in the Ice Cave demo</i>	6-134
6.7	<i>Using Early-z</i>	6-137
6.8	<i>Dirty lens effect</i>	6-138
6.9	<i>Light shafts</i>	6-141
6.10	<i>Fog effects</i>	6-145
6.11	<i>Bloom</i>	6-152
6.12	<i>Icy wall effect</i>	6-159
6.13	<i>Procedural skybox</i>	6-165
6.14	<i>Fireflies</i>	6-173
6.15	<i>The tangent space to world space normal conversion tool</i>	6-177
Chapter 7	<i>Virtual Reality</i>	
7.1	<i>Virtual reality hardware support for Unity</i>	7-185
7.2	<i>The Unity VR porting process</i>	7-186
7.3	<i>Things to consider when porting to VR</i>	7-189
7.4	<i>Reflections in VR</i>	7-191
7.5	<i>The result</i>	7-196
Chapter 8	<i>Vulkan</i>	
8.1	<i>About Vulkan</i>	8-198
8.2	<i>About Vulkan in Unity</i>	8-200
8.3	<i>Enabling Vulkan in Unity</i>	8-201
8.4	<i>Vulkan case study</i>	8-202
Chapter 9	<i>The Mali Graphics Debugger</i>	
9.1	<i>About Mali Graphics Debugger</i>	9-207
9.2	<i>MGD features</i>	9-208
9.3	<i>Enabling MGD</i>	9-210
9.4	<i>Enabling MGD with Vulkan</i>	9-214
Appendix A	<i>Revisions</i>	
A.1	<i>Revisions</i>	Appx-A-218

Preface

This preface introduces the *ARM® Guide for Unity Developers Optimizing Mobile Gaming Graphics*.

It contains the following:

- [About this book](#) on page 7.
- [Feedback](#) on page 9.

About this book

This book is designed to help you create applications and content that make the best use of Unity on mobile platforms, especially those with Mali™ GPUs.

Product revision status

The `rmpn` identifier indicates the revision status of the product described in this book, for example, `r1p2`, where:

- `rm` Identifies the major revision of the product, for example, `r1`.
- `pn` Identifies the minor revision or modification status of the product, for example, `p2`.

Intended audience

This book is for beginner to intermediate developers.

Using this book

This book is organized into the following chapters:

Chapter 1 Introduction

This chapter introduces the ARM Guide to Unity: Enhancing Your Mobile Games

Chapter 2 Optimizing applications

This chapter describes how to optimize applications in Unity.

Chapter 3 Profiling your application

This chapter describes profiling your application.

Chapter 4 Optimization lists

This chapter lists a number of optimizations for your Unity application.

Chapter 5 Global illumination in Unity with Enlighten

This chapter describes global illumination in Unity with Enlighten.

Chapter 6 Advanced graphics techniques

This chapter lists a number of advanced graphics techniques that you can use.

Chapter 7 Virtual Reality

This chapter describes the process of adapting an application or game to run on virtual reality hardware, and some differences in the implementation of reflections in virtual reality.

Chapter 8 Vulkan

This chapter describes Vulkan and how you enable it.

Chapter 9 The Mali Graphics Debugger

This chapter describes the *Mali Graphics Debugger* (MGD).

Appendix A Revisions

This appendix describes the changes between released issues of this book.

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.

Typographic conventions

italic

Introduces special terminology, denotes cross-references, and citations.

bold

Highlights interface elements, such as menu names. Denotes signal names. Also used for terms in descriptive lists, where appropriate.

monospace

Denotes text that you can enter at the keyboard, such as commands, file and program names, and source code.

monospace

Denotes a permitted abbreviation for a command or option. You can enter the underlined text instead of the full command or option name.

monospace italic

Denotes arguments to monospace text where the argument is to be replaced by a specific value.

monospace bold

Denotes language keywords when used outside example code.

<and>

Encloses replaceable terms for assembler syntax where they appear in code or code fragments.

For example:

```
MRC p15, 0, <Rd>, <CRn>, <CRm>, <Opcode_2>
```

SMALL CAPITALS

Used in body text for a few terms that have specific technical meanings, that are defined in the *ARM® Glossary*. For example, IMPLEMENTATION DEFINED, IMPLEMENTATION SPECIFIC, UNKNOWN, and UNPREDICTABLE.

Additional reading

Information published by ARM and by third parties.

See <http://infocenter.arm.com> for access to ARM documentation.

ARM publications

This book contains information that is specific to this product. See the following documents for other relevant information:

None.

Developer resources:

[https://developer.arm.com\(graphics/](https://developer.arm.com(graphics/).

Other publications

Relevant documents published by third parties:

OpenGL ES 1.1 Specification at <http://www.khronos.org>.

OpenGL ES 2.0 Specification at <http://www.khronos.org>.

OpenGL ES 3.0 Specification at <http://www.khronos.org>.

OpenGL ES 3.1 Specification at <http://www.khronos.org>.

Unity Scripting Reference at [Unity](#).

GPU Gems: Programming Techniques, Tips, and Tricks for Real-time Graphics by Randima Fernando (Series Editor).

GPU Pro: Advanced Rendering Techniques by Wolfgang Engel (Editor).

<https://developer.oculus.com/osig>.

Feedback

Feedback on this product

If you have any comments or suggestions about this product, contact your supplier and give:

- The product name.
- The product revision or version.
- An explanation with as much information as you can provide. Include symptoms and diagnostic procedures if appropriate.

Feedback on content

If you have comments on content then send an e-mail to errata@arm.com. Give:

- The title *ARM Guide for Unity Developers Optimizing Mobile Gaming Graphics*.
- The number ARM 100140_0303_01_en.
- If applicable, the page number(s) to which your comments refer.
- A concise explanation of your comments.

ARM also welcomes general suggestions for additions and improvements.

————— Note ————

ARM tests the PDF only in Adobe Acrobat and Acrobat Reader, and cannot guarantee the quality of the represented document when used with any other PDF reader.

Chapter 1

Introduction

This chapter introduces the ARM Guide to Unity: Enhancing Your Mobile Games

It contains the following sections:

- [*1.1 About Unity* on page 1-11.](#)
- [*1.2 About ARM® Mali™ GPUs* on page 1-12.](#)
- [*1.3 About optimization* on page 1-13.](#)
- [*1.4 About the Ice Cave demo* on page 1-14.](#)

1.1 About Unity

Unity is a software platform that enables you to create and distribute 2D games, 3D games, and other applications.

This book is intended to help you create applications and content that make the best use of Unity on mobile platforms, especially those with Mali™ GPUs. It describes techniques and best practices that you can use to improve the performance of your applications.

Note

Unless otherwise noted, the techniques described here also work on other platforms.

1.2 About ARM® Mali™ GPUs

ARM Mali GPUs are designed for mobile or embedded devices. ARM Mali GPUs are divided into the following families:

Bifrost GPUs

The Bifrost GPUs have unified shader cores that perform vertex, fragment, geometry, tessellation, and compute processing. They are used for graphics and compute applications with Vulkan, OpenGL ES 1.1 to OpenGL ES 3.2, and OpenCL 1.2 Full Profile.

Midgard GPUs

The Bifrost GPUs have unified shader cores that perform vertex, fragment, geometry, tessellation, and compute processing. They are used for graphics and compute applications with Vulkan, OpenGL ES 1.1 to OpenGL ES 3.2, and OpenCL 1.2 Full Profile.

Utgard GPUs

The Utgard GPUs have a vertex processor and one or more fragment processors. They are used for graphics only applications with OpenGL ES 1.1 and 2.0.

1.3 About optimization

Graphics is about making things look good. Optimization is about making things look good with the least computational effort. This is especially important for mobile devices that keep power consumption low by limiting computing power and memory bandwidth.

1.4 About the Ice Cave demo

The Ice Cave demo is a demonstration application created by ARM that uses a number of optimized techniques to produce high-quality visual content for mobile devices.

This document describes the graphical techniques used in the Ice Cave demo, and solutions to problems that were encountered during the development of the project.

The Ice Cave demo was developed for and tested on a mobile device with an ARM Cortex®-A57 MP4 processor and an ARM Mali-T760 MP8 GPU.

Chapter 2

Optimizing applications

This chapter describes how to optimize applications in Unity.

It contains the following sections:

- [2.1 The optimization process](#) on page 2-16.
- [2.2 Unity quality settings](#) on page 2-17.

2.1 The optimization process

Optimization is the process of taking an application and making it more efficient. For graphical applications, this typically means modifying the application to make it faster.

For example, a game with a low frame rate might appear jumpy. This gives a bad impression and can make a game difficult to play. You can use optimization to improve the frame rate of a game making it a better, smoother experience.

To optimize your code, use the optimization process. This is an iterative process that guides you through finding and removing performance problems.

The optimization process consists of the following steps:

1. Take measurements of your application with a profiler.
2. Analyze the data to locate the bottleneck.
3. Determine the relevant optimization to apply.
4. Verify that the optimization works.
5. If the performance is not acceptable return to step 1 and repeat the process.

The following is an example of the optimization process:

1. If you have a game that does not have the performance you require, you can use a profiler to take measurements.
2. Use the profiler to analyze the measurements so you can isolate and identify the source of the performance problem.
3. The problem with the game is, for example, rendering too many vertices.
4. Reduce the number of vertices in your meshes.
5. Execute the game again to ensure the optimization worked.

If you do this and the game still does not perform as expected, restart the process by profiling the application again to find out what else is causing problems.

Expect to repeat this process a number of times. Optimization is an iterative process where you might find performance problems in a number of different areas.

2.2 Unity quality settings

It is useful to know the Unity quality settings to ensure you select the correct settings for your application.

Unity has a number of options that alter the image quality of your game. Some of these options have a high computational cost and can have a negative impact on the performance of your game.

The following figure shows quality settings in the **Inspector**:

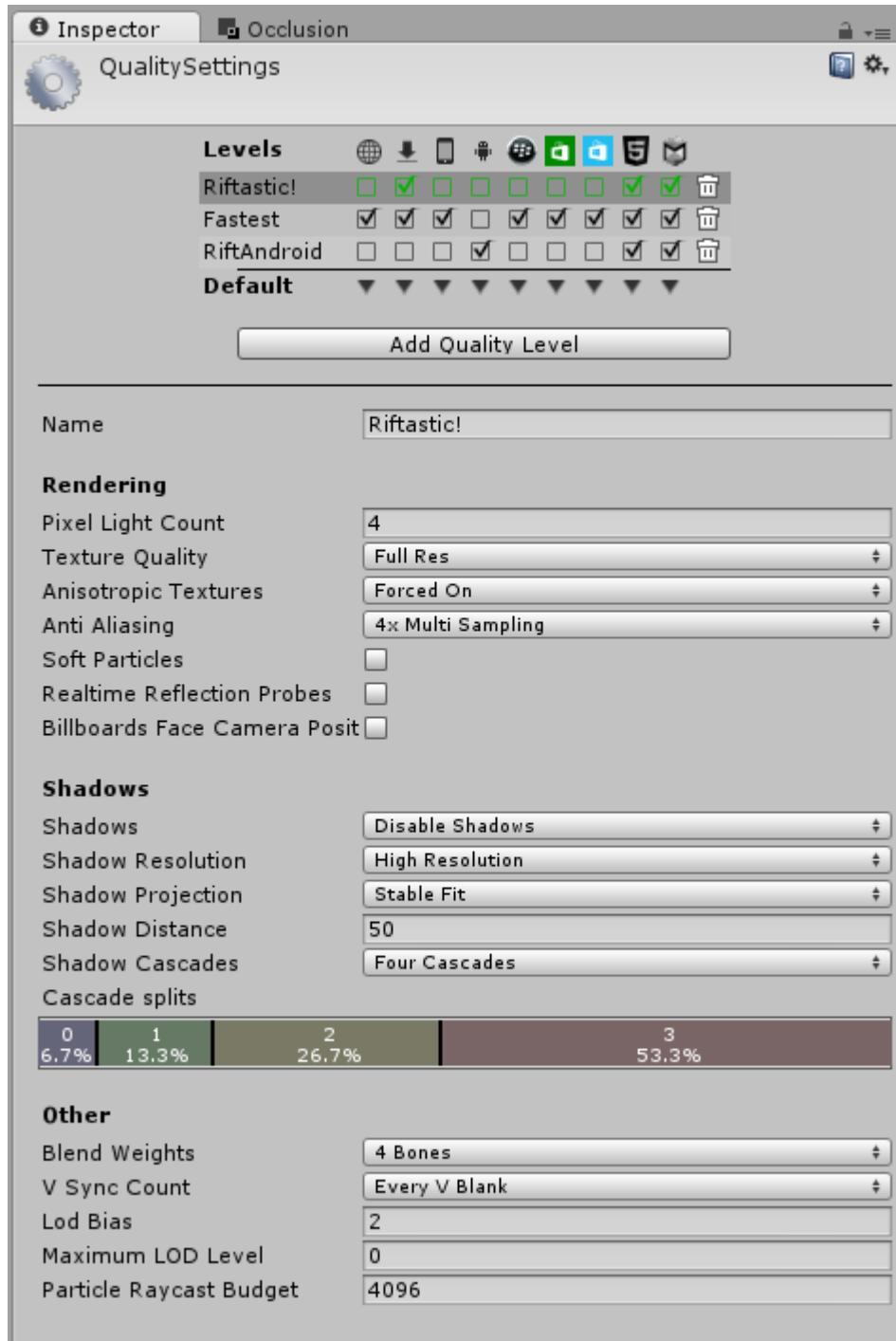


Figure 2-1 Quality settings

There are a number of options that can increase the image quality of your game with only a small trade-off in performance. For example, if the frame rate of your game is low, the GPU might be processing too much information when performing a complex graphical effect. You can perform less complex versions of graphical effects, such as shadows and lighting, for a relatively small impact on the graphical quality. Simpler effects can reduce the load on the GPU significantly, providing a higher frame rate.

The default settings for lighting can sometimes be too complex for a mobile device, so some games written for mobile platforms avoid complex techniques or use game-specific techniques. This might involve techniques such as pre-baking lighting into light maps or projecting textures instead of casting shadows.

In **Project Settings > Quality** there are a number of options that can have a large impact on the performance of your game:

Pixel Light Count

Pixel Light Count is the number of lights that can affect a given pixel. A high pixel light count requires a large number of calculations. Most games can use very few dynamic and real-time lights with minimal impact on image quality. Consider using techniques such as light maps and projected textures in your game, if lighting is causing performance problems.

Texture Quality

Texture Quality can load the GPU but it typically does not cause performance problems. Reducing texture quality can negatively impact the visual quality of your game, so only reduce the quality if you must. In the Ice Cave demo, **Texture Quality** is set to full resolution. If textures are causing performance problems, try using mipmapping. Mipmapping reduces compute and bandwidth requirements without impacting image quality.

AntiAliasing

AntiAliasing is an edge-smoothing technique that blends the pixels around triangle edges. This provides a noticeable improvement to the visual quality of your game. There are several methods of anti-aliasing, but in this case the toggle is for *Multi-Sampled Anti-Aliasing* (MSAA). 4x MSAA is very low cost operation on Mali GPUs, so always use it if possible.

Soft Particles

Soft Particles requires rendering to a depth texture or rendering in deferred mode. This increases the load on the GPU, but can be worth it in terms of achieving realistic visuals on your particles. On mobile platforms, rendering to and reading from a depth texture uses up valuable bandwidth, and rendering using a deferred path means you have no access to MSAA. Consider whether soft particles are important enough to your game to use them.

Anisotropic Textures

Anisotropic Textures is a technique that removes distortion from textures drawn at high gradients. This improves the image quality but it is an expensive technique. Avoid using this technique unless the distortion is especially noticeable.

Shadows

Shadows can be computationally intensive if they are high quality. If shadows cause performance problems, try simple shadows or switch them off. If shadows are important in your game, consider using simple dynamic shadowing techniques such as projected textures.

Realtime Reflection Probes

The **Realtime Reflection Probes** option can have a significant negative impact on the runtime performance.

When a reflection probe is rendered, every face of the cubemap is rendered separately by a camera at the origin of the probe. If inter-reflections are considered, this process takes place for every reflection bounce level. In the case of glossy reflections, the cubemap mipmaps are also used to apply a blurring process.

The following factors influence the rendering of the cubemap:

Cubemap Resolution

Higher resolution cubemaps increase rendering time. Use the lowest resolution cubemap possible for the quality you require.

Culling Mask

Use the culling mask when rendering the cubemap to avoid rendering any geometry not relevant in the reflections.

Cubemap Updating

the **Refresh Mode** option defines the update frequency for a cubemap:

- The **Every Frame** option renders the cubemap every frame. This is the most computationally expensive option, so avoid using it unless you require it.
- The **On Awake** option renders the cubemap at runtime one time, when the scene starts.
- The **Via Scripting** option provides you with control over when the cubemap is updated. With this option, you can limit the use of runtime resources by specifying the conditions when an update takes place.

Chapter 3

Profiling your application

This chapter describes profiling your application.

It contains the following sections:

- [*3.1 About profiling* on page 3-21](#).
- [*3.2 Profiling with the Unity profiler* on page 3-22](#).
- [*3.3 The Unity Frame Debugger* on page 3-24](#).

3.1 About profiling

You profile your application to find performance bottlenecks. When you have identified these, optimizing in these areas improves your application performance.

You can profile your Unity application with the following tools:

- Unity Profiler.
- Unity Frame Debugger.
- ARM Mali Graphics Debugger.
- ARM DS-5 Streamline.

For more information about the Mali Graphics Debugger, see [Chapter 9 The Mali Graphics Debugger on page 9-206](#).

3.2 Profiling with the Unity profiler

The Unity profiler provides detailed per-frame performance data in a series of charts to help you find the bottlenecks in your game.

If you click in a chart, you see a vertical slice and this selects a single frame. You can read information from this frame in the display panel at the bottom of the screen. If you click on a different chart without modifying the frame selection, the panel shows data from the profiler you have selected.

The following figure shows the Unity profiler:

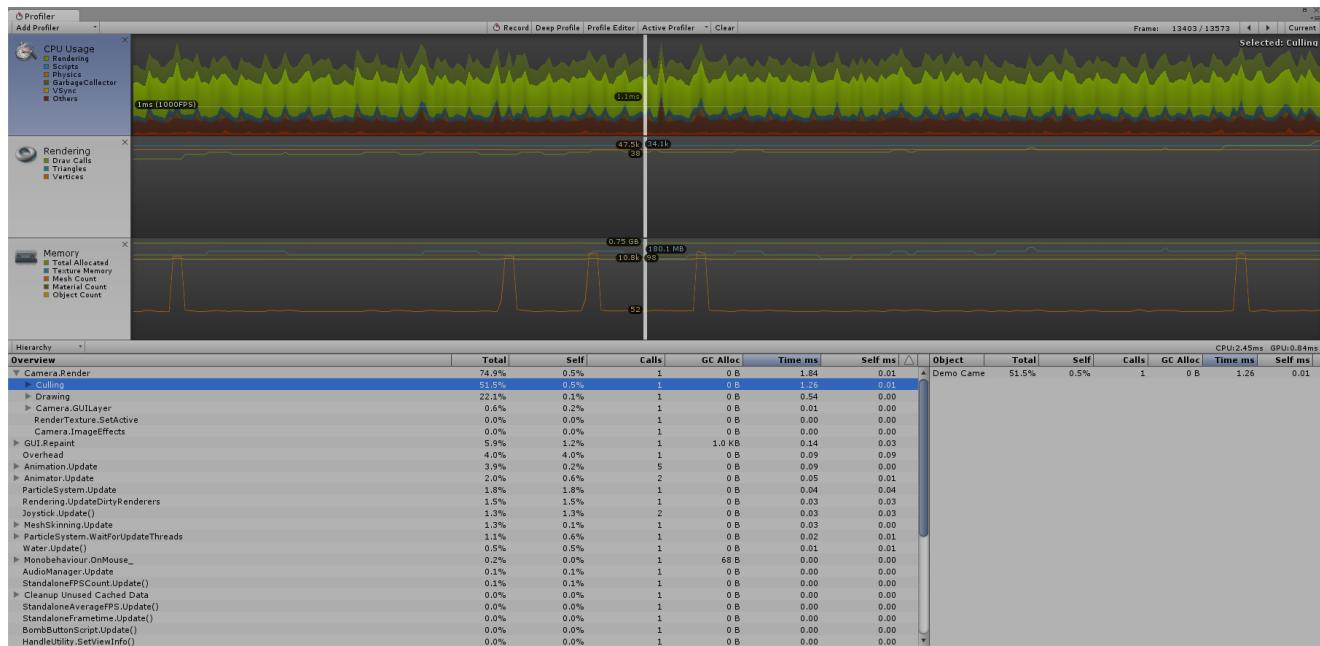


Figure 3-1 Unity profiler

The Unity profiler provides the following functions:

CPU Usage profiler

The CPU Usage profiler chart shows a breakdown of CPU utilization, highlighting different components, such as rendering, scripts, or physics. If you select a frame on the chart, the panel displays the functions that most contributed to that particular frame, in time taken, number of calls or memory allocations.

Concentrate on the functions that spend more time executing or allocate too much memory.

Note

In a multi-processor system, the values are averaged.

Rendering profiler

The Rendering profiler chart shows the number of draw calls, triangles, and vertices rendered in your scene. Selecting a frame on the chart displays more information about batching, textures, and memory consumption.

Look at the number of draw calls, triangles, and vertices that your scene renders. These are the most important numbers in mobile platforms.

Memory profiler

The Memory profiler chart displays the amount of memory allocated and the amount of resources used by the game, such as meshes or materials. Selecting a frame on the chart displays the memory consumption of your assets, the graphics and audio subsystems, or of profiler data itself.

Memory is limited on mobile platforms, so you must monitor the requirements of your game during its lifetime and check the number of resources in use. Some techniques, if applied incorrectly, can cause a large number of new objects to be created. For example, a texture atlas applied incorrectly can lead to the creation of a large number of new material objects.

The Add Profiler function

Add Profiler is an option on the drop-down menu in the top left corner of the profiler window. It enables you to add more charts to the profiler window, for example, CPU usage, rendering, or memory.

The Profiler.BeginSample() and Profiler.EndSample() methods

The Unity profiler enables you to use the `Profiler.BeginSample()` and `Profiler.EndSample()` methods. You can mark a region in your script, attach a custom label to it, and this region appears in the profiler hierarchy as a separate entry. Doing this enables you to obtain information about a specific piece of code without the compute and memory overhead of the `Deep Profile` option.

```
void Update()
{
    Profiler.BeginSample("ProfiledSection");
    [...]
    Profiler.EndSample();
}
```

The following figure shows a profiled section:

Overview	Total	Self	Calls	GC Alloc	Time ms	Self ms
WaitForTargetFPS	95.7%	95.7%	1	0 B	15.06	15.06
▼ ProfiledSectionTest.Update()	2.1%	0.0%	1	0 B	0.33	0.00
ProfiledSection	2.1%	2.1%	1	0 B	0.33	0.33

Figure 3-2 Profiled section

3.3 The Unity Frame Debugger

The Frame Debugger is a profiling tool that enables you to trace draw calls on a frame-by-frame basis.

The Frame Debugger is available from the **Window** menu.

The left pane shows the tree of draw calls issued in the frame.

In the right pane it shows additional information related to the selected draw call, such as geometry details and the shader that draws it.

The following figure shows a Phoenix object in the Frame Debugger:

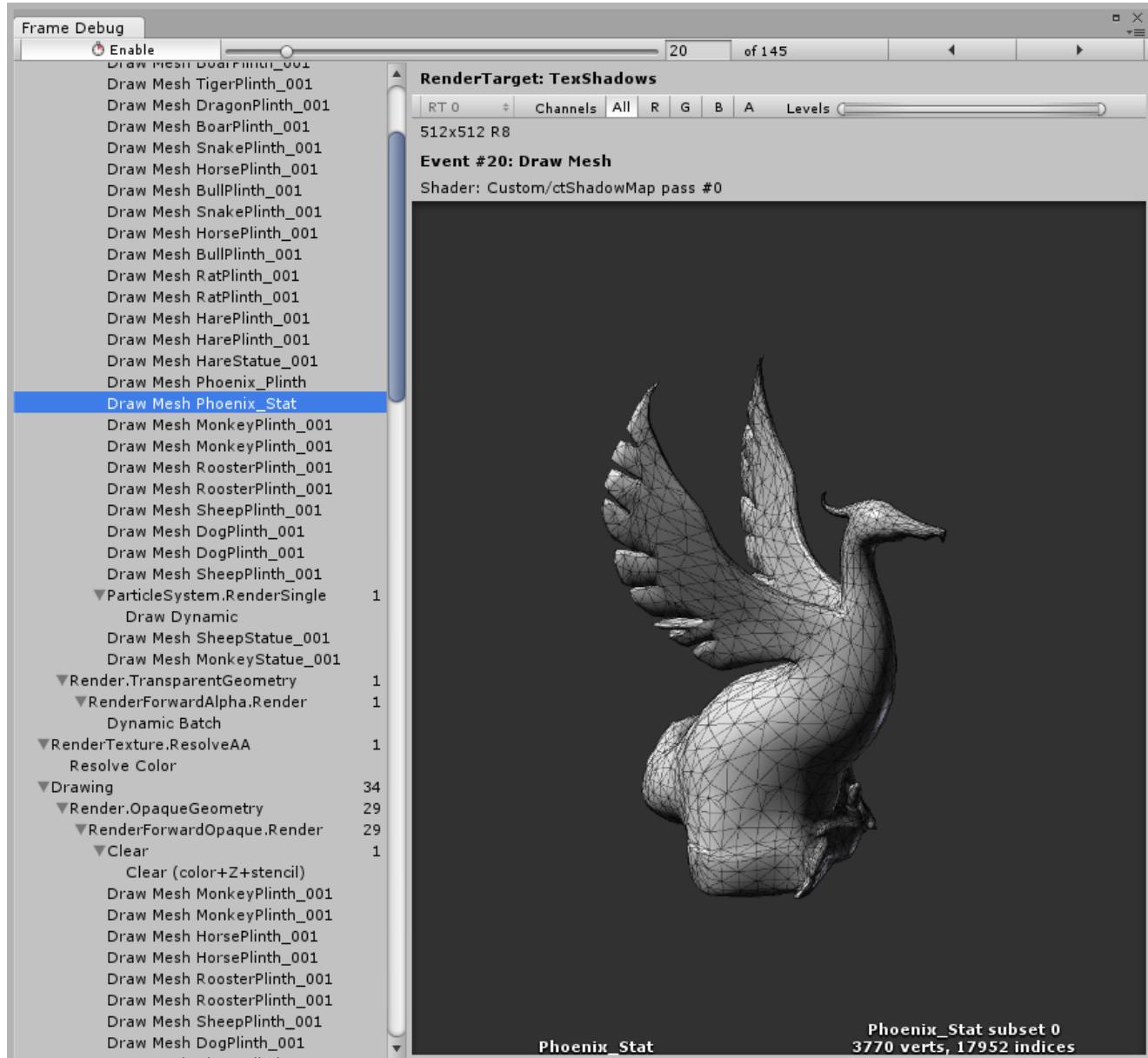


Figure 3-3 Frame Debugger

If the camera is rendering to a target at the selected draw call, you can visualize the rendered texture in the **Game View**.

The following figure shows a Phoenix object in the **Game View**:

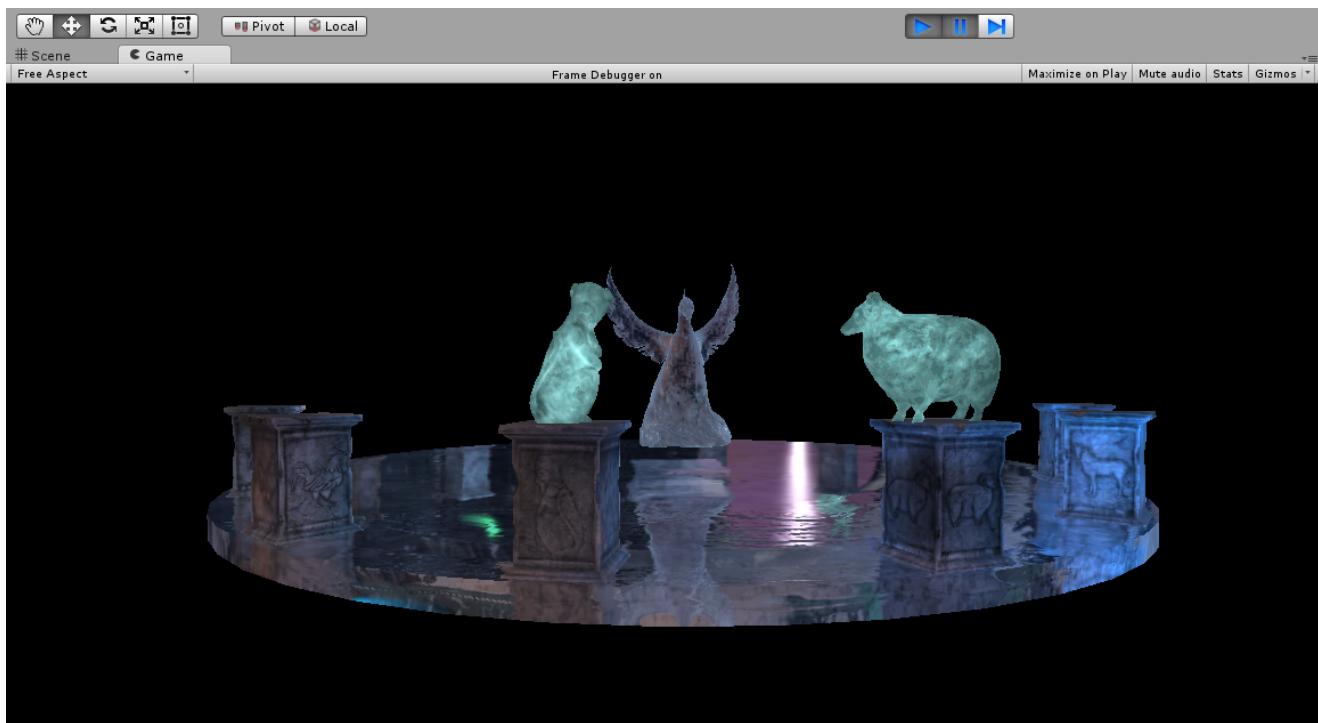


Figure 3-4 Frame Debugger Game View

For a description of how to increase the performance of a game by using the Frame Debugger to help optimize the rendering order of objects, see [4.2.9 Specify the rendering order](#) on page 4-50.

Chapter 4

Optimization lists

This chapter lists a number of optimizations for your Unity application.

It contains the following sections:

- [*4.1 Application processor optimizations* on page 4-27](#).
- [*4.2 GPU optimizations* on page 4-33](#).
- [*4.3 Asset optimizations* on page 4-53](#).
- [*4.4 Optimizing with the Mali™ Offline Shader Compiler* on page 4-55](#).

4.1 Application processor optimizations

The following list describes application processor optimizations:

Use coroutines instead of Invoke()

The `Monobehaviour.Invoke()` method is a fast and convenient way to call a method in a class with a time delay, but it has the following limitations:

- It uses reflection in C# to find the method to call, this can be slower than calling the method directly.
- There are no compile-time checks on the method signature.
- You cannot supply additional parameters.

The following code shows the `Invoke()` function:

```
public void Function()
{
    [...]
} Invoke("Function", 3.0f);
```

An alternative method is to use coroutines. A coroutine is a function of type `IEnumerator` that can return control to Unity with a special `yield return` statement. You can call the function again later and it resumes where it left off earlier.

You can call coroutines through the `MonoBehaviour.StartCoroutine()` method:

```
public IEnumerator Function(float delay)
{
    yield return new WaitForSeconds(delay);
    [...]
} StartCoroutine(Function(3.0f));
```

Changing from the `Monobehaviour.Invoke()` method to using coroutines provides more flexibility over the parameters passed to the functions dealing with animation states.

Use coroutines for relaxed updates

If your game requires an action every specific time interval, try launching a coroutine in the `MonoBehaviour.Start()` callback, instead of performing an action every frame in the `MonoBehaviour.Update()` callback. For example:

```
void Update()
{
    // Perform an action every frame
}

IEnumerator Start()
{
    while(true)
    {
        // Do something every quarter of second
        yield return new WaitForSeconds(0.25f);
    }
}
```

Note

An alternative use of this technique is to spawn enemies at irregular intervals. Use an infinite loop inside the coroutine that spawns an enemy and generates a random number. Pass the random number to the `WaitForSeconds()` function.

Avoid hard-coded strings for tags

Avoid hard-coded values for tags because they restrict the scalability and robustness of your game. For example, with tag names, if you refer to the names directly by strings, you cannot easily modify them and you are potentially exposed to spelling errors. For example:

```
if(gameObject.CompareTag("Player"))
{
    [...]
}
```

You can improve this by implementing a special class for tags that exposes public constant strings. For example:

```
public class Tags
{
    public const string Player = "Player";
    [...]
}

if(gameObject.CompareTag(Tags.Player))
{
    [...]
}
```

Note

You can use a Tags class with public constant strings to assist adding new tags in a consistent and scalable manner.

Reduce the number of physics calculations by changing the fixed time step

You can reduce the computational load of physics calculations by changing the fixed time step. Typically, most physics calculations take place at a fixed time step and you can increase or decrease the length of this step.

Increasing the time step decreases the load on the application processor but reduces the accuracy of physics calculations.

You can access the time manager from the main menu: **Edit > Project Settings > Time**.

The following figure shows the time manager:

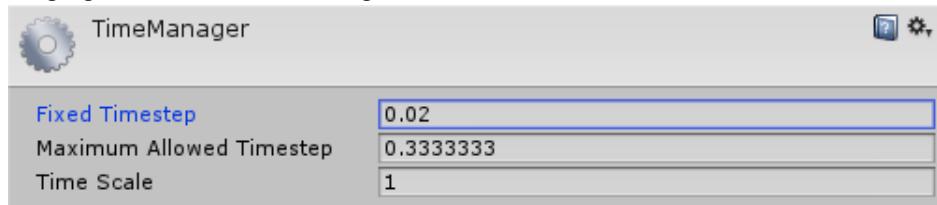


Figure 4-1 Fixed timestep settings

Remove empty callbacks

If your code includes empty definitions for functions such as `Awake()`, `Start()`, or `Update()`, remove them. There is an overhead associated with these because the engine still attempts to access them even though they are empty. For example:

```
// Remove the following empty definition

void Awake()
{
}
```

Avoid using `GameObject.Find()` in every frame.

`GameObject.Find()` is a function that iterates through every object in the scene. This can cause a significant increase in the main thread size if it is used in the incorrect part of your code. For example:

```
void Update()
{
    GameObject playerGO = GameObject.Find("Player");
    playerGO.transform.Translate(Vector3.forward * Time.deltaTime);
}
```

A better technique is to call `GameObject.Find()` on startup and cache the result, for example in the `Start()` or `Awake()` function:

```
private GameObject _playerGO = null ;

void Start()
{
    _playerGO = GameObject.Find("Player");
}

void Update()
{
    _playerGO.transform.Translate(Vector3.forward * Time.deltaTime);
}
```

Another alternative is to use `GameObject.FindWithTag()`:

```
void Update()
{
    GameObject playerGO = GameObject.FindWithTag("Player");
    playerGO.transform.Translate(Vector3.forward * Time.deltaTime);
}
```

Note

Use a dedicated class called `LocatorManager` that performs all the object retrievals immediately when the scene finishes loading. Other classes can use this as a service so objects are not retrieved multiple times.

Use the `StringBuilder` class to concatenate strings

When concatenating complex strings, use the `System.Text.StringBuilder` class. This is faster than the `string.Format()` method and uses less memory than concatenation with the plus operator:

```
// Concatenation with the plus operator
string str = "foo" + "bar";

// String.Format() method
string str = string.Format("{1}{2}", "foo", "bar");
```

The `System.Text.StringBuilder` class:

```
// StringBuilder class
using System.Text;

StringBuilder strBld = new StringBuilder();
strBld.Append("foo");
strBld.Append("bar");
string str = strBld.ToString();
```

The following figure shows string concatenations:

Overview	Total	Self	Calls	GC Alloc	Time ms	Self ms
▼ StringConcatenationTest.Start()	99.8%	0.0%	1	1.12 GB	2488.86	0.33
▼ StringConcatenationTest.StringFormatMethod()	51.0%	0.2%	1	0.56 GB	1273.01	6.78
► String.Format()	50.8%	0.3%	10000	0.56 GB	1266.22	7.51
▼ StringConcatenationTest.PlusConcatenation()	48.4%	0.2%	1	0.56 GB	1208.21	6.53
► String.Concat()	48.2%	0.5%	10000	0.56 GB	1201.68	13.15
▼ StringConcatenationTest.StringBuilderClass()	0.2%	0.2%	1	256.2 KB	7.31	6.97
► StringBuilder.Append()	0.0%	0.0%	10000	256.1 KB	0.32	0.12
► StringBuilder..ctor()	0.0%	0.0%	1	0 B	0.00	0.00
► StringBuilder.ToString()	0.0%	0.0%	1	0 B	0.01	0.00

Figure 4-2 String concatenations

Use the `CompareTag()` method instead of the `tag` property

Use the `GameObject.CompareTag()` method instead of the `GameObject.tag` property. The `CompareTag()` method is faster and does not allocate extra memory:

```
GameObject mainCamera = GameObject.Find("Main Camera");

// GameObject.tag property
if(mainCamera.tag == "MainCamera")
{
    // Perform an action
}

// GameObject.CompareTag() method
if(mainCamera.CompareTag("MainCamera"))
{
    // Perform an action
}
```

The following figure shows the use of `CompareTag()`:

Overview	Total	Self	Calls	GC Alloc	Time ms	Self ms
▼ TagComparisonTest.Start()	97.2%	0.1%	1	3.6 MB	127.48	0.26
▼ TagComparisonTest.TagProperty()	68.0%	59.6%	1	3.6 MB	89.27	78.14
▼ GameObject.get_tag()	7.9%	1.0%	100000	3.6 MB	10.44	1.32
GC.Collect	6.9%	6.9%	4	0 B	9.12	9.12
▼ String.op_Equality()	0.5%	0.3%	100000	0 B	0.69	0.50
String.Equals()	0.1%	0.1%	100000	0 B	0.18	0.18
► TagComparisonTest.CompareTagMethod()	28.9%	28.4%	1	0 B	37.94	37.27
GameObject.Find()	0.0%	0.0%	1	0 B	0.00	0.00

Figure 4-3 Compare tag

Use object pools

If your game has many objects of the same kind that are created and destroyed at runtime, you can use the design pattern *object pool*. This design pattern avoids the performance penalty of allocating and freeing many objects dynamically.

If you know the total number of objects that you require, you can create them all immediately and disable the objects that are not immediately required. When a new object is required, search the pool for the first unused one and enable it.

When an object is not required anymore, you can return it to the pool. This means resetting the object to a default starting state and disabling it.

This technique can be used with objects such as enemies, projectiles, and particles. If you do not know the exact number of objects that you require, do a test to find how many are used and create a pool slightly bigger than the number you find.

————— Note —————

Use object pools for enemies and bombs. This restricts the allocation of those objects to the loading phase of the game.

Cache component retrievals

Cache the component instance returned by `GameObject.GetComponent<Type>()`. The function call involved is quite expensive.

Properties such as `GameObject.camera`, `GameObject.renderer` or `GameObject.transform` are shortcuts to the corresponding `GameObject.GetComponent<Camera>()`,
`GameObject.GetComponent<Renderer>()`, and `GameObject.GetComponent<Transform>()`:

```
private Transform _transform = null;  
  
void Start()  
{  
    _transform = GameObject.GetComponent<Transform>();  
}  
  
void Update()  
{  
    _transform.Translate(Vector3.forward * Time.deltaTime);  
}
```

Consider caching the return value of `Transform.position`. Even if it is a C# getter property, there is overhead associated with an iteration over the transform hierarchy to calculate the global position.

————— Note —————

In Unity 5 and higher, the transform component is automatically cached.

Use `OnBecameVisible()` and `OnBecameInvisible()` callbacks

Callbacks such as `MonoBehaviour.OnBecameVisible()` and

`MonoBehaviour.OnBecameInvisible()` notify your scripts if their associated game objects become visible or invisible on screen.

These calls enable you to, for example, disable computational heavy code routines or effects when a game object is not rendered on screen.

Use `sqrMagnitude` for comparing vector magnitudes

If your application requires the comparison of vector magnitudes, use `Vector3.sqrMagnitude` instead of `Vector3.Distance()` or `Vector3.magnitude`.

`Vector3.sqrMagnitude` sums the squared components without calculating the root, but this is useful for comparisons. The other calls use a computationally expensive square root.

The following code shows the three different techniques used in comparisons of two positions in space:

```
// Vector3.sqrMagnitude property
if ((_transform.position - targetPos).sqrMagnitude < maxDistance * maxDistance)
{
    // Perform an action
}

// Vector3.Distance() method
if (Vector3.Distance(transform.position, targetPos) < maxDistance)
{
    // Perform an action
}

// Vector3.magnitude property
if ((_transform.position - targetPos).magnitude < maxDistance)
{
    // Perform an action
}
```

Use built-in arrays

If you know the size of an array in advance, use the built-in arrays.

`ArrayList` and `List` classes have more flexibility because they grow in size the more elements you insert, but they are slower than the built-in arrays.

Use planes as collision targets

If your scene only requires particle collisions with planar objects such as floors or walls, change the particle system collision mode to `Planes`. Changing the setting to use planes reduces the computations required. In this mode you can provide Unity with a list of empty GameObjects to act as the collider planes.

The following figure shows collision settings:

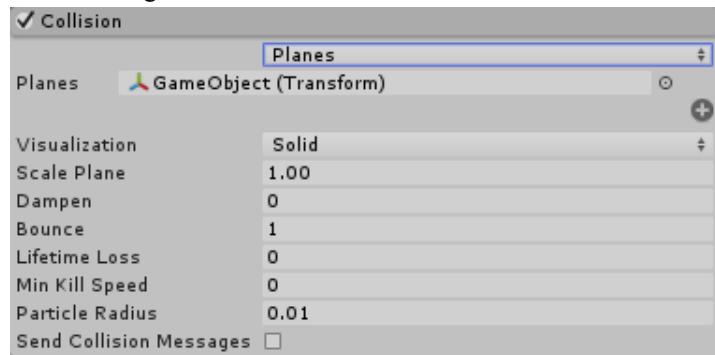


Figure 4-4 Collision settings

Use compound primitive colliders rather than mesh colliders

Mesh colliders are based on the real geometry of an object. These are very accurate for collision detection but are computationally expensive.

You can combine shapes such as boxes, capsules, or spheres into a compound collider that mimics the shape of the original mesh. This enables you to achieve similar results with much lower computational overhead.

4.2 GPU optimizations

This section lists GPU optimizations.

This section contains the following subsections:

- [4.2.1 Miscellaneous GPU optimizations on page 4-33](#).
- [4.2.2 Lightmaps and light probes on page 4-34](#).
- [4.2.3 ASTC texture compression on page 4-43](#).
- [4.2.4 Mipmapping on page 4-47](#).
- [4.2.5 Skyboxes on page 4-48](#).
- [4.2.6 Shadows on page 4-48](#).
- [4.2.7 Occlusion culling on page 4-49](#).
- [4.2.8 Use OnBecameVisible\(\) and OnBecomeInvisible\(\) callbacks on page 4-50](#).
- [4.2.9 Specify the rendering order on page 4-50](#).
- [4.2.10 Use depth pre-pass on page 4-52](#).

4.2.1 Miscellaneous GPU optimizations

The following list describes miscellaneous GPU optimizations:

Use static batching

Static batching is a common optimization technique that reduces the number of draw calls and this, in turn, reduces application processor utilization.

Dynamic Batching is performed transparently by Unity, but it cannot be applied to objects made up of a large number of vertices, because the computational overhead becomes too large.

Static Batching can work on objects made up of a large number of vertices, but the batched objects must not move, rotate, or scale during rendering.

To enable Unity to group objects for static batching, mark them as static in the Inspector.

The following figure shows static batching settings:

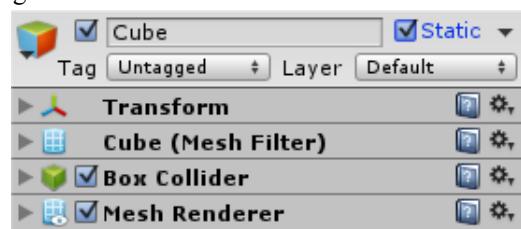


Figure 4-5 Static batching settings

Use 4x MSAA

ARM Mali GPUs can do 4x multi-sample anti-aliasing with very low computational overhead.

You can enable 4x MSAA in the Unity Quality Settings.

The following figure shows the MSAA settings:

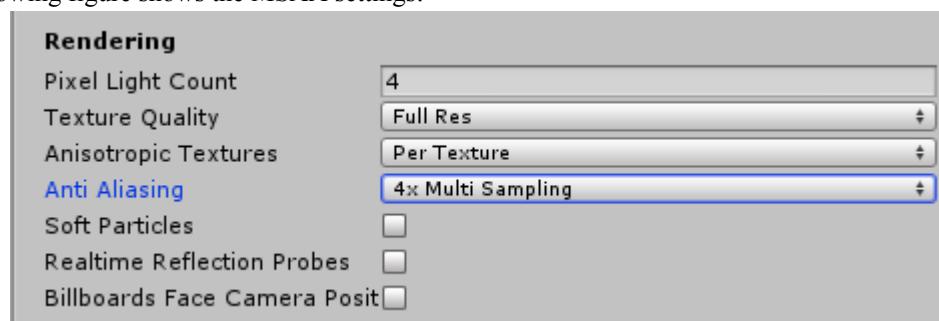


Figure 4-6 MSAA settings

Use level of detail

Level of detail (LOD) is a technique where the Unity engine renders different meshes for the same object depending on the distance from the camera.

Geometry is more detailed when the object is close to the camera. The detail level is reduced as the object moves away from the camera and at the furthest distance you can use a planar billboard.

You must set up LOD groups properly to manage the meshes to use and the associated distance ranges.

To access the setup of LOD groups, select: **Add Component > Rendering > LOD Group**.

In Unity 5 you can set a **Fade Mode** for each LOD level to blend to contiguous LODs. This smooths the transition between them. Unity calculates a blending factor according to the screen size of the object and passes it to your shader for blending. You must implement the geometry blending in a shader.

The following figure shows the LOD group settings:

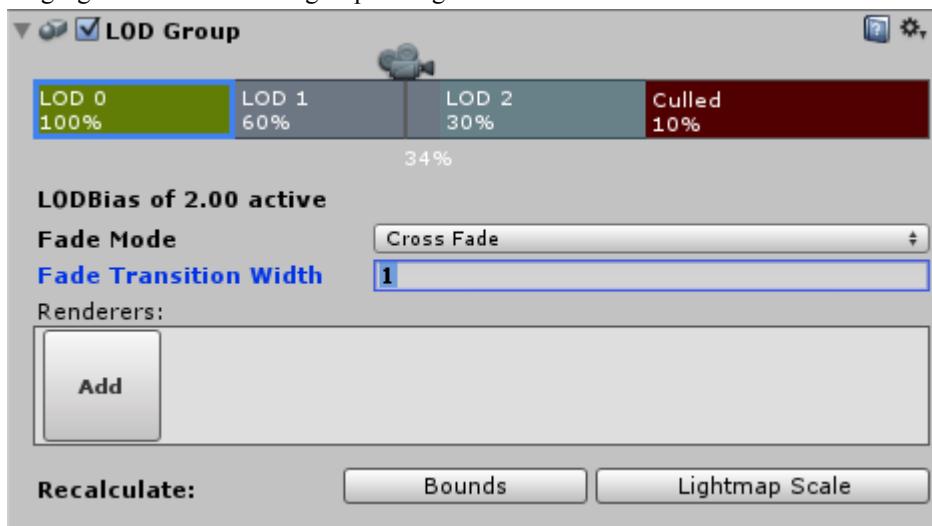


Figure 4-7 LOD group settings

Avoid mathematical functions in custom shaders

When you are writing custom shaders, try to avoid using expensive built-in mathematical functions such as:

- `pow()`.
- `exp()`.
- `log()`.
- `cos()`.
- `sin()`.
- `tan()`.

4.2.2 Lightmaps and light probes

Runtime lighting calculations are computationally expensive. A popular technique to reduce the computation requirements called lightmapping pre-computes the lighting calculations and bakes them into a texture called a lightmap.

This means you lose the flexibility of a fully dynamically lit environment, but you do get very high quality images without impacting performance.

To bake the resulting lighting in a static lightmap

- Set the geometry that receives the lighting to **static**.
- Set the **Baking** option in the light to **Baked** instead of **Realtime**.
- Tick the **Baked GI** option in Scene tab of the Lightmapping window.

To see the resulting lightmap:

- Select the geometry.
- Open the Lighting window by selecting **Window > Lighting**.
- Press the **Object** button.
- Select **Baked Intensity lightmap** in the preview options.

If the Continuous Baking option is selected Unity bakes the lightmap and updates the scene in the Editor in seconds.

A quick way to check that the lightmap has set up correctly is to run the game in the Editor and disable the light. If the lighting is still there, the lightmap has been created correctly and it is in use.

The following figure shows an intensity lightmap in the **Lighting** tab.

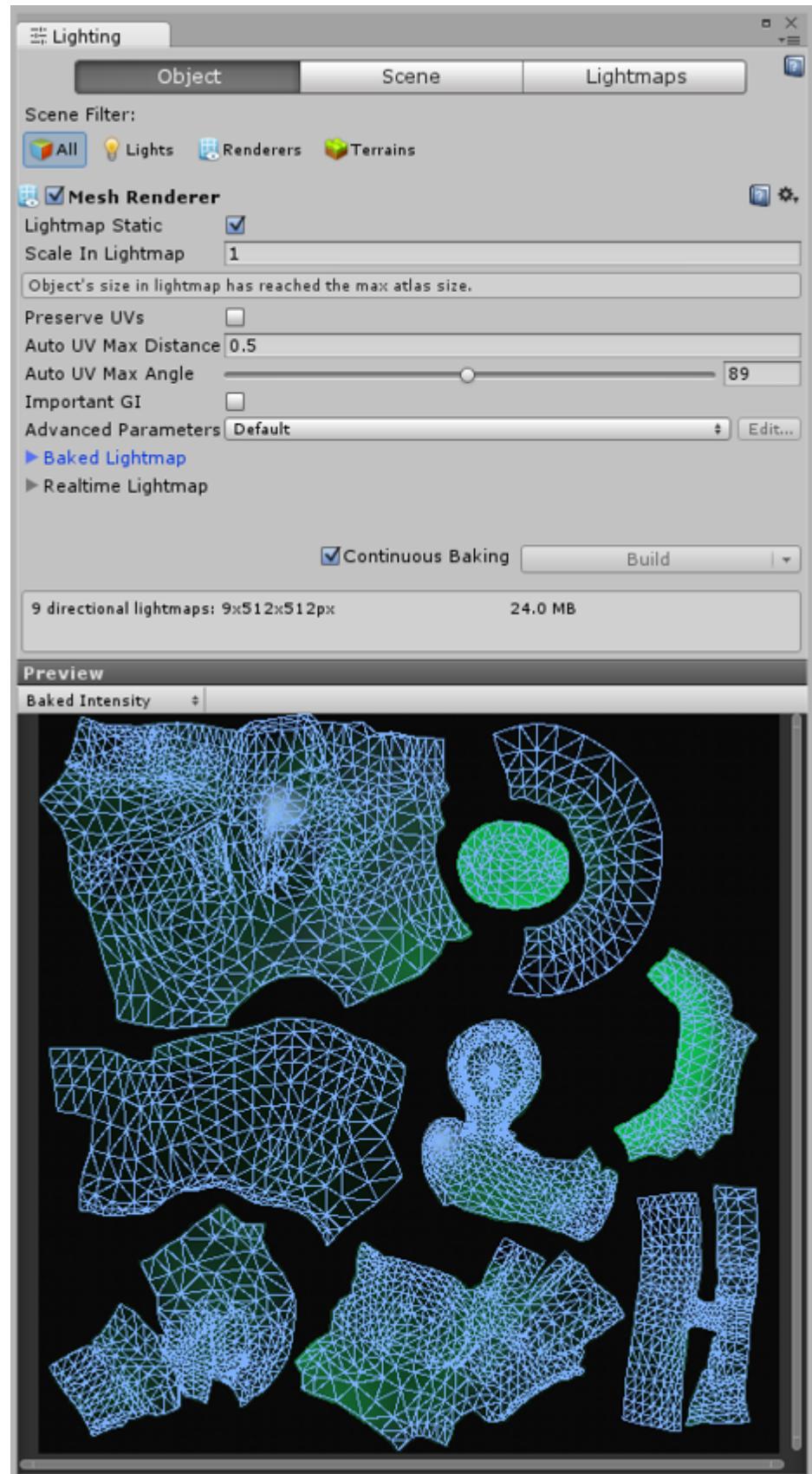


Figure 4-8 The intensity lightmap

The following figure shows the editor displaying lighting from a green light at the end of a cave. This lighting is generated with a static lightmap.

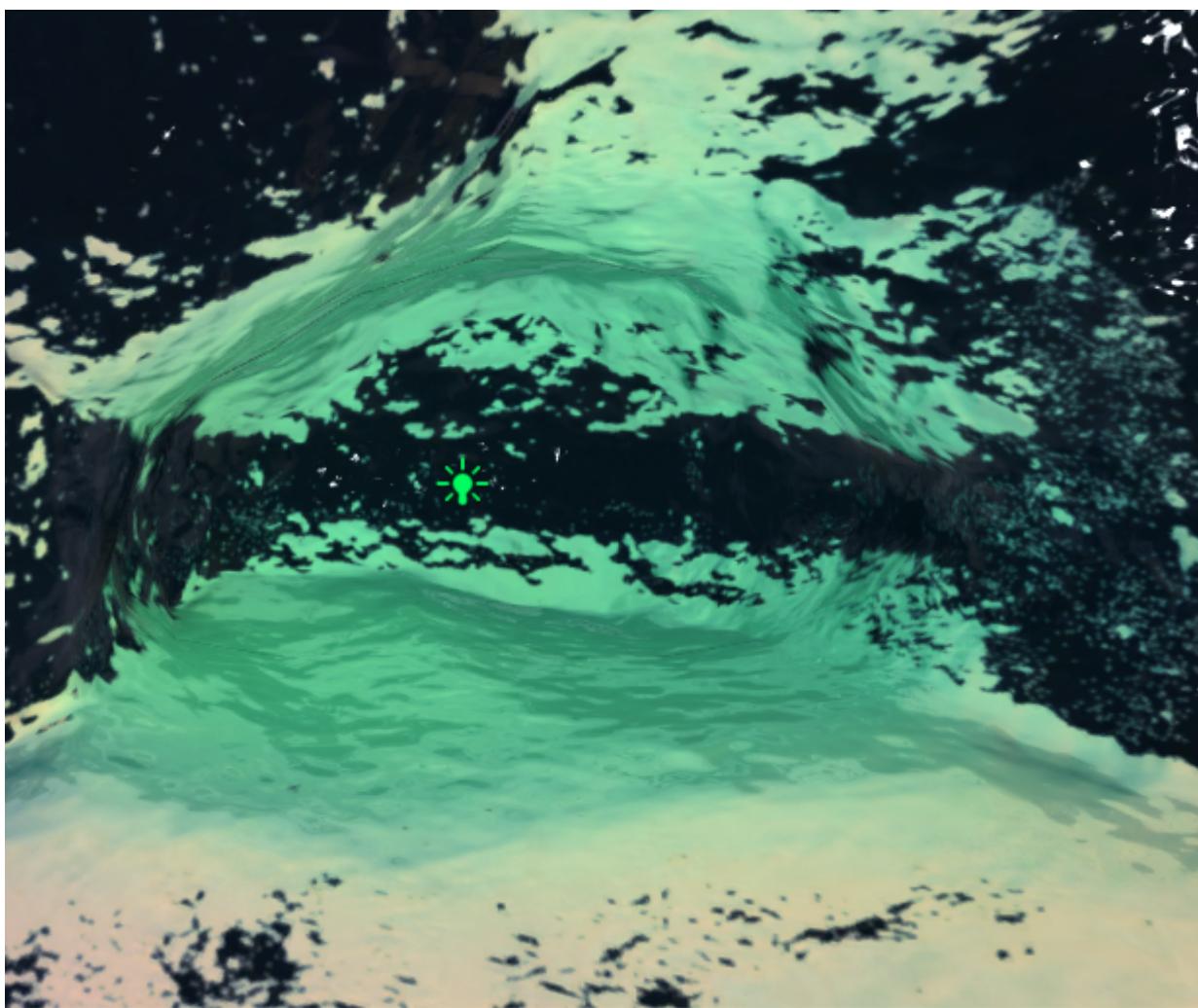


Figure 4-9 Adding a light to bake a static lightmap

The following figure shows the result of the static lightmap in the Ice Cave demo.

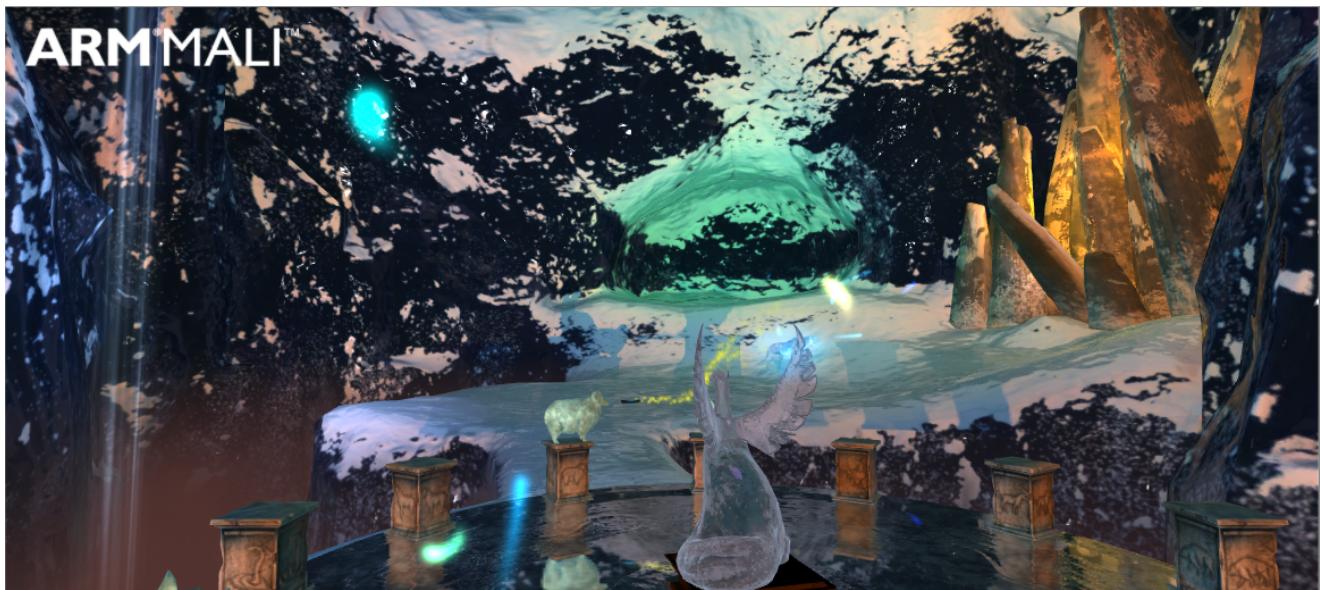


Figure 4-10 Lightmapped cave

Setting up lightmapping

To prepare an object for lightmapping you must have:

- A model in your scene with lightmap UVs.
- The model must be marked as **lightmap static**.
- There must be a light within range of the model.
- The **Baking type** of the light must be set to **Baked**.

————— Note ————

Only the static objects in your scene are lightmapped. They are not likely to be perfect, so experiment to see what works best for your game.

Objects that are not marked as static are not placed in the lightmap. Selecting a renderer provides you with a number of settings and enables you to set whether its lightmap is static or not.

Open the **Lighting** window from the main menu of the editor window and select **Window** and **Lightmapping**. There are three buttons:

- **Object**.
- **Scene**.
- **Lightmaps**.

The following figure shows lightmap options:



Figure 4-11 Lightmap options

Object

Clicking the **Object** button enables you to change settings related to lightmapping on the object you have selected in the hierarchy. This enables you to modify the object settings that impact the lightmapping process. Select a light, and you can change a number of options:

- **Baked Only** enables the light at baking time and disables it at runtime.
- **Baked If Baked GI** is selected, the light is baked.
- **Realtime** The light works for both pre-computed real-time GI and without GI.
- **Realtime Only** disables the light at baking time and enables it at runtime.
- **Mixed** The light is baked, but it is still present at runtime to give direct lighting to non-static objects.

Setting the majority of lights to **Baked** ensures the number of calculations at runtime is relatively low.

Scene

The **Scene** tab contains settings that apply to the whole scene. You can enable and disable the **Pre-computed Realtime GI** and **Baked GI** features in this tab.

In the **Environment Lighting** section, there are a number of options that allow you to define several factors influencing the environment lighting, such as the Skybox, the type of Ambient Source, and the Ambient Intensity:

- The **Reflection Bounces** option is the most important from the performance point of view. **Reflection Bounces** defines the number of inter-reflections between reflective objects, that is, the number of bake times for the probe that sees the objects. This option can have a large negative impact on performance if the reflection probes are updated at runtime. Only set the number of bounces higher than one if the reflective objects shall be clearly visible in the probes.
- In the **Precomputed Realtime GI** tab the **CPU Usage** option defines the amount of processor time that is spent evaluating GI at runtime. A higher value for **CPU Usage** results in faster reactions from the lighting, but might affect frame rate. The impact on the performance is lower in multiprocessor systems.
- The **Baked GI** tab contains an option where you can set the lightmap texture to be compressed. Compressing lightmap textures requires less storage space and less bandwidth at runtime but the compression process can add artifacts to the texture.
- In the **General GI** tab, be careful with the **Directional Mode** option. If you cannot use deferred lighting with dual lightmaps, another technique is to use directional lightmaps. These enable you to use normal mapping and specular lighting without real-time lights. Use directional lightmaps if normal mapping must be preserved but dual lightmaps are not available. This is typically the case on mobile devices. When **Directional Mode** is set to **Directional** an additional lightmap is created to store the dominant direction of incoming light. As a result this mode requires about twice as much storage space.
- In **Directional Specular** mode, additional data is stored for specular reflection and normal maps. In this case the storage requirements increase four times.
- The **Lightmaps** tab enables you to set and locate the lightmap asset file used for the scene. To access the Lightmap Snapshot box the **Continuous Baking** option must be unchecked.

The following figure shows lightmaps in the **Lighting** tab:

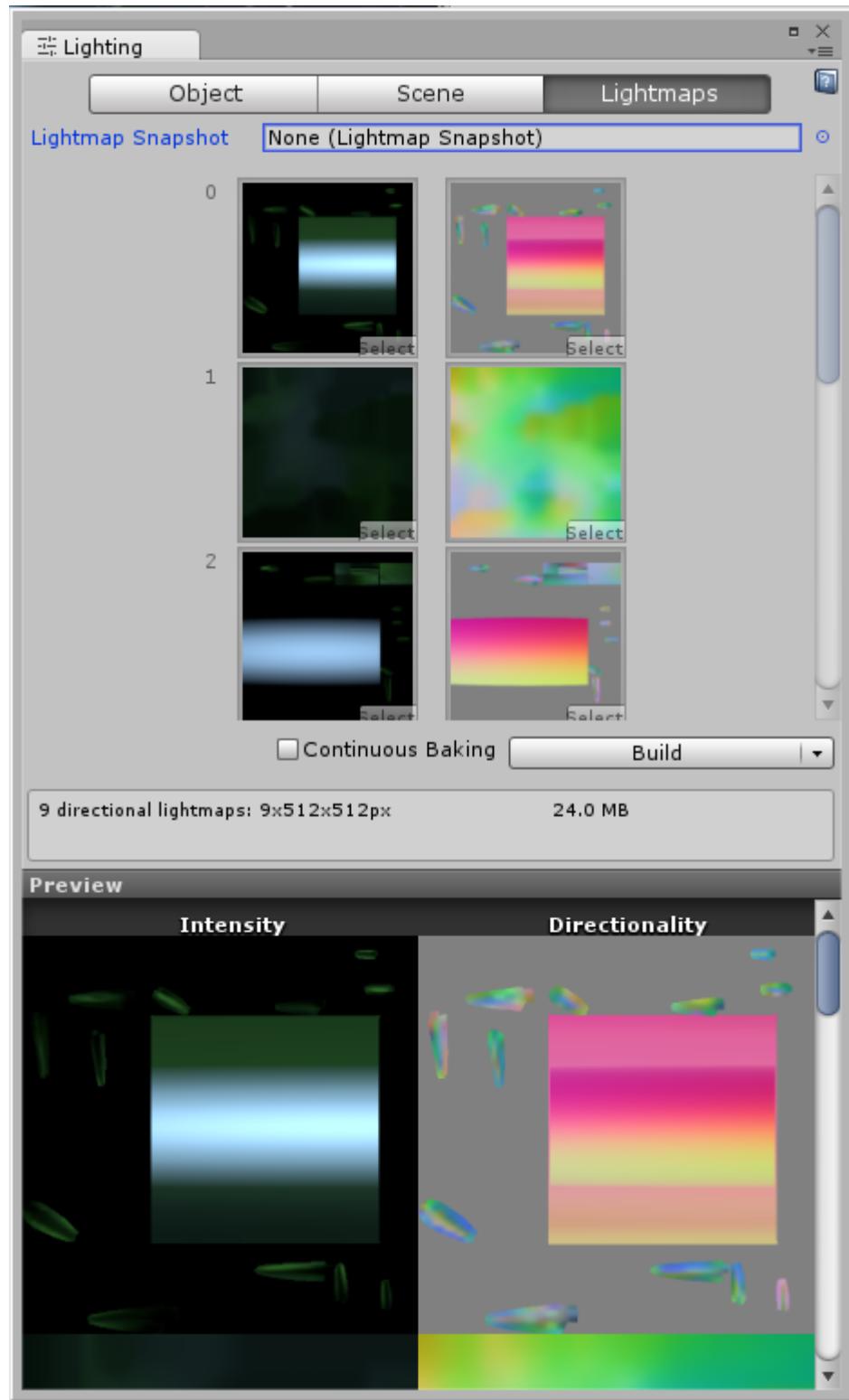


Figure 4-12 Lightmaps in the Lighting tab

Use directional lightmaps

If you cannot use deferred lighting with dual lightmaps, another technique is to use directional lightmaps. These enable you to use normal mapping and specular lighting without real-time lights.

Use directional lightmaps if normal mapping must be preserved, but dual lightmaps are not available. This is typically the case on mobile devices.

————— Note ————

This technique requires more video memory because it computes a second set of lightmaps to store directional information.

Use light probes for dynamic objects in your game

Light probes enable you to add some dynamic lighting to lightmapped scenes.

The following figure shows light probe settings:

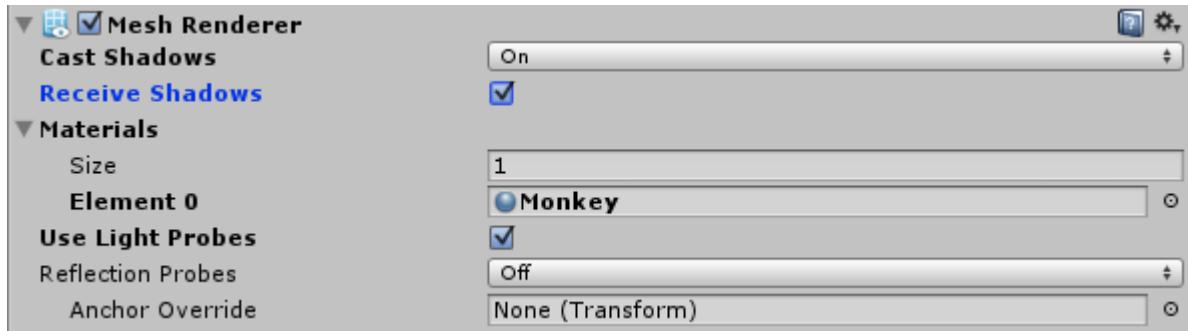


Figure 4-13 Light probe settings

Light probes take a sample, or probe, of the lighting in an area. If the probes form a volume, or cell, the lighting is interpolated between these probes depending on their position within the cell.

The following figure shows light probes:

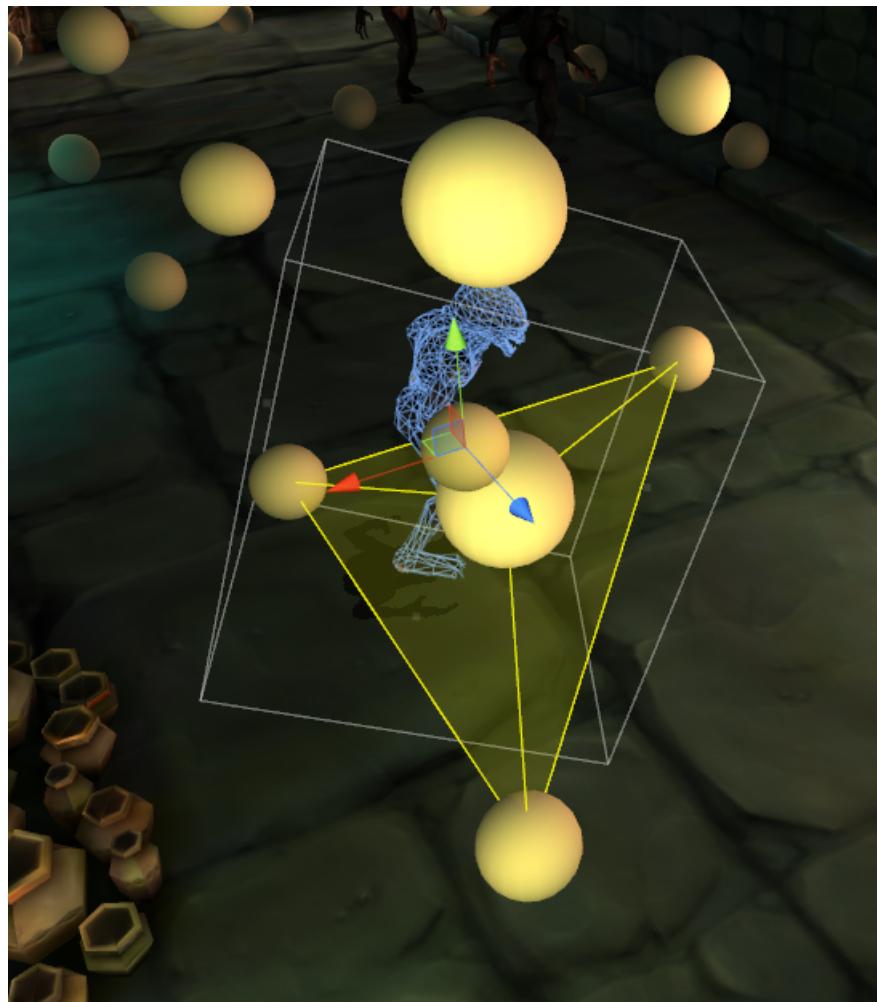


Figure 4-14 Light probes

The more probes there are, the more accurate the lighting is. You do not typically require many light probes because there is interpolation between probes. You require more light probes in areas where there are large changes in light color or intensity.

The lighting at any position can then be approximated by interpolating between the samples taken by the nearest probes.

Take care placing the light probes and mark the meshes you want to be influenced by them with the **Use Light Probes** option.

The following figure shows multiple light probes:

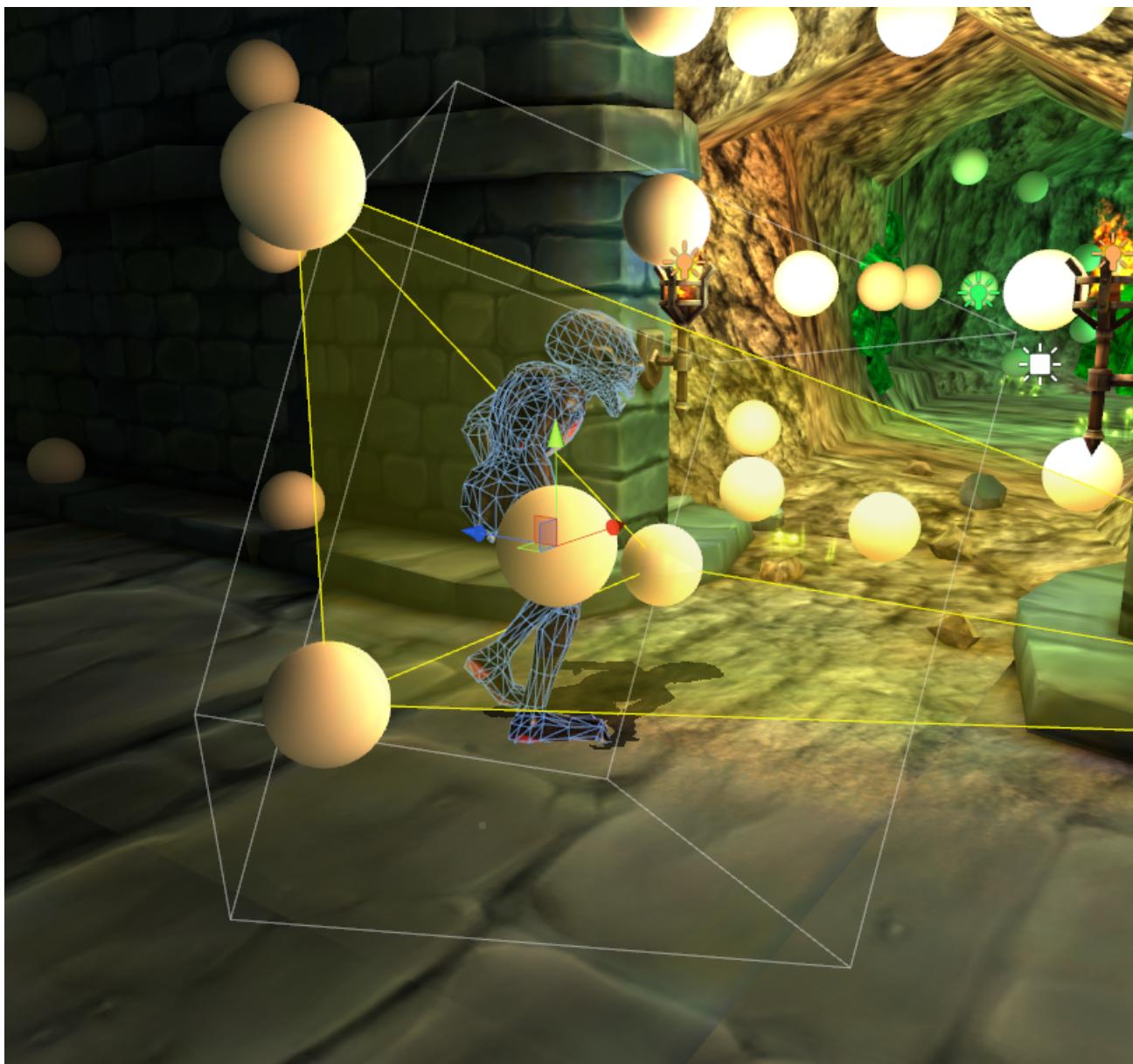


Figure 4-15 Multiple light probes

4.2.3 ASTC texture compression

ASTC texture compression is an official extension to the OpenGL and OpenGL ES graphics APIs. ASTC can reduce the memory required by your application and reduce the memory bandwidth required by the GPU.

ASTC offers texture compression with high quality, low bitrate and has many control options. It includes the following features:

- Bit rates range from *8 bits per pixel* (bpp) to less than 1bpp. This enables you to fine-tune the trade-off of file size against quality.
- Support for 1 to 4 color channels.
- Support for both *low dynamic range* (LDR) and *high dynamic range* (HDR) images.
- Support for 2D and 3D images.
- Support for selecting different combinations of features.

The following figure shows the ASTC settings window:

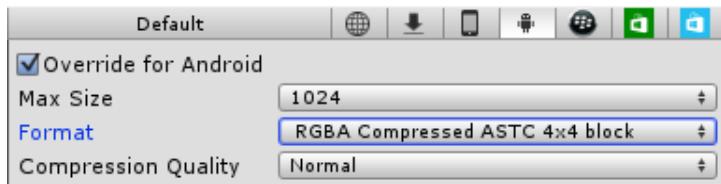


Figure 4-16 ASTC settings

There are a number of block size options available in the ASTC settings window. You can choose among these options and select the block size that best fits the assets. The larger block sizes provide higher compression. Select large block sizes for textures that are not shown in great detail, for example, objects far away from camera. Select smaller block sizes for textures that are show more detail, for example those closer to camera.

————— Note ————

- If your device supports ASTC, use it to compress the textures in your 3D content. If your device does not support ASTC, try using ETC2.
- You must differentiate between textures used in 3D content from textures used in the GUI elements. In some cases it might be best to leave the GUI textures uncompressed to avoid unwanted artifacts.

The following figure shows the block sizes available for different texture compression formats:

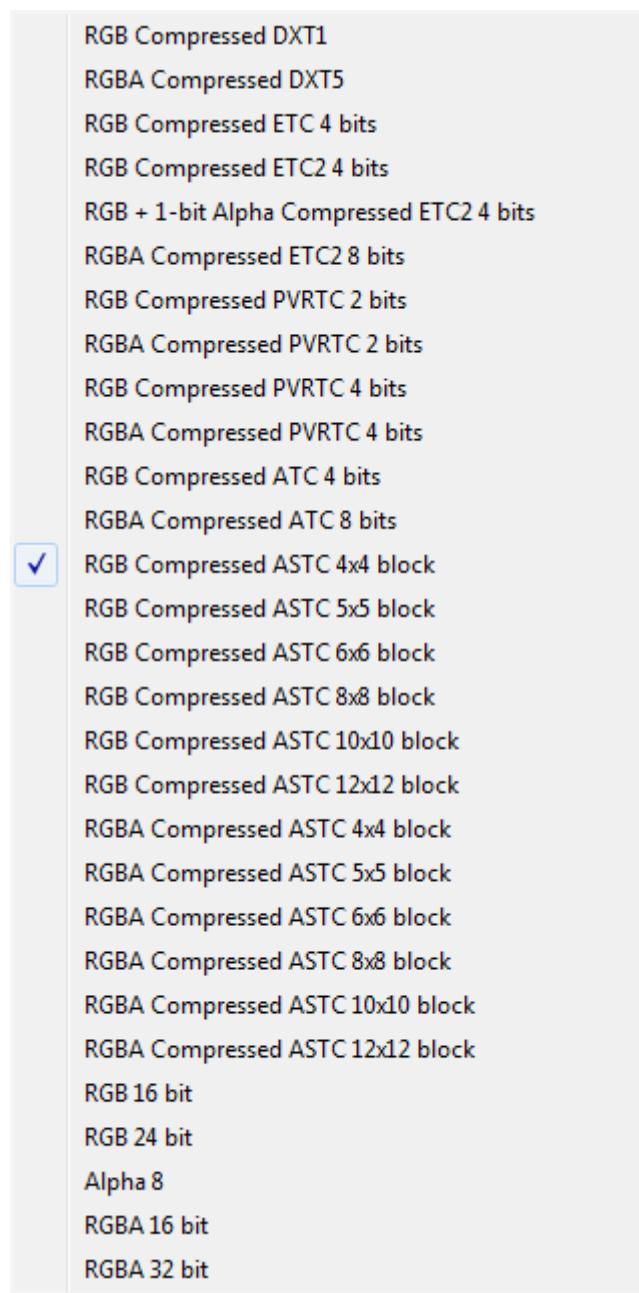


Figure 4-17 Texture compression block sizes

Selecting the correct format for ASTC textures

When compressing an ASTC texture, there are a number of options to choose from.

Texture compression algorithms have different channel formats, typically RGB and RGBA. ASTC supports several other formats, but these formats are not exposed within Unity. Each texture is typically used for a different purpose such as, standard texturing, normal mapping, specular, HDR, alpha, and look up textures. All of these texture types require a different compression format to achieve the best possible results.

The following figure shows texture settings:

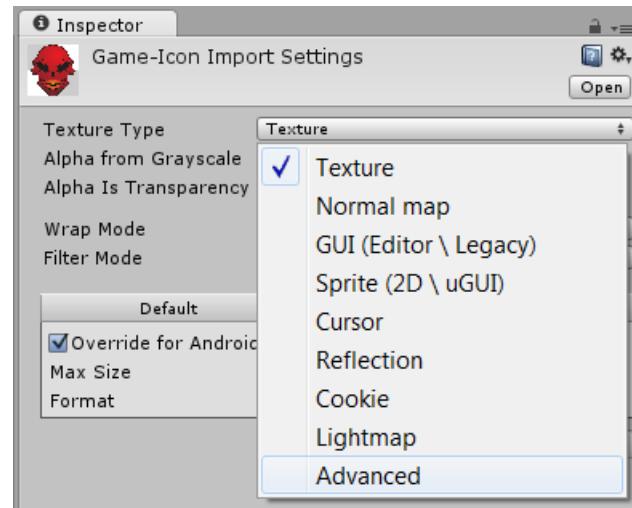


Figure 4-18 Texture settings

Do not compress all of your textures with one format in **Build Settings**. Keep texture compression as **Don't Override**.

Find your texture within the project hierarchy and bring it up in the **Inspector**. Unity typically imports your texture as the type **Texture**. This type only provides limited options for compression. Set the type to **Advanced** to show a larger choice of options.

The following diagram shows settings for a GUI texture with some transparency. The texture is for a GUI, so **sRGB** and **MipMaps** are disabled. To include transparency, you require the alpha channel. To enable this, tick the **Alpha Is Transparency** box and tick the **Override for Android** box.

The following figure shows advanced texture settings:

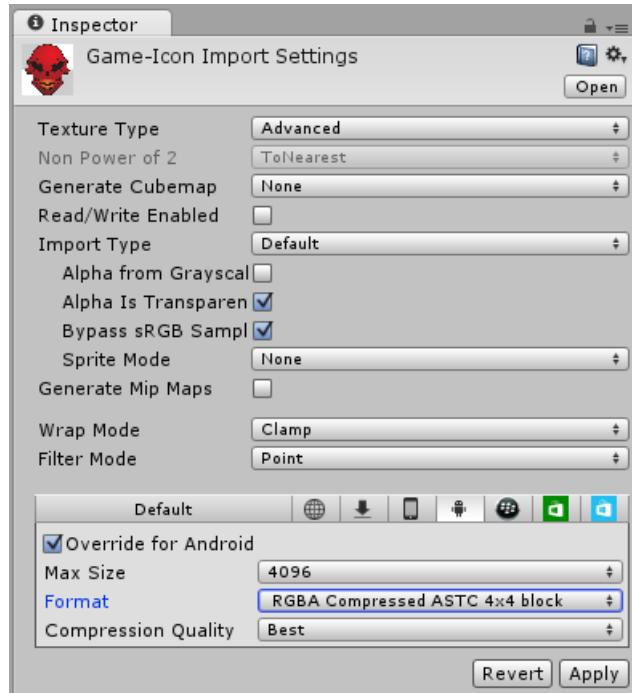


Figure 4-19 Advanced texture settings

There is an option to select a format and block size. RGBA includes the alpha channel and 4x4 is the smallest block size you can select. Set **Max texture size** to the maximum and set **Compression Quality**, this setting defines how much time is spent looking for accurate compression.

Selecting specific settings for all of your textures improves the visual quality of your project and avoids unnecessary texture data at compression time.

The following table shows the compression ratio for the available ASTC block sizes in Unity for an RGBA 8 bits per channel texture with a 1024x1024 pixel resolution at 4 MB in size.

Table 4-1 Compression ratios for the ASTC block sizes available in Unity

ASTC block size	Size	Compression ratio
4x4	1 MB	4.00
5x5	655 KB	6.25
6x6	455 KB	9.00
8x8	256 KB	16.00
10x10	164 KB	24.97
12x12	144 KB	35.93

4.2.4 Mipmapping

Mipmapping is a technique related to textures that can both enhance the visual quality and the performance of your game.

Mipmaps are pre-calculated versions of a texture at different sizes. Each texture generated is called a level, and it is half as wide and half as high as the preceding one. Unity can automatically generate the complete set of levels from the first level at the original size down to a 1x1 pixel version.

To generate the mipmaps do the following:

1. Select a texture in the **Project window**.
2. Change the **Texture Type** to **Advanced**.
3. Enable the **Generate Mip Maps** option in the **Inspector**.

The following figure shows mipmap settings:

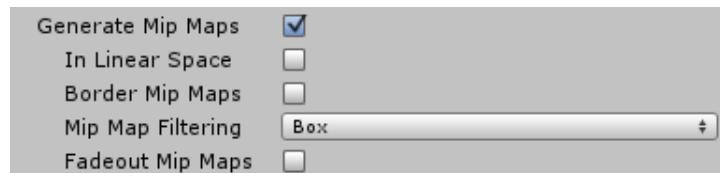


Figure 4-20 Mip map settings

If a texture does not have mipmap levels, when the area in pixels covered by a textured surface is smaller than the size of its texture, the GPU scales the texture down to fit the smaller area. However, some accuracy is lost in this process even with a filter to interpolate the pixel colors.

If a texture does have mipmap levels, the GPU fetches pixel data from the level closest to the object size to render the texture. This ensures both higher image quality and reduced bandwidth because the levels are scaled offline for better quality and only the texture data from the correct level is fetched by the GPU. The disadvantage of mipmapping is it requires 33% more memory to store the texture data.

Mipmaps and GUI textures

You do not usually require mipmapping for textures used in a 2D UI. UI textures are typically rendered on screen without scaling so they only use the first level in the mipmap chain.

To change this setting select a texture in the **Project window** then go in the **Inspector** and look at **Texture Type**. Set the type to **Editor GUI and Legacy GUI**, or set the type to **Advanced** and disable the **Generate Mip Maps** option.

4.2.5

Skyboxes

Skyboxes are often used in games and other applications, there are several methods to implement them.

You can draw a skybox by rendering the background of the camera using a single cubemap.

This requires one cubemap texture and one draw call. This uses less memory, memory bandwidth, and draw calls, compared to other methods.

To set up a skybox using this method:

1. Select the camera.
2. Ensure **Clear Flags** is set to **Skybox**.
3. Select or add a skybox component.

The skybox component has one spot for a material. This is the material that Unity uses to draw the background of your camera at the start of each frame.

The material you use is the one that contains all the information you require. Create a material using the **Mobile > Skybox** shader and fill the six images of the Skybox with the material. Your material preview displays an image.

When you have completed the material, drag it into the skybox component. Your skybox renders in the background correctly, without any obvious seams or unnecessary draw calls.

4.2.6

Shadows

Shadows help add perspective and realism to your scenes. Without them it can sometimes be difficult to tell the depth of objects, especially if they are similar to the surrounding objects.

Shadow algorithms can be very complex, especially when rendering accurate, high resolution shadows. Ensure you select an appropriate level of complexity and resolution for the shadows in your game.

Unity supports transform feedback for calculating real-time shadows.

————— Note —————

The Ice Cave demo implements custom shadows. Shadows based on local cubemaps are combined with shadows rendered at runtime.

Unity has a number of options for shadows under **Edit > Project Settings > Quality** that can impact the performance of your game:

Hard/Hard and Soft Shadows

Soft shadows look more realistic but take longer to calculate.

Shadow Distance

The **Shadow Distance** option defines the distance from the camera that shadows appear in.

Increasing the shadow distance increases the number of shadows visible, and this increases the computational load. Increasing the shadow distance also increases the number of texels available for the shadows in the shadow map, passively increasing the resolution of your shadows.

You can use hard shadows with a small shadow distance and a high resolution. This produces reasonable quality shadows within a good distance of the camera that are not too complex.

Light mapped objects do not produce realtime shadows, so the more static shadows you can bake into the scene, the fewer real-time calculations the GPU does.

The following figure shows an alien character with a shadow:



Figure 4-21 Alien with shadow

Use real-time shadows sparingly

Real-time shadows can dramatically enhance the realism of a scene but they are computationally expensive.

On mobile devices, try to limit the number of lights that include real-time only shadows and try to use lightmapping instead.

Consider the mesh renderer component of the objects in your scene. If you do not intend to use them for casting or receiving shadows disable the **Cast Shadows** and **Receive Shadows** options accordingly. This reduces the computation cost of rendering shadows.

You can find more settings for shadows in the **Quality Settings** section such as:

- **Shadow Resolution** enables you to select the balance between quality and processing time.
- **Shadow Distance** enables you to limit shadow generation to objects close to the camera.
- **Shadow Cascades** enables you to select the balance between quality and processing time. You can set this to zero, two or four. Cascaded Shadow Maps are used for directional lights to achieve very good shadow quality, especially for long viewing distances. A higher number of cascades produces better quality, but increases processing overhead.

4.2.7 Occlusion culling

Occlusion culling is a process that disables the rendering of objects when they are obscured from the view of the camera. This saves GPU processing time by rendering fewer objects.

Unity automatically performs frustum culling when objects exit the camera frustum completely however, depending on your style of application, there might still be other objects that cannot be seen and do not have to be rendered.

Unity includes an occlusion culling system called Umbra. For more information on Umbra see occlusion culling in the Unity documentation.

The settings you use for occlusion culling depends on the style of your game. You must be careful picking settings because using occlusion culling in a scene with the incorrect settings can degrade performance.

4.2.8 Use `OnBecameVisible()` and `OnBecomeInvisible()` callbacks

If you use the callbacks `MonoBehaviour.OnBecameVisible()` and `MonoBehaviour.OnBecomeInvisible()`, Unity notifies your scripts when their associated game objects moves in or out of a camera frustum. Your application can then act accordingly.

You can use `OnBecameVisible()` and `OnBecomeInvisible()` to optimize the rendering process, for example, rendering reflections on a pool with a second camera and render targets.

This involves rendering geometry and combining textures off screen before rendering to the final screen surface. This technique is relatively expensive so it is only used when it is necessary. You are only required to render a reflection when it is visible. That is when:

- The reflection surface is within the camera frustum.
- Nothing opaque is in front of the surface.

These conditions are checked with the `OnBecameVisible()` and `OnBecomeInvisible()` callbacks from the reflective surface:

```
void OnBecomeVisible()
{
    enabled = true;
}

void OnBecomeInvisible()
{
    enabled = false;
}
```

Even with these checks in place there can still be times when a reflection might be rendered off screen even though it is not visible onscreen. To avoid this you can add another condition:

For example, the camera must be inside the room of the reflective surface:

```
void OnBecomeVisible()
{
    if (inside == false)
    {
        return;
    }
    enabled = true;
}

void OnBecomeInvisible()
{
    if (inside == false)
    {
        return;
    }
    enabled = false;
}

void OnTriggerEnter()
{
    inside = true;
}

void OnTriggerExit()
{
    inside = false;
}
```

These conditions restrict the rendering of reflections to specific areas of the game. This means you can add effects in other, less compute intensive areas of the game.

4.2.9 Specify the rendering order

In a scene, the object rendering order is very important for performance.

If objects were rendered in random order, an object might be rendered and subsequently occluded by another object in front of it. This means all the computations required to render the occluded object are wasted.

Various software and hardware techniques exist to reduce the amount of wasted computation because of occluded objects, but you can guide this process because you have knowledge of how the scene is explored by a player.

One of the hardware techniques for reducing wasted computation that is available on the ARM Mali GPUs from the Mali-T600 series onwards, is Early-Z. Early-Z is a completely transparent system from your point of view that performs a Z-test before the fragment shader is actually processed. If the GPU cannot enable Early-Z optimization, the depth test is executed after the fragment shader. This can be computationally expensive and the computations can be wasted if the fragment is occluded. The Early-Z system checks that the depth of the pixel being processed is not already occupied by a nearer pixel. If it is occupied, it does not execute the fragment shader. This system provides performance benefits but it is automatically disabled in some cases, such as if the fragment shader modifies the depth by writing into the `gl_FragDepth` variable, the fragment shader calls `discard`, or if blending or alpha testing is enabled for objects such as transparent objects. To assist this system to achieve maximum efficiency, ensure that opaque objects are rendered from front to back. This helps to reduce the overdraw factor in scenes with only opaque objects.

Ordering the rendering of each frame front-to-back can be expensive and also incorrect if you render transparent objects in the same pass. ARM Mali GPUs from T620 onwards provide a mechanism called *Pixel Forward Kill* (PFK). Mali GPUs are pipelined so multiple threads can be concurrently executing for the same pixel. If a thread completes its execution, the PFK system stops all other threads for that pixel if the current one covers them. The effect is a reduction of wasted computations.

Unity provides you with **Queue** options inside the shaders, or in the material to specify the order of rendering. This can be set in the shader, so objects that have a material that uses that shader are rendered together. Inside this rendering group, the order of rendering is random except in some cases such as transparency.

By default, Unity provides some standard groups that are rendered from first to last in the following order:

Table 4-2 Queue values for specifying rendering order

Name	Value	Notes
Background	1000	-
Geometry	2000	Default, used for opaque geometry.
AlphaTest	3000	This is drawn after all opaque objects. For example, foliage.
Transparent	4000	This group is also rendered in back to front order to provide the correct results.
Overlay	5000	Overlay effects such as user interface, lens flares, dirty lens.

The integer values can be used instead of their string names. These are not the only values available. You can specify other queues using an integer value between those shown. The higher the number, the later it is rendered.

For example, you can use one of the following instructions to render a shader after the Geometry queue, but before the AlphaTest queue:

```
Tags { "Queue" = "Geometry+1" }

Tags { "Queue" = "2001" }
```

Using the rendering order to increase performance

In the Ice Cave demo, the cave covers large part of the screen and its shaders are expensive. Avoiding rendering parts of it when possible can increase performance.

Rendering order optimization was included after looking at the composition of the framebuffers using the Unity Frame Debugger and other tools such as the ARM Mali Graphics Debugger. These enable you to see the rendering order.

To open the Unity **Frame Debugger**, select the menu option **Window > Frame Debugger**. This is useful because there might be things that appear to be correct in editor mode, but might not work correctly when you execute them. This can be the case if for example, you have run time only settings, or if you are rendering to a texture from another camera. After starting the demo in play mode and positioning the camera, you can enable the Frame Debugger and get the sequence of drawings executed by Unity.

In the Ice Cave demo, scrolling down the draw calls shows that the cave is rendered first. The objects are then rendered into the scene occluding parts of the cave that are already rendered. Another example is the reflective crystals that in some scenes are occluded by the cave. In these cases, setting a higher rendering order results in a reduction in computations because fragment shaders are not executed for the occluded crystals.

4.2.10 Use depth pre-pass

Setting the rendering order for objects to avoid overdraw is useful, but it is not always possible to specify the rendering order for each object.

For example, if you have a set of objects with computationally expensive shaders and the camera can rotate around them freely, some objects that were at the back can move to the front. In this case, if there is a static rendering order set for these objects, some might be drawn last, even if they are occluded. This can also happen if an object can cover parts of itself.

In these cases, you can use a depth-prepass to reduce overdraw. The depth-prepass renders the geometry without writing colors in the framebuffer. This initializes the depth buffer for each pixel with the depth of the nearest visible object. After this pre-pass the geometry is rendered as usual but using the Early-Z technique, only the objects that contribute to the final scene are actually rendered. Extra vertex shader computations are required for this technique because the vertex shader is computed two times for each object, once for filling the depth buffer and another for the actual rendering. This technique is very useful if your game is fragment bound and you have spare capacity in the vertex shader.

In Unity, to do a rendering pre-pass for objects with custom shaders, add an extra pass to your shaders:

```
// extra pass that renders to depth buffer only
Pass {
    ZWrite On
    ColorMask 0
}
```

After adding this pass, the frame debugger shows that the objects are rendered twice. The first time they are rendered there are no changes in the color buffer.

Note

You can see the depth buffer by choosing it in the top left menu of the frame debugger.

4.3 Asset optimizations

The following list describes asset optimizations:

Disable Read/Write for static textures

If you do not dynamically modify a texture, ensure the **Read/Write Enabled** option in the **Inspector** is disabled.

Combine meshes to reduce draw calls

To reduce the number of draw calls required for rendering you can combine several meshes into one with the `Mesh.CombineMeshes()` method. If the meshes all share the same material, set the `mergeSubMeshes` argument to `true` so it generates a single submesh out of each mesh in the combine group.

Combining several meshes into a single larger mesh helps you:

- Create more effective occluders.
- Turn tile-based assets into a single large seamless solid asset.

The mesh combine script can be useful for performance optimization but this depends on the makeup of your scene. Large meshes tend to stay in view longer than smaller meshes, so experiment to get the correct size.

One way to apply this technique is to create an empty game object in the hierarchy, make it the parent of all the meshes that you want to combine and then attach it to a script.

For more information on the mesh combine script, see the Unity documentation at:

<http://unity3d.com>.

Do not import animations data on FBX mesh models that do not animate

When importing an FBX mesh that does not contain any animation data, you can set the **Animation Type** to **None** in the **Rig** tab of the import settings. If this is set, placing your mesh into the hierarchy Unity does not generate an unused animator component.

Avoid Read/Write meshes

If your model is modified at runtime, Unity keeps a second copy of the mesh data in memory to modify while preserving the original.

If your model is not modified at run-time, even to be scaled, disable the **Read/Write Enabled** option from the **Model** tab of the import settings. A second copy is no longer required so this saves memory.

Use texture atlases

You can use a texture atlas to reduce the number of draw calls required for a set of objects.

A texture atlas is a group of textures combined into one large texture. Multiple objects can reuse this texture with an appropriate set of coordinates. This helps with the automatic batching that Unity applies to objects sharing the same material.

When setting the UV texture coordinates of an object, avoid changing the `mainTextureScale` and `mainTextureOffset` properties of its material. This creates a new unique material that does not work with batching. Instead, access the mesh data through the `MeshFilter` component and change the coordinates per vertex using the `Mesh.uv` property.

The following figure shows a texture atlas:

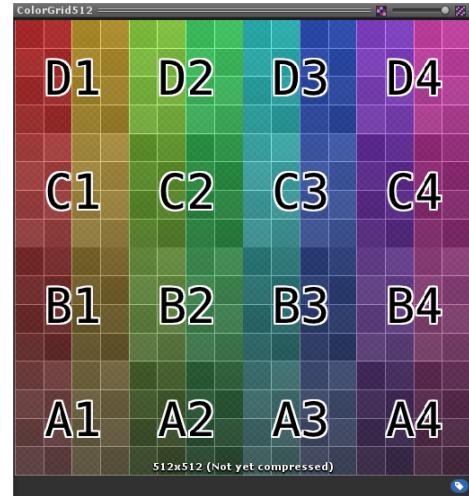


Figure 4-22 Texture atlas

4.4 Optimizing with the Mali™ Offline Shader Compiler

The Mali Offline Shader Compiler is a tool that enables you to compile vertex, fragment, and compute shaders into a binary form. You can also use the compiler as a profiling tool.

This section contains the following subsections:

- [4.4.1 About the Mali™ Offline Shader Compiler](#) on page 4-55.
- [4.4.2 Measuring Unity shaders](#) on page 4-55.
- [4.4.3 Analyzing statistics](#) on page 4-56.
- [4.4.4 Optimizing for the Mali™ GPU pipelines](#) on page 4-56.
- [4.4.5 Additional techniques for reducing pipeline cycles](#) on page 4-57.

4.4.1 About the Mali™ Offline Shader Compiler

The shaders in your application run on the GPU. This requires the GPU to spend time computing the final result of your shader, such as the vertex position or the color of the pixel.

The Mali Offline Shader Compiler provides information about the number of cycles shaders require to execute in each pipeline of the Mali GPU.

The cycle values produced are tailored for a specific GPU. You select the GPU as an option on the command line. Ensure you choose GPUs that correspond to that range of devices that you want your application to target. This ensures that the statistics you get from the tool are realistic and match your typical use case scenario.

4.4.2 Measuring Unity shaders

You write Unity shaders in the programming language *C for Graphics* (Cg). Cg is based on the C programming language with some modifications to make it more suitable for GPU programming.

Unity translates the Cg to OpenGL, OpenGL ES, or DirectX during the build process.

To retrieve the OpenGL ES shader code:

1. Select the shader you want to analyze in Unity.
2. Choose **OpenGL ES30** or **OpenGL ES20** as the Custom Platform you want to build for.
3. Click the **Compile and show** button.

The result is displayed in your development environment.

Note

- The Mali Offline Shader Compiler only supports OpenGL ES shaders.
 - If your build platform is set to Android, Unity builds **OpenGL ES30** shaders by default.
-

Vertex and fragment sessions are typically delimited by `#ifdef VERTEX` or `#ifdef FRAGMENT`. If you use an option such as `#pragma multi_compile <FEATURE_OFF> <FEATURE_ON>`, multiple shader variants are built in the file.

Typically there are multiple `VERTEX` and `FRAGMENT` sections. Each variant is compiled separately when Unity starts your application. When you enable a feature, the relevant variant is selected.

Because the code has been translated into OpenGL ES, you can copy the vertex and fragment shader code into two separate files and compile them with the Mali Offline Shader Compiler.

Compile the shaders with one of the following options:

- `-v` for vertex shaders.
- `-f` for fragment shaders.

The following figure shows the output of the Mali Offline Shader Compiler:

```

ARM Mali Offline Compiler v4.6.0
(C) Copyright 2007-2015 ARM Limited.
All rights reserved.

No driver specified, using "offline_api" as default.

No core specified, using "Mali-T880" as default.

No core revision specified, using "r2p0" as default.

8 work registers used, 13 uniform registers used, spilling not used.

          A      L/S      T      Total      Bound
Cycles:   43      11      13      67      A
Shortest Path: 16      11      13      40      A
Longest Path: 16      11      13      40      A

Note: The cycles counts do not include possible stalls due to cache misses.

Compilation succeeded.

```

Figure 4-23 Output of the Mali Offline Shader Compiler

4.4.3 Analyzing statistics

The statistics produced by the Mali Offline Shader Compiler provide a measurement of how many cycles per vertex or pixel the shader requires.

The result is subdivided into three lines:

- Total.
- Shortest path.
- Longest path.

The shortest and longest paths are measured by looking at the effect of taking or not taking branches in your code. This provides an estimate of the smallest and largest number of cycles of execution.

For Arithmetic, the measurement in the first line is divided by the number of Arithmetic pipelines. This is one, two or three, depending on the Mali GPU.

The second and third lines are for the Load/Store and Texture pipelines. These do not take into account cache misses, so it is best to multiply those numbers by 1.5 to get more realistic estimates.

4.4.4 Optimizing for the Mali™ GPU pipelines

The Mali Offline Shader Compiler provides numbers of cycles used in each pipeline. The shader is slowest in the pipeline with the highest number of cycles. Optimize your shader with optimizations that target that pipeline it is slowest in.

Mali GPUs contain three types of processing pipeline:

- Arithmetic pipeline.
- The Load/Store pipeline.
- The Texture pipeline.

The pipelines all run in parallel. Your shaders typically use all three types of pipeline.

The Arithmetic pipeline

All arithmetic operations consume cycles in the Arithmetic pipeline.

The following are a number of ways you can reduce the Arithmetic pipeline usage:

- Avoid using complex arithmetic such as:
 - The inverse matrix function.
 - Modulo operators.
 - Division.
 - Determinant.

- Sine.
- Cosine.
- For integer operands, use operations such as shifts to compute divisions, modulo, and multiplications.
- Use transpose instead of inverse for orthogonal matrices.
- To avoid computing the transpose, switch the order of operands in a matrix-vector or matrix-matrix multiplication if one of the matrices is transposed. For example:

```
Transpose(A)*Vector == Vector * A.
```

You can also reduce load on the Arithmetic pipe by moving load to the other pipelines:

- Pass matrices as uniforms instead of computing them. This uses the Load/Store pipeline.
- Use a texture to store a set of precomputed values that represent a function such as sine or cosine. This moves the load to the Texture pipeline.

The Load/Store pipeline

The Load/Store pipeline is used for reading uniforms, writing varyings, and accessing buffers in the shaders such as Uniform Buffer Objects or Shader Storage Buffer Objects.

If your application is Load/Store pipeline bound, try the following techniques:

- Use a texture instead of a buffer object to read data in the shader.
- Compute data using arithmetic operations.
- Compress or reduce uniforms and varyings.

The Texture pipeline

Texture accesses use cycles in the Texture pipeline and use memory bandwidth. Using large textures can be detrimental because cache misses are more likely and this can cause multiple threads to stall while waiting for data.

To improve the performance of the Texture pipeline try the following:

Use mipmaps

Mipmaps increase the cache hit rate because it selects the best resolution of the texture to use based on the variation of texture coordinates.

Use texture compression

This is also good for reducing the memory bandwidth and increasing the cache hit rate. Each compressed block contains more than one texel, so accessing it makes it more cacheable.

Avoid trilinear or anisotropic filtering

Trilinear and anisotropic filtering increase the number of operations required to fetch texels. Avoid using these techniques unless you absolutely require them.

4.4.5 Additional techniques for reducing pipeline cycles

There are a number of additional techniques you can use to reduce the cycles used in each pipeline.

Avoid register spilling

The Mali Offline Shader Compiler indicates if your shader spills registers. Register spilling is typically caused in a thread by a high number of variables that cannot fit entirely in the register set.

Register spilling is typically caused in a thread by a high number of:

- Input uniforms.
- Varyings.
- Temporary variables.

Register spilling can also occur if variables are high precision.

Register spilling forces the Mali GPU to read some uniforms from memory, this increases the load on the Load/Store unit and reduces performance. To solve this issue, try to reduce the number and the precision of the uniforms you supply to the shader.

In the Ice Cave demo, some of the shaders suffered from register spilling, for example:

```
8 work registers used, 16 uniform registers used, spilling used.
```

Figure 4-24 Shader with register spilling.

Reducing the number of uniforms permitted solves this problem, and the result is an increase in performance, for example:

```
8 work registers used, 13 uniform registers used, spilling not used.
```

Figure 4-25 Shader with no register spilling.

Reduce the precision of varyings and uniforms

When you write custom shaders, you can specify the floating point precision of uniforms and varyings using 32-bit floats or 16-bit half-floats. The precision determines the minimum and maximum values and the granularity of values that the variable can represent.

There are several advantages of using half-floats:

- Bandwidth usage is reduced.
- The cycles used in the arithmetic pipeline are reduced because the shader compiler can optimize your code to use more parallelization.
- The number of uniform registers required is reduced and this in turn reduces the risk of register spilling.

The following code provides examples of a simple fragment shader variant from the Ice Cave demo. The shader is compiled with the Mali Offline Shader Compiler twice.

The first code example is compiled with floats:

```
$ malisc -f -V Compiled-CaveMaliStandardFloat.frag
ARM Mali Offline Compiler v4.6.0
(C) Copyright 2007-2015 ARM Limited.
All rights reserved.

No driver specified, using "offline_api" as default.

No core specified, using "Mali-T880" as default.

No core revision specified, using "r2p0" as default.

8 work registers used, 13 uniform registers used, spilling used.

          A      L/S      T      Total      Bound
Cycles:    47      15      10      72      A
Shortest Path: 15      13       9      37      A
Longest Path: 16      15      10      41      A

Note: The cycles counts do not include possible stalls due to cache misses.

Compilation succeeded.
```

Figure 4-26 Shader compiled with Floats

The second code example is compiled with half-floats:

```
$ malisc -f -V Compiled-CaveMaliStandardHalf.frag
ARM Mali Offline Compiler v4.6.0
(C) Copyright 2007-2015 ARM Limited.
All rights reserved.

No driver specified, using "offline_api" as default.

No core specified, using "Mali-T880" as default.

No core revision specified, using "r2p0" as default.

7 work registers used, 7 uniform registers used, spilling not used.

          A      L/S      T      Total      Bound
Cycles:   47      11      10      68      A
Shortest Path: 15       9       9      33      A
Longest Path: 15      11      10      36      A

Note: The cycles counts do not include possible stalls due to cache misses.

Compilation succeeded.
```

Figure 4-27 Shader compiled with Half floats

The number of Load/Store instructions is reduced in the half-float version. The number of work and uniform registers used is reduced and there is no register spilling.

The code generated with half-floats is also smaller than code generated with floats. This improves the cache hit rate on the Mali GPU increasing performance.

Use world space normal maps for static objects

You can use Tangent space normal maps to increase the details of a model without increasing the geometric detail. You can use tangent space normal maps on animated objects without modifying them because of their locality to each triangle of the mesh.

Unfortunately these require more arithmetic operations to be performed in the shaders to achieve the correct result. For static objects, these calculations are typically unnecessary.

You can alternatively use local space normal maps or world space normal maps. Using local space normal maps reduces the number of calculations performed in the shaders but transformations on the model must be applied to the sampled normal. World space normal maps do not require any transformations but these are static and the objects cannot move. In the Ice Cave demo, the cave and other high quality objects are static and using world space normal maps reduces the number of ALU operations required by the shaders considerably. Most common 3D modeling tools can create world space normal maps or you can generate them by code in an offline process.

Chapter 5

Global illumination in Unity with Enlighten

This chapter describes global illumination in Unity with Enlighten.

It contains the following sections:

- [5.1 About Enlighten on page 5-61](#).
- [5.2 The structure of Enlighten on page 5-62](#).
- [5.3 Setting up a scene with Enlighten on page 5-71](#).
- [5.4 Example: Enlighten setup of the Ice Cave on page 5-73](#).
- [5.5 Using Enlighten in custom shaders on page 5-80](#).

5.1 About Enlighten

Unity 5 includes Enlighten, a real-time global illumination solution from Geomerics. Geomerics is an ARM company.

Unity 5 uses Enlighten for simulating real-time indirect lighting. Indirect lighting is light that bounces off surfaces back into a scene.

The core of Enlighten is a real-time radiosity engine. The engine generates a light map and light probes that contain only the indirect lighting, which is projected back into a scene. Light maps and light probes are updated in real time in the Unity 5 editor and in-game on many of the platforms that are supported by Unity. These include Windows, OS X, iOS, and Android. In addition, Unity uses Enlighten to produce baked light maps.

When Enlighten calculates the indirect light in a scene, the results are stored in a light map or in light probes and they are applied in the shader code. The standard Unity materials make full use of Enlighten, including custom material shaders. The shaders and the Unity rendering engine ensure that the direct light is properly rendered, so that the final game looks like the final composition.

The following figure shows an example of a scene with only direct lighting, just Enlighten indirect lighting, and the effect of both direct and indirect lighting combined.



Figure 5-1 Arches

The scene in the left image is illuminated by direct lighting, with no bounce light and no added ambient light term. Any surfaces in shadow are black.

The middle image shows the same scene and lighting setup with only the Enlighten indirect light contribution and no albedo textures applied. This is the lighting output that Enlighten generates in real time.

The right image shows the final composition, which includes the direct lighting and Enlighten indirect lighting.

5.2 The structure of Enlighten

Enlighten consists of a number of components that are configurable in Unity. It is useful to understand the components and how they work together.

This section contains the following subsections:

- [5.2.1 Overview on page 5-62](#).
- [5.2.2 Precompute on page 5-62](#).
- [5.2.3 Real-time solver on page 5-69](#).
- [5.2.4 Baking light maps on page 5-70](#).

5.2.1 Overview

The components of Enlighten are:

The precompute

This precalculates the light transport in a scene. The precompute only depends on the static geometry and not the materials or light.

The real-time solver

This takes the precalculated data and combines it with the material and light information to produce light maps and probes in real time.

The light map baker

This produces baked light maps for direct and indirect light and also ambient occlusion. The light map baker relies on the precompute data.

The following sections describe the components in more detail.

5.2.2 Precompute

The precompute component computes data based on the geometry of a scene. At runtime, Enlighten uses the precompute data to generate the indirect lighting for any lighting setup in real time.

The precompute is started by pressing the Build button in the Lighting window, or when changing a scene when the auto build option is selected.

Once the precompute has started, the status of the precompute running tasks is shown in the progress bar at the bottom of Unity.

The following figure shows the progress bar and the current job.

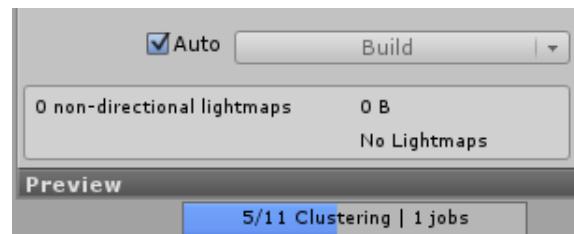


Figure 5-2 Progress

Note

The progress bar shows more precompute steps than are described in this section, but these are less important.

The precompute performs the following steps:

1. Packing.
2. Clustering
3. Computing the light transport.

To help understand the individual steps and the settings that affect them, consider the scene in the following figure. The scene consists of a floor, a wall with a corner, and a TV hanging on the wall.

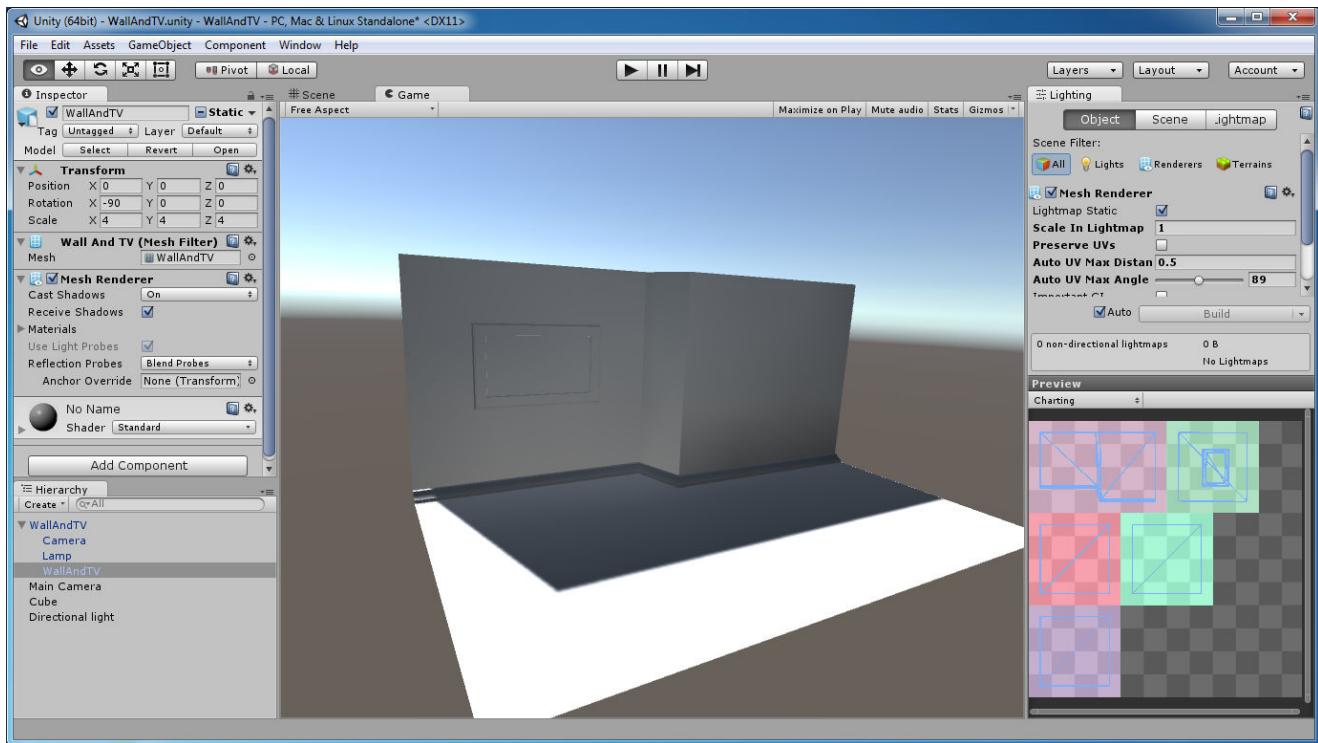


Figure 5-3 Wall and TV scene

Packing

Packing identifies charts in UVs then packs the charts into the light map that is used for the real-time lighting.

Packing uses existing UVs or the UVs Unity created for you when you select Generate Lightmap UVs in the Import settings.

Charts are groups of triangles in UVs that share vertices. After charts are identified, Enlighten repacks the UVs ensuring that there is no light leaking between charts and that UVs are packed as tightly as possible. You can see the results of the packing in the **UV Charts** render mode and also in the Charting preview. The following figures in this section show charts in the **UV Charts** render mode with the Charting preview shown in the bottom right of the window.

The pre-compute, the run-time memory usage, and the run-time performance of Enlighten all depend on the number of light map pixels. The more complex your mesh is, the more difficult it is to generate a UV unwrap with few charts. You can optimize the number of charts by using Enlighten to generate simplified UVs that remove smaller details in your mesh. Generating simplified UVs is enabled by default, but if you have carefully crafted your own UVs, you can disable this by ticking **Preserve UVs**. You can see the results in the **UV Charts** mode.

The following figure shows the **UV Charts** mode and the Charting preview with **Preserve UVs** enabled. Each chart is shown in a different color.

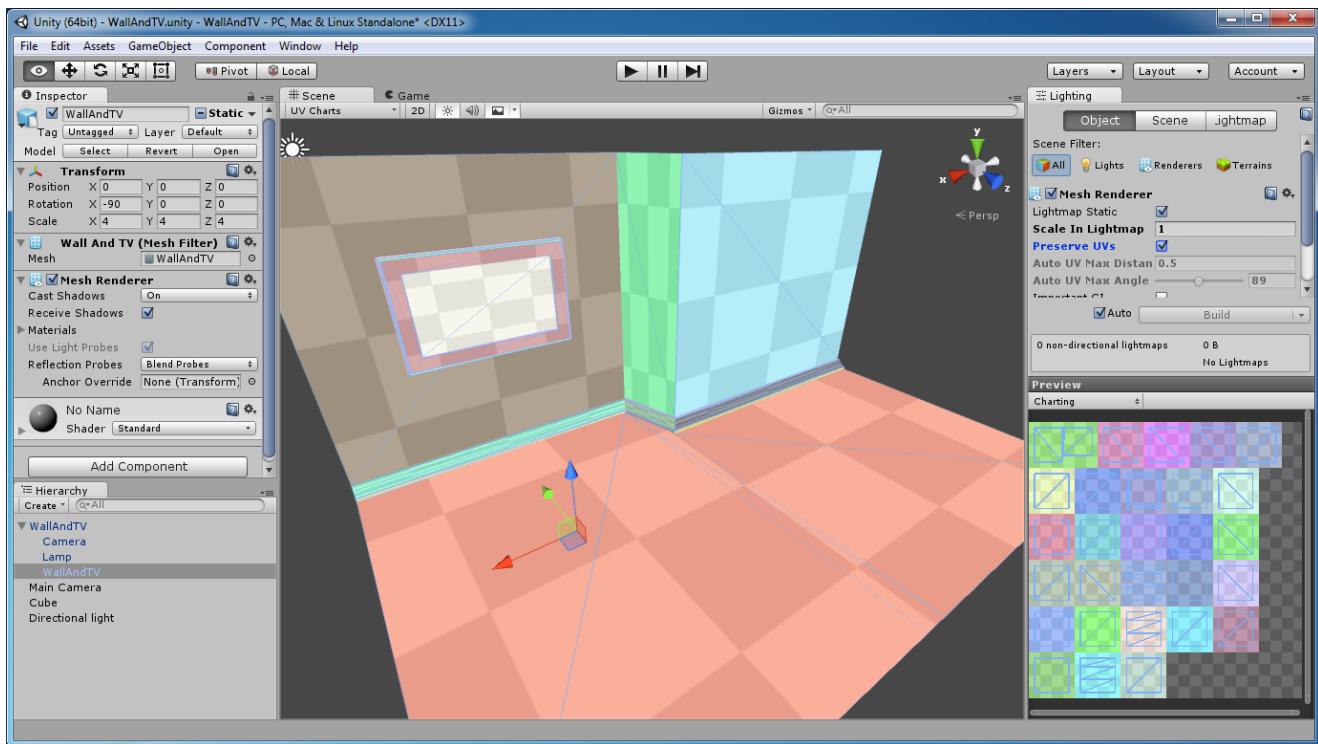


Figure 5-4 Wall and TV scene with Preserve UVs

As the image shows, parts of the TV, the wall, and the skirting board, all consist of multiple charts.

Only enable **Preserve UVs** if you have carefully authored UVs that you want to preserve in Unity. Usually it is best to not enable this flag.

The most important setting that influences the simplified UV generation is the **Auto UV Max Distance** setting. The **Auto UV Max Distance** tells Enlighten what details can be omitted from the illumination. Enlighten attempts to merge charts by creating a UV layout for the combined charts, but it does not split input charts. The new UV layout is created by projecting the vertices of all of the charts that are identified onto an optimal plane in world space. Only charts whose world space vertices are within the given distance above the plane are considered for merging.

When the value of **Auto UV Max Distance** is set to a small value, such as **0.01**, and **Preserve UVs** is disabled, the packing result looks like the following figure.

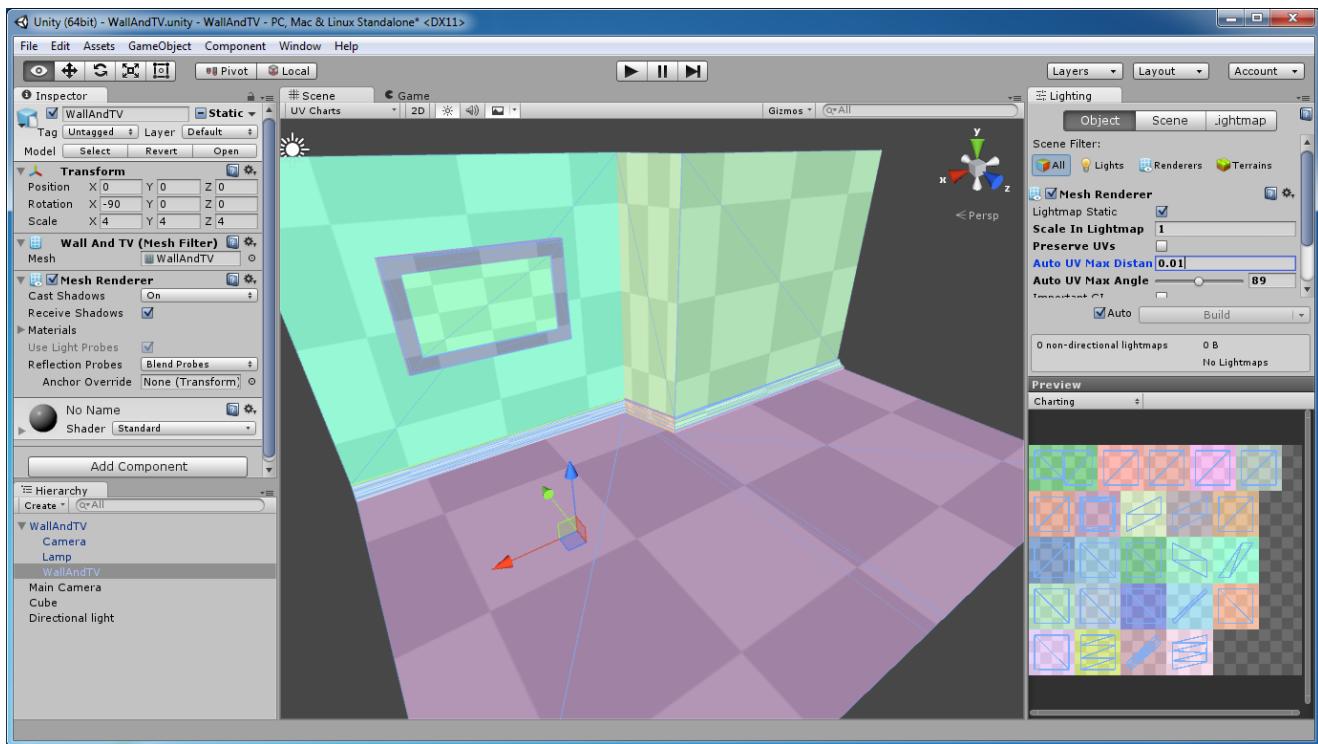


Figure 5-5 Wall and TV scene with Auto UV Max Distance at 0.01

Except for the colors of the visualization, the packing shown in [Figure 5-5 Wall and TV scene with Auto UV Max Distance at 0.01](#) on page 5-65 is similar to [Figure 5-4 Wall and TV scene with Preserve UVs](#) on page 5-64 which has **Preserve UVs** enabled. It has approximately the same number of charts, but the unwrapping is slightly different.

With the chosen light map resolution, only a few light map pixels are used for the wall and the floor. This is enough to capture the coarse global illumination of the scene. The TV and the skirting board are so small that they barely require one pixel. However, they are allocated far more pixels, because they consist of several charts, and charts have a minimum size of 4 x 4 pixels.

————— Note —————

If the value of **Min Chart Size** is set to 2, charts can have a smaller resolution of 2 x 2 pixels. However, this disables stitching and can lead to visible seams.

The following figure shows the result of setting the value of **Auto UV Max Distance** to 0.5.

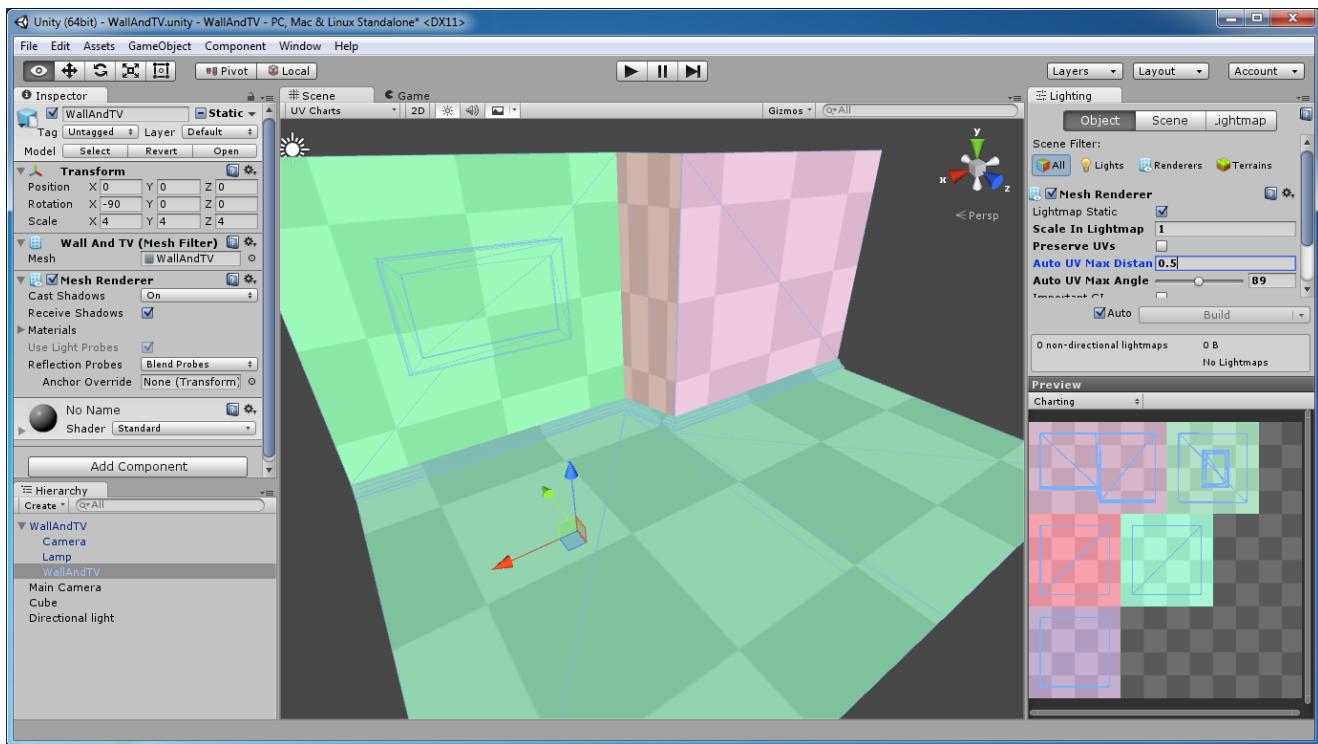


Figure 5-6 Wall and TV scene with Auto UV Max Distance at 0.5

When the value of **Auto UV Max Distance** set to 0.5, Enlighten projects the UVs of the TV and the skirting board onto the wall and floor respectively. As a result, they do not require any additional pixels in the light map. This is shown in the Charting preview. Increasing the value of **Auto UV Max Distance** results in more details being ignored.

The following figure shows the result of setting the value of **Auto UV Max Distance** to 2.

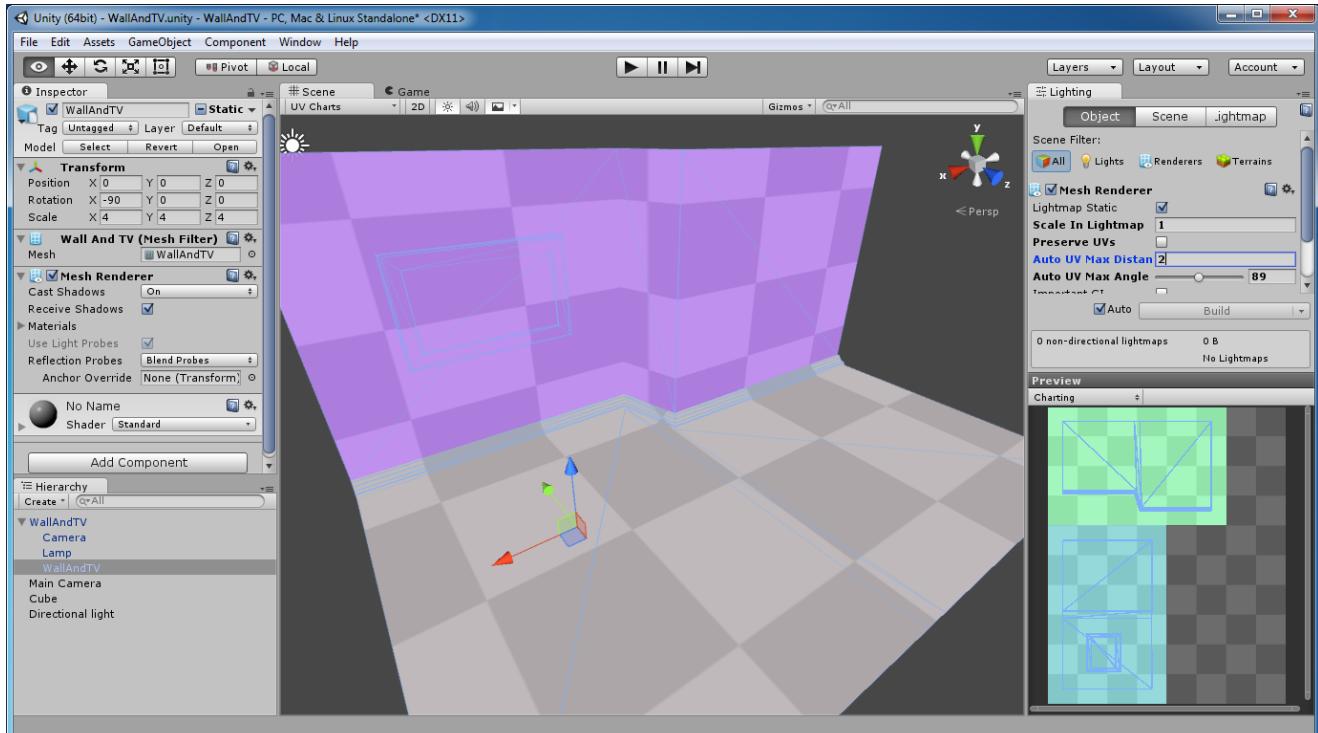


Figure 5-7 Wall and TV scene with Auto UV Max Distance at 2

When the value of **Auto UV Max Distance** set to 2, Enlighten considers the different parts of the wall, including the corner, to be one segment, with pixels spanning the different parts. As a result, the light map is smaller. Increasing the distance causes more geometry to be projected into one plane which can lead to a loss of detail.

Similar to Auto UV Max Distance, you can define what charts Enlighten considers for merging based on the angle between triangles. You can define the maximum angle by setting the value of **Auto UV Max Angle**. The default value of 89, in degrees, ensures that Enlighten does not merge charts that are set at a right angle to each other in world space.

Clustering

To model the light transport in a scene, Enlighten splits the scene geometry into clusters.

The clusters are separate from the light map pixels and generated only from the position and orientation of the triangles. Materials or UV coordinates are not used.

The following figure shows the clusters for the wall and TV example scene.

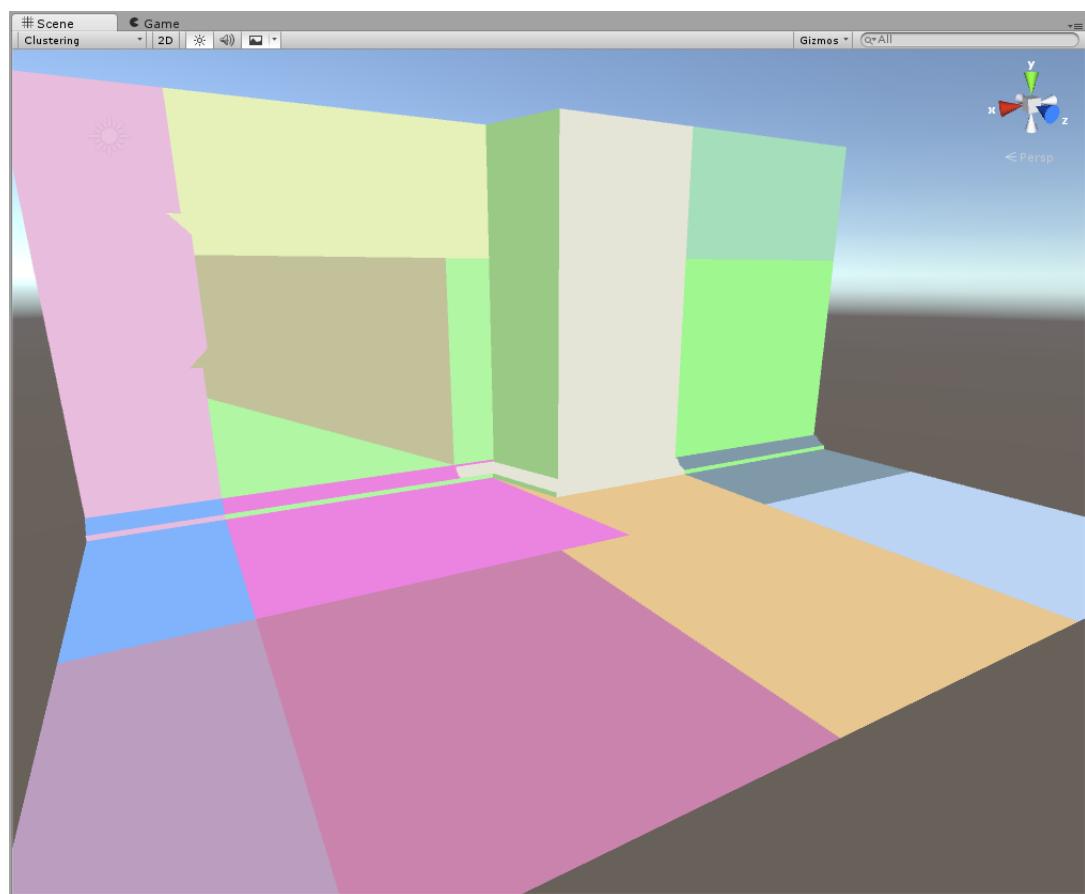


Figure 5-8 Wall and TV Clusters

Computing the light transport

Every cluster receives input lighting from lights that are placed in the scene. A cluster can also emit light.

For every pixel in a light map, and for every light probe, Enlighten casts rays in all directions of its hemisphere or in all directions for probes. It calculates the visibility of the cluster from the pixel or probe. Enlighten assumes that the scene only consists of diffuse surfaces, so that the visibility is directly proportional to the light that the pixel receives from the clusters. The visibility is often called form-factor in computer graphics. When the form-factors have been calculated, they are compressed to reduce memory usage and performance requirements.

The following figure shows a simplified view of this process using an example scene with a floor, wall, and cube. The figure shows a pixel X that sees the clusters A, B, C, and D with visibility values of 0.05, 0.1, and 0.2. At runtime, Enlighten can evaluate the lighting that this pixel receives by evaluating the following expression:

$$0.05*A+0.1*B+0.05*C+0.2*D$$

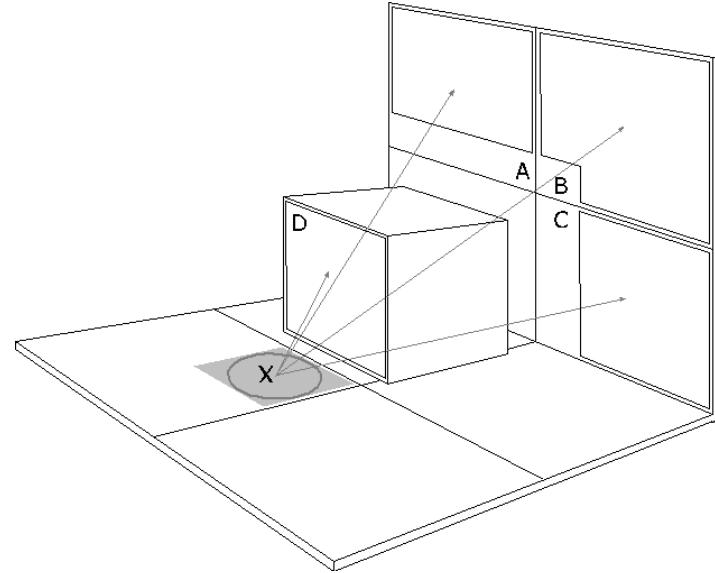


Figure 5-9 Form Factors

The following figure shows a lighting setup with a red light shining on cluster A, a green light on cluster B, a blue light on cluster C and no light onto cluster D.

Because the weight for cluster B is larger than the weight for cluster A and C, the final output value for pixel X has a green tint.

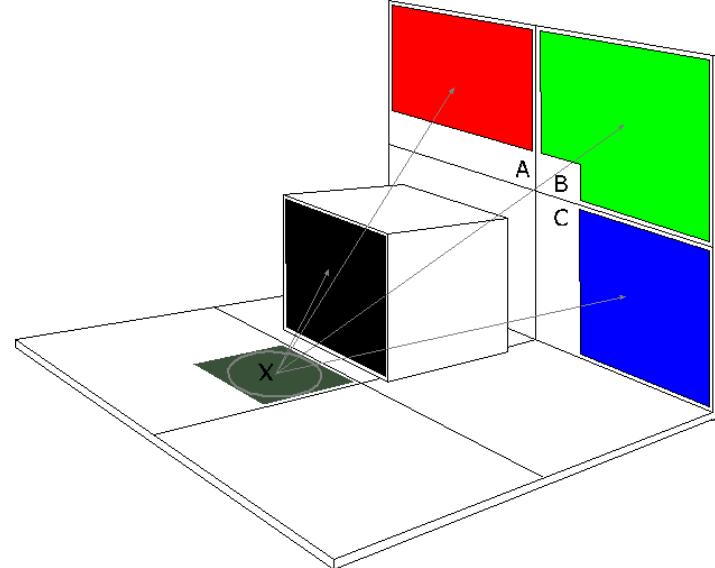


Figure 5-10 Form Factors with lighting

The precomputed data must be updated whenever the geometry is changed, because the clusters depend on the geometry of the scene. Unity can automatically trigger the precompute when it is required, or you can turn the continuous build off, and only trigger it when you are finished with the scene changes. You must precompute a scene before you can use lighting or bake lights.

Unity exposes the **Cluster Resolution** setting that controls the size of a cluster in relation to the light map resolution. The default value of 0.5 means that a cluster is twice the size of a light map pixel. There is usually no reason to change the **Cluster Resolution** setting.

The light transport calculates the form-factors. Enlighten has two important parameters that affect this step:

- The **Irradiance Budget** is the number of form-factors that Enlighten stores. Enlighten merges form-factors as much as required to stay within a given budget, but a too small budget produces oversimplified illumination. On the other hand, a larger budget consumes more runtime memory and the runtime solving cost is higher. Normally the default value of 128 is a good choice, but values as low as 64 can also give good results. Increasing this value does not increase the light transport time and therefore also not the precompute time.
- The **Irradiance Quality** is the number of rays to cast per pixel. Increasing this value results in longer precompute times. The default value of 8192 normally produces good results, but if you have chosen a high light map resolution, ARM recommends that you increase this value, because the indirect light might contain noise. This value does not affect runtime memory or performance.

To speed development, leave this option set to the default or a value that produces reasonable results when developing your game. Increase the value before shipping to get the maximum quality.

5.2.3

Real-time solver

The real-time solver component is present in the editor and in the game. It generates light maps and probes for direct feedback of indirect lighting.

The real-time solver consists of the following stages:

Input lighting stage

The input lighting stage calculates how lights illuminate the clusters, including dynamic lights. The light values are added together for all lights, and the final value is multiplied by the albedo of the cluster. This is based on the material of the geometry that the cluster is made of. Realistic rendering of lights, including input lighting, always requires shadows. Directional lights can correctly handle shadowing, that is, clusters that are in the shadow area of the light do not receive any light. Enlighten supports shadowed spot lights and point lights, but this is not implemented in Unity 5. Because of this, set an appropriate range for both spot and point lights to avoid leaking of indirect light.

Solve stage

The solve stage sums up the cluster value multiplied by the stored form-factors, and stores the results in the light map. The runtime performance is therefore directly linked to the number of form-factors that are stored per pixel or probe.

Bounce stage

The bounce stage reads back the values from the light maps and bounces light values back to the clusters. This way Enlighten simulates multiple light bounces. Light bounces are limited to the light map update rate in Enlighten. The light map update rate is typically determined by the render update rate.

Lights, materials, and the pre-compute data are inputs to this component. Both lights and materials can be changed at runtime. This enables you to fine-tune the lighting and the appearance of the scene with instant feedback.

You can control the number of threads to run the solver using the **CPU Usage** setting. The default value is Low. Increase this value if you want to increase the update rate of the indirect lighting and are not able to optimize your scene, for example, by setting appropriate simplified UV settings.

The following image shows dropdown menu for the **CPU Usage** setting.

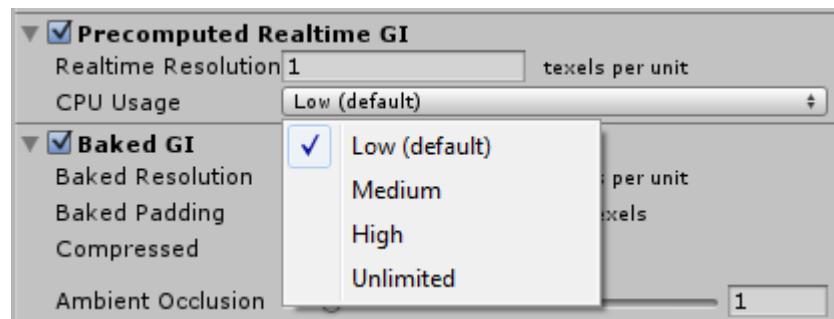


Figure 5-11 The CPU Usage setting

5.2.4 Baking light maps

Enlighten can generate baked light maps that contain direct lighting, indirect lighting, and ambient occlusion. The indirect light map is generated based on the real-time results and these are up-sampled and filtered to produce high-quality output.

You can also enable the final gather stage that uses the baked light maps as input, and Enlighten calculates another, final light bounce, with the precision bound by the bake pixel resolution. This step uses the standard baked light maps as input, so only consider switching to final gather when you are happy with the lighting in your scene.

Both the real-time and baked lighting can be combined. In Unity, this is exposed on a per-light setting. Both types of lighting can blend together to contribute to the indirect lighting.

The following figure shows light map baking settings:

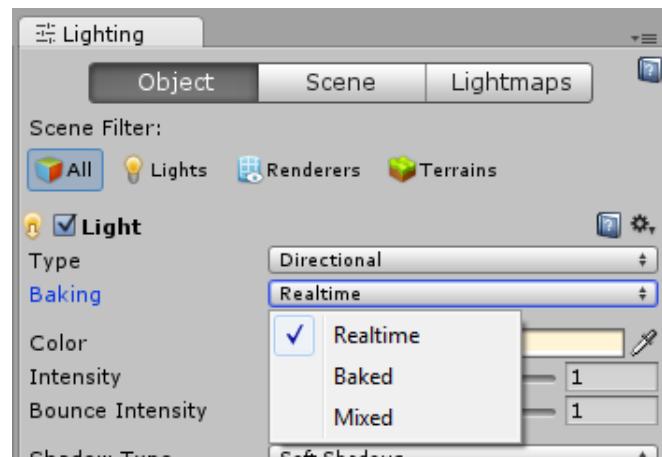


Figure 5-12 Light map baking settings

5.3

Setting up a scene with Enlighten

This section provides some general tips on how best to set up a scene with Enlighten in Unity.

Turning off lightmap baking

For real-time feedback on your lighting, and to enable a faster iteration time, set up your scene as if you are only using real-time global illumination. Do this even if you only plan on using baked light maps. An advantage of setting up a scene for only real-time global illumination is that for more powerful target platforms you can enable real-time lighting, whereas for less powerful target platforms you can use baked light maps.

When you are happy with your lighting, you can turn on light map baking for all the required lights. Enlighten is responsible for both the real-time lighting and the light map baking, so the baked light maps match the real-time rendered results, with differences only visible in soft shadows and area lights.

Setting the Enlighten light map resolution

The most important metric to optimize for is the Enlighten light map resolution. In a simulated real-world scene, with human sized characters, you typically set the texture resolution to be one pixel per meter. Smaller scale details are best handled with screen space ambient occlusion.

Optimizing the precompute

The **UV Charts** visualization mode is usually the best view mode to use for setting up a scene for Enlighten in Unity. Meshes are only rendered in this mode when a pre-compute is started. However, the packing is the first stage in the precompute and Unity updates the results when a stage is finished.

The packing stage is quick to compute. If it takes a long time for meshes to appear in the **UV Charts** render mode, then the resolution might be too high. To test whether the resolution is too high, decrease the resolution and then rerun the precompute. The precompute is triggered automatically, unless you have deactivated it.

When you have adjusted the resolution to your requirements, wait for the pre-compute to finish. If you notice lighting artifacts such as light leaking, use the **UV Charts** render mode again. Enlighten does not split input charts, and charts that go through a wall can leak. A typical example is a floor mesh that spans multiple rooms and only has one chart. In such cases, consider creating smaller charts by splitting and separating the input UVs.

Adjusting the scene elements

When the precompute has completed, you can add lights and get instant feedback about the overall scene appearance. This is not limited to light sources and their position.

You can also change material properties, such as surface color, texture, or emissive settings.

With Enlighten, any surfaces can be set up to emit light and can therefore be turned into an area light. These area lights have the benefit of having no associated render cost because all their lighting is done by Enlighten. This can be useful for lower-end mobile devices, where the number of dynamic light sources is very limited. Use the **Important GI** checkbox in the object panel for emissive surfaces, especially if they are small, because this ensures that the clusters for these surfaces are not merged with other clusters.

Lighting small objects with light probes

Consider lighting smaller objects in your scene using light probes instead of light maps. To use light probes, do not click the **Lightmap static** box.

Typically, smaller objects do not contribute much to the global illumination and setting them to be probe lit removes them from the pre-compute stage, making the pre-compute faster.

Smaller objects are often also more difficult to generate UVs for. You can also merge the meshes of smaller objects with larger objects, such as the TV and the wall in the example that is used throughout [5.2.2 Precompute](#) on page [5-62](#).

5.4 Example: Enlighten setup of the Ice Cave

Unity has several global illumination visualization modes. This section describes these modes with the Ice Cave as an example scene.

The following figure shows the different visualization modes that are listed in the Unity Global Illumination visualization menu.

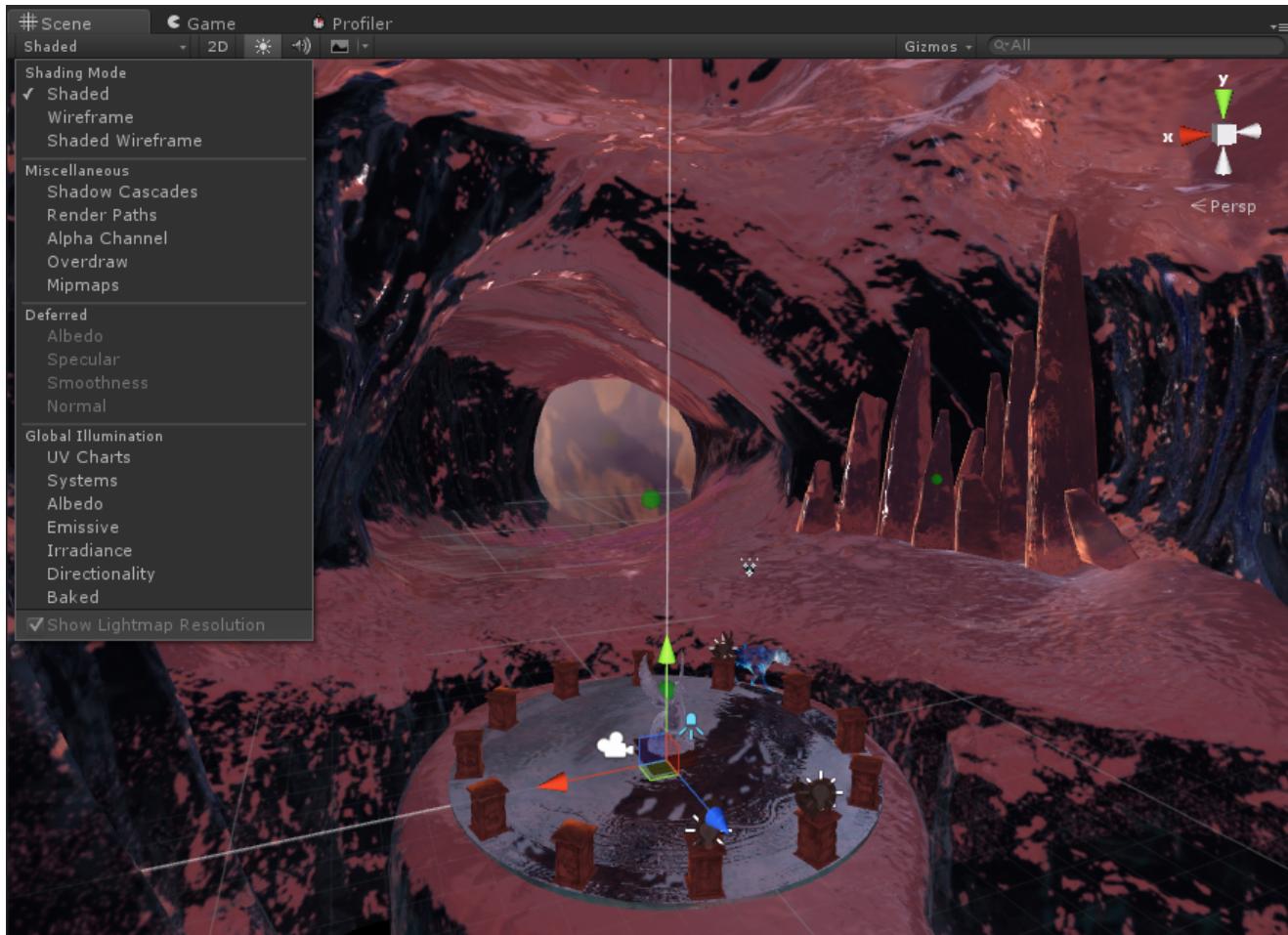


Figure 5-13 Global illumination visualization menu

UV Charts view

The **UV Charts** view is recommended for verifying both your pixel size and UV charts. Compared to the earlier example scenes, the more complex Ice Cave scene shows meshes in the **UV Charts** with multiple charts and sharp boundaries. See, for example, the cave ceiling. To avoid any visible seams in the final rendered image, Enlighten attempts to stitch these boundaries automatically.

Note

The result of the stitching can be seen in the **Irradiance** and **Directionality** views. However, in the final image the seams are much less visible because of the applied material textures. Therefore, do not overestimate visible seams in these two views. Ideally, you can hide the seams by moving them to less visible parts of the mesh.

The following figure shows Ice Cave UV parts in the **UV charts** view.

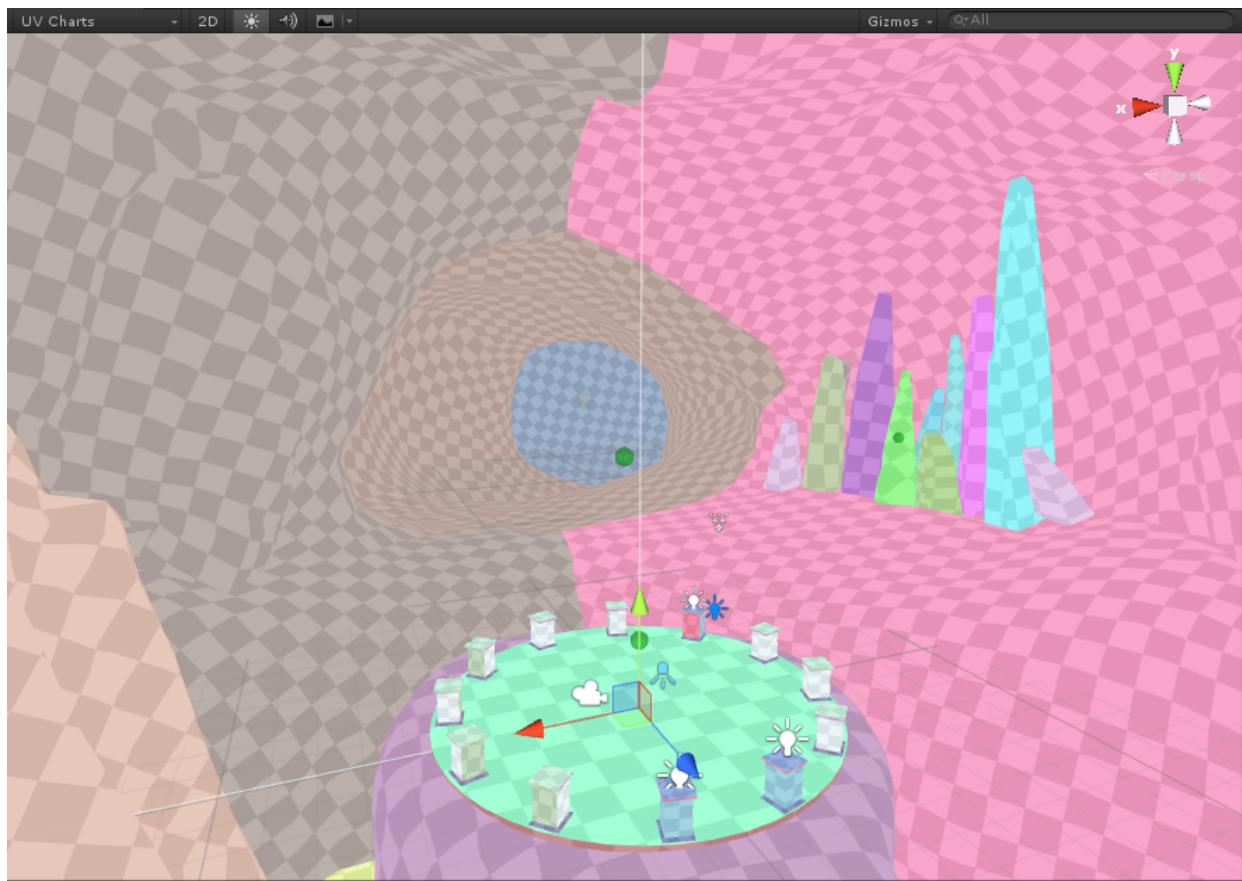


Figure 5-14 Ice Cave UV Charts

Systems view

The **Systems** view shows the systems that are automatically generated by Enlighten. Multiple systems enable the precompute and the runtime to be executed in parallel. In general, the automatic generation works fine without the need to change any settings.

In the following figure, each color represents a system. Each system can contain multiple objects.

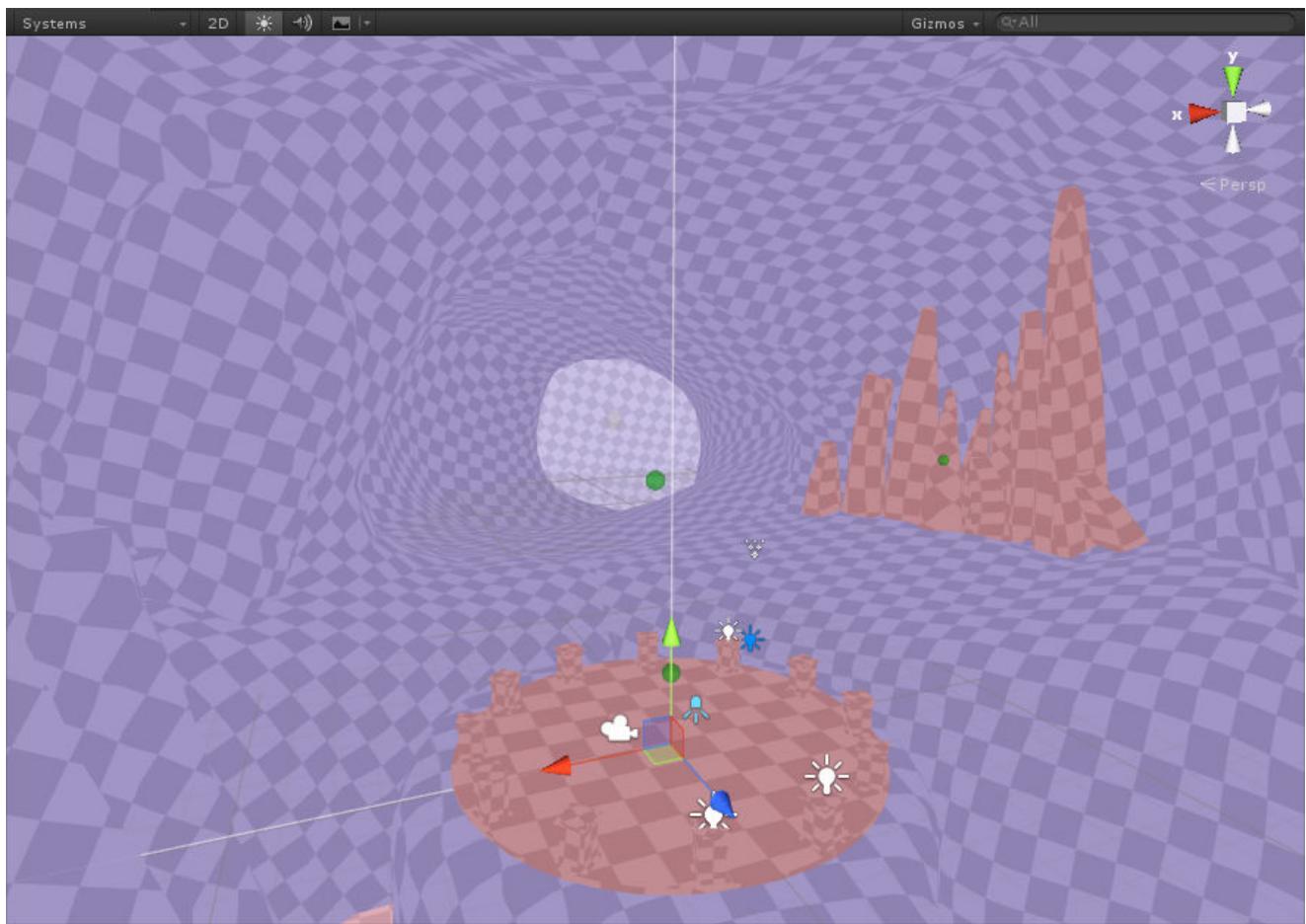


Figure 5-15 Ice Cave Systems

Albedo view

The **Albedo** view shows the light map pixels that are rendered with the diffuse color that Enlighten is using for light bounced or emitted in the scene. It effectively shows at what detail Enlighten uses the color information from your scene.

The following figure shows that the albedo is an average of the albedo map and the albedo color.

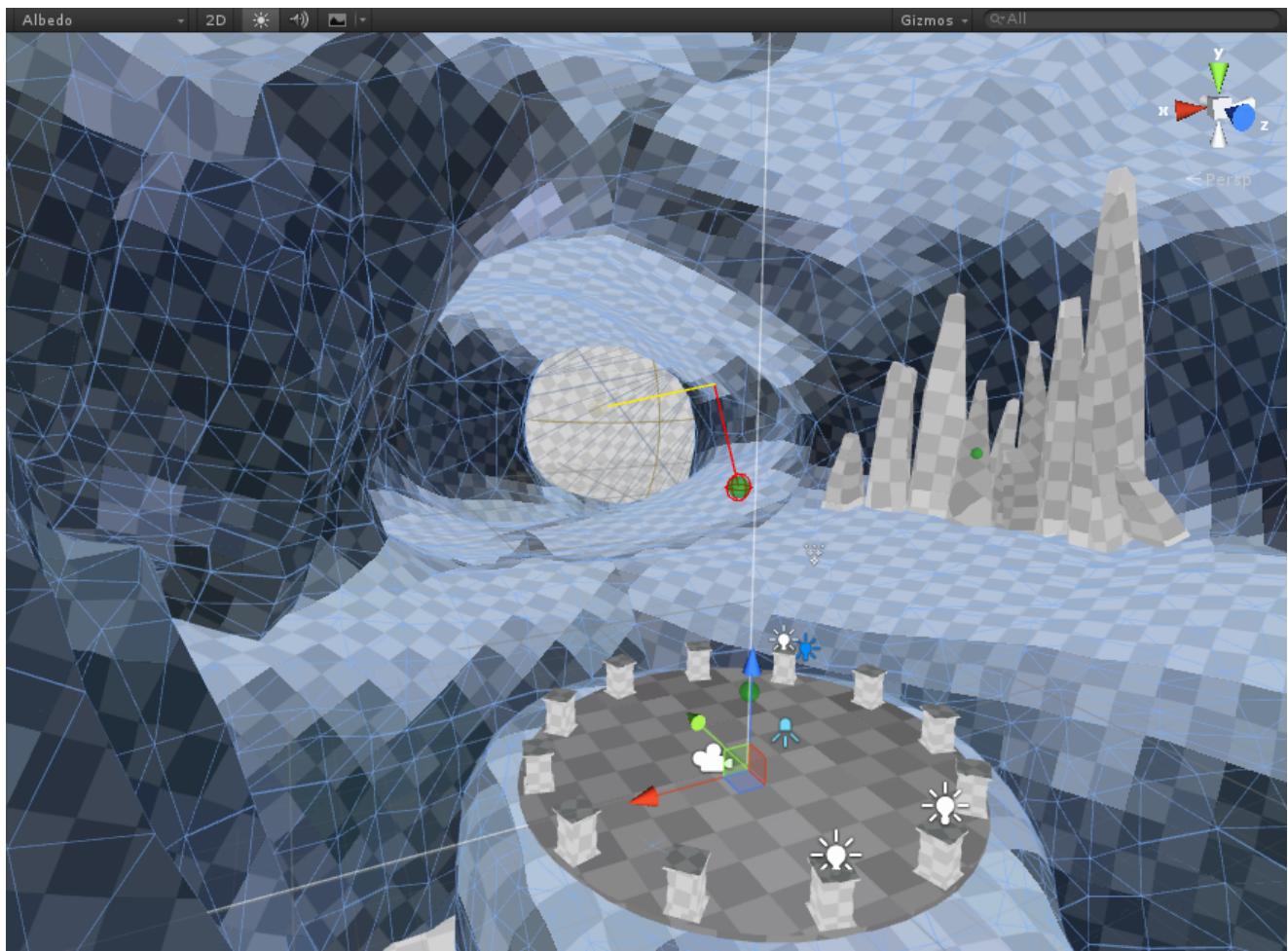


Figure 5-16 Ice Cave Albedo

Emissive view

The **Emissive** view is similar to the albedo view, but it shows the color and intensity of the emitted light from emissive objects. You can use emissive objects with Enlighten to create area lights without any rendering cost.

The following figure shows that the emissive value is the average of the emissive map and emissive color. For example, the blue crystals and green platform are emissive.

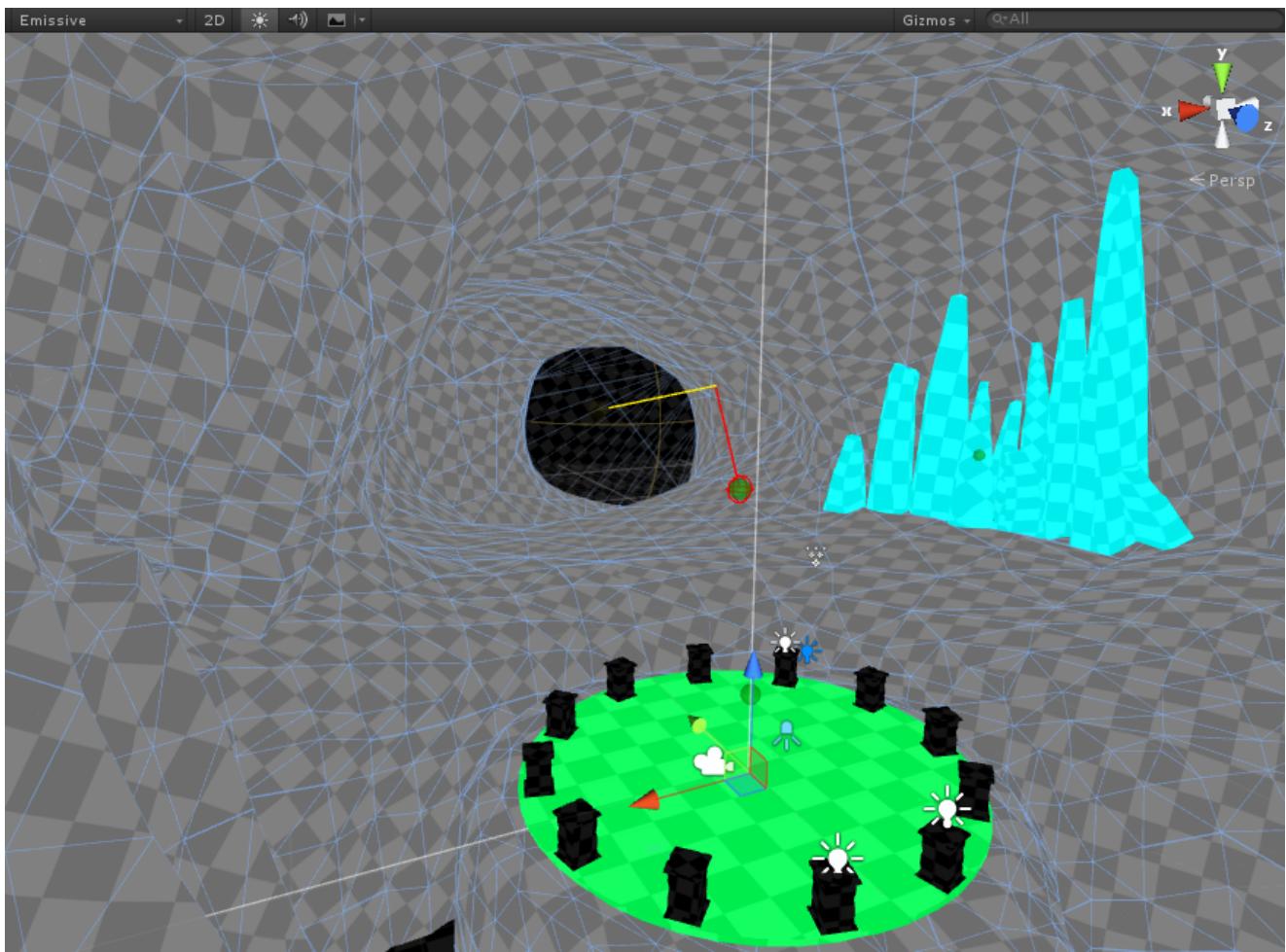


Figure 5-17 Ice Cave Emissive

Irradiance view

The **Irradiance** view only shows the irradiance that is received by surfaces using light maps. This view helps to identify lighting or scene setup issues, as often the standard render view, with textures applied, makes it difficult to identify the exact problem.

The following figure shows the **Irradiance** view with light coming from the fireflies and a light placed behind the crystals.

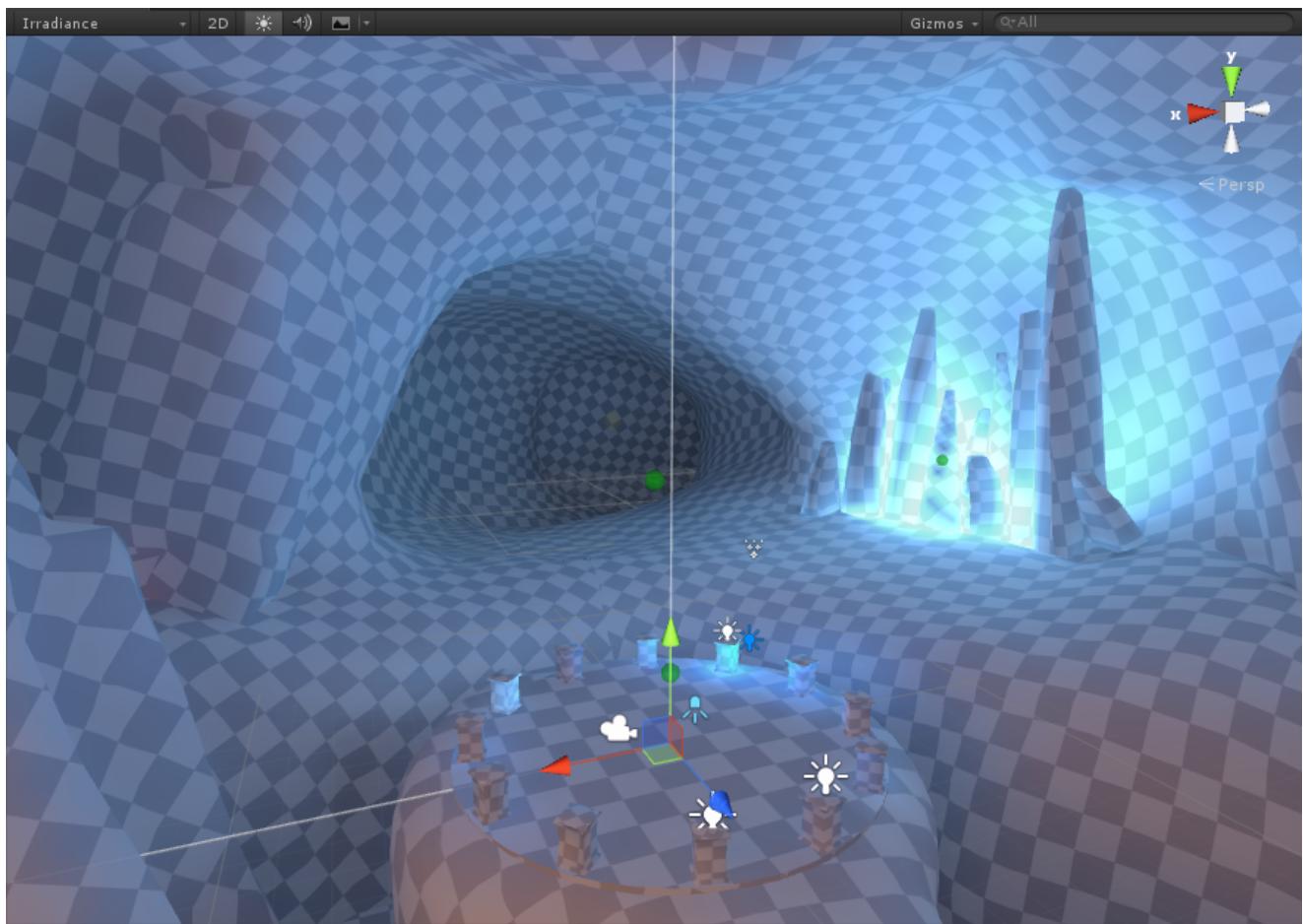


Figure 5-18 Ice Cave Irradiance

Directionality view

The **Directionality** view shows the dominant light direction for each pixel. This is used if the Directional or Directional Specular options are selected. Like the Irradiance view, changing the position of lights makes this map change.

In the following figure, the bottom of the scene mainly receives light from the sun which is above the cave. The green color indicates that this is positive on the Y axis. The directionality changes near the crystals because the main light is positioned behind them.

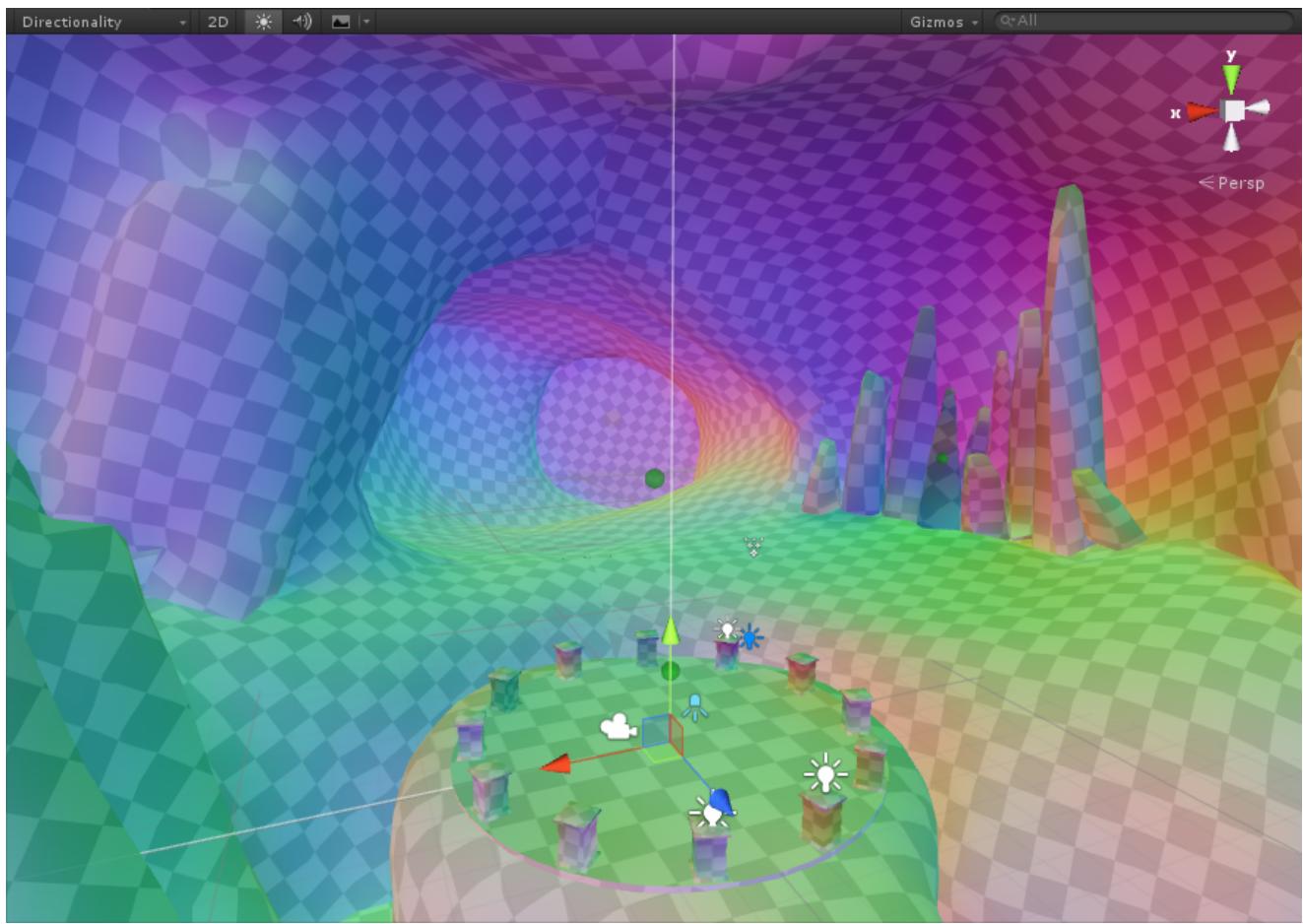


Figure 5-19 Ice Cave Directionality

You can also visualize these maps as a 2D texture for each object. Select an object, and in the Lighting tab select the Object switch. In the Preview window at the bottom of the tab, you can see that the same charts the object refers to are selected.

5.5 Using Enlighten in custom shaders

This section describes using Enlighten in custom shaders. This is an advanced topic and is only required if you want to completely replace the shaders and still make use of the Enlighten light maps and light probes output.

This section contains the following subsections:

- [5.5.1 Global illumination code flow on page 5-80](#).
- [5.5.2 Custom shader code on page 5-81](#).
- [5.5.3 *Enlighten.cginc* on page 5-82](#).

5.5.1 Global illumination code flow

This section describes the flow of the code that is executed by the forward renderer of the Standard Shader.

Vertex shader

In the vertex shader, you must set up the texture coordinates for the lightmap so that you can access the global illumination data.

Set up the texture co-ordinates using the `vertForwardBase()` function that is located in the file `UnityStandardCore.cginc`. This function computes the texture coordinates for static and dynamic lightmaps, and a per-vertex ambient color that is based on the spherical harmonics.

Fragment shader

The main function that computes global illumination is `UnityGlobalIllumination()`. `UnityGlobalIllumination()` is located in the file `UnityGlobalIllumination.cginc`.

The following code shows the sequence of calls for the Unity Standard Shader, rendered with forward rendering:

```
Pass Forward:Standard.shader ->
    fragForwardBase (UnityStandardCore.cginc) ->
        FragmentGI (UnityStandardCore.cginc)->
            UnityGlobalIllumination (UnityGlobalIllumination.cginc)
```

`fragForwardBase()` is the starting point for understanding how the elements are combined.

The function `fragForwardBase()` initializes the `FragmentCommonData` structure. This structure includes the correct values for normal, position, eye vectors, and other values. These values are based on the settings of the materials such as parallax maps and normal maps.

The `fragForwardBase()` function computes a shadow attenuation factor and an occlusion factor that is based on the shadows and occlusion maps. The value that is used to calculate this depends on the shader model of the building platform.

The `fragForwardBase()` function accesses the global illumination data using the following function:

```
UnityGI gi = FragmentGI (s.posWorld, occlusion, i.ambientOrLightmapUV, atten,
s.oneMinusRoughness, s.normalWorld, s.eyeVec, mainLight);
```

The `FragmentGI()` function initializes:

- The `UnityGIInput` structure with the information passed as parameters.
- The UV coordinates to access the lightmaps. Alternatively, if lightmaps are not used, it initializes the ambient color.
- The data that is required to access the reflection cubemaps.

The `FragmentGI()` function calls the `UnityGlobalIllumination()` function.

The `UnityGlobalIllumination()` function checks the current object, and if necessary, it samples the indirect lighting from the light probes.

The probes store data as spherical harmonics. The Shader Model defines the maximum spherical harmonic order used. Higher orders are better quality but they are more computationally expensive.

Probes are typically used to light dynamic objects with indirect light.

The result is stored as an indirect diffuse color. The static or dynamic lightmaps are ignored for objects that sample from light probes.

If the object does not use any lightmaps, then the code stores the color of the direct light that is received by the object.

If the object uses static lightmaps, they are sampled. The type of light maps depends on the type of global illumination:

- Non-directional.
- Directional.
- Directional Specular.

These types of global illumination are also available for dynamic lightmaps. These lightmaps are updated at runtime by Enlighten.

Both static and dynamic lightmaps modify the value of the indirect diffuse color that is returned by the function.

The `UnityGlobalIllumination()` function samples the reflections from a single specular cubemap, or interpolates the reflections from two specular cubemaps, to set the indirect specular color.

The `fragForwardBase` code calculates the color at the pixel location using `UNITY_BRDF_PBS`.

This macro is defined by various functions that depend on the shader model in use.

The macro computes the color of the pixel based on the material property, the direct light, and the indirect light that is computed in the `fragForwardBase()` function.

The `UNITY_BRDF_GI()` function computes the contribution of the light coming from the static and the dynamic lightmaps.

The `UNITY_BRDF_GI()` function reads the `gi.light2` and `gi.light3` fields of the `UnityGI` structure. These fields contain information about the position and intensity of the lights that are retrieved from the lightmaps for Directional and Directional Specular global illumination.

The `fragForwardBase()` function applies the Emission factor to the pixel.

The `fragForwardBase()` function applies the fog color.

5.5.2 Custom shader code

You can modify the Unity Standard Shader to include directional global illumination.

This section describes an example of the Unity Standard Shader that is modified to include directional global illumination. The example shader does not include specular lighting or global illumination of dynamic objects using spherical harmonics. It is based on Unity 5.0.

Enlighten texture coordinates

Enlighten and Unity generate an extra set of UV coordinates for all objects that are marked as static, to access the lightmaps they update.

It passes these coordinates as `TEXCOORD1` for the static lightmaps, and `TEXCOORD2` for the dynamic lightmaps.

You must copy these from the vertex shader to the fragment shader.

The following code shows how to add this input to your custom shader:

```
struct vertexInput {
    half4 vertex : POSITION;
    half2 texcoord : TEXCOORD0;
    half2 uv1 : TEXCOORD1; //Static Lightmap texture coordinates
    half2 uv2 : TEXCOORD2; //Dynamic Lightmap texture coordinates
    half3 normal : NORMAL;
};
```

Use the correct names for the variables

Enlighten constructs the albedo and emissive maps that are used at run-time. When Enlighten is doing this, it accesses the attributes of the materials used.

These attributes must be named correctly. If not, the maps are constructed with incorrect colors.

Ensure the names of your custom shader attributes match the names in the Standard Shader.

To download the Standard Shader file, go to: <https://unity3d.com/get-unity/download/archive>.

Download the appropriate *Built-in shaders* file for your OS and version of Unity.

The Standard Shader is called *Standard.shader*.

The following code shows the name of the texture as declared in the standard shader. This texture name is used by Enlighten to update its lightmaps with the materials albedo information:

```
Color("Color", Color) = (1, 1, 1, 1)
_MainTex("Albedo", 2D) = "white" {}
```

Insert the correct macros

At run-time, Unity defines `multi_compile` macros. These macros specify if the object uses Enlighten and how to sample the data that is generated. For example, they can specify that dynamic objects access the light probes instead of the enlighten data.

The following code shows the macro definitions:

```
//Those multi_compile options are set automatically by Unity.
#pragma multi_compile LIGHTMAP_OFF LIGHTMAP_ON
#pragma multi_compile DIRLIGHTMAP_OFF DIRLIGHTMAP_COMBINED DIRLIGHTMAP_SEPARATE
#pragma multi_compile DYNAMICLIGHTMAP_OFF DYNAMICLIGHTMAP_ON
```

Vertex shader

The vertex shader uses a utility function to retrieve data:

```
output.enlightenData = GetEnlightenInput(output.vertexInWorld, output.normalInWorld,
    input.uv1, input.uv2);
```

This function is described in the code, see [5.5.3 Enlighten.cginc on page 5-82](#).

Fragment shader

The fragment shader samples the indirect lighting contribution and uses it as ambient color:

```
//If the shader uses enlighten, the ambient term is embedded in the lightmap. Do not
//add it again
ambient = texColorTerm * SampleIndirectContribution(input.enlightenData,
normalInWorld, 1.0);
```

5.5.3 Enlighten.cginc

Enlighten.cginc is a utility file that integrates global illumination with custom shaders.

The following code defines the `GetEnlightenInput()` and `SampleIndirectContribution()` functions that you can use in your shaders:

```
#ifndef ENLIGHTEN_INCLUDED
#define ENLIGHTEN_INCLUDED

#include "UnityShaderVariables.cginc"
#include "UnityStandardBRDF.cginc"
#include "UnityPBSLighting.cginc"

struct enlightenInput
{
    half4 posWorld;
    half4 ambientOrLightmapUV;
    half3 eyeVec;
    half3 normalWorld;
};

enlightenInput GetEnlightenInput(half4 posWorld, half3 normalWorld, half2 uvCoord1,
half2 uvCoord2)
{
    enlightenInput o;

    // Static lightmap texture coordinates
    o.ambientOrLightmapUV.xy = uvCoord1.xy * unity_LightmapST.xy + unity_LightmapST.zw;

    // Dynamic lightmap texture coordinates
    o.ambientOrLightmapUV.zw = uvCoord2.xy * unity_DynamicLightmapST.xy +
unity_DynamicLightmapST.zw;

    o.posWorld = posWorld;
    o.eyeVec = normalize(posWorld.xyz - _WorldSpaceCameraPos);
    o.normalWorld = normalWorld;

    return o;
}

half3 SampleIndirectContribution(enlightenInput i, half3 normalWorld, half
shadowAttenuation)
{
    half3 output = half3(0.0,0.0,1.0);
    #ifdef ENLIGHTEN_ON
    //First Light

    // Do not consider any attenuation or occlusion
    half atten = 0.0;
    half occlusion = 1.0;
    UnityGI gi = SimplifiedFragmentGI (i.posWorld, occlusion, i.ambientOrLightmapUV,
atten, normalWorld, i.eyeVec);

    return half3(gi.indirect.diffuse);
    #elif defined (ENLIGHTEN_OFF)
    return output;
    #else
        #error You must include "#pragma multi_compile ENLIGHTEN_OFF ENLIGHTENS_ON" in your
shader to allow projectors to be enabled and disabled
    #endif
    }
#endif
```

`SimplifiedFragmentGI()` is a simplified version of the `FragmentGI()` function inside the `UnityStandardCore.cginc` file.

This version does not use the data that is required for the cubemaps and the light probes, so those parts have been removed. Within this function, the call to `UnityStandardGlobalIllumination()` is modified to point to a simplified version of it. `UnityStandardGlobalIllumination()` is contained in `UnityPBSLighting.cginc` and contains the code for sampling the various lightmaps, light probes, and cubemaps.

Chapter 6

Advanced graphics techniques

This chapter lists a number of advanced graphics techniques that you can use.

It contains the following sections:

- [6.1 Custom shaders](#) on page 6-85.
- [6.2 Implementing reflections with a local cubemap](#) on page 6-98.
- [6.3 Combining reflections](#) on page 6-114.
- [6.4 Dynamic soft shadows based on local cubemaps](#) on page 6-120.
- [6.5 Refraction based on local cubemaps](#) on page 6-128.
- [6.6 Specular effects in the Ice Cave demo](#) on page 6-134.
- [6.7 Using Early-z](#) on page 6-137.
- [6.8 Dirty lens effect](#) on page 6-138.
- [6.9 Light shafts](#) on page 6-141.
- [6.10 Fog effects](#) on page 6-145.
- [6.11 Bloom](#) on page 6-152.
- [6.12 Icy wall effect](#) on page 6-159.
- [6.13 Procedural skybox](#) on page 6-165.
- [6.14 Fireflies](#) on page 6-173.
- [6.15 The tangent space to world space normal conversion tool](#) on page 6-177.

6.1 Custom shaders

This section describes custom shaders.

This section contains the following subsections:

- [6.1.1 About custom shaders on page 6-85](#).
- [6.1.2 Shader structure on page 6-86](#).
- [6.1.3 Compilation Directives on page 6-88](#).
- [6.1.4 Includes on page 6-88](#).
- [6.1.5 The OpenGL ES 3.0 graphics pipeline on page 6-89](#).
- [6.1.6 Vertex shaders on page 6-90](#).
- [6.1.7 Vertex shader input on page 6-91](#).
- [6.1.8 Vertex shader output and varyings on page 6-91](#).
- [6.1.9 Fragment Shaders on page 6-93](#).
- [6.1.10 Providing data to shaders on page 6-93](#).
- [6.1.11 Debugging shaders in Unity on page 6-95](#).

6.1.1 About custom shaders

Unity 5 and higher includes a Physically Based Shading (*PBS*) model that simulates the interactions between material and light. This provides a high level of realism and makes it possible to achieve a consistent look under different lighting conditions.

You can easily use PBS with the standard shader. If you create your own material, it is automatically assigned the standard shader.

You can easily access the Standard Shader. If you create your own material the Standard Shader is assigned to it. There are a number of other built-in shaders that are very useful for beginners. You can see all the available built-in shaders divided into families by clicking on the **shader** drop down menu in the **Inspector**.

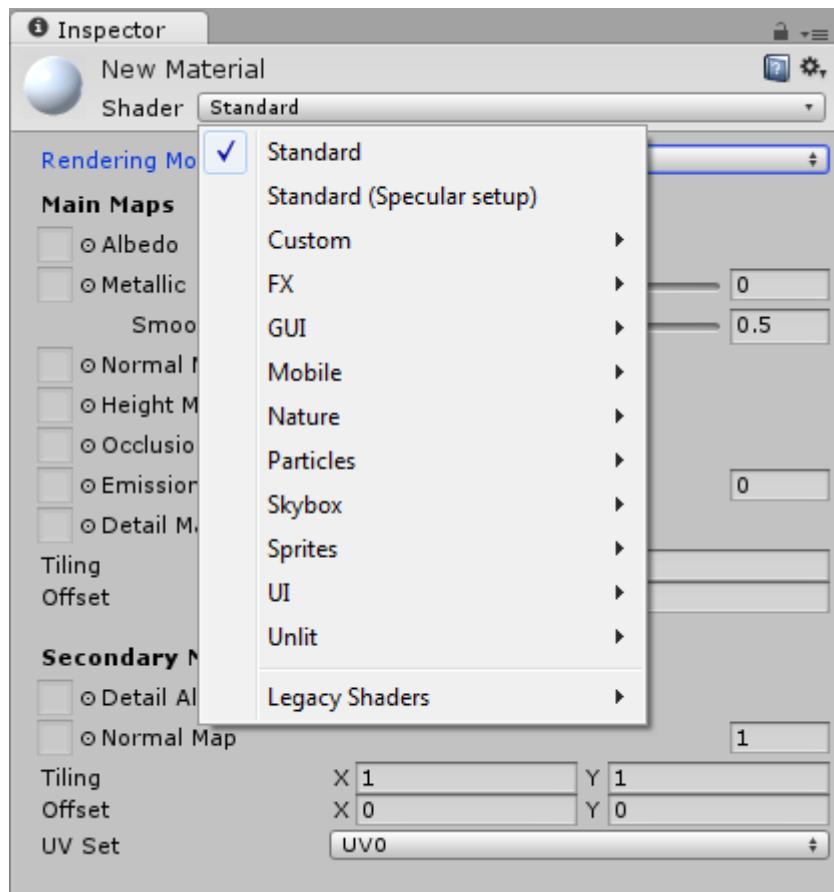


Figure 6-1 Unity built-in shaders

The source code of built-in shaders is available in the Unity download archive <http://unity3d.com/> that contains more than 120 shaders. You can learn a lot from reading and trying to understand the code of these shaders.

In addition to these, there are many effects that cannot be achieved by using existing shaders. For example, shaders that implement reflections based on local cubemaps. For more information see [6.2 Implementing reflections with a local cubemap](#) on page 6-98.

In Unity there are typically two ways of writing shaders:

Surface shaders

These are commonly used when shaders are affected by lights and shadows. Unity does the work related to the lighting model for you, enabling you to write more compact shaders.

Vertex and fragment shaders

These are the most flexible shaders but you must implement everything. The Unity ShaderLab does more than vertex and fragment shaders but these are in the main programmable part of the graphics pipeline where all the shading is done so it is important to know how to write custom vertex and fragment shaders.

6.1.2 Shader structure

The following code shows a very simple vertex and fragment shader that contains most of the elements required in vertex or fragment shader.

The shader example is written in Cg. Unity also supports the HLSL language for shader snippets.

```
Shader "Custom/ctTextured"
{
    Properties
```

```

        _AmbientColor ("Ambient Color", Color) = (0.2,0.2,0.2,1.0)
        _MainTex ("Base (RGB)", 2D) = "white" {}
    }

SubShader
{
    Pass
    {
        CGPROGRAM
        #pragma target 3.0
        #pragma glsl
        #pragma vertex vert
        #pragma fragment frag

        #include "UnityCG.cginc"

        // User-specified uniforms
        uniform float4 _AmbientColor;
        uniform sampler2D _MainTex;

        struct vertexInput
        {
            float4 vertex : POSITION;
            float4 texCoord : TEXCOORD0;
        };
        struct vertexOutput
        {
            float4 pos : SV_POSITION;
            float4 tex : TEXCOORD0;
        };

        // Vertex shader.
        vertexOutput vert(vertexInput input)
        {
            vertexOutput output;

            output.tex = input.texCoord;
            output.pos = mul(UNITY_MATRIX_MVP, input.vertex);
            return output;
        }

        // Fragment shader.
        float4 frag(vertexOutput input) : COLOR
        {
            float4 texColor = tex2D(_MainTex, float2(input.tex));
            return _AmbientColor + texColor;
        }

        ENDCG
    }
}
Fallback "Diffuse"
}

```

The first key word is **Shader** followed by the *path/name* of the shader. The path defines the category where the shader is displayed in the drop down menu when you are setting a material. The shader from the example is displayed under the category of **Custom** shaders in the drop down menu.

The **Properties{}** block lists the shader parameters that are visible in the inspector and what parameters you can interact with.

Each shader in Unity consists of a list of **subshaders**. When Unity renders a mesh, it looks for the shader to use, and selects the first **subshader** that can run on the graphics card. This way shaders are executed correctly on different graphics cards that support different shader models. This feature is important because GPU hardware and APIs are constantly evolving. For example, you can write your main shader targeting a Mali Midgard GPU to make use of the latest features of OpenGL ES 3.0, while in a separate subshader, write a replacement shader for graphics cards supporting OpenGL ES 2.0 and below.

The **Pass** block causes the geometry of an object to be rendered one time. A shader can contain one or more passes. You can use multiple passes on old hardware, or to achieve special effects.

If Unity cannot find a **subshader** in the body of the shader that can render the geometry correctly it rolls back to another shader defined after the **Fallback** statement. In the example this is the Diffuse built-in shader.

Cg program snippets are written between CGPROGRAM and ENDCG.

6.1.3 Compilation Directives

You pass compilation directives as `#pragma` statements. The compilation directives indicate the shader functions to be compiled.

Each compilation directive must contain at least the directives to compile the vertex and the fragment shader: `#pragma vertex name`, `#pragma fragment name`.

By default, Unity compiles shaders into shader model 2.0. The directive `#pragma target` enables shaders to be compiled into other capability levels. If the shader becomes large you get an error of the following type:

```
Shader error in 'Custom/MyShader': Arithmetic instruction limit of 64 exceeded; 83 arithmetic instructions needed to compile program;
```

If this is the case you must change from shader model 2.0 to shader model 3.0 by adding the `#pragma target 3.0` statement. Shader model 3.0 has a much higher instruction limit.

If you pass several varyings from vertex shader to fragment shader you might get the following error:

```
Shader error in 'Custom/MyShader': Too many interpolators used (maybe you want #pragma glsl?) at line 75.
```

If this is the case add the compilation directive `#pragma glsl`. This directive converts Cg or HLSL code into GLSL.

The `#pragma only_renderers` directive.

Unity supports several rendering platforms such as `gles`, `gles3`, `opengl`, `d3d11`, `d3d11_9x`, `xbox360`, `ps3` and `flash`. By default, shaders are compiled to all these platforms unless you explicitly limit this number using the `#pragma only_renderers` followed by the render APIs you want leaving a blank space between them.

If you are targeting mobile devices only limit shader compilations to `gles` and `gles3`. You must also add the `opengl` and `d3d11` renderers used by Unity Editor:

```
#pragma only_renderers gles gles3 [opengl, d3d11]
```

6.1.4 Includes

It is possible to add include files in the shader to make use of Unity predefined variables and helper functions.

You can see the available includes in `C:\Program Files\Unity\Editor\Data\CGIncludes`. For example, in the include `UnityCG.cginc` you can find several useful helper functions and macros used in many standard shaders. To use them, declare the includes in your shader.

A number of Unity built-in variables are available to shaders. They are located in the include `UnityShaderVariables.cginc`. You are not required to include this file in your shader because Unity does this automatically. Several useful transformation matrices and magnitudes are directly available in the shaders. It is important to know all of these to avoid duplicating the work. For example, before considering how to pass a matrix to the shader, camera position or projection parameters or light parameters, check if an include already provides this.

To improve performance, it is sometimes preferable to execute an operation in the CPU and pass the result to the GPU instead of executing it in the vertex shader for every vertex. For example, this is the case of multiplications of matrix uniforms. This is the reason why Unity made available for us as built-in uniforms several compound matrices. Some of the important Unity shader built-in values are shown in the following table:

Table 6-1 Important Unity shader built-in values

Built-in Uniform	Description
UNITY_MATRIX_V	Current view matrix
UNITY_MATRIX_P	Current projection matrix
Object2World	Current model matrix
_World2Object	Inverse of current world matrix
UNITY_MATRIX_VP	Current view * projection matrix
UNITY_MATRIX_MV	Current model * view matrix
UNITY_MATRIX_MVP	Current model * view * projection matrix
UNITY_MATRIX_IT_MV	Invert transpose of current model * view matrix
_WorldSpaceCameraPos	Camera position in world space
_ProjectionParams	Near and far planes and 1/farPlane as components of a vector
_Time	Current time and fractions in a vector (t/20, t, t*2, t*3)

6.1.5 The OpenGL ES 3.0 graphics pipeline

It is important to know where in the graphic pipeline the programmable vertex and fragment shaders are located.

The following figure shows a schematic view of the OpenGL ES 3.0 graphic pipeline flow. OpenGL ES 3.0 is a major step in the evolution of embedded graphics and is derived from the OpenGL 3.3 specification.

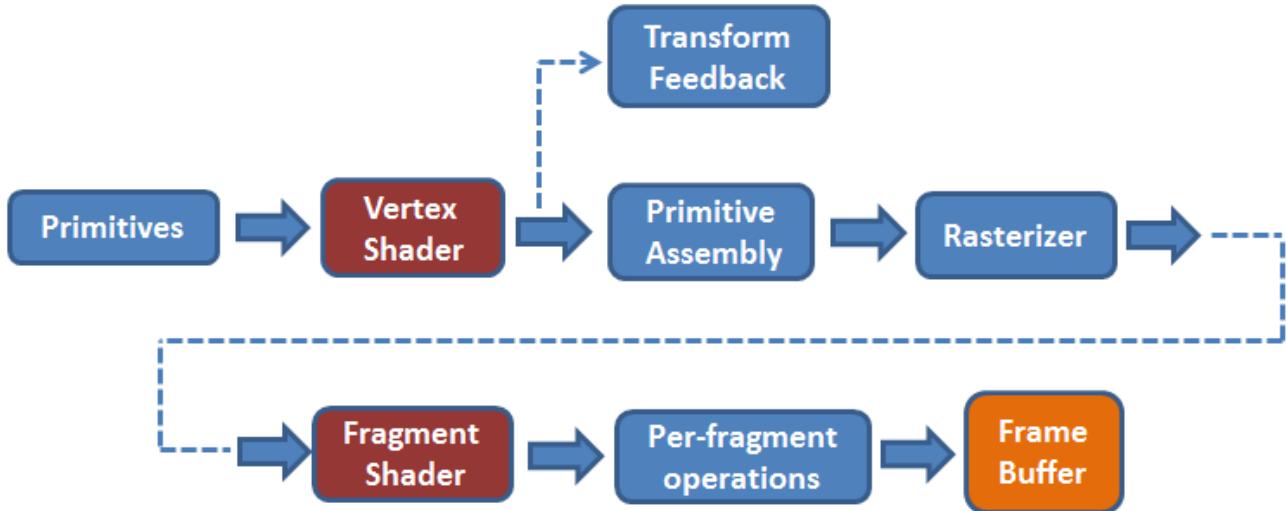


Figure 6-2 OpenGL ES 3.0 Programmable Pipeline

Primitives

In the primitives stage the pipeline operates on the geometric primitives described by vertices, points, lines and polygons.

Vertex Shader

The vertex shader implements a general-purpose programmable method for operating on vertices. The vertex shader transforms and lights vertices.

Primitive assembly

In primitive assembly the vertices are assembled into geometric primitives. The resulting primitives are clipped to a clipping volume and sent to the rasterizer.

Rasterization

Output values from the vertex shader are calculated for every generated fragment. This process is known as interpolation. During rasterization, the primitives are converted into a set of two-dimensional fragments that are then sent to the fragment shader.

Transform feedback

Transform feedback, enables writing selective writing to an output buffer that the vertex shader outputs and is later sent back to the vertex shader. This feature is not exposed by Unity but it is used internally, for example, to optimize the skinning of characters.

Fragment shader

The fragment shader implements a general-purpose programmable method for operating on fragments before they are sent to the next stage.

Per-fragment operations

In Per-fragment operations several functions and tests are applied on each fragment: pixel ownership test, scissor test, stencil and depth tests, blending and dithering. As a result of this per-fragment stage either the fragment is discarded or the fragment color, depth or stencil value is written to the frame buffer in screen coordinates.

6.1.6 Vertex shaders

The vertex shader example runs once for each vertex of the geometry. The purpose of the vertex shader is to transform the 3D position of each vertex, given in the local coordinates of the object, to the projected 2D position in screen space and calculate the depth value for the Z-buffer.

The vertex shader example sample code is in [6.1.2 Shader structure on page 6-86](#).

The transformed position is expected in the output of the vertex shader. If the vertex shader does not return a value the console displays the following error:

```
Shader error in 'Custom/ctTextured': '' : function does not return a value: vert at
line 36
```

In the example, the vertex shader receives as input, the vertex coordinates in local space and the texture coordinates. Vertex coordinates are transformed from local to screen space using the Model View Projection matrix `UNITY_MATRIX_MVP` that is a Unity built-in value:

```
output.pos = mul(UNITY_MATRIX_MVP, input.vertex);
```

Texture coordinates are passed to fragment shaders as a varying but this it does not mean that they are not transformed.

Normals are transformed from object space to world space in a different manner. To guarantee that the normal is still normal to the triangle after a non-uniform scaling operation, it must be multiplied by the transpose of the inverse of the transformation matrix. To apply the transpose operation you flip the order of factors in the multiplication. The inverse of the local to world matrix is the built-in `_World2Object` Unity matrix. It is a 4x4 matrix so you must build a 4 component vector from the 3 component normal input vector.

```
float4 normalWorld = mul(float4(input.normal, 0.0), _World2Object);
```

When building the four component vector you add a zero as the fourth component. This is necessary to handle vector transformation correctly in the fourth dimensional space while for coordinates the fourth component must be a unit.

You can skip the process of normal transformation if normals are supplied already in world coordinates. This saves work in the vertex shader. Avoid this hint if the object mesh could potentially be handled by any Unity built in shader because in this case normals are expected in object coordinates.

Most of the graphics effects are implemented in the fragment shader but you can also do some effects in the vertex shader. Vertex Displacement Mapping, also known as Displacement Mapping is a well-known

technique enabling you to deform a polygonal mesh using a texture to add surface detail, for example, in terrain generation using height maps. To have access in the vertex shader to this texture, also known as displacement map, you must add the pragma directive `#pragma target 3.0` because it is only available in shader model 3.0. According to the shader model 3.0 at least 4 texture units must be accessible inside the vertex shader. If you force the editor to use the OpenGL renderer then you must also add the `#pragma glsl` directive. If you do not declare this directive the error message produced suggests it:

```
Shader error in 'Custom/ctTextured': function "tex2D" not supported in this profile
(maybe you want #pragma glsl?) at line 57
```

In the vertex shader you also can animate vertices using “procedural animation” techniques. You can use the time variable in shaders enabling you to modify the vertex coordinates as a function of time. Mesh skinning is another type of functionality implemented in the vertex shader. Unity uses this to animate the vertices of the meshes associated with character skeletons.

6.1.7 Vertex shader input

The input and output of the vertex shader are defined by means of structures. In the input structure of the example you declare only the vertex attributes position and texture coordinates.

You can define more attributes as input, for example a second set of texture coordinates, normals in object coordinates, colors and tangents, using the following semantics.

```
struct vertexInput
{
    float4 vertex : POSITION;
    float4 tangent : TANGENT;
    float3 normal : NORMAL;
    float4 texcoord : TEXCOORD0;
    float4 texcoord1 : TEXCOORD1;
    fixed4 color : COLOR;
};
```

A semantic is a string attached to a shader input or output that provides information about the use of a parameter. You must specify a semantic for all variables passed between shader stages.

If you use incorrect semantics such as `float3 tangent2 : TANGENTIAL`, you get an error of the following type:

```
Shader error in 'Custom/ctTextured': unknown semantics "TANGENTIAL" specified for
"tangent2" at line 32
```

For performance, only specify the parameters in the input structure that you strictly require. Unity has some predefined input structures for the most common cases of input parameter combinations: `appdata_base`, `appdata_tan` and `appdata_full`. These are described in the `UnityCG.cginc` include file. The previous vertex input structure example corresponds to `appdata_full`. In this case you are not required to declare the structure, only declare the include file.

6.1.8 Vertex shader output and varyings

Vertex shader output is defined in an output structure that must contain the vertex transformed coordinates. In the following example, the output structure is very simple but you can add other magnitudes.

The following code lists the semantics supported by Unity:

```
struct vertexOutput
{
    float4 pos : SV_POSITION;
    float4 tex : TEXCOORD0;
    float4 texSpecular : TEXCOORD1;
    float3 vertexInWorld : TEXCOORD2;
    float3 viewDirInWorld : TEXCOORD3;
    float3 normalInWorld : TEXCOORD4;
    float3 vertexToLightInWorld : TEXCOORD5;
    float4 vertexInScreenCoords : TEXCOORD6;
    float4 shadowsVertexInScreenCoords : TEXCOORD7;
};
```

The transformed vertex coordinates are defined with the semantic `SV_POSITION`. Two textures, several vectors, and coordinates in different spaces calling the semantic `TEXCOORDn` are also passed to the fragment shader.

`TEXCOORD0` is typically reserved for UVs and `TEXCOORD1` for lightmap UVs, but technically you can send anything from `TEXCOORD0` to `TEXCOORD7` to the fragment shader. It is important to notice that each interpolator, that is each semantic, can only process a maximum of 4 floats. Put larger variables such as matrices into multiple interpolators. This means that if you define a matrix to be passed as a varying:
`float4x4 myMatrix : TEXCOORD2`, Unity uses the interpolators from `TEXCOORD2` to `TEXCOORD5`.

Everything you send from the vertex shader to the fragment shader is linearly interpolated by default. For every pixel in the triangle defined by the vertices V_1 , V_2 and V_3 the rasterizer, located in the graphic pipeline between vertex and fragment shaders, calculates the pixel coordinates as a linear interpolation of the vertices coordinates using the barycentric coordinates λ_1 , λ_2 and λ_3 .

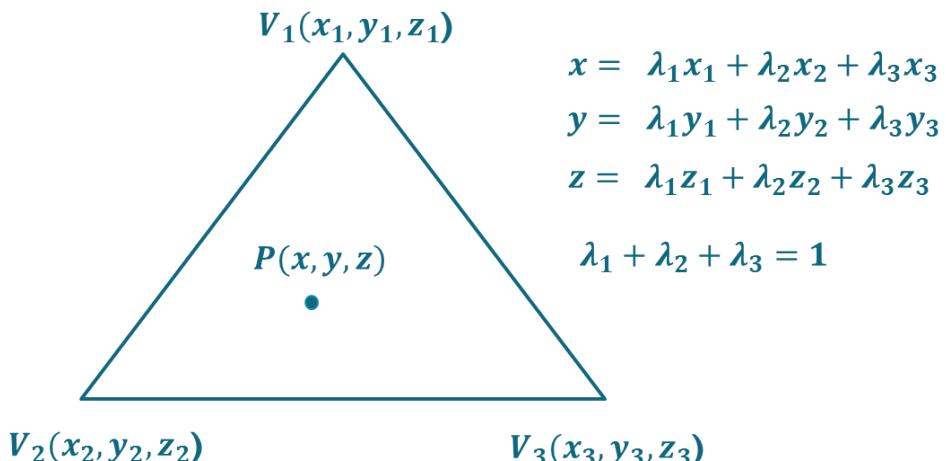


Figure 6-3 Linear interpolation using barycentric coordinates

The following diagram shows the results of color interpolation in a triangle with vertex colors red, green and blue.

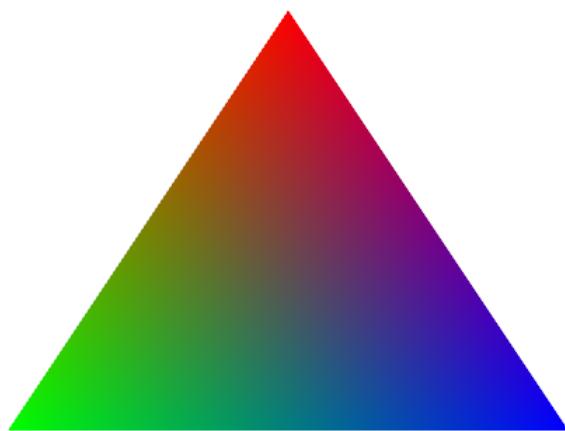


Figure 6-4 Color Interpolation

The same interpolation is applied to any varying passed from the vertex to the fragment shader. This is a very powerful tool because there is a hardware linear interpolator. For example, if you have a plane and you want to apply a color as a function of the distance to the center C, you pass the coordinate of the center C to the vertex shader, calculate the squared distance from the vertex to C and pass that magnitude to the fragment shader. The value of the distance is automatically interpolated for you in every pixel of every triangle.

Values are linearly interpolated so it is possible to perform per-vertex computations and reuse them in the fragment shader, that is, a value that can be linearly interpolated in the fragment shader can be calculated in the vertex shader instead. This can provide a substantial performance boost because the vertex shader runs on a much smaller data set than the fragment shader.

You must be careful with the use of varyings especially in mobiles where performance and memory bandwidth consumption are critical to the success of many games. The more varyings there are, the more vertex accesses and fragment shader varying reads bandwidth. Aim for a reasonable balance when using varyings.

6.1.9 Fragment Shaders

The fragment shader is the graphics pipeline stage after primitive rasterization.

For each sample of the pixels covered by a primitive, a fragment is generated. The fragment shader code is executed for each generated fragment. There are many more fragments than vertices so you must take care about the number of operations performed in the fragment shader.

In the fragment shader you can access the fragment coordinates in the windows space among other values that contains all interpolated per-vertex output values from the vertex shader.

In the shader example in [6.1.2 Shader structure on page 6-86](#), the fragment shader receives the interpolated texture coordinates from the vertex shader and performs a texture lookup to obtain the color at these coordinates. It combines this color with the ambient color to produce the final output color. From the declaration of the fragment shader `float4 frag(vertexOutput input) : COLOR` it is clear that it is expected to produce the fragment color. The fragment shader is where you do the operations to achieve the required effect. This ultimately consists of assigning the correct color to a fragment.

6.1.10 Providing data to shaders

Any data declared as a uniform in the Pass block is available to both vertex and fragment shaders.

A uniform can be considered as a type of global constant variable because it cannot be modified inside the shader.

You can supply this uniform to the shader in the following ways:

- Using the **Properties** block.
- Programmatically from a script.

The **Properties** block enables you to define uniforms interactively in the **Inspector**. Any variable declared in the **Properties** block appears in the material inspector listed with the variable name.

The following code shows the **Properties** block of the shader example associated to material `ctSphereMat`:

```
Properties
{
    _AmbientColor ("Ambient Color", Color) = (0.2,0.2,0.2,1.0)
    _MainTex ("Base (RGB)", 2D) = "white" {}
}
```

The variables `_AmbientColor` and `_MainTex` declared in the **Properties** block with the names **Ambient Color** and **Base (RGB)** respectively appear in the **Material Inspector** with those names.

The following figure shows Properties in Material Inspector:

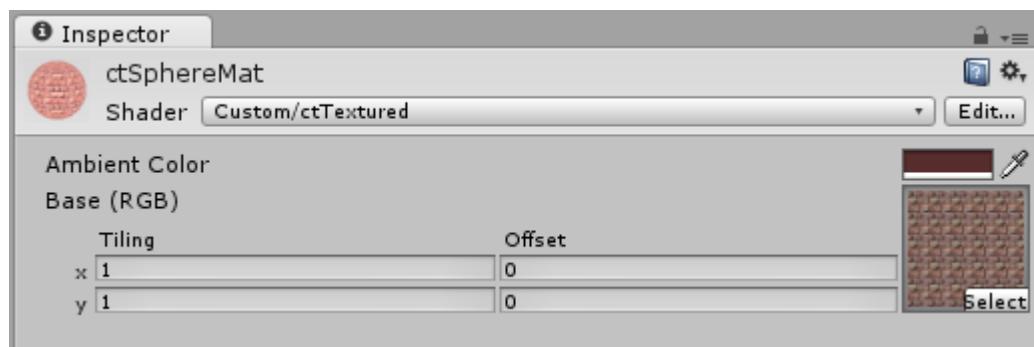


Figure 6-5 Properties in Material Inspector

Passing data to the shader by means of the **Properties** block is very useful especially when you are in the development stage of the shader because you can change the data interactively and see the effect at run time.

You can put the following types of variables in the **Properties** block:

- Float.
- Color.
- Texture 2D.
- Cubemap.
- Rectangle.
- Vector.

The **Properties** block is not a useful way of passing data if for example, data is required from a previous calculation or data is required to be passed at specific point in time.

An alternative method of passing data to the shaders is programmatically from a script.

The material class exposes several methods that you can use to pass data associated with a material to a shader. The following table lists the most common methods:

Table 6-2 Common methods for passing data associated with a material to a shader

Method
<code>SetColor (propertyName: string, color: Color);</code>
<code>SetFloat (propertyName: string, value: float);</code>
<code>SetInt (propertyName: string, value: int);</code>
<code>SetMatrix (propertyName: string, matrix: Matrix4x4);</code>
<code>SetVector (propertyName: string, vector: Vector4);</code>
<code>SetTexture (propertyName: string, texture: Texture);</code>

In the following code, immediately before the main camera renders the scene, a secondary camera `shwCam` renders the shadows to a texture to be combined with the main camera render pass.

For the shadow texture projection process the vertices must be transformed in a convenient manner. The shadow camera projection matrix (`shwCam.projectionMatrix`), world to local transformation matrix (`shwCam.transform.worldToLocalMatrix`), and the rendered shadow texture (`m_ShadowsTexture`) are passed to the shader.

These values are available in the shader as uniforms with the names `_ShwCamProjMat`, `_ShwCamViewMat` and `m_ShadowsTexture`.

The following code shows how matrices and textures are sent to the shader by means of materials contained in the list `shwMats`.

```
// Called before object is rendered.

public void OnWillRenderObject()
{
    // Perform different checks.
    ...
    CreateShadowsTexture();
    // Set up shadows camera shwCam.
    ...
    /// Pass matrices to the shader
    for(int i = 0; i < shwMats.Count; i++)
    {
        shwMats[i].SetMatrix("_ShwCamProjMat", shwCam.projectionMatrix);
        shwMats[i].SetMatrix("_ShwCamViewMat", shwCam.transform.worldToLocalMatrix);
    }
    // Render shadows texture
    shwCam.Render();
    for(int i = 0; i < shwMats.Count; i++)
    {
        shwMats[i].SetTexture( "_ShadowsTex", m_ShadowsTexture );
    }
    s_RenderingShadows = false;
}
```

6.1.11 Debugging shaders in Unity

In Unity it is not possible to debug shaders in the same way as you do with traditional code. You can however use the output of the fragment shader to visualize the values you want to debug. You then have to interpret the image produced.

The following figure shows the output of the shader `ctRef1LocalCubemap.shader` applied to the reflective surface of the floor from [6.2 Implementing reflections with a local cubemap on page 6-98](#):



Figure 6-6 Chess room with reflections

In following fragment shader, the output color has been replaced by the normalized local corrected reflected vector:

```
return float4(normalize(localCorrRef1DirWS), 1.0);
```

Instead of the reflected image, it visualizes the components of the reflected vector normalized as colors.

The reddish color zone in the floor indicates that the reflected vector has a strong X component, that is, it is mostly oriented to toward the X Axis. The reddish part shows the reflection coming from that direction, that is, from the windowed wall.

The bluish zone indicates a predominance of reflected vectors oriented to Z axis, that is, the reflection from the right wall.

In the black zone the vectors are mainly oriented to $-Z$ but the colors can only have positive components because the negative components are clamped to 0.

The following figure shows the result of replacing the output color of the fragment by the normalized local reflected vector:



Figure 6-7 Shader debugging with multiple colors

It might initially be difficult to interpret the meaning of the colors while debugging so try to focus on a single color component. For example, you can return only the Y component of the normalized local corrected reflected vector:

```
float3 normLocalCorrRef1DirWS = normalize(localCorrRef1DirWS);
return float4(0, normLocalCorrRef1DirWS.y, 0, 1);
```

In this case, the output is only the reflections coming mainly from the roof above the camera. That is, the part of the room oriented to the Y axis. The reflections from the walls of the room are coming from X, Z and $-Z$ directions, so they are rendered in black.

The following figure shows shader debugging with a single color:



Figure 6-8 Shader debugging with a single color

Check that the magnitude you are debugging with color, is between 0 and 1 because any other value is automatically clamped. Any negative value is assigned zero and any value greater than 1 is assigned 1.

6.2 Implementing reflections with a local cubemap

Reflections based on a local cubemap is a useful technique for rendering reflections on mobile devices.

Unity version 5 and higher implements reflections based on local cubemaps as **Reflection Probes**. You can combine these with other types of reflections, such as reflections rendered at runtime in your own custom shader.

This section contains the following subsections:

- [6.2.1 History of reflection implementations](#) on page 6-98.
- [6.2.2 Generating correct reflections with a local cubemap](#) on page 6-100.
- [6.2.3 Shader Implementation](#) on page 6-102.
- [6.2.4 Filtering cubemaps](#) on page 6-104.
- [6.2.5 Ray-box intersection algorithm](#) on page 6-108.
- [6.2.6 Source code for editor script to generate cubemaps](#) on page 6-111.

6.2.1 History of reflection implementations

Graphics developers have always tried to find computationally cheap methods to implement reflections.

One of the first solutions is spherical mapping. This technique simulates reflections or lighting on objects without using expensive ray-tracing or lighting calculations.

The following figure shows an environment map on a sphere:

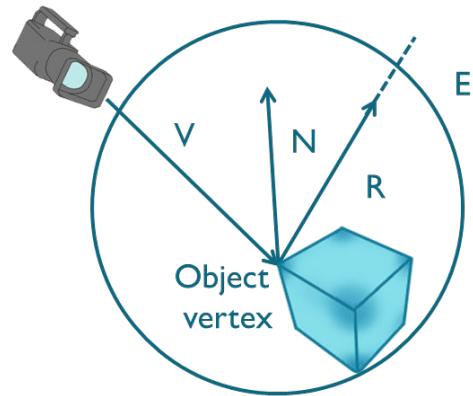


Figure 6-9 Environment map on a sphere

The following figure shows the equation for mapping a spherical surface into two dimensions:

The spherical surface is mapped into 2D:

$$u = \frac{R_x}{m} + \frac{1}{2}$$

$$v = \frac{R_y}{m} + \frac{1}{2}$$

$$m = 2 \cdot \sqrt{R_x^2 + R_y^2 + (R_z + 1)^2}$$

Figure 6-10 Spherical surface 2D mapping equation

This approach has several disadvantages, but the main problem is the distortions that occur when mapping a picture onto a sphere. In 1999, it became possible to use cubemaps with hardware acceleration. Cubemaps solved the problems of image distortions, viewpoint dependency and computational inefficiency related to spherical mapping.

The following figure shows an unfolded cube:

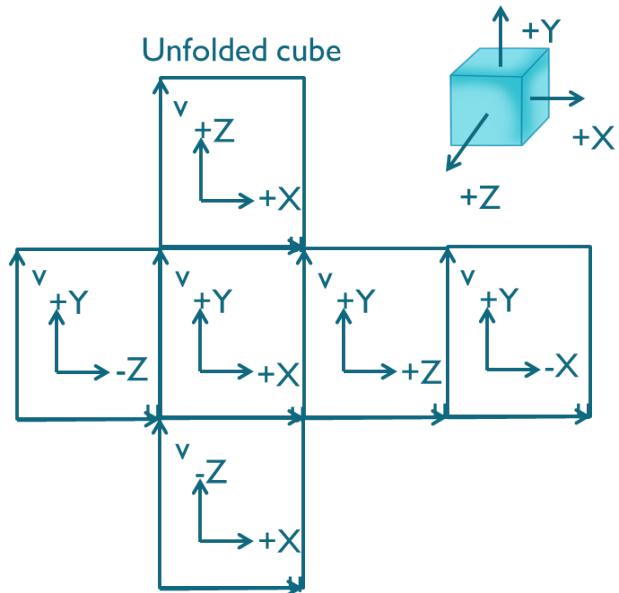


Figure 6-11 Unfolded cube

Cubemapping uses the six faces of a cube as the map shape. The environment is projected onto each side of a cube and stored as six square textures, or unfolded into six regions of a single texture. The cubemap is generated by rendering the scene from a given position with six different camera orientations with a 90 degree view frustum representing each a cube face. Source images are sampled directly. No distortion is introduced by resampling into an intermediate environment map.

The following figure shows infinite reflections:

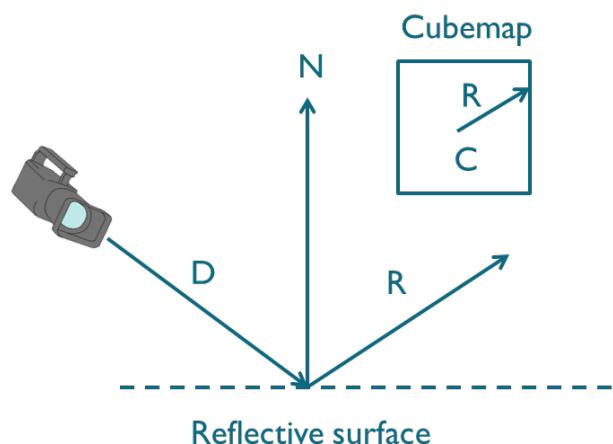


Figure 6-12 Infinite reflections

To implement reflections based on cubemaps, evaluate the reflected vector R and use it to fetch the texel from the cubemap _Cubemap using the available texture lookup function `texCUBE()`:

```
float4 color = texCUBE(_Cubemap, R);
```

The normal \mathbf{N} and view vector \mathbf{D} are passed to fragment shader from the vertex shader. The fragment shader fetches the texture color from the cubemap:

```
float3 R = reflect(D, N);
float4 color = texCUBE(_Cubemap, R);
```

This approach can only reproduce reflections correctly from a distant environment where the cubemap position is not relevant. This simple and effective technique is mainly used in outdoor lighting, for example, to add reflections of the sky.

The following figure shows incorrect reflections:

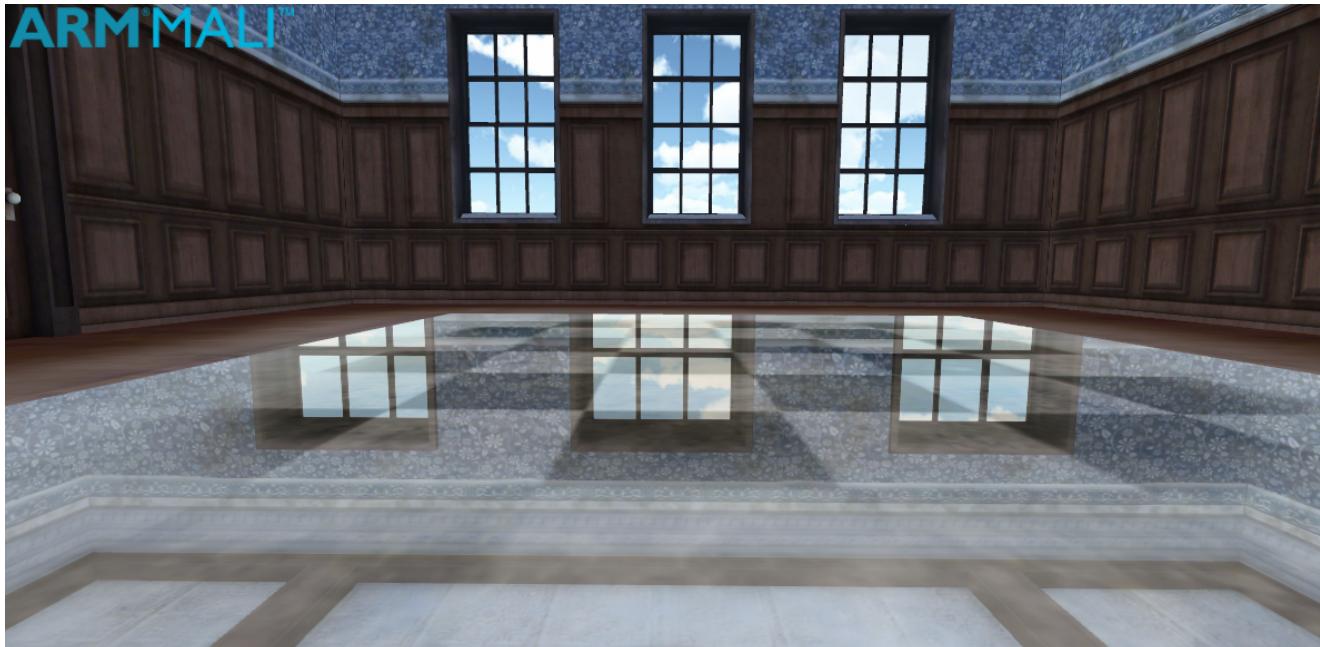


Figure 6-13 Incorrect reflections

If you use this technique in a local environment it produces inaccurate reflections. The reason why the reflections are incorrect is that in the expression `float4 color = texCUBE(_Cubemap, R);` there is no binding to the local geometry. For example, if you walk across a reflective floor looking at it from the same angle you always see the same reflection. The reflected vector is always the same and the expression always produces the same result. This is because the direction of the view vector does not change. In the real world reflections depend on both viewing angle and viewing position.

6.2.2 Generating correct reflections with a local cubemap

A solution to this problem involves binding to the local geometry in the procedure to calculate the reflection.

This solution is described in *GPU Gems: Programming Techniques, Tips, and Tricks for Real-time Graphics* by Randima Fernando (Series Editor).

The following figure shows a local correction using a bounding sphere:

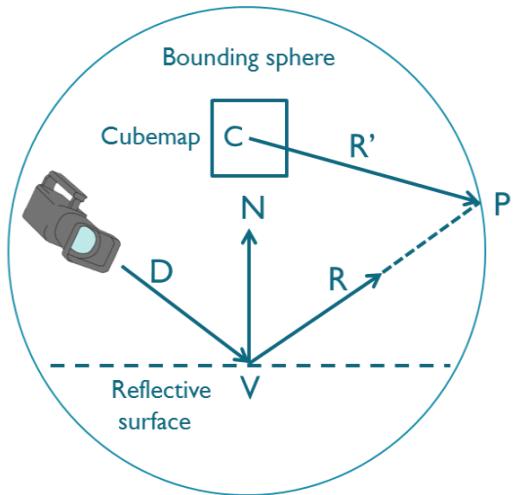


Figure 6-14 Local correction using a bounding sphere

A bounding sphere is used as a proxy volume that delimits the scene to be reflected. Instead of using the reflected vector R to fetch the texture from the cubemap a new vector R' is used. To build this new vector you find the intersection point P in the bounding sphere of the ray from the local point V in the direction of the reflected vector R . Create a new vector R' from the center of the cubemap C , where the cubemap was generated, to the intersection point P . Use this vector to fetch the texture from the cubemap.

```
float3 R = reflect(D, N);
Find intersection point P
Find vector R' = CP
float4 col = texCUBE(_Cubemap, R');
```

This approach produces good results in the surfaces of objects with a near spherical shape but reflections in plane reflective surfaces are deformed. Another drawback of this method is that the algorithm to calculate the intersection point with the bounding sphere solves a second degree equation and this is relatively complex.

In 2010 a developer proposed a better solution in a forum at <http://www.gamedev.net>. This approach replaces the previous bounding sphere by a box and solves the deformations and complexity problems of the previous method. For more information see: <http://www.gamedev.net/topic/568829-box-projected-cubemap-environment-mapping/?&p=4637262>.

The following figure shows a local correction using a bounding box:

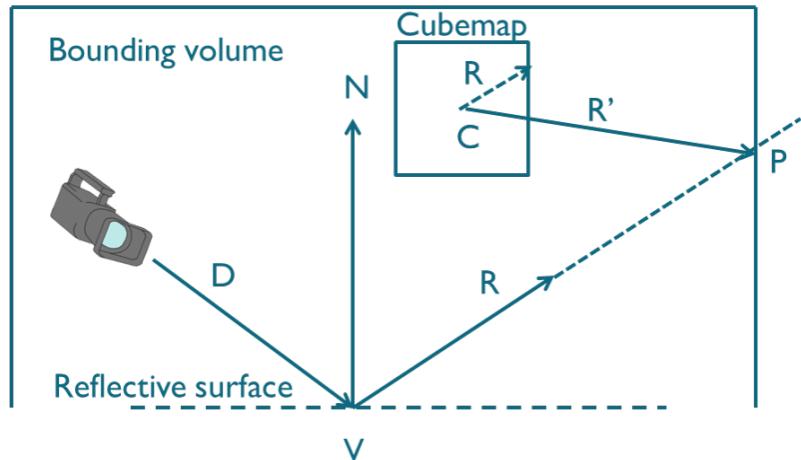


Figure 6-15 Local correction using a bounding box

A more recent work in 2012 by Sébastien Lagarde uses this new approach to simulate more complex ambient specular lighting using several cubemaps and uses an algorithm to evaluate the contribution of each cubemap and efficiently blend on the GPU. See <http://seblagarde.wordpress.com>

Table 6-3 Differences between infinite and local cubemaps

Infinite Cubemaps	Local Cubemaps
<ul style="list-style-type: none"> Mainly used outdoors to represent the lighting from a distant environment. Cubemap position is not relevant. 	<ul style="list-style-type: none"> Represents the lighting from a finite local environment. Cubemap position is relevant. The lighting from these cubemaps is right only at the location where the cubemap was created. Local correction must be applied to adapt the intrinsic infinite nature of cubemaps to local environment.

The following figure shows the scene with correct reflections generated with local cubemaps.



Figure 6-16 Correct reflections

6.2.3 Shader Implementation

This section describes shaders that implement reflections using local cubemaps.

The vertex shader calculates three magnitudes that are passed to the fragment shader as interpolated values:

- The vertex position.
- The view direction.
- The normal.

These values are in world coordinates.

The following code shows a shader implementation of reflections using local cubemaps, for Unity.

```
vertexOutput vert(vertexInput input)
{
    vertexOutput output;
    output.tex = input.texcoord;
    // Transform vertex coordinates from local to world.
    float4 vertexWorld = mul(_Object2World, input.vertex);
    // Transform normal to world coordinates.
    float4 normalWorld = mul(float4(input.normal, 0.0), _World2Object);
    // Final vertex output position. output.pos = mul(UNITY_MATRIX_MVP, input.vertex);
```

```
// ----- Local correction -----
output.vertexInWorld = vertexWorld.xyz;
output.viewDirInWorld = vertexWorld.xyz - _WorldSpaceCameraPos;
output.normalInWorld = normalWorld.xyz;
return output;
}
```

The intersection point in the volume box and the reflected vector are computed in the fragment shader. You build new local corrected reflection vector and use it to fetch the reflection texture from the local cubemap. You then combine the texture and reflection to produce the output color:

```
float4 frag(vertexOutput input) : COLOR
{
    float4 reflColor = float4(1, 1, 0, 0);
    // Find reflected vector in WS.
    float3 viewDirWS = normalize(input.viewDirInWorld);
    float3 normalWS = normalize(input.normalInWorld);
    float3 reflDirWS = reflect(viewDirWS, normalWS);
    // Working in World Coordinate System.
    float3 localPosWS = input.vertexInWorld;
    float3 intersectMaxPointPlanes = (_BBoxMax - localPosWS) / reflDirWS;
    float3 intersectMinPointPlanes = (_BBoxMin - localPosWS) / reflDirWS;
    // Looking only for intersections in the forward direction of the ray.
    float3 largestParams = max(intersectMaxPointPlanes, intersectMinPointPlanes);
    // Smallest value of the ray parameters gives us the intersection.
    float distToIntersect = min(min(largestParams.x, largestParams.y), largestParams.z);
    // Find the position of the intersection point.
    float3 intersectPositionWS = localPosWS + reflDirWS * distToIntersect;
    // Get local corrected reflection vector.
    float3 localCorrRefDirWS = intersectPositionWS - _EnviCubeMapPos;
    // Lookup the environment reflection texture with the right vector.
    reflColor = texCUBE(_Cube, localCorrRefDirWS);
    // Lookup the texture color.
    float4 texColor = tex2D(_MainTex, float2(input.tex));
    return _AmbientColor + texColor * _ReflAmount * reflColor;
}
```

In the previous code for the fragment shader, the magnitudes `_BBoxMax` and `_BBoxMin` are the maximum and minimum points of the bounding volume. The variable `_EnviCubeMapPos` is the position where the cubemap was created. Pass these values to the shader from the following script:

```
[ExecuteInEditMode]
public class InfoToReflMaterial : MonoBehaviour
{
    // The proxy volume used for local reflection calculations.
    public GameObject boundingBox;

    void Start()
    {
        Vector3 bboxLength = boundingBox.transform.localScale;
        Vector3 centerBBox = boundingBox.transform.position;

        // Min and max BBox points in world coordinates
        Vector3 BMin = centerBBox - bboxLength/2;
        Vector3 BMax = centerBBox + bboxLength/2;

        // Pass the values to the material.
        gameObject.renderer.sharedMaterial.SetVector("_BBoxMin", BMin);
        gameObject.renderer.sharedMaterial.SetVector("_BBoxMax", BMax);
        gameObject.renderer.sharedMaterial.SetVector("_EnviCubeMapPos", centerBBox);
    }
}
```

Pass the values for `_AmbientColor`, `_ReflAmount`, the main texture, and cubemap texture to the shader as uniforms from the properties block:

```
Shader "Custom/ctReflLocalCubemap"
{
    Properties
    {
        _MainTex ("Base (RGB)", 2D) = "white" { }
        _Cube("Reflection Map", Cube) = "" {}
        _AmbientColor("Ambient Color", Color) = (1, 1, 1, 1)
        _ReflAmount("Reflection Amount", Float) = 0.5
    }
    SubShader
```

```
{
    Pass
    {
        CGPROGRAM
        #pragma glsl
        #pragma vertex vert
        #pragma fragment frag
        #include "UnityCG.cginc"

        // User-specified uniforms
        uniform sampler2D _MainTex;
        uniform samplerCUBE _Cube;
        uniform float4 _AmbientColor;
        uniform float _ReflAmount;
        uniform float _ToggleLocalCorrection;
        // ---Passed from script InfoRoReflmaterial.cs -----
        uniform float3 _BBoxMin;
        uniform float3 _BBoxMax;
        uniform float3 _EnvCubeMapPos;

        struct vertexInput
        {
            float4 vertex : POSITION;
            float3 normal : NORMAL;
            float4 texcoord : TEXCOORD0;
        };
        struct vertexOutput
        {
            float4 pos : SV_POSITION;
            float4 tex : TEXCOORD0;
            float3 vertexInWorld : TEXCOORD1;
            float3 viewDirInWorld : TEXCOORD2;
            float3 normalInWorld : TEXCOORD3;
        };

        Vertex shader    { }
        Fragment shader { }

        ENDCG
    }
}
}
```

The algorithm to calculate the intersection point in the bounding volume is based on the use of the parametric representation of the reflected ray from the local position or fragment. For a description of the ray-box intersection algorithm, see [6.2.5 Ray-box intersection algorithm on page 6-108](#).

6.2.4 Filtering cubemaps

One of the advantages of implementing reflections using local cubemaps is the fact that the cubemap is static. That is, it is generated during development rather than at run-time. This provides an opportunity to apply any filtering to the cubemap images to achieve an effect.

CubeMapGen is a tool by AMD that applies filtering to cubemaps. You can obtain CubeMapGen from the AMD developer web site at: <http://developer.amd.com>.

To export cubemap images from Unity to CubeMapGen you must save each cubemap image separately. For the source code of a tool that saves the images, see [6.2.6 Source code for editor script to generate cubemaps on page 6-111](#). This tool can create a cubemap and can optionally save each cubemap image separately.

You must place the script for this tool in a folder called `Editor` in the Asset directory.

To use the cubemap editor tool:

1. Create the cubemap.
2. Launch the Bake CubeMap Tool from GameObject menu.
3. Provide the cubemap and the camera render position.
4. Optionally save individual images if you plan to apply filtering to the cubemap.

The following figure shows the Bake CubeMap tool interface:

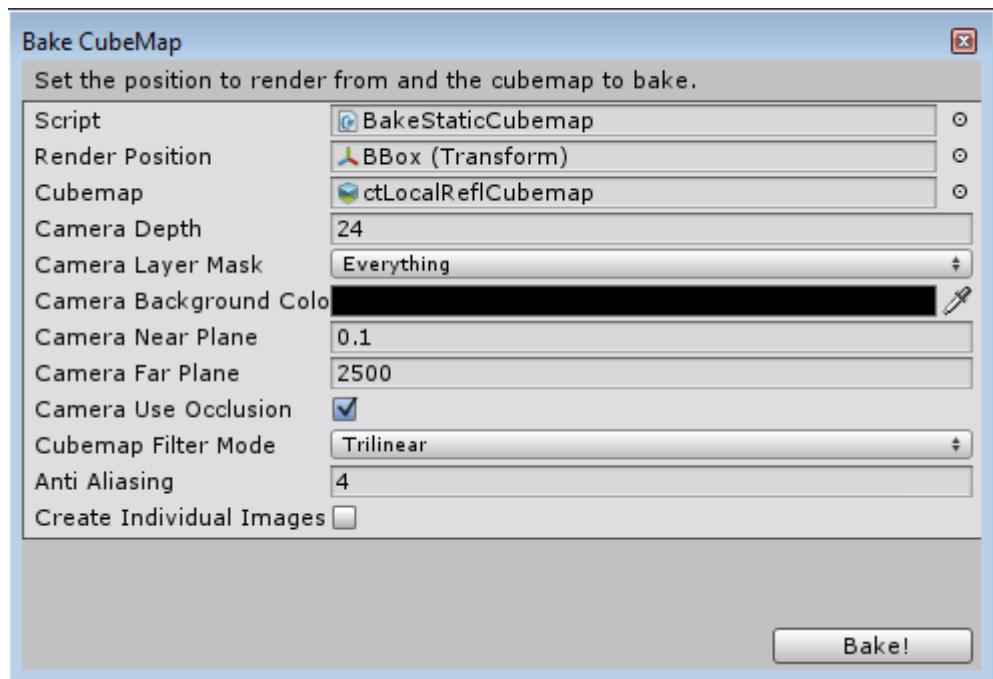


Figure 6-17 Bake CubeMap tool interface

You can load each of the images for the cubemap separately with CubeMapGen.

Select what face to load from the **Select Cube Face** drop down menu and then press **Load Cubemap Face** button. When all faces have been loaded it is possible to rotate the cubemap and check that it is correct.

CubeMapGen has a number of different filtering options in the **Filter Type** drop down menu. Select the filter settings you require and press **Filter Cubemap** to apply the filter. The filtering can take up to several minutes depending on the size of the cubemap. There is no undo option so save the cubemap as a single image before applying any filtering. If the result of the filtering is not what you expect you can reload the cubemap and try adjusting the parameters.

Use the following procedure for importing cubemap images into CubeMapGen:

1. Check the box to save individual images when baking the cubemap.
2. Launch the CubeMapGen tool and load cubemap images following the relations shown in the following table.
3. Save the cubemap as a single dds or cube cross image. Undo is not available so this enables you to reload the cubemap if you are experimenting with filters.
4. Apply filters to cubemap as required until the results are satisfactory.
5. Save the cubemap as individual images.

The following table shows the equivalence of cubemap face index between CubeMapGen and Unity.

Table 6-4 Equivalence of cubemap face index between CubeMapGen and Unity

AMD CubeMapGen	Unity
X+	-X
X-	+X
Y+	+Y
Y-	-Y
Z+	+Z
Z-	-Z

The following figure shows CubeMapGen after loading the six cubemap images:



Figure 6-18 CubeMapGen

The following figure shows the result of CubeMapGen applying a Gaussian filtering to achieve a *frosty* effect:

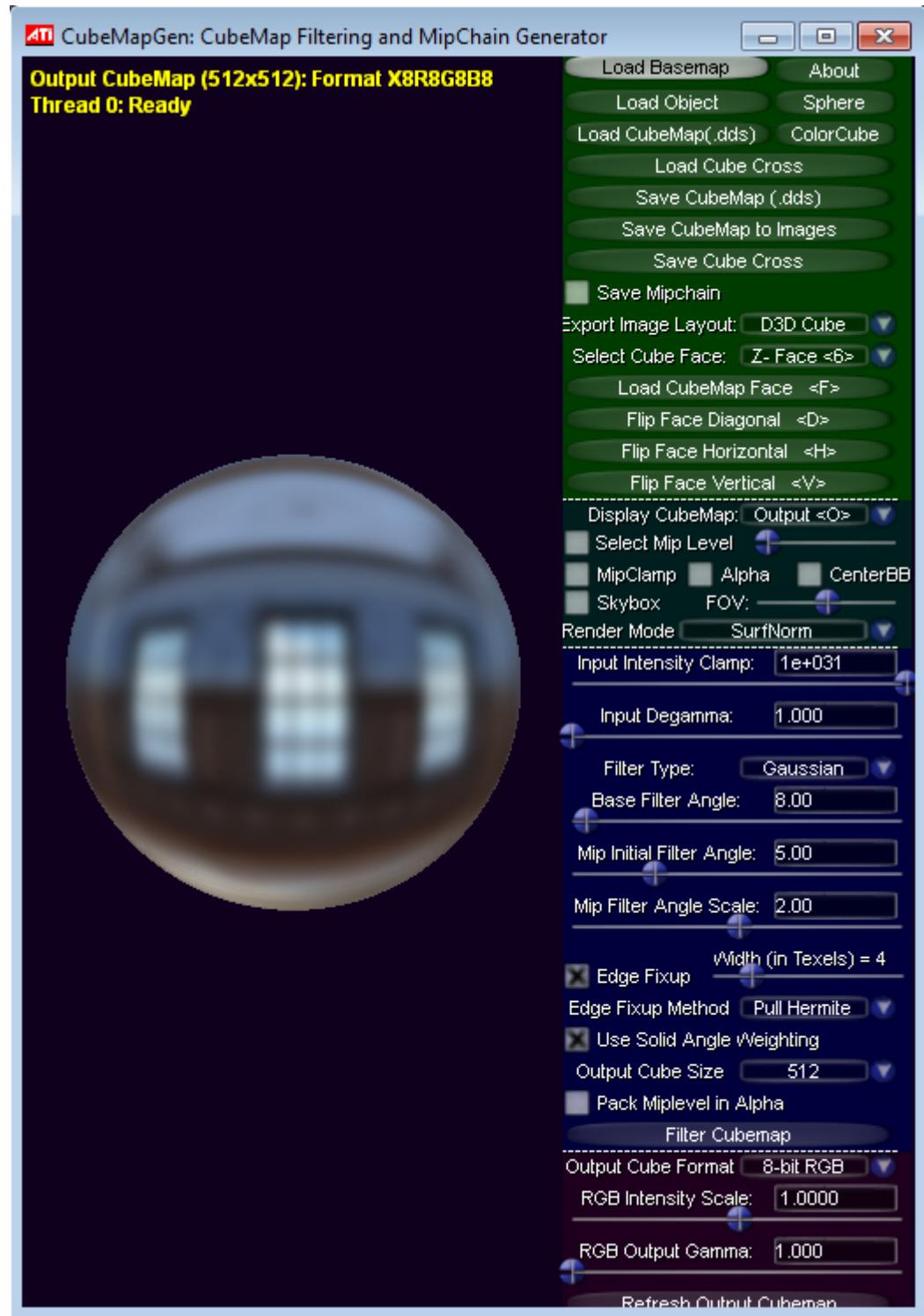


Figure 6-19 CubeMapGen showing frosty effect

The following table shows the filter parameters used with the Gaussian filter to achieve the frosty effect.

Table 6-5 Parameters used in CubeMapGen to produce a frosty effect in the reflections.

Filter settings	Value
Type	Gaussian
Base Filter Angle	8
Mip Initial Filter Angle	5
Mip Filter Angle Scale	2.0
Edge Fixup	Checked
Edge Fixup Width	4

The following figure shows a reflection generated with a cubemap with a frosty effect:



Figure 6-20 Reflection with frosty effect

The following figure summarizes the work flow to apply filtering to Unity cubemaps with the CubeMapGen tool.

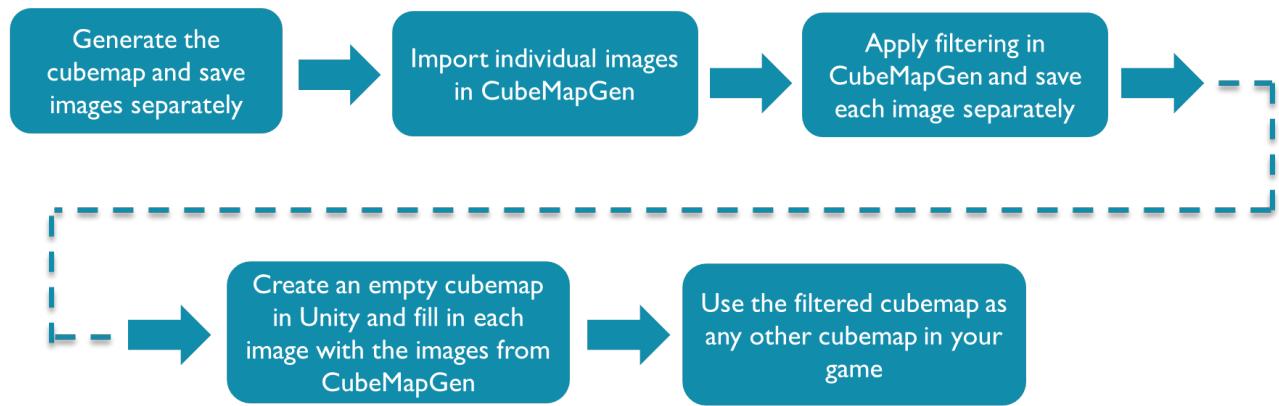


Figure 6-21 Cubemap filtering workflow

6.2.5 Ray-box intersection algorithm

This section describes the ray-box intersection algorithm.

The following figure shows a graph with line equations:

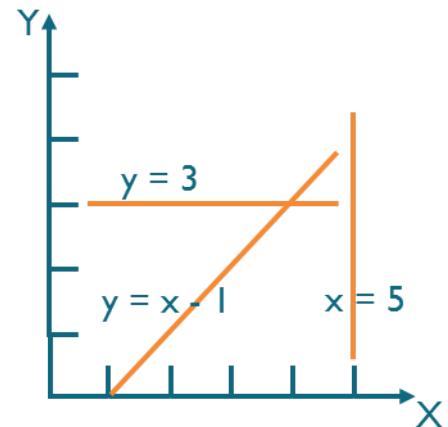


Figure 6-22 Graph with line equations

Equation of a line

$$y = mx + b$$

The vector form of this equation is:

$$\mathbf{r} = \mathbf{o} + t\mathbf{d}$$

Where:

\mathbf{o} is the origin point

\mathbf{d} is the direction vector

t is the parameter

The following figure shows an axis aligned bounding box:

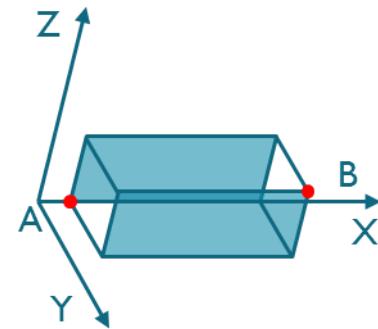


Figure 6-23 Axis aligned bounding box

An axis aligned bounding box AABB can be defined by its min and max points A and B.

AABB defines a set of lines parallel to coordinate axis. Each component of line can be defined by the following equation:

$$x = A_x; y = A_y; z = A_z; x = B_x; y = B_y; z = B_z$$

To find where a ray intersects one of those lines, equal both equations. For example:

$$\mathbf{o}_x + t_x \mathbf{d}_x = A_x$$

You can write the solution as:

$$t_{Ax} = (A_x - o_x) / d_x$$

Obtain the solution for all components of both intersection points in the same manner:

$$\begin{aligned}t_{Ax} &= (Ax - Ox) / Dx \\t_{Ay} &= (Ay - Oy) / Dy \\t_{Az} &= (Az - Oz) / Dz \\t_{Bx} &= (Bx - Ox) / Dx \\t_{By} &= (By - Oy) / Dy \\t_{Bz} &= (Bz - Oz) / Dz\end{aligned}$$

In vector form these are:

$$\begin{aligned}\mathbf{t}_A &= (\mathbf{A} - \mathbf{O}) / D \\ \mathbf{t}_B &= (\mathbf{B} - \mathbf{O}) / D\end{aligned}$$

This finds where the line intersects the planes defined by the faces of the cube but it does not guarantee that the intersections lie on the cube.

The following figure shows a 2D representation of ray-box intersection:

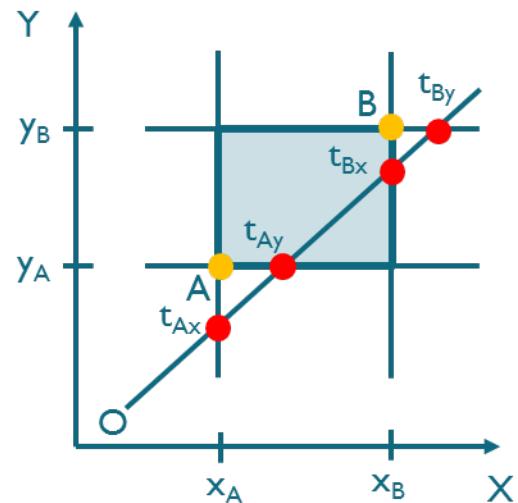


Figure 6-24 Ray-box intersection 2D representation

To find what solution is really an intersection with the box, you require the greater value of the t parameter for the intersection at the min plane.

$$t_{\min} = (t_{Ax} > t_{Ay}) ? t_{Ax} : t_{Ay}$$

You require the smaller value of the parameter t for the intersection at the max plane.

$$t_{\min} = (t_{Ax} > t_{Ay}) ? t_{Ax} : t_{Ay}$$

You also must consider those cases when you get no intersections.

The following figure shows a ray-box with no intersection:

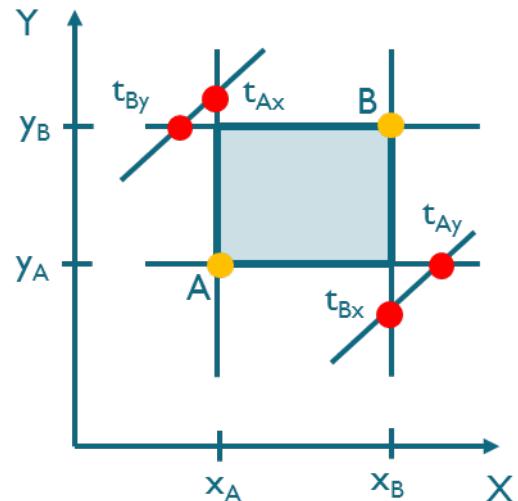


Figure 6-25 Ray-box with no intersection

If you guarantee that the reflective surface is enclosed by the BBox, that is, the origin of the reflected ray is inside the BBox, then there are always two intersections with the box, and the handling of different cases is simplified.

The following figure shows a ray-box intersection in BBox:

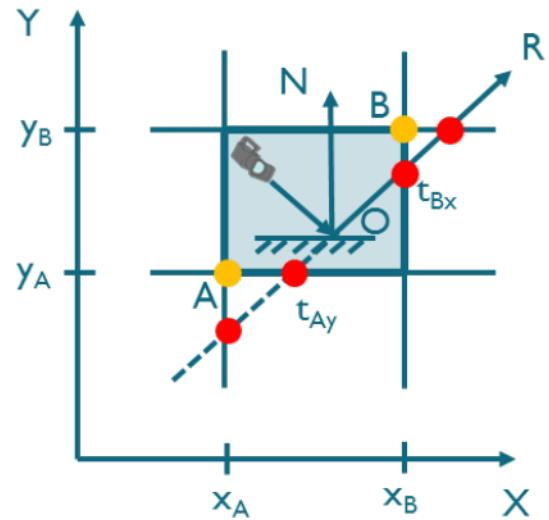


Figure 6-26 Ray-box intersection in BBox

6.2.6 Source code for editor script to generate cubemaps

This section provides the source code for editor script to generate cubemaps.

```
/*
 * This confidential and proprietary software may be used only as
 * authorised by a licensing agreement from ARM Limited
 * (C) COPYRIGHT 2014 ARM Limited
 * ALL RIGHTS RESERVED
 * The entire notice above must be reproduced on all authorised
 * copies and copies may only be made to the extent permitted
 * by a licensing agreement from ARM Limited.
 */
using UnityEngine;
using UnityEditor;
using System.IO;
/**
```

```

* This script must be placed in the Editor folder.
* The script renders the scene into a cubemap and optionally
* saves each cubemap image individually.
* The script is available in the Editor mode from the
* Game Object menu as "Bake Cubemap" option.
* Be sure the camera far plane is enough to render the scene.
*/
public class BakeStaticCubemap : ScriptableWizard
{
    public Transform renderPosition;
    public Cubemap cubemap;
    // Camera settings.
    public int cameraDepth = 24;
    public LayerMask cameraLayerMask = -1;
    public Color cameraBackgroundColor;
    public float cameraNearPlane = 0.1f;
    public float cameraFarPlane = 2500.0f;
    public bool cameraUseOcclusion = true;
    // Cubemap settings.
    public FilterMode cubemapFilterMode = FilterMode.Trilinear;
    // Quality settings.
    public int antiAliasing = 4;

    public bool createIndividualImages = false;

    // The folder where individual cubemap images will be saved
    static string imageDirectory = "Assets/CubemapImages";
    static string[] cubemapImage
        = new string[]{"front+Z", "right+X", "back-Z", "left-X", "top+Y", "bottom-Y"};
    static Vector3[] eulerAngles = new Vector3[]{new Vector3(0.0f, 0.0f, 0.0f),
        new Vector3(0.0f, -90.0f, 0.0f), new Vector3(0.0f, 180.0f, 0.0f),
        new Vector3(0.0f, 90.0f, 0.0f), new Vector3(-90.0f, 0.0f, 0.0f),
        new Vector3(90.0f, 0.0f, 0.0f)};

    void OnWizardUpdate()
    {
        helpString = "Set the position to render from and the cubemap to bake.";
        if(renderPosition != null && cubemap != null)
        {
            isValid = true;
        }
        else
        {
            isValid = false;
        }
    }

    void OnWizardCreate ()
    {

        // Create temporary camera for rendering.
        GameObject go = new GameObject( "CubemapCam", typeof(Camera) );
        // Camera settings.
        go.camera.depth = cameraDepth;
        go.camera.backgroundColor = cameraBackgroundColor;
        go.camera.cullingMask = cameraLayerMask;
        go.camera.nearClipPlane = cameraNearPlane;
        go.camera.farClipPlane = cameraFarPlane;
        go.camera.useOcclusionCulling = cameraUseOcclusion;
        // Cubemap settings
        cubemap.filterMode = cubemapFilterMode;
        // Set antialiasing
        QualitySettings.antiAliasing = antiAliasing;

        // Place the camera on the render position.
        go.transform.position = renderPosition.position;
        go.transform.rotation = Quaternion.identity;

        // Bake the cubemap
        go.camera.RenderToCubemap(cubemap);

        // Rendering individual images
        if(createIndividualImages)
        {
            if (!Directory.Exists(imageDirectory))
            {
                Directory.CreateDirectory(imageDirectory);
            }
            RenderIndividualCubemapImages(go);
        }
    }
}

```

```
// Destroy the camera after rendering.  
DestroyImmediate(go);  
}  
  
void RenderIndividualCubemapImages(GameObject go)  
{  
    go.camera.backgroundColor = Color.black;  
    go.camera.ClearFlags = CameraClearFlags.Skybox;  
    go.camera.fieldOfView = 90;  
    go.camera.aspect = 1.0f;  
  
    go.transform.rotation = Quaternion.identity;  
  
    //Render individual images  
    for (int camOrientation = 0; camOrientation < eulerAngles.Length ; camOrientation++)  
    {  
        string imageName = Path.Combine(imageDirectory, cubemap.name + "_"  
            + cubemapImage[camOrientation] + ".png");  
        go.camera.transform.eulerAngles = eulerAngles[camOrientation];  
        RenderTexture renderTex = new RenderTexture(cubemap.height,  
            cubemap.height, cameraDepth);  
        go.camera.targetTexture = renderTex;  
        go.camera.Render();  
        RenderTexture.active = renderTex;  
  
        Texture2D img = new Texture2D(cubemap.height, cubemap.height,  
            TextureFormat.RGB24, false);  
        img.ReadPixels(new Rect(0, 0, cubemap.height, cubemap.height), 0, 0);  
  
        RenderTexture.active = null;  
        GameObject.DestroyImmediate(renderTex);  
  
        byte[] imgBytes = img.EncodeToPNG();  
        File.WriteAllBytes(imageName, imgBytes);  
  
        AssetDatabase.ImportAsset(imageName, ImportAssetOptions.ForceUpdate);  
    }  
  
    AssetDatabase.Refresh();  
}  
  
[MenuItem("GameObject/Bake Cubemap")]  
static void RenderCubemap ()  
{  
    ScriptableWizard.DisplayWizard("Bake CubeMap", typeof(BakeStaticCubemap), "Bake!");  
}
```

6.3 Combining reflections

This section describes combining reflections.

This section contains the following subsections:

- [6.3.1 About combining reflections](#) on page 6-114.
- [6.3.2 Combining reflections shader implementation](#) on page 6-116.
- [6.3.3 Combining reflections from a distant environment](#) on page 6-118.

6.3.1 About combining reflections

The reflections based on local cubemaps technique enables rendering high quality, efficient reflections based on static local cubemaps. However, if the objects are dynamic, the static local cubemap is no longer valid and the technique does not work.

You can solve this by combining static reflections with dynamically generated reflections.

The following figure shows combining reflections from static and dynamic geometry:

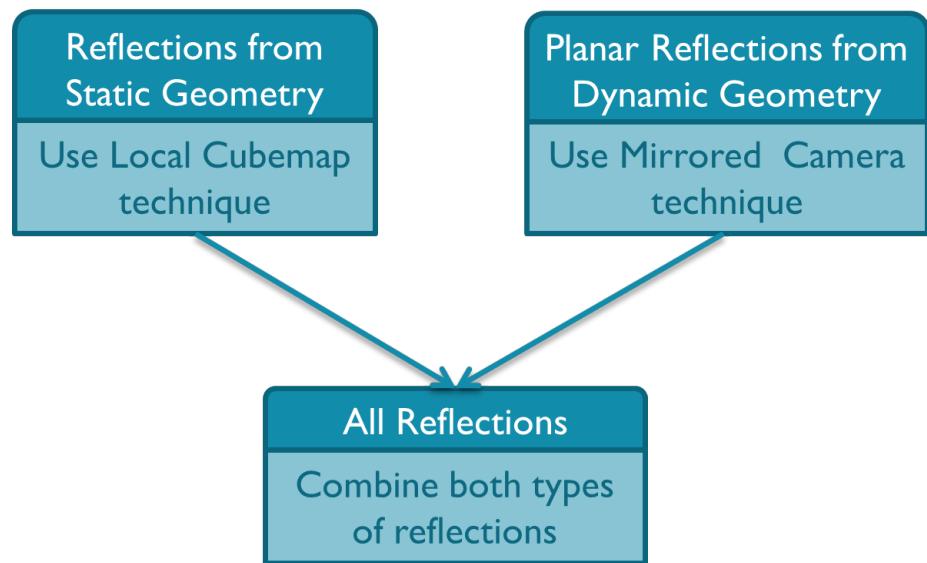


Figure 6-27 Combining reflections from static and dynamic geometry

If the reflective surface is planar, you can generate dynamic reflections with a mirrored camera.

To create a mirrored camera, calculate the position and orientation of the main camera that renders the reflections at runtime.

Mirror the position and orientation of the main camera, relative to the reflective plane.

The following figure shows the mirrored camera technique:

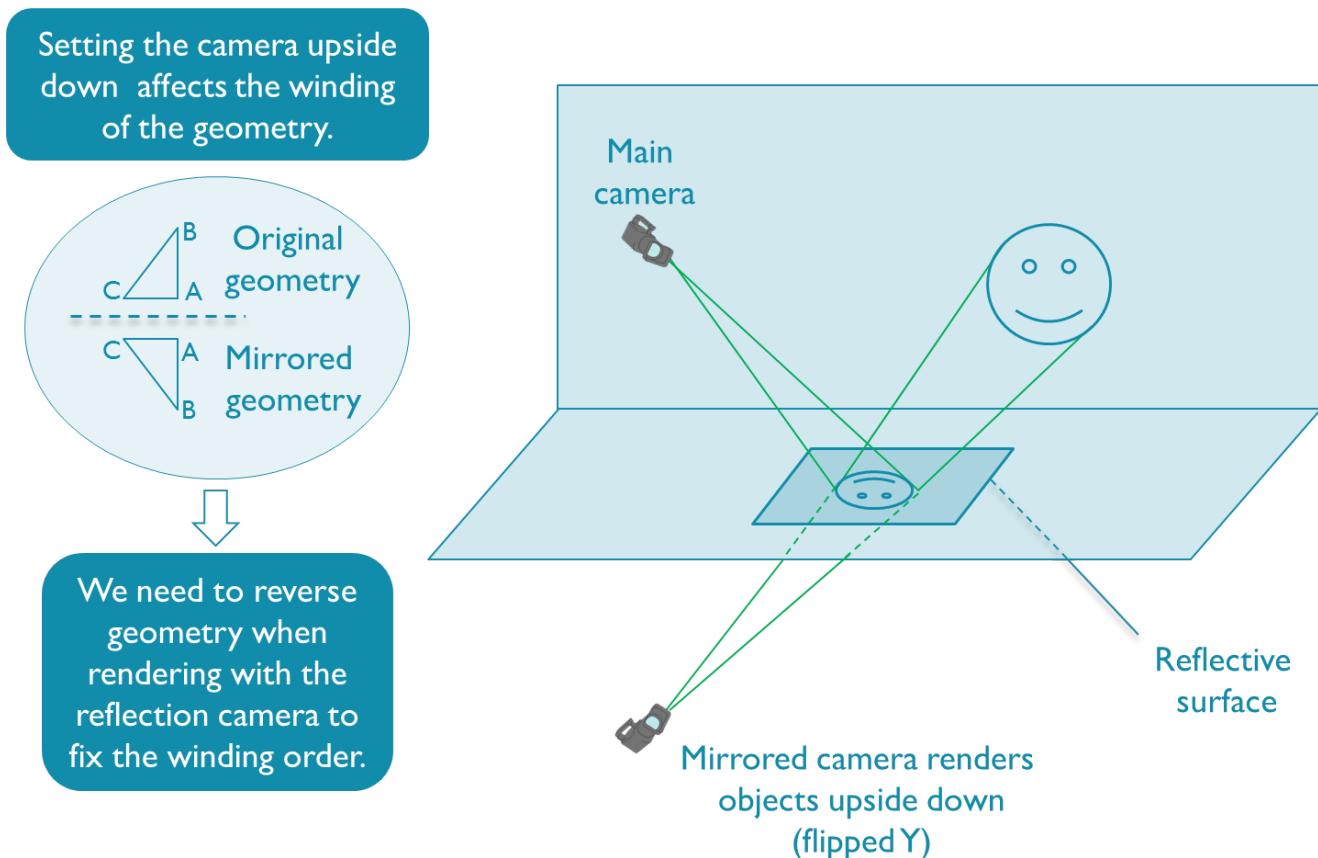


Figure 6-28 The mirrored camera technique for rendering planar reflections

In the mirroring process, the new reflection camera ends up with its axis in the opposite orientation. In the same manner as a physical mirror, reflections from the left and right are inverted. This means the reflection camera renders the geometry with an opposite winding.

To render the geometry correctly you must invert the winding of the geometry before rendering the reflections. When you have finished rendering the reflections, restore the original winding.

The following figure shows the sequence of steps required to set up the mirrored camera and render the reflections:

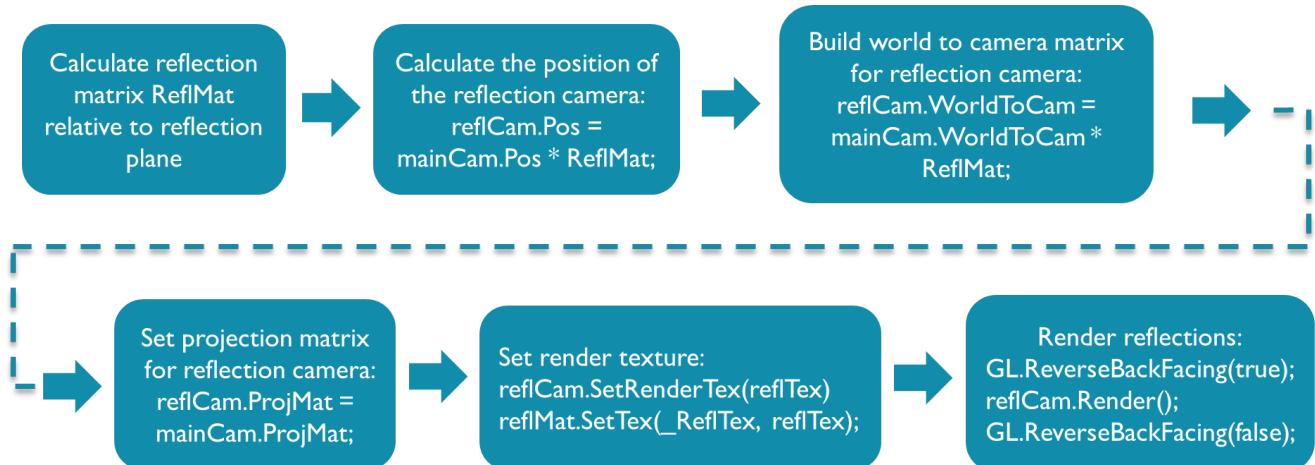


Figure 6-29 Main steps for setting up the mirrored camera and rendering the reflections

Build the mirror reflection transformation matrix. Use this matrix to calculate the position, and the world-to-camera transformation matrix, of the reflection camera.

The following figure shows the mirror reflection transformation matrix:

$$R = \begin{bmatrix} 1 - 2n_x^2 & -2n_xn_y & -2n_xn_z & -2n_xn_w \\ -2n_xn_y & 1 - 2n_y^2 & -2n_yn_z & -2n_yn_w \\ -2n_xn_z & -2n_yn_z & 1 - 2n_z^2 & -2n_zn_w \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\begin{aligned} n_x &= \text{planeNormal}.x \\ n_y &= \text{planeNormal}.y \\ n_z &= \text{planeNormal}.z \\ n_w &= -\text{dot}(\text{planeNormal}, \text{planePos}) \end{aligned}$$

Figure 6-30 The mirror reflection transformation matrix

Apply the reflection matrix transformation to the position and world-to-camera matrix of the main camera. This provides you with the position and world-to-camera matrix of the reflection camera.

The projection matrix of the reflection camera must be the same as the projection matrix of the main camera.

The reflection camera renders reflections to a texture.

For good results, you must set up this texture properly before rendering:

- Use Mipmaps.
- Set the filtering mode to trilinear.
- Use multisampling.

Ensure the texture size is proportional to the area of the reflective surface. The larger the texture is, the less pixelated the reflections are.

An example script for a mirrored camera is located at: [http://wiki.unity3d.com/index.php/
File:MirrorReflection.png](http://wiki.unity3d.com/index.php/File:MirrorReflection.png).

6.3.2 Combining reflections shader implementation

You can combine static environment reflections, with dynamic planar reflections in a shader.

To combine reflections in the shaders, you must modify the shader code provided in [6.2.3 Shader Implementation on page 6-102](#).

The shader must incorporate the planar reflections, rendered at runtime with the reflection camera. To achieve this, the texture `_ReflectionTex` from the reflection camera, is passed to the fragment shader as a uniform, and combined with the planar reflection result using a `lerp()` function.

In addition to the data related to the local correction, the vertex shader also calculates the screen coordinates of the vertex using the built-in function `ComputeScreenPos()`. It passes these coordinates to the fragment shader:

```
vertexOutput vert(vertexInput input)
{
    vertexOutput output;
    output.tex = input.texcoord;
    // Transform vertex coordinates from local to world.
```

```

float4 vertexWorld = mul(_Object2World, input.vertex);

// Transform normal to world coordinates.
float4 normalWorld = mul(float4(input.normal, 0.0), _World2Object);

// Final vertex output position.
output.pos = mul(UNITY_MATRIX_MVP, input.vertex);

// ----- Local correction -----
output.vertexInWorld = vertexWorld.xyz;
output.viewDirInWorld = vertexWorld.xyz - _WorldSpaceCameraPos;
output.normalInWorld = normalWorld.xyz;

// ----- Planar reflections -----
output.vertexInScreenCoords = ComputeScreenPos(output.pos);
return output;
}

```

The planar reflections are rendered to a texture, so the fragment shader is able to access the screen coordinates of the fragment. To enable this, pass the vertex screen coordinates to the fragment shader as a varying.

In the fragment shader:

- Apply the local correction to the reflection vector.
- Retrieve the color of the environment reflections `staticReflColor` from the local cubemap.

The following code shows how to combine static environment reflections, using the local cubemap technique, with dynamic planar reflections, that are rendered at runtime using the mirrored camera technique:

```

float4 frag(vertexOutput input) : COLOR
{
    float4 staticReflColor = float4(1, 1, 1, 1);

    // Find reflected vector in WS.
    float3 viewDirWS = normalize(input.viewDirInWorld);
    float3 normalWS = normalize(input.normalInWorld);
    float3 reflDirWS = reflect(viewDirWS, normalWS);

    // Working in World Coordinate System.
    float3 localPosWS = input.vertexInWorld;
    float3 intersectMaxPointPlanes = (_BBoxMax - localPosWS) / reflDirWS;
    float3 intersectMinPointPlanes = (_BBoxMin - localPosWS) / reflDirWS;

    // Look only for intersections in the forward direction of the ray.
    float3 largestParams = max(intersectMaxPointPlanes, intersectMinPointPlanes);

    // Smallest value of the ray parameters gives us the intersection.
    float distToIntersect = min(min(largestParams.x, largestParams.y), largestParams.z);

    // Find the position of the intersection point.
    float3 intersectPositionWS = localPosWS + reflDirWS * distToIntersect;

    // Get local corrected reflection vector.
    float3 localCorrReflDirWS = intersectPositionWS - _EnvCubeMapPos;

    // Lookup the environment reflection texture with the right vector.
    float4 staticReflColor = texCUBE(_Cube, localCorrReflDirWS);

    // Lookup the planar runtime texture
    float4 dynReflColor = tex2DProj(_ReflectionTex,
        UNITY_PROJ_COORD(input.vertexInScreenCoords));

    // Revert the blending with the background color of the reflection camera
    dynReflColor.rgb /= (dynReflColor.a < 0.00392)?1:dynReflColor.a;

    // Combine static environment reflections with dynamic planar reflections
    float4 combinedRefl = lerp(staticReflColor.rgb, dynReflColor.rgb, dynReflColor.a);

    // Lookup the texture color.
    float4 texColor = tex2D(_MainTex, float2(input.tex));
    return _AmbientColor + texColor * _ReflAmount * combinedRefl;
}

```

Extract the texture color `dynReflColor`, from the planar runtime reflection texture `_ReflectionTex`.

Declare `_ReflectionTex` as a uniform in the shader.

Declare `_ReflectionTex` in the Property Block. This enables you to see how it looks at runtime and this assists you with debugging while you are developing your game.

For the texture lookup, project the texture coordinates, that is, divide the texture coordinates by the last component of the coordinate vector. You can use the Unity built-in function `UNITY_PROJ_COORD()` for this.

Use the `lerp()` function to combine the static environment reflections and the dynamic planar reflections. Combine the following:

- The reflection color.
- The texture color of the reflective surface.
- The ambient color component.

6.3.3 Combining reflections from a distant environment

When you render reflections from static and dynamic objects, you might also have to consider reflections from a distant environment. For example, reflections from the sky that are visible through a window in your local environment.

In this case you must combine three different types of reflections:

- Reflections from the static environment using the local cubemap technique.
- Planar reflections from dynamic objects using the mirrored camera technique.
- Reflections from the skybox using the standard cubemap technique. The reflection vector does not require a correction before fetching the texture from the cubemap.

To incorporate the reflections from a skybox, use the reflection vector `reflDirWS` to fetch the texel from the skybox cubemap. Pass the skybox cubemap texture as a uniform to the shaders.

Note

Do not apply a local correction.

To ensure the skybox is only visible from the windows, render the transparency of the scene in the alpha channel when you are baking the static cubemap for reflections.

Assign a value of one to opaque geometry and a value of zero where there is no geometry, or the geometry is fully transparent. For example, render the pixels corresponding to the windows with zero in the alpha channel.

Pass the skybox cubemap `_Skybox` to the shaders, as a uniform.

In the fragment shader code in [6.3.2 Combining reflections shader implementation](#) on page 6-116, look for this comment:

```
// Lookup the planar runtime texture
```

Insert the following lines before the comment:

```
float4 skyboxReflColor = texCUBE(_Skybox, reflDirWS);
staticReflColor = lerp(skyboxReflColor.rgb, staticReflColor.rgb, staticReflColor.a);
```

This code combines the static reflections with reflections from the skybox.

The following figure shows combining different types of reflections:



Figure 6-31 Combining different types of reflections

6.4 Dynamic soft shadows based on local cubemaps

This technique uses a local cubemap to hold a texture that represents transparency of the static environment. This is a very efficient technique that generates high-quality soft shadows.

This section contains the following subsections:

- [6.4.1 About dynamic soft shadows based on local cubemaps on page 6-120](#).
- [6.4.2 Generating shadow cubemaps on page 6-120](#).
- [6.4.3 Rendering shadows on page 6-121](#).
- [6.4.4 Combining cubemap shadows with a shadow map on page 6-124](#).
- [6.4.5 Results of the cubemap shadow technique on page 6-125](#).

6.4.1 About dynamic soft shadows based on local cubemaps

In your scene there are moving objects and static environments such as rooms. By using this technique, you are not required to render static geometry to a shadow map every frame. This enables you to use a texture to represent the shadows.

Cubemaps can be a good approximation of many kinds of a static local environment including irregular shapes such as the cave in the Ice Cave demo. The alpha channel can also represent the amount of light entering the room.

The objects that move are typically everything except the room. Objects such as:

- The sun.
- The camera.
- Dynamic objects.

With the whole room represented by a cube texture you can access arbitrary texels of the environment within a fragment shader. For example, this means the sun can be in any arbitrary position and you can calculate the amount of light reaching a fragment based on the value fetched from the cubemap.

The alpha channel or transparency, represents the amount of light entering the local environment. In your scene, attach the cubemap texture to the fragment shaders that render the static and dynamic objects that you want to add shadows to.

6.4.2 Generating shadow cubemaps

Start with a local environment to which you want to apply shadows from light sources outside of the environment. For example, a room, a cave, or a cage.

This technique is similar to reflections based on a local cubemap. For more information see [6.2 Implementing reflections with a local cubemap on page 6-98](#).

Create the shadow cubemap in the same manner as the reflection cubemap but you must also add an alpha channel. The alpha channel or transparency, represents the amount of light entering the local environment.

Work out the position from where to render the six faces of the cubemap. In most cases this is in the center of the bounding box of the local environment. You require this position for the generation of the cubemap. You must also pass this position to the shaders to calculate a local correction vector to fetch the right texel from the cubemap.

When you have decided where to position the center of the cubemap, you can render all the faces to the cubemap texture and record the transparency or alpha of the local environment. The more transparent an area is, the more light that comes into the environment. If there is no geometry it is fully transparent. If required, you can use the RGB channels to store the color of the environment for colored shadows such as stained glass, reflections, or refractions.

6.4.3 Rendering shadows

Build a vector P_iL in world space from a vertex or fragment, to the light or lights, and fetch the cubemap shadow using this vector.

Before fetching each texel, you must apply local correction to the P_iL vector. ARM recommends you make the local correction in the fragment shader to obtain more accurate shadows.

To compute the local correction, you must calculate an intersection point of the fragment-to-light vector with the bounding box of the environment. Use this intersection point to build another vector from the cubemap origin position C to the intersection point P . This gives you the final vector CP that you use to fetch the texel.

You require the following input parameters to calculate the local correction:

- $_EnviCubeMapPos$ The cubemap origin position.
- $_BBoxMax$ The maximum point of the bounding box of the environment.
- $_BBoxMin$ The minimum point of the bounding box of the environment.
- P_i The fragment position in world space.
- PiL The normalized fragment-to-light vector in world space.

Compute the output value CP . This is the corrected fragment-to-light vector that you use to fetch a texel from the shadow cubemap.

The following figure shows the local correction of the fragment-to-light vector.

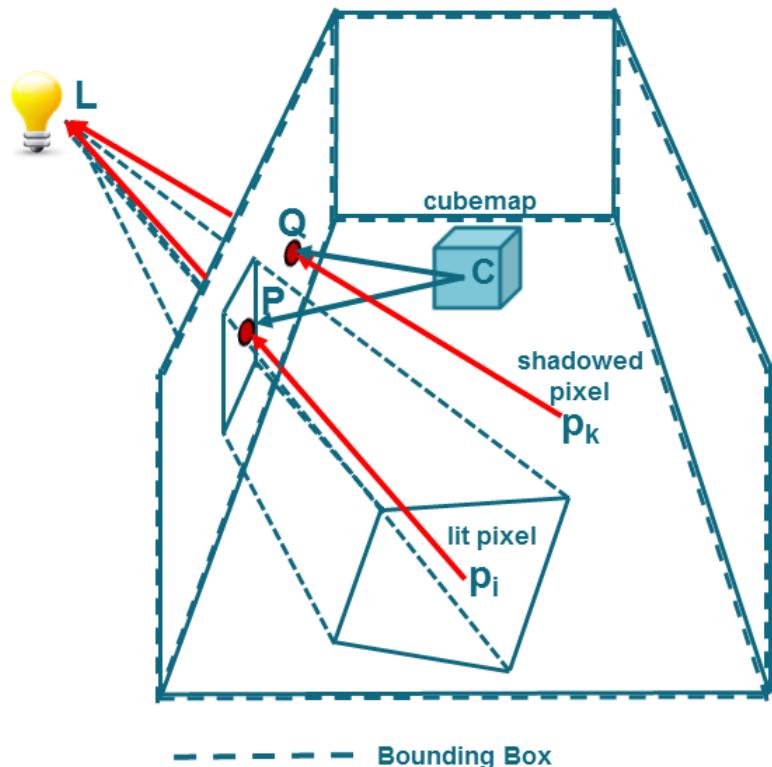


Figure 6-32 Local correction of the fragment-to-light vector

The following example code shows how to calculate the correct CP vector:

```
// Working in World Coordinate System.
vec3 intersectMaxPointPlanes = (_BBoxMax - Pi) / PiL;
vec3 intersectMinPointPlanes = (_BBoxMin - Pi) / PiL;

// Looking only for intersections in the forward direction of the ray.
vec3 largestRayParams = max(intersectMaxPointPlanes, intersectMinPointPlanes);

// Smallest value of the ray parameters gives us the intersection.
```

```
float dist = min(min(largestRayParams.x, largestRayParams.y), largestRayParams.z);

// Find the position of the intersection point.
vec3 intersectPositionWS = Pi + PiL * dist;

// Get the local corrected vector.
CP = intersectPositionWS - _EnvCubeMapPos;
```

Use the CP vector to fetch a texel from the cubemap. The alpha channel of the texel provides information about how much light or shadow you must apply to a fragment:

```
float shadow = texCUBE(cubemap, CP).a;
```

The following figure shows a chess room with hard shadows:



Figure 6-33 Chess room with hard shadows

This technique generates working shadows in your scene, but you can improve the quality of the shadows with two more steps:

- Back faces in shadow
- Smoothness

Back faces in shadow

The cubemap shadow technique does not use depth information to apply shadows. This means that some faces are incorrectly lit when they are meant to be in shadow.

The problem only occurs when a surface is facing in the opposite direction to the light. To fix this problem, check the angle between the normal vector and the fragment-to-light vector, P_iL . If the angle, in degrees, is outside of the range -90 to 90, the surface is in shadow.

The following code snippet does this check:

```
if (dot(PiL, N) < 0)
    shadow = 0.0;
```

The previous code causes each triangle into a hard switch from light to shade. For a smooth transition use the following formula:

```
shadow *= max(dot(PiL, N), 0.0);
```

Where:

- shadow is the alpha value fetched from the shadow cubemap.
- P_iL is the normalized fragment-to-light vector in world space.
- N is the normal vector of the surface in world space.

The following figure shows a chess room with back faces in shadow:



Figure 6-34 Chess room with back faces in shadow

Smoothness

This shadow technique can provide realistic soft shadows in your scene.

1. Generate mipmaps and set trilinear filtering for the cubemap texture.
2. Measure the length of a fragment-to-intersection-point vector.
3. Multiply the length by a coefficient.

The coefficient is a normalizer of a maximum distance in your environment to the number of mipmap levels. You can calculate this automatically against bounding volume and mipmap levels. You must customize the coefficient to your scene. This enables you to tweak the settings to suit the environment, improving visual quality. For example, the coefficient used in the Ice Cave project is **0.08**.

You can reuse the results of calculations that you did for the local correction. Reuse `dist` from the code snippet for local correction as the length of the segment from the fragment position to the intersection point of the fragment-to-light vector with the bounding box:

```
float texLod = dist;
```

Multiply `texLod` by the distance coefficient:

```
texLod *= distanceCoefficient;
```

To implement softness, fetch the correct mipmap level of the texture using the Cg function `texCUBEElod()` or the GLSL function `textureLod()`.

Construct a `vec4` where `XYZ` represents a direction vector and the `W` component represents the LOD.

```
CP.w = texLod;
shadow = texCUBEElod(cubemap, CP).a;
```

This technique provides high-quality, smooth shadows for your scene.

The following figure shows a chess room with smooth shadows:



Figure 6-35 Smooth shadows

6.4.4 Combining cubemap shadows with a shadow map

To get shadows complete with dynamic content, you must combine cubemap shadows with a traditional shadow map technique. This is more work but it is still worth it because you are only required to render dynamic objects to the shadow map.

The following figure shows the chess room with smooth shadows only:



Figure 6-36 Smooth shadows

The following figure shows the chess room with smooth shadows combined with dynamic shadows:



Figure 6-37 Smooth shadows combined with dynamic shadows

6.4.5 Results of the cubemap shadow technique

In traditional techniques, rendering shadows can be quite expensive because it involves rendering the whole scene from the perspective of each shadow-casting light source. The cubemap shadow technique described here delivers improved performance because it is mostly pre-baked.

This technique is also independent of output-resolution. It produces the same visual quality at 1080p, 720p and other resolutions.

The softness filtration is calculated in hardware so the smoothness comes at almost no computational cost. The smoother the shadow, the more efficient the technique is. This is because the smaller mipmap levels result in less data than traditional shadow map techniques. Traditional techniques require a large kernel to make shadows smooth enough to be more visually appealing. This requires high memory bandwidth and this reduces performance.

The quality you get with the cubemap shadow technique is higher than you might expect. It provides realistic softness and stable shadows with no shimmering at the edges. Shimmering edges can be observed when using traditional shadow map techniques because of rasterization and aliasing effects. However, none of the anti-aliasing algorithms can fix this problem entirely.

The cubemap shadow technique does not have a shimmering problem. The edges are stable even if you use a much lower resolution than that used in the render target. You can use four times lower resolution than the output and there are neither artifacts nor unwanted shimmering. Using four times lower resolution also saves memory bandwidth and this improves performance.

This technique can be used on any device on the market that supports shaders such as those with OpenGL ES 2.0 or higher. If you already know where and when to use the reflections based on local cubemaps technique, then you can easily apply the shadow technique to your implementation.

————— Note ————

The technique cannot be used for everything in your scene. Dynamic objects, for instance, receive shadows from the cubemap but they cannot be pre-baked to the cubemap texture. For dynamic objects, use shadow maps for generating shadows, blended with the cubemap shadow technique.

The following figure shows the Ice Cave with shadows:



Figure 6-38 Ice Cave with shadows

The following figure shows the Ice Cave with smooth shadows:



Figure 6-39 Ice Cave with smooth shadows

The following figure shows the Ice Cave with smooth shadows:

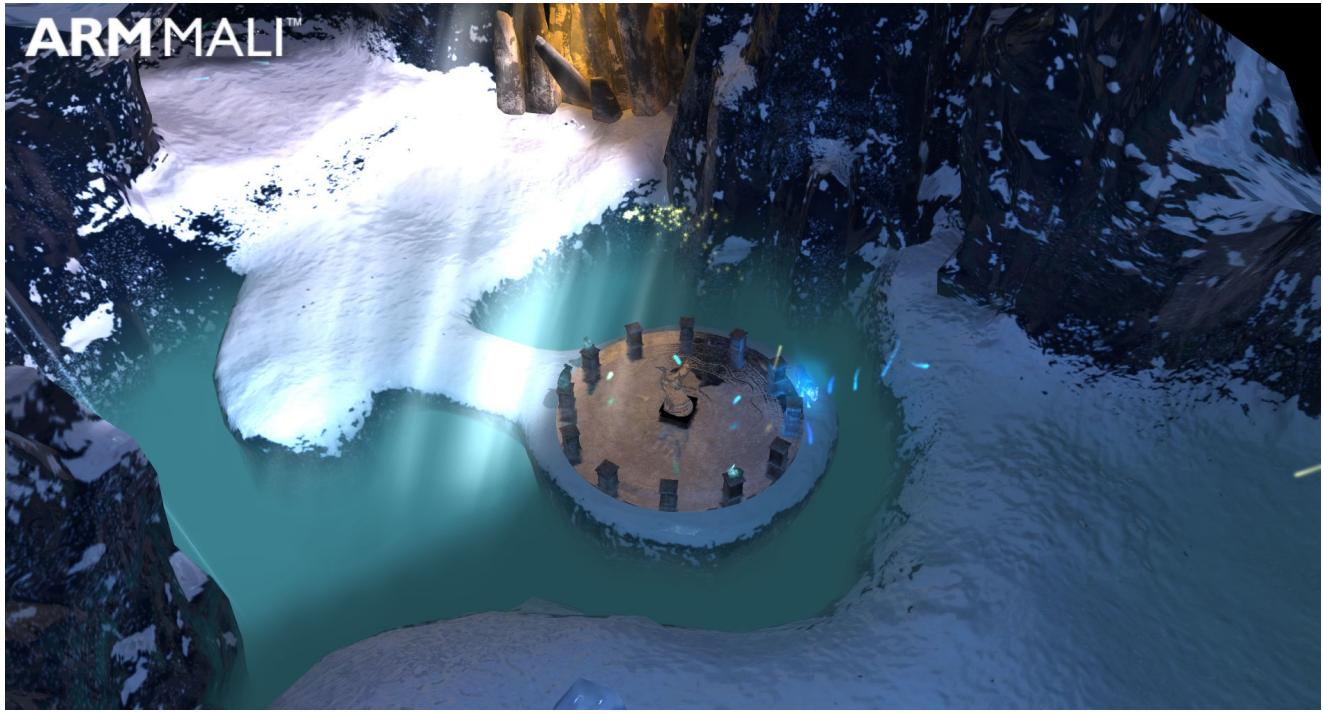


Figure 6-40 Ice Cave with smooth shadows

6.5 Refraction based on local cubemaps

You can use local cubemaps to implement high quality refractions. You can combine these with reflections at runtime.

This section contains the following subsections:

- [6.5.1 About refractions on page 6-128](#).
- [6.5.2 Refraction implementations on page 6-128](#).
- [6.5.3 About refractions based on local cubemaps on page 6-129](#).
- [6.5.4 Preparing the cubemap on page 6-129](#).
- [6.5.5 Shader implementation on page 6-131](#).

6.5.1 About refractions

Game developers are regularly looking for efficient methods to implementing visually impressive effects in their games. This is especially important when targeting mobile platforms because you must carefully balance resources to achieve maximum performance.

Refraction is the change in direction of a light wave because of a change in the medium it is passing through. If you want extra realism with semi-transparent geometry, refraction is an important effect to consider.

The refractive index determines how much light is bent, or refracted, when entering a material. Refraction is defined as the bending of light as it passes from one medium, with refractive index n_1 to another medium with refractive index n_2 .

You use Snell's Law to calculate the relationship between the refractive indices and the sine of the incident and refracted angles.

The following figure shows Snell's Law and refraction:

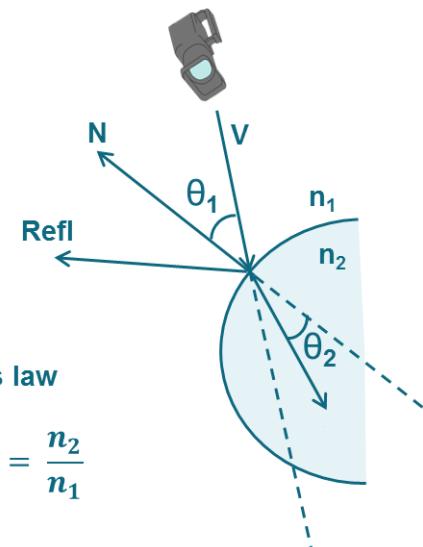


Figure 6-41 Refraction of light as it passes through one medium to another

6.5.2 Refraction implementations

Developers have tried to render refraction since they started to render reflections because these processes take place together in any semi-transparent surface. There are several techniques for rendering reflections but not many for refraction.

Existing methods for implementing refraction at runtime differ depending on the specific type of refraction. Most of the techniques render the scene behind the refractive object to a texture at runtime, and then apply a texture distortion in a second pass to achieve the refracted appearance. Depending on

the texture distortion, you can use this approach to render refraction effects such as water, heat haze, glass, and other effects.

Some of these techniques can achieve good results, but the texture distortion is not physically based so the results are not always correct. For example, if you render a texture from the point of view of the refraction camera, there might be areas that are not directly visible to the camera but are visible in a physically based refraction.

The main limitation of using render-to-texture methods is quality. When the camera is moving, pixel shimmering or pixel instability is often visible.

6.5.3 About refractions based on local cubemaps

Local cubemaps are an excellent technique for rendering reflections and developers have used static cubemaps to implement both reflections and refractions since they became available.

However, if you use static cubemaps to implement reflections or refractions in a local environment, the results are incorrect if you do not apply a local correction.

In the technique described here a local correction is applied to ensure correct results. This technique is highly optimized. It is especially useful for mobile devices where runtime resources are limited so must be carefully balanced.

6.5.4 Preparing the cubemap

You must prepare the cubemap to be used in the refraction implementation:

To prepare the cubemap, do the following:

1. Place a camera in the center of the refractive geometry.
2. Hide the refractive object and render the surrounding static environment to a cubemap in the six directions. You can use this cubemap for implementing both refraction and reflection.
3. Bake the environment surrounding the refractive object into a static cubemap.
4. Determine the direction of the refraction vector, and find where it intersects with the bounding box of the local environment.
5. Apply the local correction in the same way as [6.4 Dynamic soft shadows based on local cubemaps on page 6-120](#).
6. Build a new vector from the position where the cubemap was generated, to the intersection point. Use this final vector to fetch the texel from the cubemap, to render what is behind the refractive object.

Instead of fetching the texel from the cubemap using the refracted vector R_{rf} , you find the point P where the refracted vector intersects the bounding box and build a new vector R'_{rf} from the center of the cubemap C to the intersection point P. Use this new vector to fetch the texture color from the cubemap.

```
float eta=n2/n1;  
float3Rrf = refract(D,N,eta);  
Find intersection point P  
Find vector R'rf = CP;  
Float4 col = texCube(Cubemap, R'rf);
```

The following figure shows a scene with a cubemap and the refraction vectors:

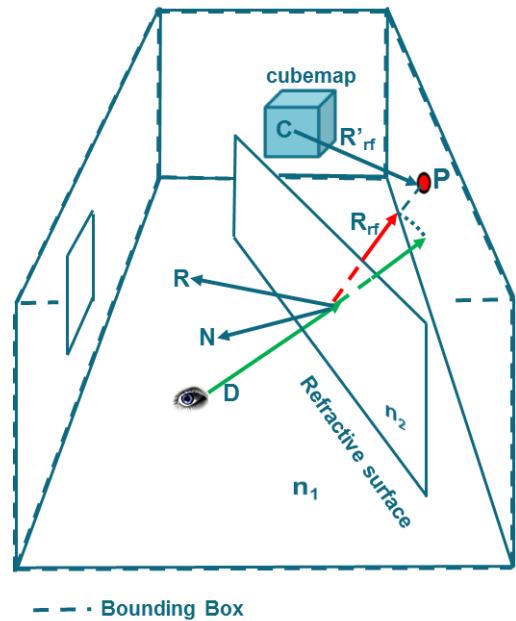


Figure 6-42 The local correction to refraction vector

The refraction produced by this technique is accurately physically based, because the direction of the refraction vector is calculated using Snell's Law.

There is also a built-in function that you can use in your shader to find the refraction vector R strictly according to the Snell's law:

```
R = refract( I, N, eta);
```

Where:

- I is the normalized view or incident vector.
- N is the normalized normal vector.
- eta is the ratio of indices of refractions n_1/n_2 .

The following figure shows the flow of shaders that implement refraction based on a local cubemap:

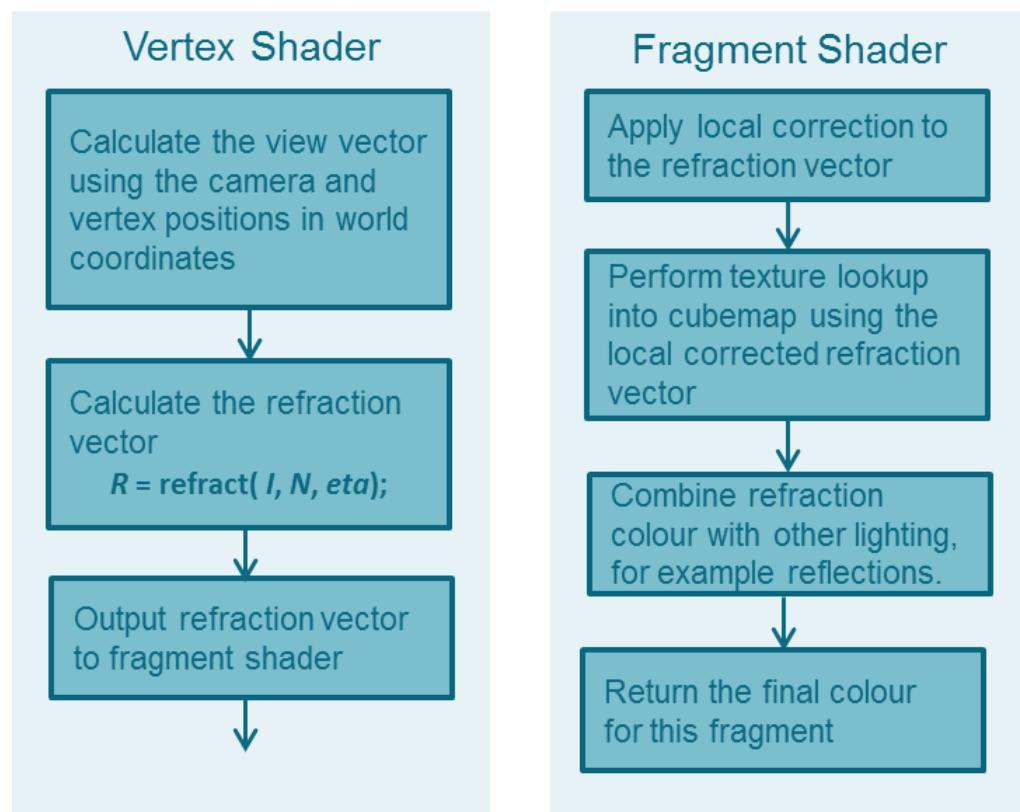


Figure 6-43 Shader implementations of refraction based on local cubemap

6.5.5 Shader implementation

When you fetch the texel corresponding to the locally-corrected refraction direction, you might want to combine the refraction color with other lighting. For example, reflections that take place simultaneously with refraction.

To combine the refraction color with other lighting, you must pass an additional view vector to the fragment shader and apply the local correction to it. Use the result to fetch the reflection color from the same cubemap.

The following code snippet shows how to combine reflection and refraction to produce the final output color:

```

// ----- Environment reflections -----
float3 newReflDirWS = LocalCorrect(input.reflDirWS, _BBoxMin, _BBoxMax, input.posWorld,
_EnviCubeMapPos);
float4 staticReflColor = texCUBE(_EnviCubeMap, newReflDirWS);
// ----- Environment refractions -----
float3 newRefractDirWS = LocalCorrect(RefractDirWS, _BBoxMin, _BBoxMax, input.posWorld,
_EnviCubeMapPos);
float4 staticRefractColor = texCUBE(_EnviCubeMap, newRefractDirWS);
// ----- Combined reflections and refractions -----
float4 combinedReflRefract = lerp(staticReflColor, staticRefractColor, _ReflAmount);

float4 finalColor = _AmbientColor + combinedReflRefract;

```

The coefficient `_ReflAmount` is passed as a uniform to the fragment shader. Use this coefficient to adjust the balance between reflection and refraction contributions. You can manually adjust `_ReflAmount` to achieve the visual effect you require.

You can find the implementation of the `LocalCorrect` function in the reflections blog at: <http://community.arm.com/groups/arm-mali-graphics/blog/2014/08/07/reflections-based-on-local-cubemaps>.

When the refractive geometry is a hollow object, refractions and reflections take place in both the front and back surfaces.

The following figure shows refraction on a glass bishop based on a local cubemap:



Figure 6-44 Refraction on a glass bishop based on a local cubemap

The image on the left shows the first pass that renders only back faces with local refraction and reflections.

The image on the right shows the second pass renders only front faces with local refraction and reflections and alpha blending with the first pass.

- In the first pass, render the semi-transparent object in the same manner that you render opaque geometry. Render the object last with front-face culling on, that is, only render the back faces. You do not want to occlude other objects so do not write to the depth buffer.

The color of the back face is obtained by mixing the colors calculated from the reflection, refraction, and the diffuse color of the object itself.

- In the second pass, render the front faces with back face culling, do this last in the rendering queue. Ensure depth writing is off. Obtain the front-face color by mixing the refraction and reflection textures with the diffuse color. The refraction in the second pass adds more realism to the final rendering. You can skip this step if the refraction on the back faces is enough to highlight the effect.
- In the final pass you alpha-blend the resulting color with the first pass.

The following figure shows the result of implementing refractions based on a local cubemap, on a semi-transparent phoenix in the Ice Cave demo.



Figure 6-45 Semi-transparent phoenix refractions

The following figure shows a semi-transparent phoenix wing:



Figure 6-46 Semi-transparent phoenix wing

6.6 Specular effects in the Ice Cave demo

Specular effects in the Ice Cave demo are implemented using the Blinn technique. This is a very efficient technique that produces good results.

The following code shows how to implement specular effects with the Blinn technique:

```
// Returns intensity of a specular effect without taking into account shadows
float SpecularBlinn(float3 vert2Light, float3 viewDir, float3 normalVec, float4 power)
{
    float3 floatDir = normalize(vert2Light - viewDir);
    float specAngle = max(dot(floatDir, normalVec), 0.0);
    return pow(specAngle, power);
}
```

A downside of the Blinn technique is that it can produce incorrect results in certain circumstances. For example, specular effects can appear in regions that are in shadow.

The following figure shows an example of a shadowed region with no specular effects:



Figure 6-47 Shadowed region

The following figure shows an example of specular effects in a shadowed region. These are incorrect:



Figure 6-48 Shadowed region with incorrect specular effects

The following figure shows an example of specular effects in a lit region. This is correct:



Figure 6-49 Lit region with correct specular effects

The shadows should make the specular effects intensity either stronger or weaker, depending on the light reaching the specular surface. All the information required to correct the intensity of the specular effects is already in the Ice Cave demo, so fixing it is relatively simple. The environment cubemap texture that is used for reflection and shadow effects contains two types of information. The RGB channels contain environment colors that are used for reflections. The alpha channel contains opacity and this is used for shadows. You can use the alpha channel to determine the specular intensity because the alpha channel

represents holes in the cave that let light in to the environment. The alpha channel is therefore used to ensure the specular effect is applied only to surfaces that are lit.

To do this, calculate a corrected reflection vector in a fragment shader to make the reflection effect. For more information about creating the corrected reflection vector see [6.2 Implementing reflections with a local cubemap on page 6-98](#). Use this vector to fetch the RGBA texel from a cubemap texture:

```
// Locally corrected static reflections
const half4 reflColor = SampleCubemapWithLocalCorrection(
    ReflDirectionWS,
    _ReflBoundingBoxMinWorld,
    _ReflBoundingBoxMaxWorld,
    input.vertexInWorld,
    _ReflCubePosWorld,
    _ReflCube);
```

For more information about what the `SampleCubemapWithLocalCorrection()` function, see [6.2 Implementing reflections with a local cubemap on page 6-98](#).

`reflColor` is in RGBA format where the RGB components contain color data for reflections and the alpha channel contains the intensity of the specular effect. In the Ice Cave demo, the alpha channel is multiplied by the specular color, that is calculated with the Blinn technique:

```
half3 specular = _SpecularColor.rgb *
    SpecularBlinn(
        input.vertexToLight01InWorld,
        viewDirInWorld,
        normalInWorld,
        _SpecularPower) *
    reflColor.a;
```

The `specular` value represents the final specular color. You can add this to your lighting model.

6.7 Using Early-z

Mali GPUs include the ability to do an Early-Z algorithm. Early-Z improves performance by removing overdrawn fragments.

The Mali GPU typically executes the Early-Z algorithm on most content, but there are cases where, to preserve correctness, it is not executed. This can be difficult to control from within Unity because it depends on both the Unity engine and the code generated by the compiler, however there are some signs you can look for.

Compile your shader for mobile and have a look at the code. Ensure your shader does not fall into one of the following categories:

Shader has side effects

This means a shader thread modifies global state during its execution, so executing the shader a second time might produce different results. Typically this means that your shader writes to a shared read/write memory buffer such as shader storage buffer objects or images. For example, if you create a shader that increments a counter to measure performance, this shader has side effects.

The following are not classed as side effects:

- Read-only memory accesses.
- Writes to write-only buffers.
- Purely local memory accesses.

Shader calls discard()

If the fragment shader can call `discard()` during its execution, the Mali GPU cannot enable Early-Z. This is because the fragment shader can discard the current fragment, but the depth value was previously modified by the Early-Z test and this cannot be reverted.

Alpha-to-coverage is enabled

If Alpha-to-coverage is enabled, the fragment shader computes data that is later accessed to define the alpha.

For example, when rendering the leaves of a tree, they are typically represented as a plane and the texture defines what region of the leaf that is transparent or opaque. If early-z is enabled you get incorrect results because part of the scene can be occluded by a transparent part of the plane.

Depth source not fixed function

The depth value used for depth testing does not come from the vertex shader. If your fragment shader writes to `gl_FragDepth`, the Mali GPU cannot perform the Early-Z test.

6.8 Dirty lens effect

You can use a dirty lens effect to invoke a sense of drama. It is often used together with a lens flare effect.

You can implement a dirty lens effect in a very light and simple way that is suitable for mobile devices.

In the Ice Cave demo, the dirty lens effect is implemented in a minimal shader that renders a variable intensity, full screen quad on top of the scene. The intensity of the quad is passed to it by a script.

The shader renders the quad using additive alpha blending, at the very end, after all transparent geometry has been rendered.

This section contains the following subsections:

- [6.8.1 Shader implementation](#) on page 6-138.
- [6.8.2 Script implementation](#) on page 6-140.
- [6.8.3 Dirty lens shader code](#) on page 6-140.

6.8.1 Shader implementation

This section describes the dirty lens shader.

For the full source code of the dirty lens shader, see [6.8.3 Dirty lens shader code](#) on page 6-140.

The following sub-shader tag instructs the full screen quad to be rendered after all opaque geometry, and after nine other transparent objects:

```
Tags {"Queue" = "Transparent+10"}
```

The shader uses the following command to deactivate depth buffer writing. This prevents the quad from occluding the geometry behind it:

```
ZWrite Off
```

At the blending stage, the output of the fragment shader is blended with the pixel color already in the frame buffer.

The shader specifies the additive blending type `Blend One One`. In this blending type source and destination factors are both `float4 (1.0, 1.0, 1.0, 1.0)`.

This type of blending is often used for particle systems when they represent effects like fire that are transparent and emit light.

The shader deactivates culling and ZTest to ensure the dirty lens effect is always rendered.

The vertices of the quad are defined in viewport coordinates, so no vertex transformations take place in the vertex shader.

The fragment shader only fetches the texel from the texture and applies a factor that is proportional to the intensity of the effect.

The following figure shows the image used as the full screen quad texture for the dirty lens effect:



Figure 6-50 Texture used for the Dirty Lens effect in Ice Cave demo

The following figure shows how the Dirty Lens effect looks in the Ice Cave demo when the camera is oriented to the sun light coming from the entrance of the cave:



Figure 6-51 Dirty Lens effect as implemented in the Ice Cave demo

6.8.2 Script implementation

A simple script creates the full screen quad and calculates the intensity factor passed to the fragment shader.

The following function creates the mesh of the quad in the Start function:

```
void CreateQuadMesh()
{
    Mesh mesh = GetComponent<MeshFilter>().mesh;
    mesh.Clear();
    mesh.vertices = new Vector3[] {new Vector3(-1, -1, 0), new Vector3(1, -1, 0),
        new Vector3(1, 1, 0), new Vector3(-1, 1, 0)};
    mesh.uv = new Vector2[] {new Vector2(0, 0), new Vector2(1, 0),
        new Vector2(1, 1), new Vector2(0, 1)};
    mesh.triangles = new int[] {0, 2, 1, 0, 3, 2};
    mesh.RecalculateNormals();

    //Increase bounds to avoid frustum clipping.
    bigBounds.SetMinMax(new Vector3(-100, -100, -100), new Vector3(100, 100, 100));
    mesh.bounds = bigBounds;
}
```

When the mesh is created, its bounds are incremented to guarantee that it is never frustum clipped. Set the size of the bounds according to the size of your scene.

The script calculates the intensity factor and passes it to the fragment shader. This is based on the relative orientation of the camera-to-sun vector and the camera forward vector. The maximum intensity of the effect takes place when the camera is looking directly at the sun.

The following code shows the calculation of the intensity factor:

```
Vector3 cameraSunVec = sun.transform.position - Camera.main.transform.position;
cameraSunVec.Normalize();
float dotProd = Vector3.Dot(Camera.main.transform.forward, cameraSunVec);
float intensityFactor = Mathf.Clamp(dotProd, 0.0f, 1.0f);
```

6.8.3 Dirty lens shader code

This is the code of DirtyLensEffect.shader:

```
half3 normalInWorld = half3(0.0,0.0,0.0);
half3 bumpNormal = UnpackNormal(tex2D(_BumpMapGlobal, input.tc));
half3x3 local2WorldTranspose = half3x3(input.tangentWorld,
    input.bitangentWorld,
    input.normalInWorld);
normalInWorld = normalize(mul(bumpNormal, local2WorldTranspose));
normalInWorld = normalInWorld*0.5 + 0.5;
return half4(normalInWorld,1.0);
```

6.9 Light shafts

Light shafts simulate the effect of crepuscular rays, atmospheric scattering or shadowing. They are used to add depth and realism to a scene.

This section contains the following subsections:

- [6.9.1 About light shafts on page 6-141](#).
- [6.9.2 Fading out the cone edges on page 6-143](#).

6.9.1 About light shafts

In the Ice Cave demo, the light shaft simulates the scattering of the sun rays coming into the cave from the opening at the top of the cave.

The following figure shows the light shafts in the Ice Cave demo:

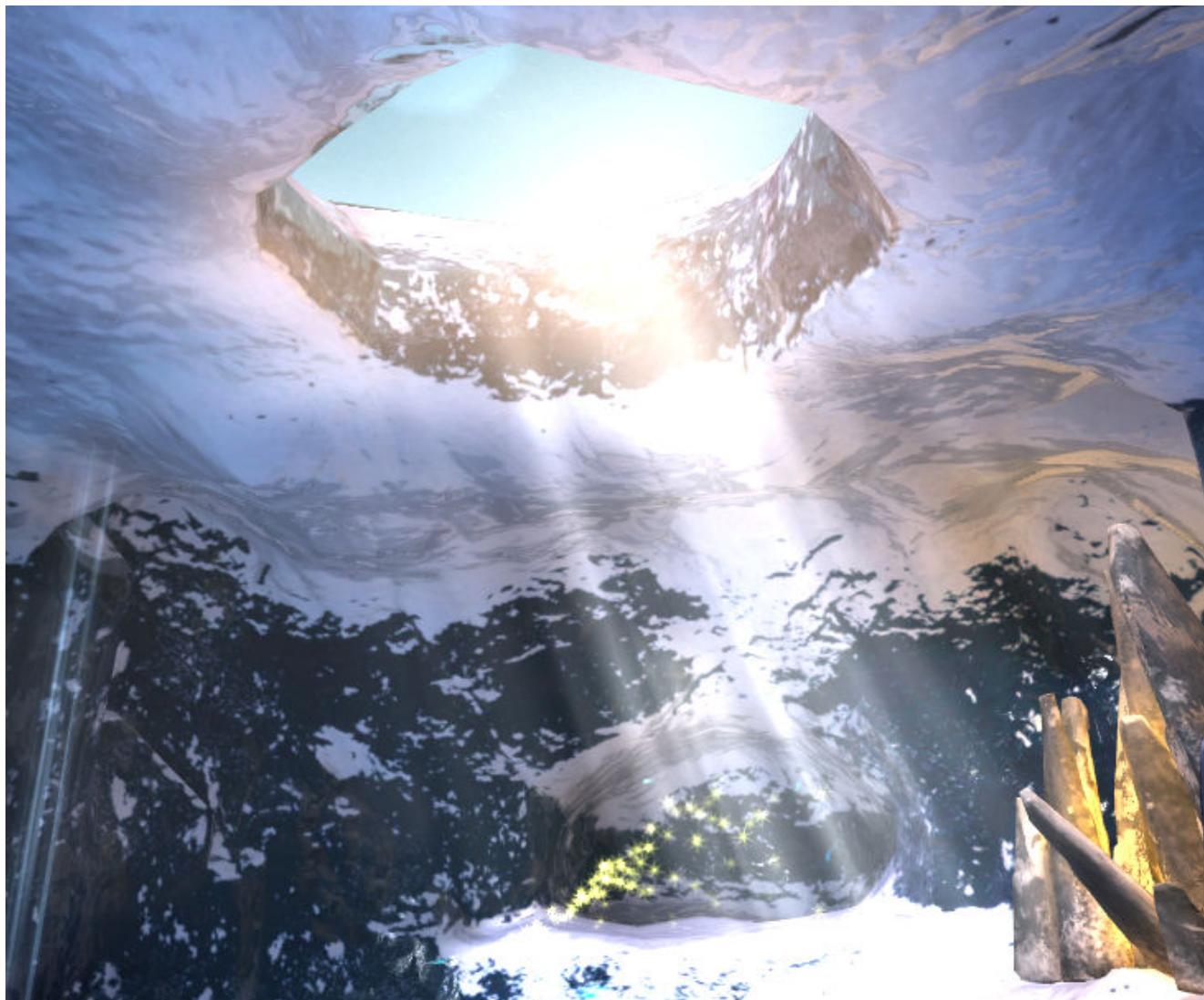


Figure 6-52 Light shaft in the Ice Cave demo

The light shaft is based on a truncated cone. The cone follows the light rays orientation in a manner that ensures the upper section of the cone is always fixed.

The following figure shows the geometry of the cone where:

- a shows that the basic geometry of the light shaft is a cylinder with two cross-sections at the top and bottom.
- b shows the lower cross-section is expanded to achieve a truncated cone geometry based on the angle θ that you define.
- c shows that the upper cross-section is fixed, and the lower cross-section is shifted horizontally according to the direction of the sun rays.

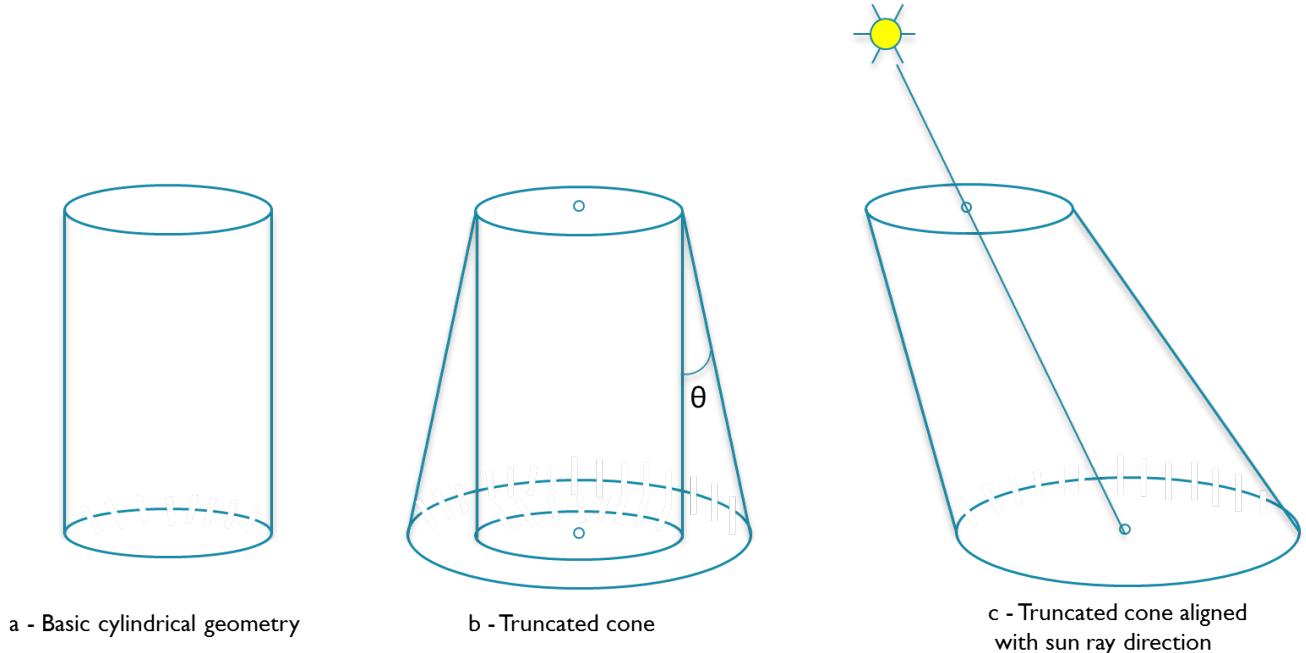


Figure 6-53 Light shaft geometry

A script uses the position of the Sun to calculate the following values:

- The magnitude of the lower cross-section cone expansion based on the value of the angle θ provided as input.
- The direction and magnitude of the cross-section shift.

The vertex shader applies a transformation based on this data. The transformation is applied to the original vertices of the cylindrical geometry, in the local coordinates of the light shaft.

When rendering the light shaft avoid rendering any hard edges that reveal its geometry. You can achieve this by using a texture mask that fades out smoothly the top and the bottom.

The following figure shows the light shaft textures:



Figure 6-54 Light shaft textures. Left mask texture. Right beam texture

6.9.2 Fading out the cone edges

In the plane parallel to the cross section, fade out the light shaft intensity depending on the relative camera-to-vertex orientation.

In the vertex shader, project the camera position on the cross section.

Build a new vector from the center of the cross section to the projection, and normalize it.

Calculate the dot product of this vector with the vertex normal, then raise the result to a power exponent.

Pass the result to the fragment shader as a varying. This is used to modulate the intensity of the light shaft.

The vertex shader performs these calculations in the *Local Coordinate System* (LCS).

The following code shows the vertex shader:

```
// Project camera position onto cross section
float3 axisY = float3(0, 1, 0);
float dotWithYAxis = dot(camPosInLCS, axisY);
float3 projOnCrossSection = camPosInLCS - (axisY * dotWithYAxis);
projOnCrossSection = normalize(projOnCrossSection);

// Dot product to fade the geometry at the edge of the cross section
float dotProd = abs(dot(projOnCrossSection, input.normal));
output.overallIntensity = pow(dotProd, _FadingEdgePower) * _CurrLightShaftIntensity;
```

The coefficient `_FadingEdgePower` enables you to fine-tune the fading of the light shaft edges.

The script passes the coefficient `_CurrLightShaftIntensity`. This makes the light shaft fade out when the camera gets close to it.

The following code shows a final touch that is added to the light shaft by slowly scrolling down the texture from a script:

```
void Update()
{
    float localOffset = (Time.time * speed) + offset;
    localOffset = localOffset % 1.0f;
    GetComponent<Renderer>().material.SetTextureOffset("_MainTex", new Vector2(0,
    localOffset));
}
```

The fragment shader fetches the beam and mask textures, and then combines them with the intensity factor:

```
float4 frag(vertexOutput input) : COLOR
{
    float4 textureColor = tex2D(_MainTex, input.tex.xy);
    float textureMask = tex2D(_MaskTex, input.tex.zw).a;
    textureColor *= input.overallIntensity * textureMask;
    textureColor.rgb = clamp(textureColor.a, 0, 1);
    return textureColor;
}
```

The light shaft geometry is rendered in the transparent queue after all the opaque geometry.

It uses additive blending to combine the fragment color with the corresponding pixel in the frame buffer.

The shader also disables culling and writing to depth buffer so it does not occlude other objects.

The settings for the pass are:

```
Blend One One
Cull Off
ZWrite Off
```

6.10 Fog effects

A fog effect adds atmosphere to a scene. You do not require an advanced implementation to generate fog, a simple fog effect can be effective.

This section contains the following subsections:

- [6.10.1 About fog effects on page 6-145](#).
- [6.10.2 Procedural linear fog on page 6-145](#).
- [6.10.3 Linear fog with height on page 6-146](#).
- [6.10.4 Non-uniform fog on page 6-147](#).
- [6.10.5 Pre-baked fog on page 6-147](#).
- [6.10.6 Volumetric fog using particles on page 6-147](#).

6.10.1 About fog effects

In the real world, the further you look in the distance the more colors fade. It does not have to be foggy for you to see this effect, you can also see it on sunny days and it is especially visible when looking at mountains.

This is very common in real life, so adding this effect to your game adds to the sense of realism.

This section describes two versions of the fog effect:

- Procedural linear fog.
- Particle-based fog.

You can apply both of these effects simultaneously. The Ice Cave demo uses both techniques.

6.10.2 Procedural linear fog

Ensure that the further away the object is, the more its colors fade to the defined fog color. To achieve this in your fragment shader, you can use simple linear interpolation between the fragment color and the fog color.

The following example code shows how to calculate the fog color in the vertex shader, based on the distance to the camera. This color is passed to fragment shader as a varying.

```
output.fogColor = _FogColor * clamp(vertexDistance * _FogDistanceScale, 0.0, 1.0);
```

Where:

- `vertexDistance` is the vertex to camera distance.
- `_FogDistanceScale` is a factor passed to the shader as a uniform.
- `_FogColor` is the base fog color you define and passed as a uniform.

In the fragment shader the interpolated `input.fogColor` is combined with the fragment color `output.Color`.

```
outputColor = lerp(outputColor, input.fogColor.rgb, input.fogColor.a);
```

The following figure shows the result you get in your scene:



Figure 6-55 Linear fog based on distance

————— Note ————

Calculating the fog in a shader means you are not required to do any extra post processing to produce the fog effect.

Try to merge as many effects as possible into one shader. However, ensure you check performance because if you exceed the cache, performance might reduce. You might have to split the shader into two or more passes.

You can do the fog color calculation in the vertex or fragment shader. The calculation in the fragment shader is more accurate but requires more compute performance because it is calculated for every fragment.

The vertex shader calculation is lower precision but is also higher performance because it is only computed once per vertex.

6.10.3 Linear fog with height

Fog is applied uniformly across your scene. You can make it more realistic by changing the density according to height.

Make the fog more dense lower down and thinner higher up. You can expose the height level, to adjust it manually.

The following figure shows linear fog based on distance and height:



Figure 6-56 Linear fog based on distance and height

6.10.4 Non-uniform fog

The fog does not have to be uniform. You can make fog more visually interesting by introducing some noise.

You can create non-uniform fog by applying a noise texture.

For a more complex effect, you can also apply more than one noise texture and make them slide with different speeds. For instance the noise texture that is further away slides slower than the texture that is closer to the camera.

You can apply more than one texture in one a single pass in one shader, just blend the noise textures according to the distance.

6.10.5 Pre-baked fog

You can pre-bake fog into textures if you know the camera does not get close to them. This reduces the compute power required.

6.10.6 Volumetric fog using particles

You can simulate volumetric fog with particles. This can provide high quality results.

Keep the number of particles to minimum. Particles increase overdraw because more shaders are executed per fragment. Try using larger particles instead of creating more particles.

In the Ice Cave demo, a maximum of 15 particles are used at a time.

Geometry instead of billboard rendering

Set the Render Mode of your particle system to Mesh instead of Billboard. This is because, to achieve the volumetric effect, the individual particles must have random rotations.

The following figure shows the particle geometry without textures:



Figure 6-57 Particles without texture

The following figure shows the particle geometry with textures:



Figure 6-58 Particles with texture

The following figure shows the texture that is applied to each particle:

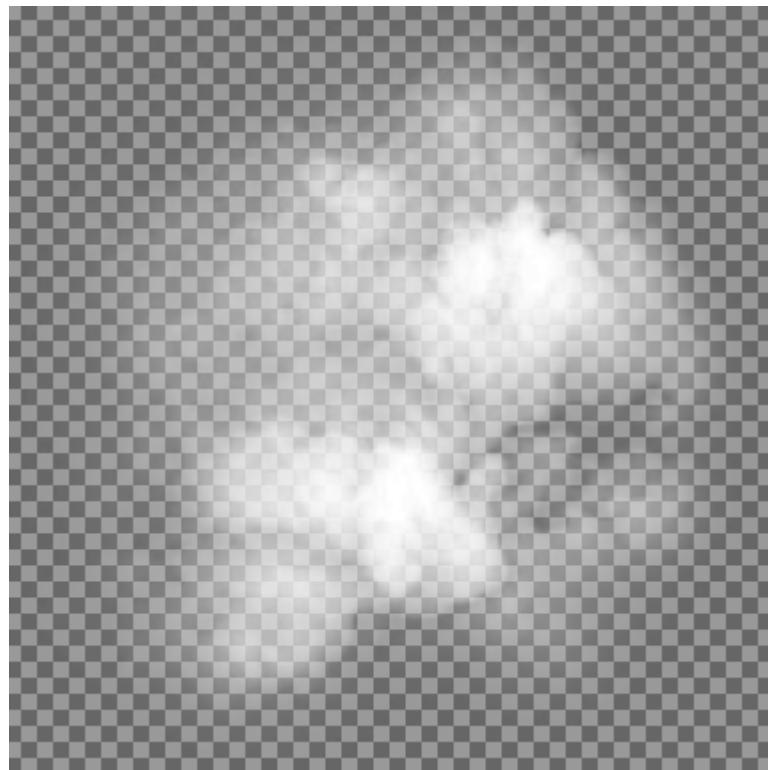


Figure 6-59 Texture of an individual particle

Angle fade effect

Fade in and out the individual particles according to their orientation compared to the camera position. If you do not fade the particles in and out, the sharp edges of the particles are visible.

The following example code shows a vertex shader that performs the fading:

```
half4 vertexInWorld = mul(_Object2World, input.vertex);
half3 normalInWorld = (mul(half4(input.normal, 0.0), _World2Object).xyz);
const half3 viewDirInWorld = normalize(vertexInWorld - _WorldSpaceCameraPos);
output.visibility = abs(dot(-normalInWorld, viewDirInWorld));
output.visibility *= output.visibility; // instead of power of 2
```

The varying parameter `output.visibility` is interpolated across the particle polygon. Read this value in a fragment shader and apply an amount of transparency.

The following code shows how to do this:

```
half4 diffuseTex = _Color * tex2D(_MainTex, half2(input.texCoord));
diffuseTex *= input.visibility;
return diffuseTex;
```

Rendering particles

To render for particles use the following procedure:

1. Render particles as the last primitives within the frame.

```
Tags { "Queue" = "Transparent+10" }
```

The +10 is present in this example because the Ice Cave demo renders 9 other transparent objects before the particles.

2. Set the appropriate blending mode.

At the beginning of a shader pass add the following line:

```
Blend SrcAlpha One
```

3. Disable writing to z-buffer.

Add the following line:

```
ZWrite Off
```

Particle system settings

The following figure shows the settings used in the Ice Cave demo for the particle system:

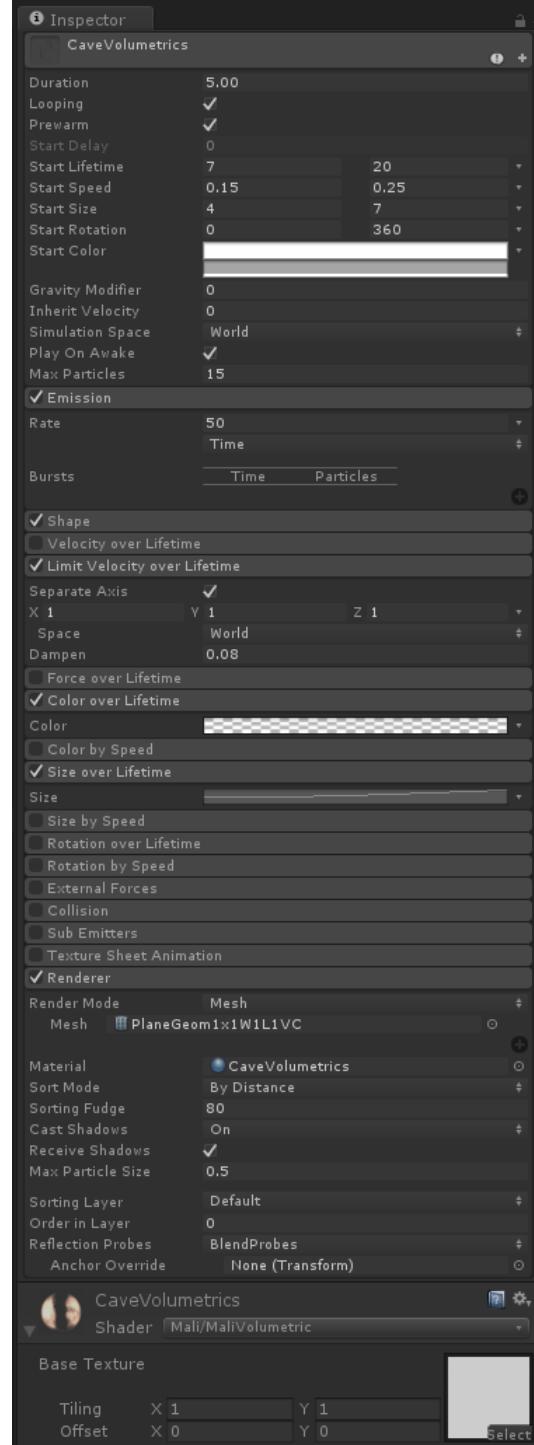


Figure 6-60 Particle system parameters from the Ice Cave demo

The following figure shows the area where the particles are spawned. This is indicated by the box in the image. The box is defined by the built-in option **Shape** in the Unity Particle System:



Figure 6-61 The particles box definition where the particles are spawned

6.11 Bloom

Bloom is used to reproduce the effects that occur in real cameras when taking pictures in a bright environment. The bloom effect simulates fringes of light extending from the borders of bright areas, creating the illusion of a bright light overwhelming the camera.

The following figure shows bloom at the entrance of the cave in the Ice Cave demo:

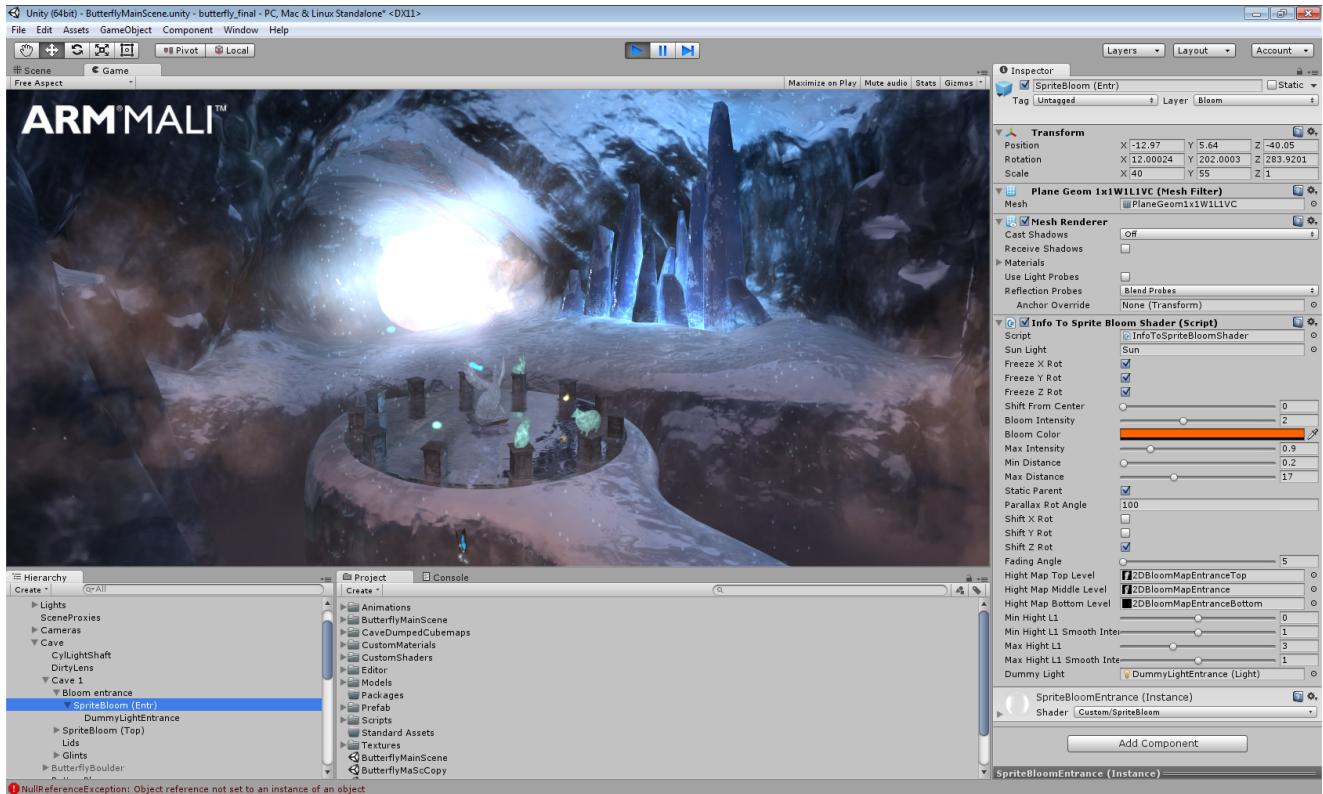


Figure 6-62 Bloom at the entrance of the cave as implemented in the Ice Cave demo

The bloom effect occurs on real lenses because they cannot focus perfectly. When the light passes through the circular aperture of the camera, some light is diffracted creating a bright disk around the image. This effect is not typically noticeable under most conditions, but it is visible under intense lighting.

This section contains the following subsections:

- [6.11.1 Implementing bloom on page 6-152](#).
- [6.11.2 Creating a bloom modulation factor script on page 6-153](#).
- [6.11.3 The vertex shader on page 6-153](#).
- [6.11.4 The fragment shader on page 6-154](#).
- [6.11.5 Correcting the bloom effect on page 6-154](#).
- [6.11.6 Interpolation between occlusion maps on page 6-157](#).

6.11.1 Implementing bloom

Bloom effects are typically implemented as a post-processing effect. Generating effects this way can use a large amount of computing power, making them unsuitable for use in mobile games. This is especially true if your game uses many complex effects like the Ice Cave demo.

An alternative approach uses a simple plane. You place the plane between the camera and the light source. The plane must be oriented with the normal pointing to the part of the scene where the camera is supposed to move around.

To create a bloom effect this way, place a plane where you want the bloom effect to occur. Modulate the intensity of the effect using a factor based on the alignment of the view vector, with the plane normal and the light source. This method is implemented in the Ice Cave demo.

The following figure shows a bloom effect implemented using a plane:



Figure 6-63 Plane at the entrance of the cave implementing bloom

6.11.2 Creating a bloom modulation factor script

The alignment factor can be calculated using a script. This factor alters the bloom effect depending on the angle that the camera is viewing the effect from.

For example, the following script calculates the alignment factor for a plane it is attached to:

```
// Light-plane
normal-camera alignmentVector3 planeToCamVec = Camera.main.transform.position
    - gameObject.transform.position;
planeToCamVec.Normalize();
Vector3 sunLightToPlainVec = origPlainPos - sunLight.transform.position;
sunLightToPlainVec.Normalize();
float sunLightPlainCameraAlignment = Vector3.Dot(planeToCamVec, sunLightToPlainVec);
sunLightPlainCameraAlignment = Mathf.Clamp (sunLightPlainCameraAlignment, 0.0f, 1.0f);
```

The alignment factor `sunLightPlainCameraAlignment` is passed to the shader to modulate the intensity of the rendered color.

6.11.3 The vertex shader

The vertex shader receives a `sunLightPlainCameraAlignment` factor and uses it to modulate the intensity of the rendered bloom color.

The vertex shader applies the MVP matrix, passes the texture coordinates to the fragment shader, and outputs the vertex coordinates.

The following code shows how this is achieved:

```
vertexOutput vert(vertexInput input)
{
    vertexOutput output;
    output.tex = input.texcoord;
    output.pos = mul(UNITY_MATRIX_MVP, input.vertex);
    return output;
}
```

6.11.4 The fragment shader

The fragment shader fetches the texture color and increments in using a `CurrBloomColor` tint color, which is passed as a uniform variable. The alignment factor modulates the color before it is used.

The following code shows how this process is performed:

```
half4 frag(vertexOutput input) : COLOR
{
    half4 textureColor = tex2D(_MainTex, input.tex.xy);
    textureColor += textureColor * _CurrBloomColor;
    return textureColor * _AlignmentFactor;
}
```

The bloom plane is rendered in the transparent queue after all opaque geometry. In the shader, use the following queue tag to set the rendering order:

```
Tags { "Queue" = "Transparent" + 1}
```

The bloom plane is applied using additive one-to-one blending, to combine the fragment color with the corresponding pixel which is already stored in the frame buffer. The shader also disables writing to the depth buffer using the instruction, `ZWrite Off`, so that existing objects are not occluded.

The following figure shows the texture that the Ice Cave demo uses for its bloom effect, and the result it creates:



Figure 6-64 Sequence showing the camera entering the occluding zone behind the opaque plinth

6.11.5 Correcting the bloom effect

Using a plane to produce a bloom effect is simple and suitable for mobile devices. However, it does not produce the expected results when an opaque object is placed between the light source producing the bloom and the camera. You can correct this with an occlusion map.

The incorrect appearance of the bloom effect behind an opaque object occurs because the effect is rendered in the transparent queue. Therefore, the blending occurs on top of any objects in front of the bloom plane.

The following figure shows an example of the incorrect bloom effect that is produced without any fixes:



Figure 6-65 Incorrect bloom effect showing over an opaque plinth

To prevent these errors from occurring, you can alter the script which calculates the alignment factor so that it processes an additional occlusion map. This occlusion map describes the areas where the camera normally applies the bloom effect incorrectly. These conflicting areas are assigned an occlusion map value of zero. This factor combines with the alignment factor to set the intensity of factor in these places to zero. This change sets the bloom to zero so that it no longer renders over object which must occlude it. The following code implements this adjustment:

```
IntensityFactor = alignmentFactor * occlusionMapFactor
```

Because the Ice Cave demo camera can move freely in three dimensions, three gray scale maps are used, each one covering a different height. Black means zero, so the occlusion map factor multiplied by the alignment factor makes the final intensity zero, occluding the bloom. White means one, so the intensity is equal to the alignment factor and the bloom is not occluded.

The following figure shows the `occlusionMapFactor` at ground level:

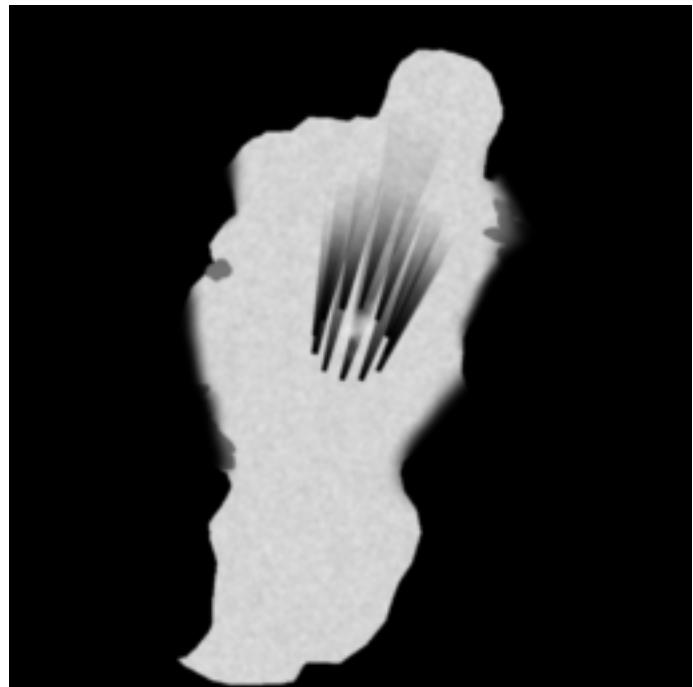


Figure 6-66 Ground level bloom occlusion map

————— Note ————

The black radiating areas in the middle of the map are zones behind the opaque plinths that occlude the bloom effect from the entrance of the cave.

There are no occluding objects at a height H_{\max} above ground level. This means that the bloom occlusion map above ground level is white. The following figure shows the `occlusionMapFactor` above ground level:

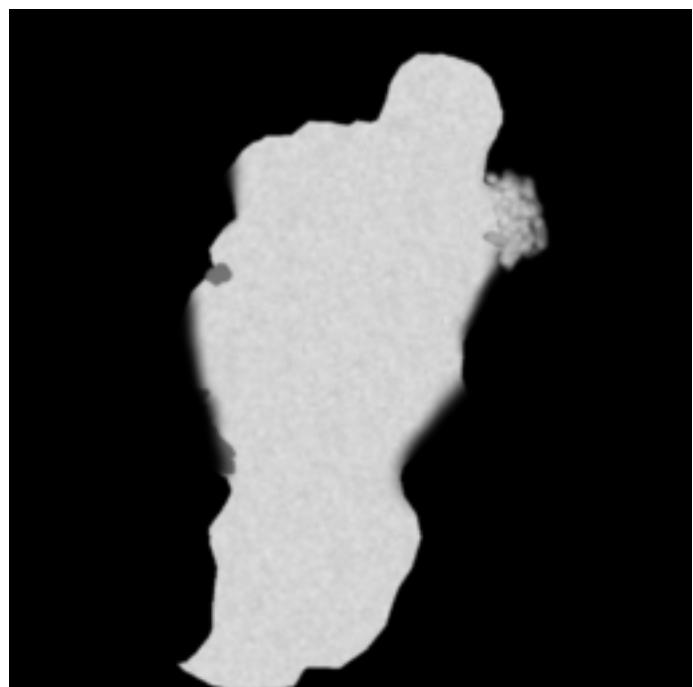


Figure 6-67 Above ground level bloom occlusion map

Below the height H_{\min} , below ground level the bloom effect is completely occluded. This means that the bloom occlusion map below ground level is completely black. The white contour is added so that the map is visible. The following figure shows the `occlusionMapFactor` below ground level:



Figure 6-68 Below ground level bloom occlusion map

6.11.6 Interpolation between occlusion maps

In the Ice Cave demo, the script calculates the XZ projection of the camera position on the 2D cave map every frame, and normalizes the result between zero and one.

If the camera height is below H_{\min} or above H_{\max} , the normalized coordinates are used to fetch the color value from a single map. If the camera height is between H_{\min} and H_{\max} , the color is retrieved from both maps and interpolated.

This interpolation creates a smooth transition between the effects at different heights.

The Ice Cave demo uses the `GetPixelBilinear()` function to retrieve the colors from the map or maps. This function returns a filter color value using the following code:

```
float groundOcclusionFactor = groundOcclusionMap.GetPixelBilinear(camPosXZNormalized.x,  
camPosXZNormalized.y)).r;
```

Using the bloom occlusion maps prevents the blending of the bloom on top of occluding opaque object. The following figure shows the resulting effect. When the camera is entering and leaving, the occluding black is behind the opaque plinth, preventing the incorrect bloom.

The following figure is sequence showing the camera entering the occluding zone behind an opaque plinth:



Figure 6-69 Entering the occluding zone behind an opaque plinth

6.12 Icy wall effect

The Ice Cave demo uses a subtle reflection effect in the icy walls of the cave. This effect is subtle, but adds extra realism and atmosphere to the scene.

This section contains the following subsections:

- [6.12.1 About the icy wall effect on page 6-159](#).
- [6.12.2 Modifying and combining normal maps to affect reflections on page 6-160](#).
- [6.12.3 Creating the reflections from the different normal maps on page 6-162](#).
- [6.12.4 Applying local correction to the reflections on page 6-164](#).

6.12.1 About the icy wall effect

Ice can be a difficult material to replicate because light scatters off it in different ways, depending on the small details of its surface. The reflection can be completely clear, completely distorted, or anywhere in between. The Ice Cave demo shows this effect and includes a parallax effect for greater realism.

The following figure shows the Ice Cave demo, demonstrating this effect:



Figure 6-70 The Ice Cave demo

The following figure shows a close up of this effect:

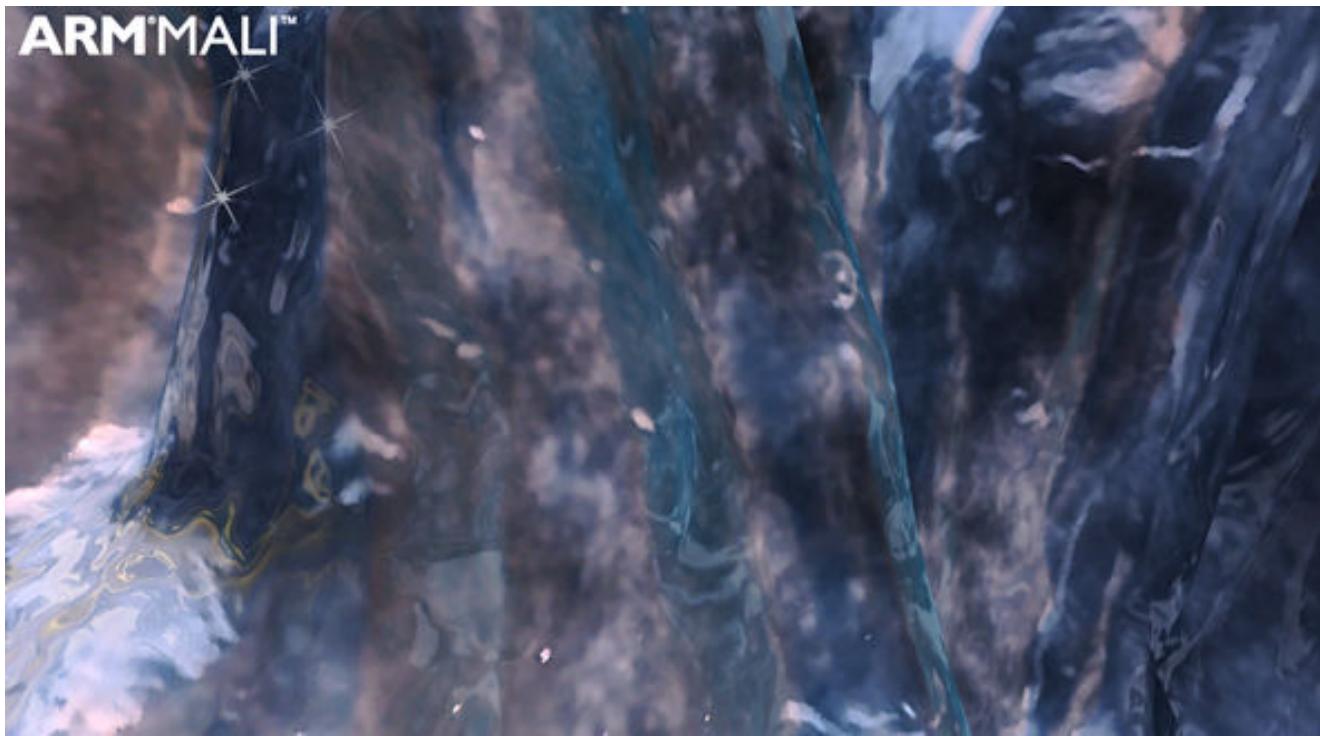


Figure 6-71 Close up of the reflective icy walls

6.12.2 Modifying and combining normal maps to affect reflections

The reflection effect in the Ice Cave demo uses the tangent space normal maps and calculated gray scale fake normal maps. Combining both of these maps with some modifiers creates the effect in the demo.

The gray scale fake normal maps are a gray scale of the tangent space normal maps. In the Ice Cave demo, most values in the gray scale maps are in the range 0.3-0.8. The following figure shows the tangent space normal maps and gray scale fake normal maps that the Ice Cave demo uses:

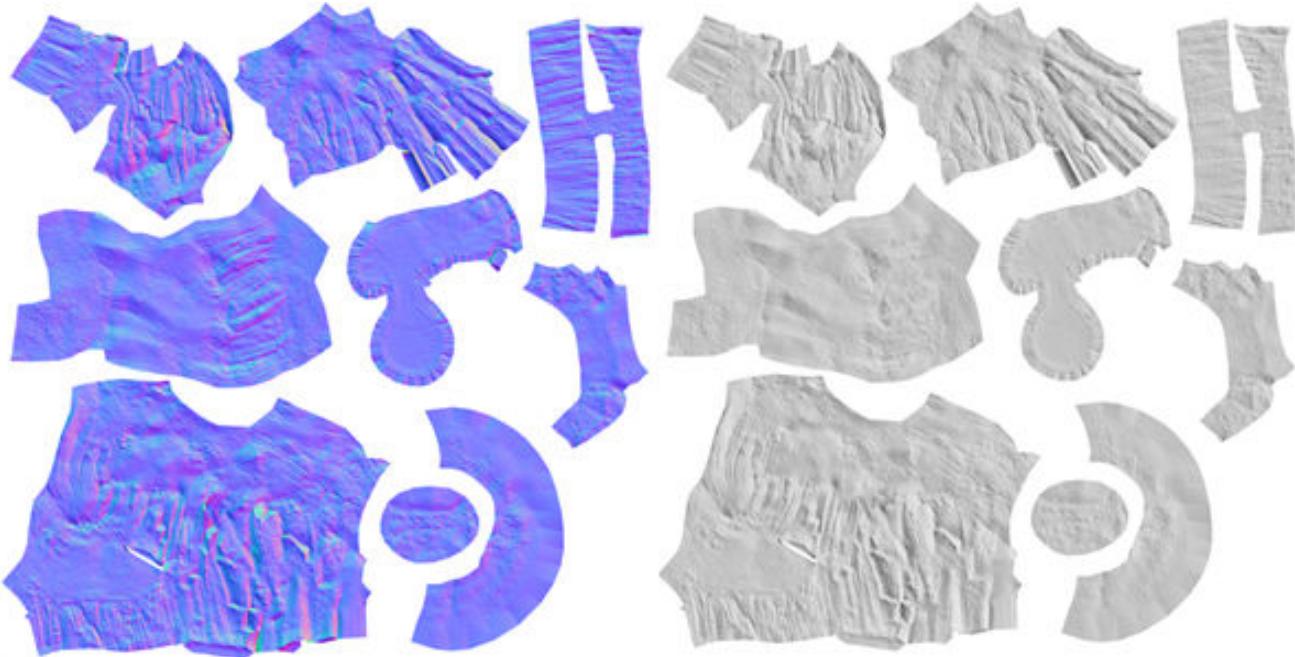


Figure 6-72 Tangent space normal maps and grayscale fake normal maps

Why the tangent space normal maps are used

The tangent space normal maps are used in the Ice Cave demo, because they have a small range of values in the resulting gray scale. This means that tangent space normal maps work well in the later stages of this rendering process.

An alternative option, is to use the object space normal maps. These maps show the same details as the tangent space normal maps, but also show where light hits. Therefore the range of values in the resulting object space normal map gray scale is too large to work well in the later stages of the rendering process. The following figure shows this effect:

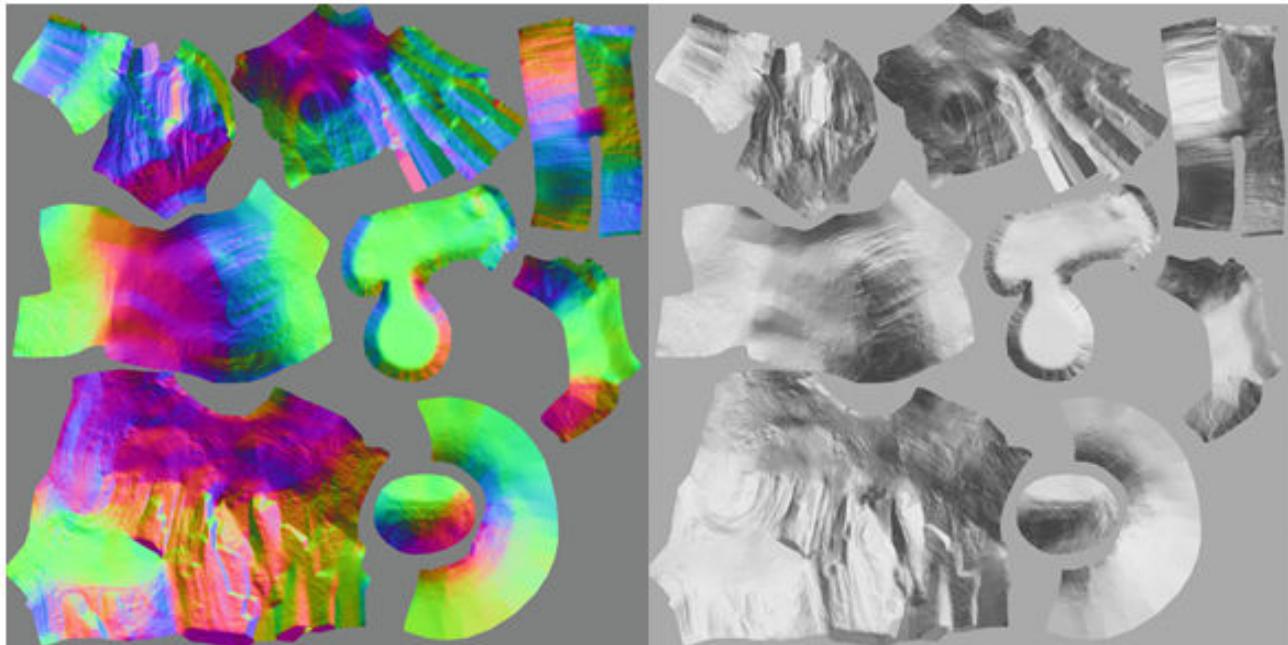


Figure 6-73 Object space normal maps and grayscale effect

Applying transparency to parts of the normal maps

The gray scale fake normal maps are only applied to the areas without snow. To prevent them from being applied to the snowy areas, the gray scale fake normal maps are modified using the diffuse texture maps for the same surfaces. This modification increases the alpha component of the gray scale fake normal maps where snow appears in the texture maps.

The following figure shows the diffuse texture maps for the Ice Cave demo static surfaces and the result when they are applied to the gray scale fake normal maps:

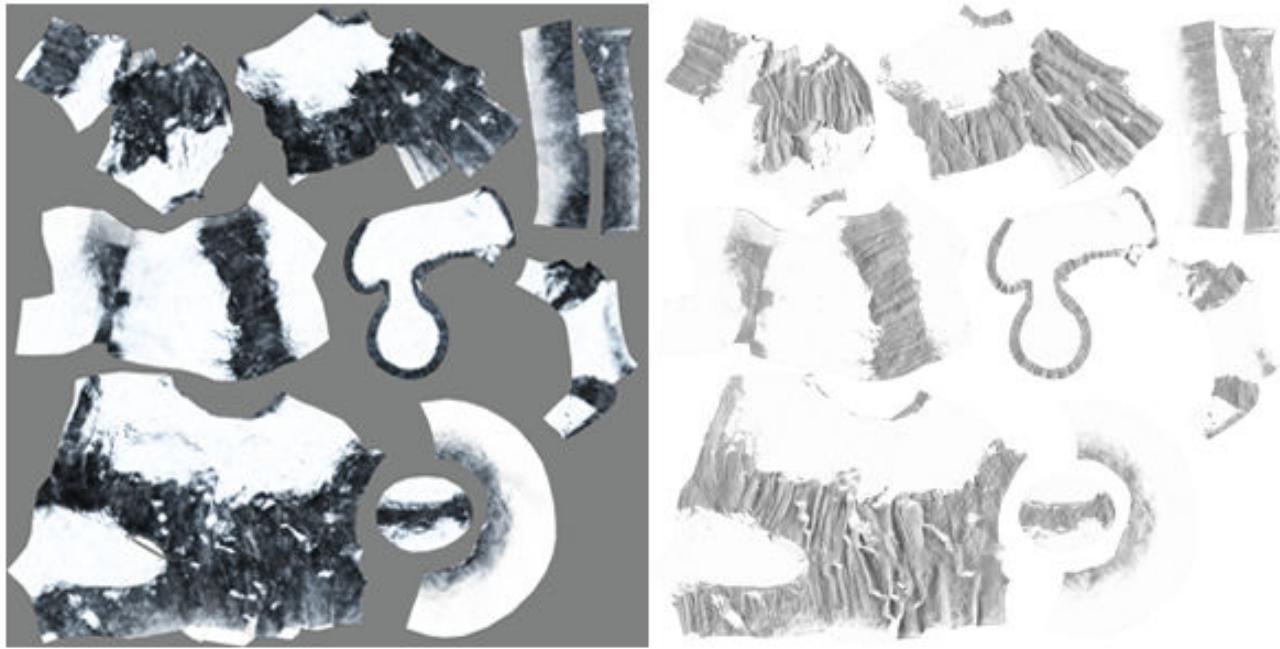


Figure 6-74 Diffuse texture and final fake normal maps with transparent areas

6.12.3 Creating the reflections from the different normal maps

The combination of the transparency adjusted gray scale fake normal maps and the true normal maps creates the reflection effects that are shown in the Ice Cave demo.

The transparency adjusted gray scale fake normal maps, `bumpFake`, and the true normal maps, `bumpNorm`, are combined proportionally. The combination uses the following function:

```
half4 bumpNormalFake = lerp(bumpNorm, bumpFake, amountFakeNormalMap);
```

This code means that in the darker parts of the cave, the main reflection contribution comes from the gray scale fake normals. In the snowy parts of the cave, the effect comes from the object space normals.

To apply the effect using the gray scale fake normal maps, the gray scale must be converted into normal vectors. To start this process, the three components of the normal vectors are set equal to the gray scale value. In the Ice Cave demo, this means that the components vectors are between (0.3, 0.3, 0.3) and (0.8, 0.8, 0.8). Because all the components are set so that they are the same, all the normal vectors point in the same direction.

The shader applies a transformation to normal components. It uses the transform that is usually used to transform values in the range 0-1 to the range -1 to 1. The equation which performs this change is, $\text{result} = 2 * \text{value} - 1$. This equation changes the normal vectors so that they either point in the direction they pointed before, or the opposite direction. For example, if the original has the components (0.3, 0.3, 0.3), then the resulting normal is (-0.4, -0.4, -0.4). If the original has the components (0.8, 0.8, 0.8), then the resulting normal is (0.6, 0.6, 0.6).

After the transformation to the range -1 to 1, the vectors are fed to the `reflect()` function. This function is designed to work using normalized normal vectors, however in this case the non-normalized normal is passed to the function. The following code shows how the shader built-in function `reflect()` works:

```
R = reflect(I, N) = I - 2 * dot(I, N) * N
```

Using this function with a non-normalized input normal with a length less than one, causes the reflection vector to deviate more than expected from the normal, according to the reflection law.

When the value of the components of the normal vector goes below 0.5, the reflection vector switches to the opposite direction. When this happens, a different part of the cubemap is read. This switching between different parts of the cubemap creates the effect of uneven spots that reflect the white parts of

the cubemap next to reflections of the rocky parts of the cubemap. Because the areas in the gray scale fake normal maps that cause the switches between positive and negative normals are also the areas that produce the most distorted angles on the reflected vectors, this creates an appealing swirl effect. The following figure shows this effect:



Figure 6-75 Swirl effect in the ice reflection

If the shader is programmed to use the fake normal values without the non-normalized clamped stage, then the result features diagonal bands. This is a significant difference and shows that these stages are important. The following figure shows the resulting effect without the clamp stage:

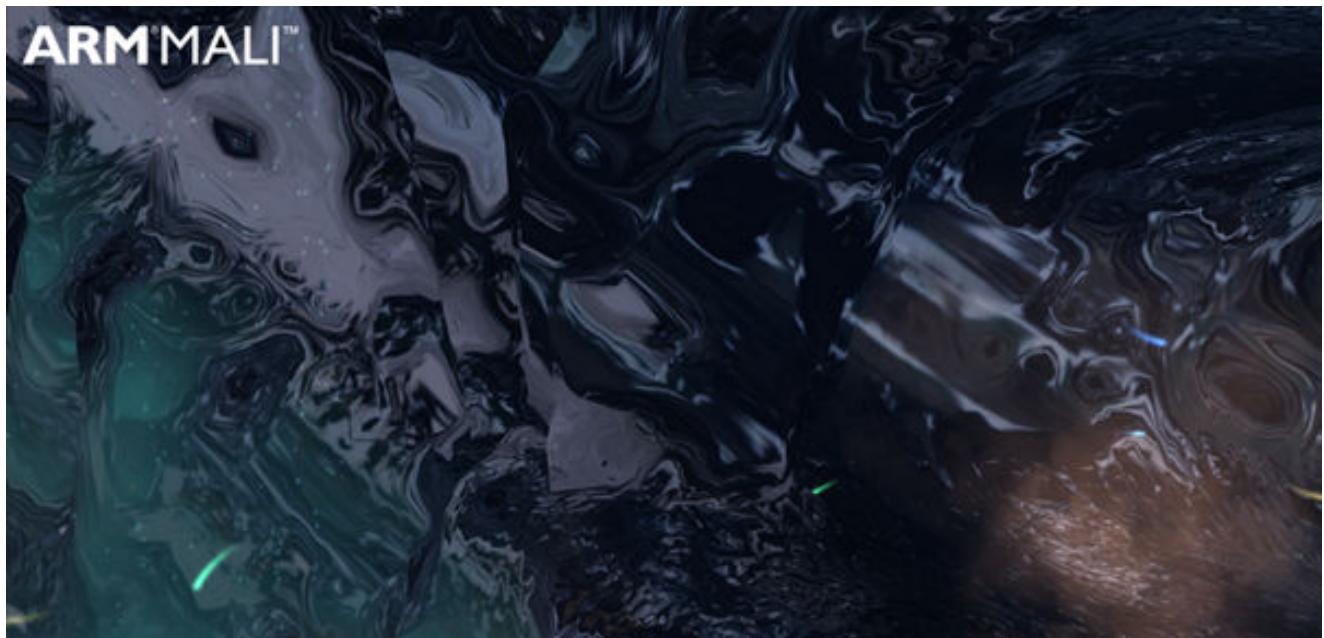


Figure 6-76 Ice reflection without the clamp stage

6.12.4 Applying local correction to the reflections

Applying local correction to the reflection vector improves the realism of the effect. Without this correction, the reflection does not change with the position of the camera like reflections in real life do. This is especially true for lateral movements of the camera.

Related information

[6.2.2 Generating correct reflections with a local cubemap](#) on page 6-100.

6.13 Procedural skybox

The Ice Cave demo uses a time-of-day system to show the dynamic shadowing effects that you can achieve with local cubemaps.

This section contains the following subsections:

- [6.13.1 About the procedural skybox on page 6-165](#).
- [6.13.2 Managing the time of day on page 6-166](#).
- [6.13.3 Rendering the Sun on page 6-167](#).
- [6.13.4 Fading the Sun behind mountains on page 6-168](#).
- [6.13.5 Subsurface scattering on page 6-170](#).

6.13.1 About the procedural skybox

To achieve a dynamic time of day effect, the following elements are combined:

- A procedurally generated sun.
- A series of fading skybox background cubemaps that represent the day to night cycle.
- A skybox clouds cubemap.

The procedural sun and the separated cloud textures also create a computationally cheap subsurface scattering effect.

The following figure shows a view of the skybox:



Figure 6-77 Sun rendering with subsurface scattering in the Ice cave demo

The Ice Cave demo uses the view direction to sample from a cubemap. This enables the demo to avoid rendering a hemisphere for the skybox. Instead it just uses planes near the holes in the cave.

Doing this increases the performance compared to rendering a hemisphere that is then mostly occluded by other geometry.

The following figure shows the skybox rendered using a plane.

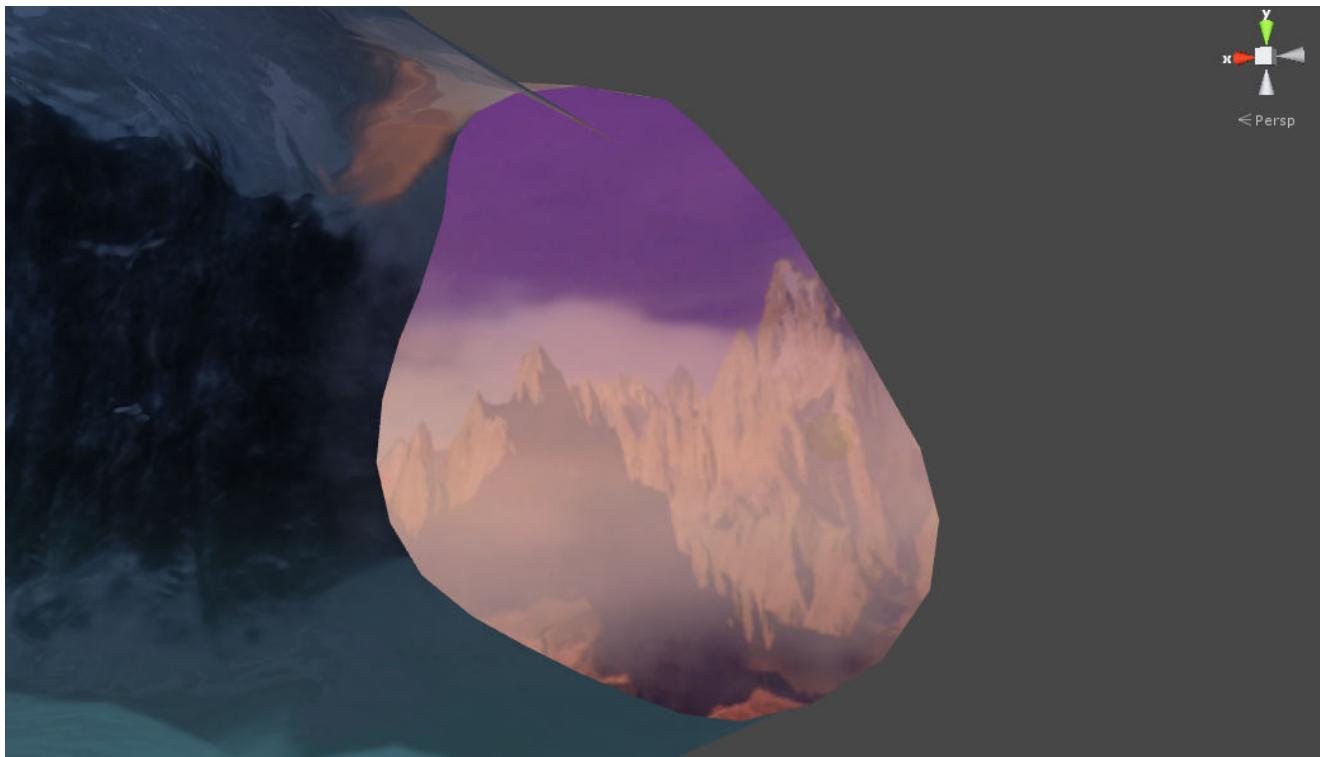


Figure 6-78 Skybox rendered using a plane

6.13.2 Managing the time of day

This effect uses a C# script to manage the mathematics for the time of day, and the animation of the day-night cycle. A shader then combines the sun and skymaps.

You must specify the following values for the script:

- The number of fading skybox background phases.
- The maximum duration of the day-night cycle.

For each frame, the script selects the skybox cubemaps to blend together. The selected skyboxes are set as textures for the shader and it blends the textures together while it is rendering.

To set the textures, the Ice Cave demo uses the Unity Shader global features. This enables you to set a texture in one place that can be used by all the shaders in your application. The following code shows this:

```
Shader.SetGlobalTexture (ShaderCubemap1, _phasesCubemaps [idx1]);
Shader.SetGlobalTexture (ShaderCubemap2, _phasesCubemaps [idx2]);
Shader.SetFloat (ShaderAlpha, blendAlpha);
Shader.SetGlobalVector (ShaderSunPosition, normalizedSunPosition);
Shader.SetGlobalVector (ShaderSunParameters, _sunParameters);
Shader.SetGlobalVector (ShaderSunColor, _sunColor);
Shader.SetGlobalTexture (ShaderCloudsCubemap, _CloudsCubemap);
```

Note

The name of the sampler used must not conflict with one locally defined by the shader.

The following is set in the script code:

- The two cubemaps that are interpolated. The values `idx1` and `idx2` are computed based on the elapsed time.
- A `blendAlpha` factor, used in the shader to blend the two cubemaps.
- A normalized sun position, used to render the sun sphere.

- A number of parameters for the sun.
- The color of the sun.
- The cloud cubemap. `ShaderCubemap1` and `ShaderCubemap2` are two strings containing the unique sampler names, in this case `_SkyboxCubemap1` and `_SkyboxCubemap2`.

To access these textures in the shader, you must declare them with the following code:

```
samplerCUBE _SkyboxCubemap1;
samplerCUBE _SkyboxCubemap2;
```

The script selects the sun color and an ambient color based on a list you specify. These are interpolated for each phase.

The sun color is passed to the shader for rendering the sun with the correct color.

The ambient color is used to dynamically set the Unity variable `RenderSettings.ambientSkyColor`:

```
RenderSettings.ambientSkyColor = _ambientColor;
```

Setting this variable enables all the materials to receive the correct ambient color and also enables Enlighten to get the correct ambient color when it updates the lightmaps.

In the Ice Cave demo, this effect causes a gradual variation of the overall color of the scene based on the phase of the day it is.

6.13.3 Rendering the Sun

To render the sun, in the fragment shader, for each pixel of the skybox you must check if it is inside the circumference of the sun.

To do this, the shader computes the dot product of the normalized sun position vector and the normalized view direction vector in world coordinates of the pixel being rendered.

The following figure shows the rendering of the sun. The dot product of the normalized view vector and the sun position is used to determine if the fragment is rendered as sky as P2, or rendered as a sun as P1.

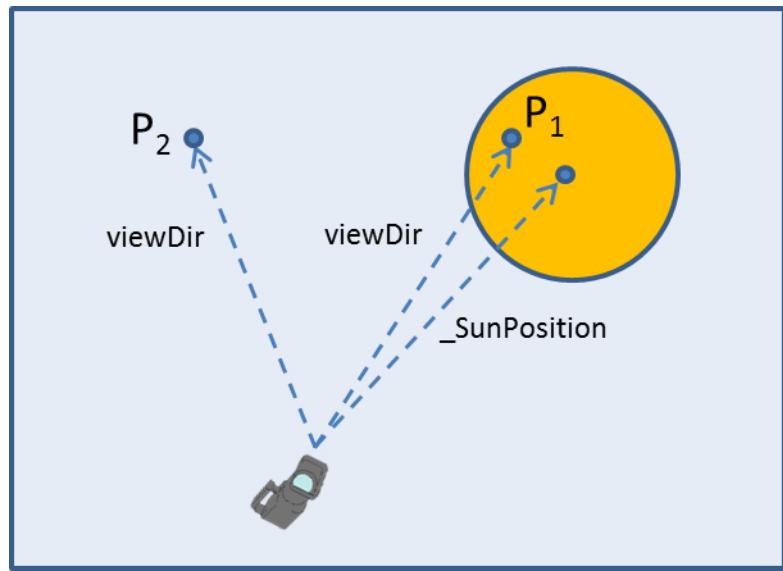


Figure 6-79 Rendering the sun

The normalized sun position vector is passed to the shader by the C# script.

- If the result is greater than a specific threshold the pixel is colored with the sun color.
- If the result is less than the threshold, the pixel is colored with sky color.

The dot product also creates a fading effect towards the edges of the sun:

```
half _sunContribution = dot(viewDir,_SunPosition);
```

The following figure shows a clear view of the sun:



Figure 6-80 Procedural Sun rendering in clear sky conditions

6.13.4 Fading the Sun behind mountains

There is a problem if the sun is low in the sky. In real life it would disappear behind the mountains.

To create this effect, the alpha channel of the cubemap is used to store a value of 0 if the texel represents the sky, and a 1 if the texel represents the mountain.

While rendering the sun, the texture is sampled and the alpha is used to make the sun fade behind the mountains. This sampling is effectively free because the texture is already sampled to render the mountains.

You can also gradually fade the alpha near the edges or the areas of the mountain where there is snow. This produces an effect of the sun bouncing off the snow for little computational effort.

A similar technique is used to create a computationally cheap subsurface scattering effect for the clouds.

The original phase cubemaps are split into two separate groups.

- One set of cubemaps contain the sky and the mountains. The alpha value is set to 0 for the sky and set to 1 for the mountains.
- The other set of cubemaps contain the clouds. The alpha value is set to 0 in the areas without clouds and gradually increases to 1 in areas where the clouds are denser.

The following figure shows a mountains texture. The mountains have an alpha value of one, the sky has a value of zero:



Figure 6-81 Mountains texture

The clouds skybox has alpha 0 for the empty areas and gradually fading to 1 when occupied by a cloud. The alpha channel fades smoothly to ensure the clouds do not look like they have been artificially put on top of the sky. The following figure shows a cloud texture with alpha:

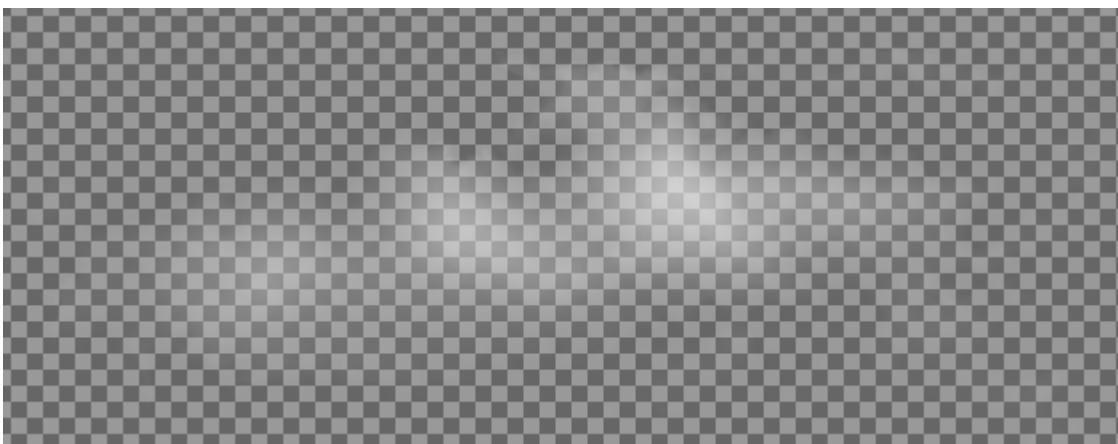


Figure 6-82 Cloud texture

The shader does the following:

1. Sample the two skyboxes that represent the current time of day phase.
2. Blend the two colors based on a blending factor calculated by the C# script.
3. Sample the cloud skybox.
4. Blend the color from point 2 with the one from the clouds using the cloud alpha.
5. Add together the cloud alpha and the skybox alpha.
6. Call a function to compute the sun color contribution for the current pixel.
7. Add together the result of point 6 with the previously blended skybox and cloud color.

————— **Note** ————

You can optimize this sequence by putting the sky, mountains, and clouds into a single skybox. These are separate in the Ice Cave demo so the artist can easily modify the skyboxes and the clouds separately.

6.13.5 Subsurface scattering

You can add a subsurface scattering effect that uses the alpha information to increase the radius of the sun.

In real life, when the sun is in a clear sky, its size appears relatively small because there are no clouds that can deviate, or scatter, the light toward your eyes. What you see are the rays directly from the sun, spread only very slightly by the atmosphere.

When the sun is occluded by clouds but they do not completely obscure it, part of the light is scattered around the cloud. This light can reach your eyes from directions that are some distance away from the sun. This can make the sun appear larger than it really is.

The following code is used to achieve this effect in the Ice Cave demo:

```
half4 sampleSun(half3 viewDir, half alpha);
{
    half _sunContribution = dot(viewDir,_SunPosition);
    half _sunDistanceFade = smoothstep(_SunParameters.z -(0.025*alpha), 1.0, _sunContribution);
    half _sunOcclusionFade = clamp( 0.9-alpha, 0.0, 1.0);

    half3 _sunColorResult = _sunDistanceFade * _SunColor * _sunOcclusionFade;
    return half4( _sunColorResult.xyz, 1.0 );
}
```

The parameters of the function are `viewDirection` and the alpha computed, when the cloud alpha and the skybox alpha are added together.

The dot product of the sun position and view directions are used to compute a scaling factor that represents the distance of the current pixels from center of the sun.

The `_sunDistanceFade` calculation uses the `smoothstep()` function to provide a more gradual fade from the center of the sun to the sky near the edges. This also enables the radius of the sun to be increased near the clouds to simulate the subsurface scattering effect.

This function has a variable domain based on the alpha, in case of clear sky, the alpha is 0 and the range is within `_SunParameters.z` and `1.0`. In this case `_SunParameters.z` is initialized to `0.995` in the C# script, this corresponds to a sun of diameter of 5 degrees, $\cos(5 \text{ degrees}) = 0.995$.

If the pixel being processed contains a cloud, the radius of the sun is increased to 13 degrees to enable an elongated scattering effect when approaching the clouds.

The `_sunOcclusionFade` factor is used to fade down the contribution of the sun based on the occlusion received from the mountains and the clouds.

The following figure shows the sun that is not occluded by the clouds:



Figure 6-83 Sun not occluded by clouds

The following figure shows the sun occluded by clouds:



Figure 6-84 Sun occluded by clouds

6.14 Fireflies

Fireflies are bright flying insects that are used in the Ice Cave demo to add more dynamism, and show the benefits of using Enlighten for real-time global illumination.

This section contains the following subsections:

- [6.14.1 About fireflies on page 6-173](#).
- [6.14.2 The firefly generator prefab on page 6-175](#).

6.14.1 About fireflies

The fireflies are composed of the following components:

- A prefab object that is instantiated at runtime.
- A box collider that is used to limit the area the fireflies can fly around in.

These two components are combined together by a C# script that manages the movement of the fireflies and defines the path they follow.

The following figure shows a firefly:



Figure 6-85 Firefly

The firefly prefab uses the Unity standard particle system to generate the trail of the firefly.

The following figure shows the Unity Particle System settings used by the fireflies:

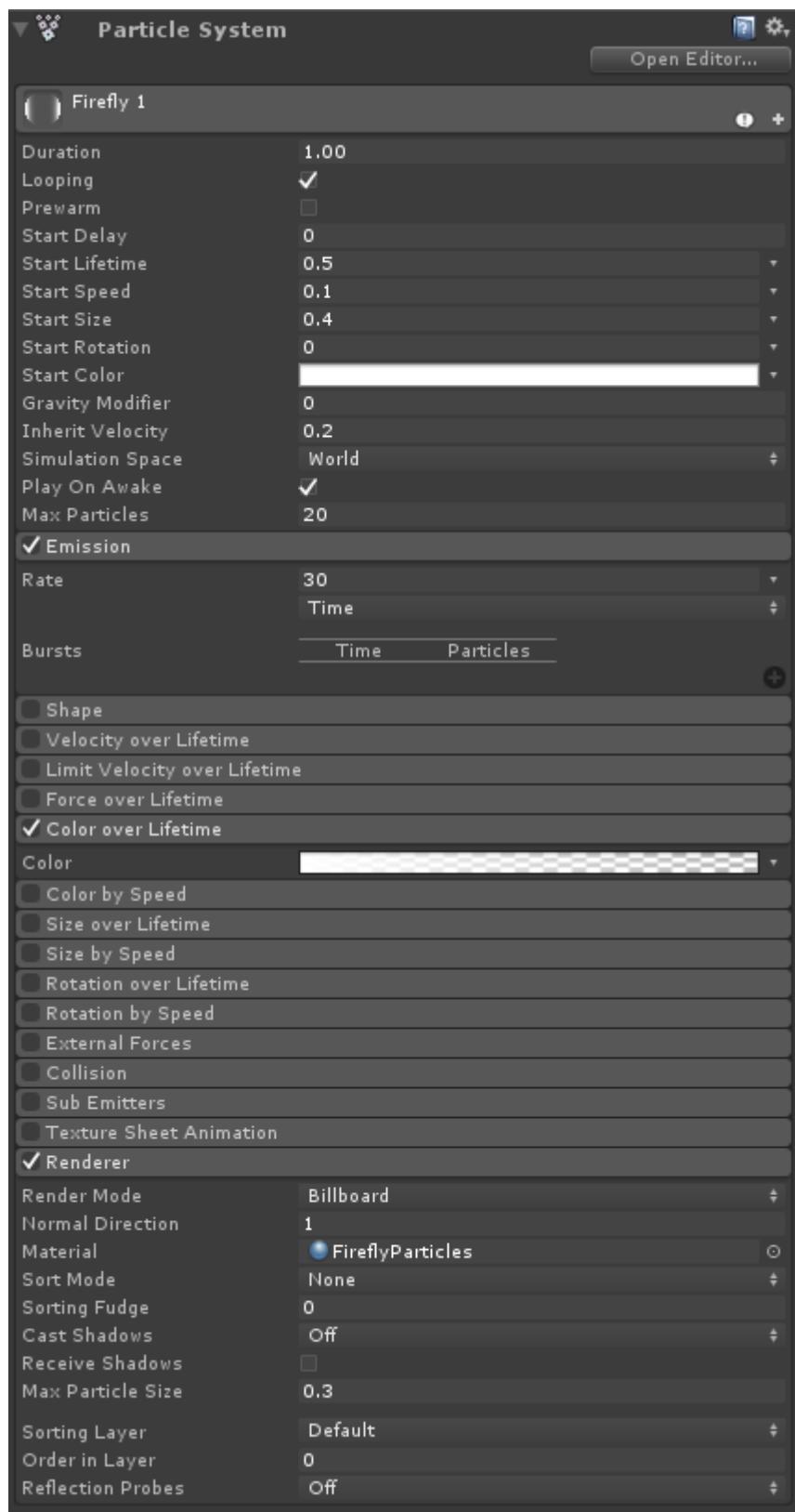


Figure 6-86 Unity Particle System settings

Depending on the effect you want to achieve, you can use a Unity Trail Renderer to provide a more continuous look.

The Trail Renderer generates a large number of triangles for each trail. You can change the number of triangles by modifying the Trail Renderer setting **Min Vertex Distance**, but a high value can cause a jerky movement of the trail if the source is moved too fast.

The **Min Vertex Distance** option defines the minimum distance between the vertices that form the trail. A high number can be fine for straight trails but it does not look smooth for curved trails.

The trail generated is always facing the camera, because of this, any abrupt movement of the source can cause the trail to overlap with itself. This causes artifacts generated by the blending of the triangles that form the trail.

The following figure shows artifacts caused by an overlapping trail:



Figure 6-87 Overlapping trail artifacts

The final component added to the prefab, is a point light that casts light in the scene while it is moving.

The light affects the Enlighten global illumination providing a bouncing light effect, especially in narrow parts of the scene where light bounces are more visible because of less light dispersion.

6.14.2 The firefly generator prefab

The firefly generator prefab manages the creation of the fireflies and updates them each frame. It is composed of a C# script for the update and a box collider that encloses the volume that each firefly can move in.

The script takes the number of fireflies to be generated as a parameter, and the prefab to instantiate the firefly object. Since the movement of the fireflies is random within the bounding box, it limits changes to a specific angle, away from the direction the firefly is moving. This ensures the firefly does not make sudden changes of direction.

To generate a random movement, a piecewise cubic Hermite interpolation is used to create the control points. The Hermite interpolation provides a smooth continuous function that behaves correctly even when different paths are connected together. The first derivative at the end points are also continuous so there are no sudden velocity changes.

This interpolation requires one control point each for the start and the end, and two tangents for each of the control points. Because these are generated randomly, the script can store three control points and two tangents. It uses the position of the first and second control points to define the first point tangent, and the second and third control points to define the tangent of the second control point.

At loading time, the script generates the following for each firefly:

- An initial position.
- An initial direction using the Unity function `Random.onUnitSphere()`.

The following code shows how the control points are initialized:

```
_fireflySpline[i*_controlPoints] = initialPosition;
Vector3 randDirection = Random.onUnitSphere;
_fireflySpline[i*_controlPoints+1] = initialPosition + randDirection;
_fireflySpline[i*_controlPoints+2] = initialPosition + randDirection * 2.0f;
```

The initial control points lie on a straight line. The tangents are generated from these control points:

```
//The tangent for the first point is in the same direction as the initial direction vector
_fireflyTangents[i*_controlPoints] = randDirection;

//This code computes the tangent from the control point positions. It is shown here for
// reference because it can be set to randDirection at initialization.
_fireflyTangents[i*_controlPoints+1] = (_fireflySpline[i*_controlPoints+2] -
_fireflySpline[i*_controlPoints+1])/2 + (_fireflySpline[i*_controlPoints+1] -
_fireflySpline[i*_controlPoints])/2;
```

To complete a firefly initialization, you must set the color of the firefly and the duration of the current path interval.

On each frame, the script updates the position of each firefly using the following code to compute the Hermite interpolation:

```
// t is the parameter that defines where in the curve the firefly is placed. It represents
// the ratio of the time the firefly has traveled along the path to the total time.
float t = _fireflyLifetime[i].y / _fireflyLifetime[i].x;

//Hermite interpolation parameters
Vector3 A = _fireflySpline[i*_controlPoints];
Vector3 B = _fireflySpline[i*_controlPoints+1];
float h00 = 2*Mathf.Pow(t,3) - 3*Mathf.Pow(t,2) + 1;
float h10 = Mathf.Pow(t,3) - 2*Mathf.Pow(t,2) + t;
float h01 = -2*Mathf.Pow(t,3) + 3*Mathf.Pow(t,2);
float h11 = Mathf.Pow(t,3) - Mathf.Pow(t,2);
//Firefly updated position
_fireflyObjects[i].transform.position = h00 * A + h10 * _fireflyTangents[i*_controlPoints]
+ h01 * B + h11 * _fireflyTangents[i*_controlPoints+1];
```

If the firefly completed the whole piece of randomly generated path, the script creates a new random piece starting from the end of the current one:

```
//t > 1.0 indicates the end of the current path
if( t >= 1.0 )
{
    //Update the new position
    //Shift the second point to the first as well as the tangent
    _fireflySpline[i*_controlPoints] = _fireflySpline[i*_controlPoints+1];
    _fireflyTangents[i*_controlPoints] = _fireflyTangents[i*_controlPoints+1];

    //Shift the third point to the second, this point doesn't have a tangent
    _fireflySpline[i*_controlPoints+1] = _fireflySpline[i*_controlPoints+2];

    //Get new random control point within a certain angle from the current fly direction
    _fireflySpline[i*_controlPoints+2] = GetNewRandomControlPoint();

    //Compute the tangent for the central point
    _fireflyTangents[i*_controlPoints+1] = (_fireflySpline[i*_controlPoints+2] -
_fireflySpline[i*_controlPoints+1])/2 + (_fireflySpline[i*_controlPoints+1] -
_fireflySpline[i*_controlPoints])/2;

    //Set how long should take to navigate this part of path
    _fireflyLifetime[i].x = _fireflyMinLifetime;

    //Timer used to check how much we traveled along the path
    _fireflyLifetime[i].y = 0.0f;
}
```

6.15 The tangent space to world space normal conversion tool

The tangent space to world space normal conversion tool is composed of a C# script and a shader. The tool runs offline from inside the Unity editor and it does not affect the run-time performance of your game.

This section contains the following subsections:

- [6.15.1 About the tangent space to world space conversion tool](#) on page 6-177.
- [6.15.2 The C# script](#) on page 6-177.
- [6.15.3 The WorldSpaceNormalCreator shader](#) on page 6-180.
- [6.15.4 The WorldSpaceNormalsCreators C# script](#) on page 6-181.
- [6.15.5 The WorldSpaceNormalCreator shader code](#) on page 6-183.

6.15.1 About the tangent space to world space conversion tool

Profiling the Ice Cave demo showed that there was a bottleneck in the Arithmetic pipeline. To reduce the load, the Ice Cave demo uses world space normal maps for the static geometry, instead of tangent space normal maps.

Tangent space normal maps are useful for animated and dynamic objects but require additional computations to correctly orient the sampled normal.

Because most of the geometry in the Ice Cave demo is static, the normal maps were converted into world space normal maps. This ensures that normal sampled for textures are already correctly oriented in world space. This change is possible because the Ice Cave demo lighting is computed in a custom shader, whereas the Unity standard shader uses tangent space normal maps.

The conversion tool is composed of:

- A C# script that adds a new option in the editor.
- A shader that performs the conversion.

The tool runs offline from inside the Unity editor and it does not affect the run-time performance of your game.

6.15.2 The C# script

You must place the C# script in the Unity Assets/Editor directory. This enables the script to add a new option to the **GameObject** menu in the Unity editor. Create this directory if it does not already exist.

The following code shows how to add a new option in the Unity editor:

```
[MenuItem("GameObject/World Space Normals Creator")]
static void CreateWorldSpaceNormals ()
{
    ScriptableWizard.DisplayWizard("Create World Space Normal",
        typeof(WorldSpaceNormalsCreator), "Create");
}
```

The following image shows the **GameObject** menu option added by the script.

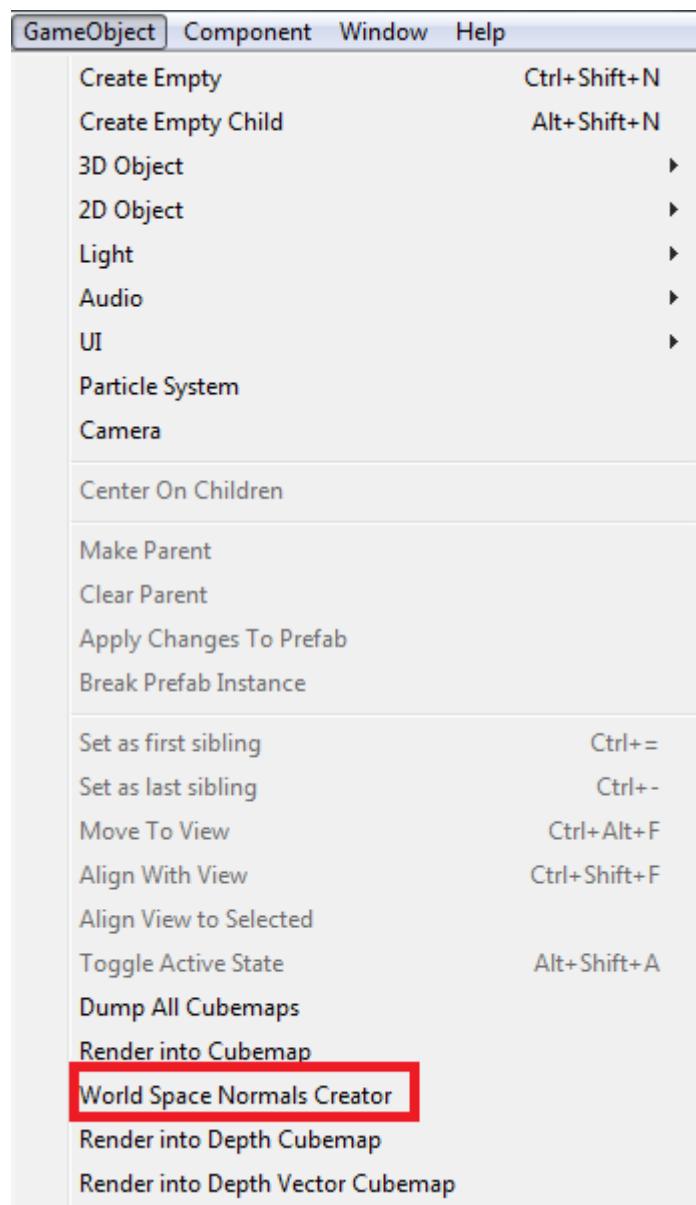


Figure 6-88 GameObject menu option added by the script

The class defined in the C# script derives from the Unity `ScriptableWizard` class and has access to some of its members. Derive from this class to create an editor wizard. Editor wizards are typically opened using a menu item.

In the `OnWizardUpdate` code, the `helpString` variable holds a help message that is displayed in the window that the wizard creates.

The `isValid` member is used to define when all the correct parameters are selected and the `Create` button is available. In this case the `_currentObj` member is checked to ensure it points to a valid object.

The fields of the Wizard windows are the public members of the class. In this case only the `_currentObj` is public so the Wizard window only has one field.

The following image shows the Custom Wizard window:

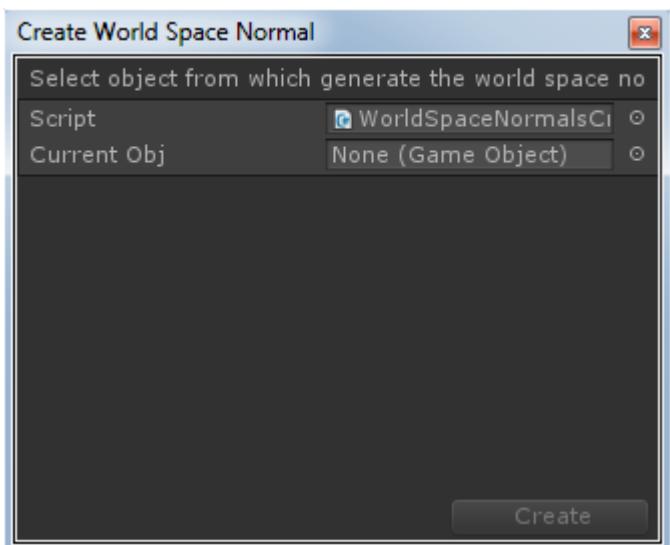


Figure 6-89 Custom Wizard window

When an object is selected and the **Create** button is clicked, the `OnWizardCreate()` function is called.

The `OnWizardCreate()` function performs the main work of the conversion.

To convert the normal, the tool creates a temporary camera that renders the new World space normal to a `RenderTargetTexture`. To do this, the camera is set to orthographic mode and the layer of the object is changed to an unused level. This means it can render the object on its own, even if it is already part of the scene.

The following code shows how the camera is set up:

```
// Set antialiasing
QualitySettings.antiAliasing = 4;
Shader wns = Shader.Find ("Custom/WorldSpaceNormalCreator");
GameObject go = new GameObject( "WorldSpaceNormalsCam", typeof(Camera) );
_renderCamera = go.GetComponent<Camera> ();
_renderCamera.orthographic = true;
_renderCamera.nearClipPlane = 0.0f;
_renderCamera.farClipPlane = 10f;
_renderCamera.orthographicSize = 1.0f;
int prevObjLayer = _currentObj.layer;
_currentObj.layer = 30; //0x40000000
```

The script sets a replacement shader that executes the conversion:

```
_renderCamera.SetReplacementShader (wns,null);
_renderCamera.useOcclusionCulling = false;
```

The camera is pointed at the object. This prevents the object being removed from rendering during frustum culling:

```
_renderCamera.transform.rotation = Quaternion.LookRotation (_currentObj.transform.position -
_renderCamera.transform.position);
```

For each material assigned to the object, the script locates the `_BumpMap` texture. This texture is set as source texture for the replacement shader using the shader global functions.

The clear color is set to `(0.5,0.5,0.5)` because normals pointing at negative directions must be represented.

```
foreach (Material m in materials)
{
    Texture t = m.GetTexture("_BumpMap");
    if( t == null )
    {
        Debug.LogError("the material has no texture assigned named Bump Map");
        continue;
    }
    Shader.SetGlobalTexture ("_BumpMapGlobal", t);
    RenderTexture rt = new RenderTexture(t.width,t.height,1);
```

```

_renderCamera.targetTexture = rt;
_renderCamera.pixelRect = new Rect(0,0,t.width,t.height);
_renderCamera.backgroundColor = new Color( 0.5f, 0.5f, 0.5f);
_renderCamera.clearFlags = CameraClearFlags.Color;
_renderCamera.cullingMask = 0x40000000;
_renderCamera.Render();
Shader.SetGlobalTexture ("_BumpMapGlobal", null);

```

After the camera renders the scene, the pixels are read back and saved as a PNG image.

```

Texture2D outTex = new Texture2D(t.width,t.height);
RenderTexture.active = rt;
outTex.ReadPixels(new Rect(0,0,t.width,t.height), 0, 0);
outTex.Apply();
RenderTexture.active = null;
byte[] _pixels = outTex.EncodeToPNG();
System.IO.File.WriteAllBytes("Assets/Textures/GeneratedWorldSpaceNormals/" + t.name
+ "_WorldSpace.png", _pixels);
}

```

The camera culling mask uses a binary mask represented in hexadecimal format, to specify what layers to render.

In this case layer 30 was used:

```
_currentObj.layer = 30;
```

The hexadecimal is 0x40000000 because its 30th bit is set to 1.

6.15.3 The WorldSpaceNormalCreator shader

The shader code that implements the conversion is quite straightforward. Instead of using the actual vertex position, it uses the texture coordinate of the vertex as its position. This causes the object to be projected onto a 2D plane, the same as when texturing.

To make the OpenGL pipeline work correctly, the UV coordinates are moved from the standard [0,1] range to the [-1,1] range and invert the Y coordinate. The Z coordinate is not used so it can be set to 0 or any value within the near and far clip planes:

```

output.pos = half4(input.tex.x*2.0 - 1.0, ((1.0 - input.tex.y)*2.0 - 1.0), 0.0, 1.0);
output.tc = input.tex;

```

The normal, tangent and bitangent are computed in the vertex shader and passed to the fragment shader to execute the conversion:

```

output.normalInWorld = normalize(mul(half4(input.normal, 0.0), _World2Object).xyz);
output.tangentWorld = normalize(mul(_Object2World, half4(input.tangent.xyz, 0.0)).xyz);
output.bitangentWorld = normalize(cross(output.normalInWorld, output.tangentWorld)
* input.tangent.w);

```

The fragment shader:

1. Converts the normal from tangent space to world space.
2. Scales the normal to the [0,1] range.
3. Outputs the normal to the new texture.

The following code shows this:

```

half3 normalInWorld = half3(0.0,0.0,0.0);
half3 bumpNormal = UnpackNormal(tex2D(_BumpMapGlobal, input.tc));
half3x3 local2WorldTranspose = half3x3(
    input.tangentWorld,
    input.bitangentWorld,
    input.normalInWorld);
normalInWorld = normalize(mul(bumpNormal, local2WorldTranspose));
normalInWorld = normalInWorld*0.5 + 0.5;
return half4(normalInWorld,1.0);

```

The following image shows a tangent space normal map before processing:



Figure 6-90 Original tangent space normal map

The following image shows a world space normal map generated by the tool:

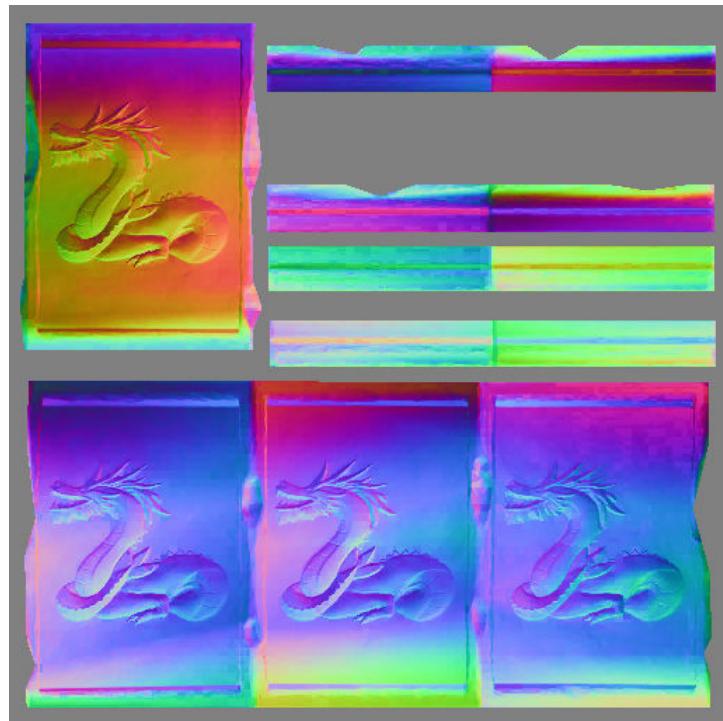


Figure 6-91 Generated world space normal map

6.15.4 The WorldSpaceNormalsCreators C# script

This is the code of the WorldSpaceNormalsCreators C# script:

```
using UnityEngine;
using UnityEditor;
using System.Collections;

public class WorldSpaceNormalsCreator : ScriptableWizard
```

```

{
    public GameObject _currentObj;
    private Camera _renderCamera;
    void OnWizardUpdate()
    {
        helpString = "Select object from which generate the world space normals";
        if(_currentObj != null)
        {
            isValid = true;
        }
        else
        {
            isValid = false;
        }
    }
    void OnWizardCreate ()
    {
        // Set antialiasing
        QualitySettings.antiAliasing = 4;
        Shader wns = Shader.Find ("Custom/WorldSpaceNormalCreator");
        GameObject go = new GameObject( "WorldSpaceNormalsCam", typeof(Camera) );
        //Set the new camera to perform orthographic projection
        _renderCamera = go.GetComponent<Camera> ();
        _renderCamera.orthographic = true;
        _renderCamera.nearClipPlane = 0.0f;
        _renderCamera.farClipPlane = 10f;
        _renderCamera.orthographicSize = 1.0f;
        //Save the current object layer and set it to a unused one
        int prevObjLayer = _currentObj.layer;
        _currentObj.layer = -30; //0x40000000
        //Set the replacement shader for the camera
        _renderCamera.SetReplacementShader (wns,null);
        _renderCamera.useOcclusionCulling = false;
        //Rotate the camera to look at the object to avoid frustum culling
        _renderCamera.transform.rotation = Quaternion.LookRotation
        (_currentObj.transform.position - _renderCamera.transform.position);
        MeshRenderer mr = _currentObj.GetComponent<MeshRenderer> ();
        Material[] materials = mr.sharedMaterials;
        foreach (Material m in materials)
        {
            Texture t = m.GetTexture("_BumpMap");
            if( t == null )
            {
                Debug.LogError("the material has no texture assigned named Bump Map");
                continue;
            }
            //Render the world space normal maps to a texture
            Shader.SetGlobalTexture ("_BumpMapGlobal", t);
            RenderTexture rt = new RenderTexture(t.width,t.height,1);
            _renderCamera.targetTexture = rt;
            _renderCamera.pixelRect = new Rect(0,0,t.width,t.height);
            _renderCamera.backgroundColor = new Color( 0.5f, 0.5f, 0.5f);
            _renderCamera.clearFlags = CameraClearFlags.Color;
            _renderCamera.cullingMask = 0x40000000;
            _renderCamera.Render();
            Shader.SetGlobalTexture ("_BumpMapGlobal", null);
            Texture2D outTex = new Texture2D(t.width,t.height);
            RenderTexture.active = rt;
            outTex.ReadPixels(new Rect(0,0,t.width,t.height), 0, 0);
            outTex.Apply();
            RenderTexture.active = null;
            //Save it to PNG
            byte[] _pixels = outTex.EncodeToPNG();
            System.IO.File.WriteAllBytes("Assets/Textures/GeneratedWorldSpaceNormals/"
            +t.name+"_WorldSpace.png",_pixels);
        }
        _currentObj.layer = prevObjLayer;
        DestroyImmediate(go);
    }
    [MenuItem("GameObject/World Space Normals Creator")]
    static void CreateWorldSpaceNormals ()
}

```

```

    {
        ScriptableWizard.DisplayWizard("Create World Space Normal",
        typeof(WorldSpaceNormalsCreator),"Create");
    }
}

```

6.15.5 The WorldSpaceNormalCreator shader code

The following is the code of the WorldSpaceNormalCreator shader.

```

Shader "Custom/WorldSpaceNormalCreator" {
    Properties {
    }
    SubShader {
        Cull off
        Pass
        {
            CGPROGRAM
            #pragma target 3.0
            #pragma glsl
            #pragma vertex vert
            #pragma fragment frag

            #include "UnityCG.cginc"

            uniform sampler2D _BumpMapGlobal;

            struct vin
            {
                half4 tex : TEXCOORD0;
                half3 normal : NORMAL;
                half4 tangent : TANGENT;
            };

            struct vout
            {
                half4 pos : POSITION;
                half2 tc : TEXCOORD0;
                half3 normalInWorld : TEXCOORD1;
                half3 tangentWorld : TEXCOORD2;
                half3 bitangentWorld : TEXCOORD3;
            };

            vout vert (vin input )
            {
                vout output;
                output.pos = half4(input.tex.x*2.0 - 1.0,((1.0-input.tex.y)*2.0 - 1.0),
                    0.0, 1.0);
                output.tc = input.tex;
                output.normalInWorld = normalize(mul(half4(input.normal, 0.0),
                    _World2Object).xyz);
                output.tangentWorld = normalize(mul(_Object2World,
                    half4(input.tangent.xyz, 0.0)).xyz);
                output.bitangentWorld = normalize(cross(output.normalInWorld,
                    output.tangentWorld) * input.tangent.w);

                return output;
            }

            float4 frag( vout input ) : COLOR
            {
                half3 normalInWorld = half3(0.0,0.0,0.0);
                half3 bumpNormal = UnpackNormal(tex2D(_BumpMapGlobal, input.tc));
                half3x3 local2WorldTranspose = half3x3(
                    input.tangentWorld,
                    input.bitangentWorld,
                    input.normalInWorld);
                normalInWorld = normalize(mul(bumpNormal, local2WorldTranspose));
                normalInWorld = normalInWorld*0.5 + 0.5;

                return half4(normalInWorld,1.0);
            }
        ENDCG
    }
}

```

Chapter 7

Virtual Reality

This chapter describes the process of adapting an application or game to run on virtual reality hardware, and some differences in the implementation of reflections in virtual reality.

It contains the following sections:

- [*7.1 Virtual reality hardware support for Unity* on page 7-185.](#)
- [*7.2 The Unity VR porting process* on page 7-186.](#)
- [*7.3 Things to consider when porting to VR* on page 7-189.](#)
- [*7.4 Reflections in VR* on page 7-191.](#)
- [*7.5 The result* on page 7-196.](#)

7.1 Virtual reality hardware support for Unity

Several types of *Virtual Reality* (VR) hardware have Unity support. Unity natively supports some devices. Plugins can enable support of some other devices.

If a device has native Unity VR support, then it produces better performance than devices that are supported using plugins. This is because Unity implements internal optimizations for the devices with native Unity VR support.

The following devices have native Unity VR support:

- Oculus Rift.
- Samsung Gear VR.
- PlayStation VR.
- Microsoft Hololens.

The following devices have Unity VR support using plugins:

- Google Cardboard.
- Moverio.
- HTC Vive and other devices the SteamVR platform supports.

Several other devices can be supported using the *Open Source Virtual Reality* (OSVR) plugin for Unity.

7.2 The Unity VR porting process

Porting an application or game to native Unity VR is a multi-stage process.

The stages required to port an application to Unity VR are:

1. Install Unity version 5.1 or later. Unity version 5.1 and later natively support VR.
2. If required, obtain the signature file for your device from the appropriate website, and place it in the required folder on your device.

For the Samsung Gear VR that runs the Ice Cave demo, this is the Oculus developer website <https://developer.oculus.com/osig>. The signature file for a Samsung device must go in the Plugins/Android/assets folder.

3. In Unity, select **Open File > Build Settings > Player Settings**. The **Player Settings** window appears. Turn on the **Virtual Reality Supported** option in the Other Settings section.

The following figure shows a screenshot of this window.

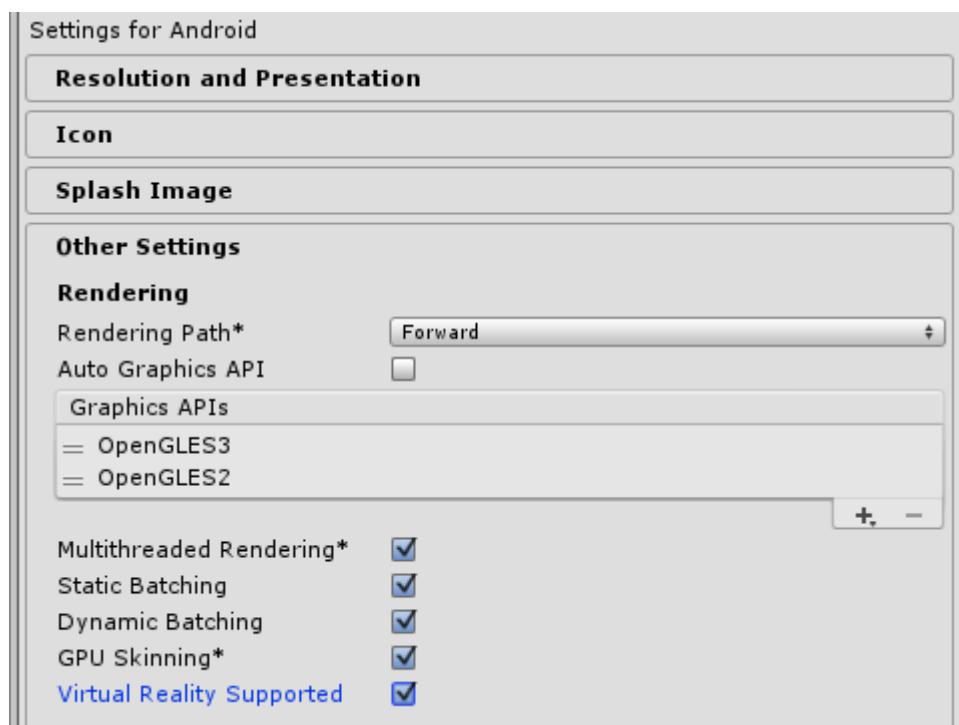


Figure 7-1 The Player Settings window

4. Set the parent for the camera. All camera controls must set the camera position and orientation to this camera parent.
5. Associate the camera control with the VR headset touch pad, if necessary.
6. For Android devices, enable the **Developer options** menu and turn on **USB debugging**.

The following figure shows the **Developer options** menu.

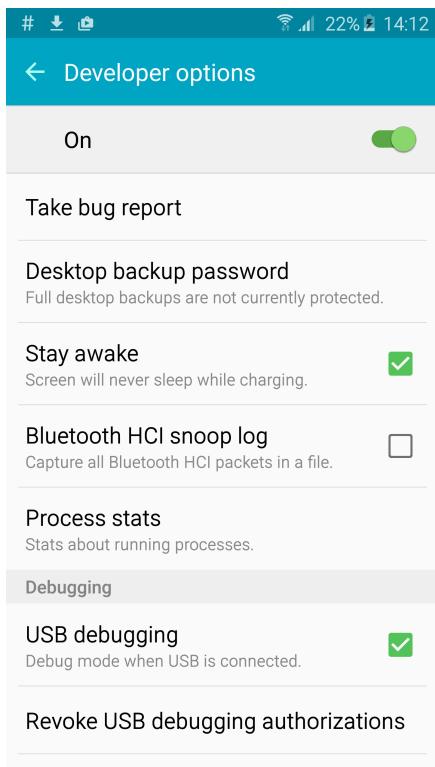


Figure 7-2 The Developer options menu

7. Build the application and install it on the device.
8. Launch the application.

When you launch the application on a Samsung device, it prompts you to insert the device into the headset. If the device is not ready for VR, it prompts you to connect to the network to download the Samsung VR software.

This section contains the following subsection:

- [7.2.1 Enabling Samsung Gear VR developer mode on page 7-187](#).

7.2.1 Enabling Samsung Gear VR developer mode

Developer mode can help you to visualize the VR application running, without inserting the device into the VR headset.

Samsung Gear VR Developer mode can be enabled if you have installed a signed VR application before. If you have not installed a signed VR application before, install a signed VR application so that you can enable this mode.

To enable developer mode:

1. On the Samsung device, select **Settings > Application Manager > Gear VR Service**.
2. Select **Manage storage**.
3. Tap on the VR service version six times.
4. Wait for the scan process to complete. A developer mode toggle appears.

Note

Using developer mode reduces the battery life of the phone, because it overrides all the settings that turn off the headset when it is not in use.

The following figure shows an example Ice Cave screenshot from the Samsung Gear VR developer mode view.



Figure 7-3 An example screenshot from the VR application running in Samsung Gear VR developer mode

7.3 Things to consider when porting to VR

VR creates a very different user experience compared to other application types. This means that some things that work for a non-VR application or game do not work for VR.

Test the application on a few different users to determine their comfort levels and adjust the code so that they find the experience comfortable.

Camera animations that are comfortable in a non-VR game might be uncomfortable in a VR version of the game. For example, the Ice Cave demo without VR has an animation mode for the camera. Some users find this uncomfortable and suffer motion sickness, especially when the camera moves backwards. Removing this mode prevents this unsettling experience.

A non-VR application can be controlled using the touch screen of the phone or by tilting the phone. These control mechanisms might not be possible in a VR application. For example, the original Ice Cave demo uses two virtual joysticks to control the camera, this does not work on a VR device because the touchscreen is not accessible. The Ice Cave VR demo is designed to run on Samsung Gear VR, which has a touchpad on the side of the headset. Using the touchpad instead of the touch screen solves this problem.

The following figure shows the touchpad on a Samsung Gear VR headset.



Figure 7-4 A Samsung Gear VR headset showing the touchpad

Users can find some visual effects that work in a non-VR application look wrong in a VR application. For example, the non-VR Ice Cave demo uses a dirty lens effect that changes its intensity based on the camera alignment with the Sun. Users testing this effect in VR found it looked wrong, so it was removed.

This section contains the following subsection:

- [7.3.1 External device control of the camera on page 7-189](#).

7.3.1 External device control of the camera

VR can benefit from control methods that are not normally associated with phones and mobile devices. Some of these methods can be worth considering, depending on the target audience for the application.

Controllers are available, that you can connect to the VR device using Bluetooth. To implement this, the Ice Cave demo uses a custom plugin that extends the Unity functionality so that it can interpret the Android Bluetooth events. These events trigger movement of the camera.

The following figure shows a Bluetooth controller that works with the Ice Cave demo.



Figure 7-5 Bluetooth controller controlling the Ice Cave demo

7.4 Reflections in VR

Reflections are important in any VR application. The real world contains many reflections, so when they are missing from a game or application, people notice.

Reflections in VR can use the same techniques that traditional games use, but they must be modified to work with the stereo visual output that a user sees.

A good implementation of reflections can make an application or game feel more realistic and immersive. However, there are a few extra issues that you must consider when you implement reflections in VR.

This section contains the following subsections:

- [7.4.1 Reflections using local cubemaps on page 7-191](#).
- [7.4.2 Combining different types of reflection on page 7-191](#).
- [7.4.3 Stereo reflections on page 7-191](#).

7.4.1 Reflections using local cubemaps

You can create reflections using local cubemaps. This method avoids creating the reflection texture each frame, instead it fetches the reflection texture from a prerendered cubemap. This method applies a local correction to the reflection vector based on where the cubemap was generated and the scene bounding box, and uses that information to fetch the correct texture.

Reflections that are generated using local cubemaps do not suffer from pixel instability or pixel shimmering that can occur when a reflection is generated at runtime for each frame.

For more information on reflections using local cubemaps, see [6.2 Implementing reflections with a local cubemap on page 6-98](#).

7.4.2 Combining different types of reflection

Different reflection generation techniques are required to achieve the best performance and effect. Which technique is required depends on the shape of the reflective surface, and whether the reflective surface and the object being reflected are static or dynamic. The result from the different reflection generation techniques must be combined to produce the result that the user sees.

For information on combining different reflection types, see [6.3 Combining reflections on page 6-114](#).

7.4.3 Stereo reflections

In a non-VR game, there is only one camera viewpoint. In VR, there is one camera viewpoint for each eye. This means that the reflection must be computed for each eye individually.

If both eyes are shown the same reflection, then users quickly notice that there is no depth in the reflections. This is inconsistent with their expectations and can break their sense of immersion, negatively affecting the quality of the VR experience.

To correct this problem, two reflections must be calculated and shown with the correct adjustment for the position of each eye as the user looks around in the game.

To implement these reflections in the Ice Cave demo, it uses two reflection textures for planar reflections from dynamic objects, and two different local corrected reflection vectors to fetch the texture from a single local cubemap for static object reflections.

Reflections can be of either dynamic or static objects. Each type of reflection requires a different set of changes to work in VR.

Implementing stereo planar reflections in Unity VR

Implementing stereo reflections in your VR game requires a few adjustments to the non-VR game code.

Before starting, ensure that you have enabled support for virtual reality in Unity. To do this, select **Build Settings > Player Settings > Other Settings** and select the checkbox for **Virtual Reality Supported**.

Dynamic stereo planar reflections

Dynamic reflections require some changes to produce a correct result for two eyes.

You must create two cameras, and a target texture for each camera to render to. Disable both cameras so that their rendering is executed programmatically. Then, attach the following script to both.

```
void OnPreRender(){
    SetUpReflectionCamera();
    // Invert winding
    GL.invertCulling = true;
}
void OnPostRender(){
    // Restore winding
    GL.invertCulling = false;
}
```

This script places and orients the reflection camera using the position and orientation of the main camera. To do this, it calls the `SetUpReflectionCamera()` function just before the left and right reflection cameras render. The following code shows how this function is implemented.

```
public GameObject reflCam;
public float clipPlaneOffset ;
...
private void SetUpReflectionCamera(){
    // Find out the reflection plane: position and normal in world space
    Vector3 pos = gameObject.transform.position;

    // Reflection plane normal in the direction of Y axis
    Vector3 normal = Vector3.up;
    float d = -Vector3.Dot(normal, pos) - clipPlaneOffset;
    Vector4 reflPlane = new Vector4(normal.x, normal.y, normal.z, d);
    Matrix4x4 reflection = Matrix4x4.zero;
    CalculateReflectionMatrix(ref reflection, reflPlane);

    // Update reflection camera considering main camera position and orientation
    // Set view matrix
    Matrix4x4 m = Camera.main.worldToCameraMatrix * reflection;
    reflCam.GetComponent<Camera>().worldToCameraMatrix = m;

    // Set projection matrix
    reflCam.GetComponent<Camera>().projectionMatrix = Camera.main.projectionMatrix;
}
```

This function calculates the view and projection matrices of the reflection camera. It determines the reflection transformation to apply to the view matrix of the main camera, `worldToCameraMatrix`.

To set the position of the cameras for each eye, add the following code after the line
`Matrix4x4 m = Camera.main.worldToCameraMatrix * reflection;`:

For the left eye

```
m[12] += stereoSeparation;
```

For the right eye

```
m[12] -= stereoSeparation;
```

The shift value `stereoSeparation` is 0.011. The `stereoSeparation` value is half the eye separation value.

Attach another script to the main camera to control the rendering of the left and right reflection cameras. The following code shows the Ice Cave implementation of this script.

```
public class RenderStereoReflections : MonoBehaviour
{
    public GameObject reflectiveObj;
    public GameObject leftReflCamera;
    public GameObject rightReflCamera;
    int eyeIndex = 0;

    void OnPreRender(){
```

```

if (eyeIndex == 0){
    // Render Left camera
    leftReflCamera.GetComponent<Camera>().Render();
    reflectiveObj.GetComponent<Renderer>().material.SetTexture(
        "_DynReflTex", leftReflCamera.GetComponent<Camera>().targetTexture);
}
else{
    // Render right camera
    rightReflCamera.GetComponent<Camera>().Render();
    reflectiveObj.GetComponent<Renderer>().material.SetTexture(
        "_DynReflTex", rightReflCamera.GetComponent<Camera>().targetTexture);
}
eyeIndex = 1 - eyeIndex;
}
}

```

This script handles the rendering of the left and right reflection cameras in the `OnPreRender()` callback function of the main camera. This script is called once for the left eye and then once for the right eye. The `eyeIndex` variable assigns the correct render order for each reflection camera and applies the correct reflection to each eye of the main camera. The first time the callback function is called it assumes that it is for the left eye. This is the order that Unity calls the `OnPreRender()` method.

Checking that different textures are in use for each eye

Checking whether the script is producing a different render texture for each eye correctly is important.

To test whether the correct texture is being shown for each eye:

Procedure

1. Change the script so that it passes the `eyeIndex` value to the shader as a uniform.
2. Use two colors for the reflection textures, one for each `eyeIndex` value.

If your script is working correctly, the output is similar to the following figure that shows a screenshot where the two different stable reflections are visible.

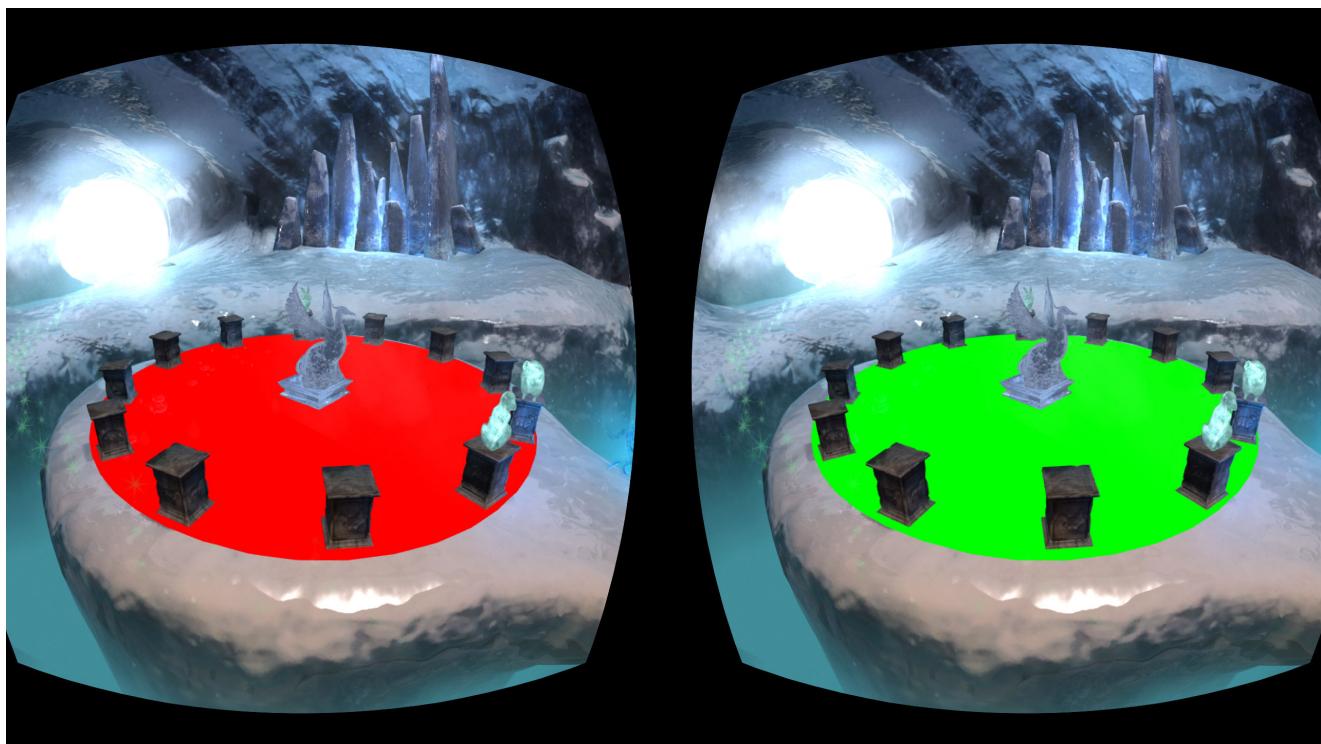


Figure 7-6 An example of a correct reflection texture output check

Static stereo reflections

You can create stereo reflections from static objects efficiently, by using cubemaps. The only difference is that you must use two reflections vectors to fetch the texels from the cubemap, one for each eye.

Unity provides a built-in value for accessing the camera position in world coordinates, in the shader:

```
_WorldSpaceCameraPos.
```

However, in VR, the position of the left and right cameras are required. `_WorldSpaceCameraPos` cannot provide the positions of the left and right cameras. So you must use a script to calculate the position of the left and right cameras, and pass the results to the shader as a single uniform.

Declare a new uniform in the shader that can pass the information for the camera positions:

```
uniform float3 _StereoCamPosWorld;
```

The best place to calculate the left and right camera positions is in the script that is attached to the main camera, because this gives easy access to main camera view matrix. The following code shows how to do this for the `eyeIndex = 0` case.

The code modifies the view matrix of the main camera to set the position of the left eye in local coordinates. The left eye position is required in world coordinates, so the inverse matrix is found. The left eye camera position is passed to the shader through the uniform `_StereoCamPosWorld`.

```
Matrix4x4 mWorldToCamera = gameObject.GetComponent<Camera>().worldToCameraMatrix;
mWorldToCamera[12] += stereoSeparation;
Matrix4x4 mCameraToWorld = mWorldToCamera.inverse;
Vector3 mainStereoCamPos = new Vector3(mCameraToWorld[12], mCameraToWorld[13],
                                         mCameraToWorld[14]);
reflectiveObj.GetComponent<Renderer>().material.SetVector("_StereoCamPosWorld",
    new Vector3 (mainStereoCamPos.x, mainStereoCamPos.y, mainStereoCamPos.z));
```

The code is the same for the right eye, except the stereo separation is subtracted from `mWorldToCamera[12]` instead of added.

In the vertex shader you must find the following line, it is responsible for calculating the view vector:

```
output.viewDirInWorld = vertexWorld.xyz - _WorldSpaceCameraPos;
```

Replace this with the following line that uses the new left and right eye camera positions in world coordinates:

```
output.viewDirInWorld = vertexWorld.xyz - _StereoCamPosWorld;
```

When the stereo reflection is implemented it is visible when the application runs in editor mode, because the reflection texture flickers as it repeatedly changes from the left eye to the right eye. This flickering is not visible in the VR device because a different texture is used for each eye.

Optimizing stereo reflections

Without further optimizations, the stereo reflection implementations run all the time. This means that processing time is wasted on reflections when they are not visible.

Insert code that checks whether a reflective surface is visible, before any work is performed on the reflections themselves. To do this, attach code similar to the following code example to the reflective object.

```
public class IsReflectiveObjectVisible : MonoBehaviour
{
    public bool reflObjIsVisible;

    void Start(){
        reflObjIsVisible = false;
    }

    void OnBecameVisible(){
        reflObjIsVisible = true;
    }

    void OnBecameInvisible(){}
```

```
        reflObjIsVisible = false;  
    }  
}
```

After defining this class, use the following `if` statement in the script attached to the main camera so that the calculations for stereo reflections are only executed when the reflective object is visible.

```
void OnPreRender(){  
    if (reflectiveObjetc.GetComponent<IsReflectiveObjectVisible>().reflObjIsVisible){  
        ...  
    }  
}
```

The rest of the code goes inside this `if` statement. This `if` statement uses the class `IsReflectiveObjectVisible` to check whether the reflective object is visible. If it is not visible, then the reflection is not calculated.

7.5 The result

This work creates a VR version of your game that implements stereo reflections that contribute to the sense of immersion, this improves the whole VR user experience. Changing to stereo reflections significantly improves the user experience because when they are not implemented, people notice that depth is missing from the reflections.

The following figure shows an example screenshot from the Ice Cave demo running in developer mode, showing the stereo reflections that are implemented.

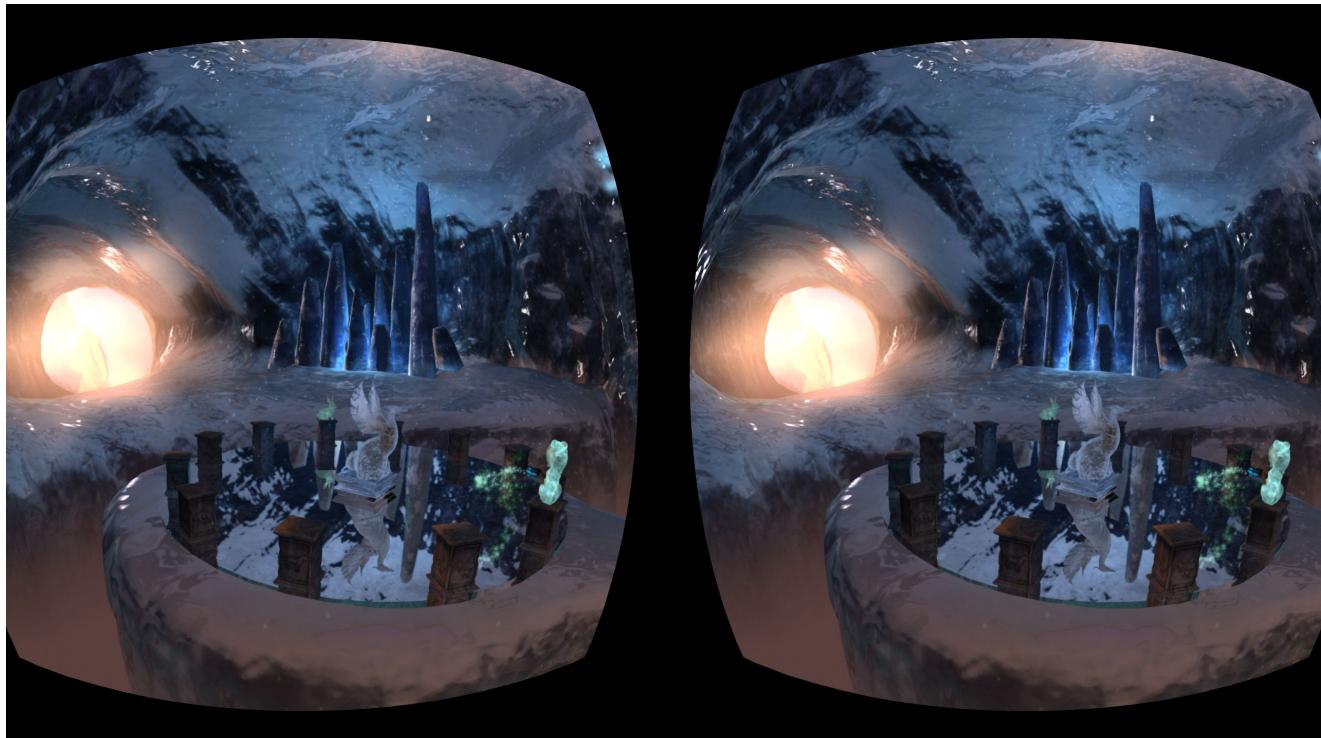


Figure 7-7 A screenshot from the Ice Cave demo

Chapter 8

Vulkan

This chapter describes Vulkan and how you enable it.

Vulkan is a cross-platform graphics and compute API from the Khronos group that offers many benefits over OpenGL and OpenGL ES. These benefits include:

- Providing one unified API framework for mobile, desktop, console, server, and embedded systems.
- Support for hardware with a wide range of capabilities.
- Minimal driver overhead enables high performance on ARM Mali GPU hardware.
- Applications gain more low-level access over GPU and compute resources.
- Reduced application processor bottlenecks.
- Support for multi-threading and multi-processing.
- Efficient use of multiple application processors.
- Use of the SPIR-V intermediate language for shaders, reducing runtime kernel compilation time.
- No requirement to ship your shader source code.
- Reduced energy consumption because of lower application processor overheads, simplified drivers, and using more on-chip memory.

It contains the following sections:

- [8.1 About Vulkan on page 8-198](#).
- [8.2 About Vulkan in Unity on page 8-200](#).
- [8.3 Enabling Vulkan in Unity on page 8-201](#).
- [8.4 Vulkan case study on page 8-202](#).

8.1 About Vulkan

Vulkan is a cross-platform graphics and compute API from the Khronos group.

Vulkan has many advantages over the older OpenGL and OpenGL ES standards:

Unified and portable

Vulkan provides one unified API framework for mobile, desktop, console, and embedded systems.

It has portability across a wide range of implementations and is useful for a wide range of applications.

Simpler drivers

Vulkan uses simpler drivers to minimize driver overhead. The lower latency and better efficiency mean that an application can achieve better performance using Vulkan than OpenGL ES 3.1.

The simpler drivers reduce application processor bottlenecks.

Your application performs resource management, and has direct low-level control over your GPU.

Multi-threading and multi-processing

Vulkan supports multi-threading across multiple application processors. This feature enables you to use multiple application processors efficiently, lowering processing load, and power consumption.

Your application controls thread management and synchronization.

Command Buffers

You can use multi-threading to do command creation for command buffers in parallel. You can also use a separate submission thread to place command buffers into command queues.

You can add graphics, compute, and DMA command buffers.

Different graphics, DMA, and compute queues provide flexibility for job dispatch.

Multi-threaded command creation enables your code to run across multiple application processor cores, increasing performance. Using multiple application processors running at a lower clock rate, rather than one higher clocked processor also reduces power consumption.

SPIR-V

Vulkan uses the SPIR-V intermediate language that enables you to use a common language front end.

SPIR-V is a multi-API, intermediate language for parallel compute and graphics. It includes flow control, graphics, and parallel compute constructs.

It provides native representation for the Vulkan shader and OpenCL kernel source languages. Multiple platforms can use the same SPIR-V front-end compiler to generate pre-compiled shaders.

Using SPIR-V means that there is no front-end compiler in the Vulkan driver, so the driver is simpler and shader compilation is faster.

Using an intermediate language means that you are not required to ship shader source code with your application. It also provides flexibility for using different shading languages in the future.

Loadable Layers

Vulkan enables you to load software layers in development for testing and debugging.

You can remove the additional software layers for production, so your shipping products do not have the testing overhead.

Multi-pass rendering

Multi-pass rendering is a technique where you declare everything ahead of time in the render pass. You can specify different outputs for each sub pass, and chain them together.

Multi-pass enables the driver to make optimizations when pixels in one sub pass access the results of the previous sub pass at the same pixel location. In this way, data can be contained on the fast on-chip memory, saving bandwidth and power. This process is more efficient on tile-based GPUs such as Mali GPUs.

Example use-cases for multi-pass rendering are:

- Deferred rendering
- Soft-particles
- Tone-mapping

The following diagram shows the structure of Vulkan.

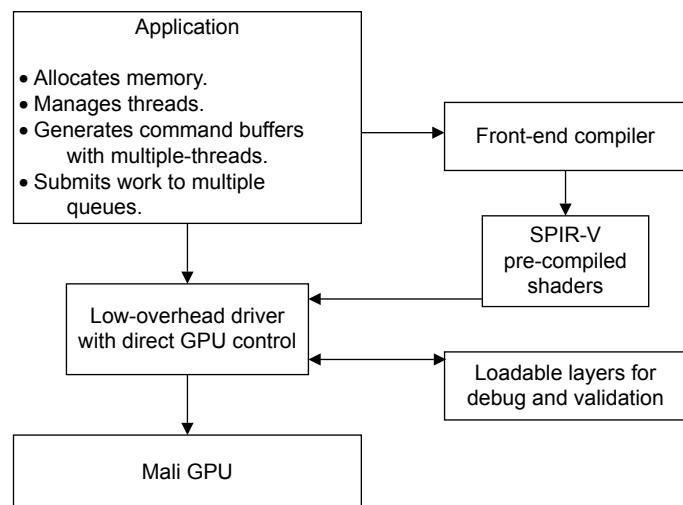


Figure 8-1 Vulkan structure

8.2 About Vulkan in Unity

The advantages of using Vulkan in Unity.

Many of the advantages of Vulkan come from the fact that Vulkan's driver is simpler compared to OpenGL and OpenGL ES.

However, there is the downside that much of the management work must be handled within your application and this low-level access makes applications more complex to write.

Unity handles most of this work for you, so your application automatically benefits from the enhanced performance of Vulkan.

In Unity, all you are required to do is add Vulkan to the list of graphics APIs that Unity can use to build your application.

Vulkan is more efficient than OpenGL ES, it reduces power consumption and increases battery life so it is typically a better choice if it is available. Consider using Vulkan if:

- You want maximum performance for your application.
- Your application is bound to the application processor.
- You suspect an OpenGL or OpenGL ES driver is causing problems.

8.3 Enabling Vulkan in Unity

This describes the procedure to enable Vulkan in Unity 5.6.

To enable Vulkan:

1. Select **File > Build Settings**
2. In the new window dialog, press the **Player Settings ...** button.
3. In the **Player Settings** window, find the **Other Settings** section.
4. Ensure the **Auto Graphics API** check box is not checked. This enables you to manually select the API.
5. To add Vulkan as an API, press **+** to add a new API to the list, and add Vulkan. Vulkan is added as the last option.
6. To make Vulkan the primary API, select **Vulkan** and move to the top of the list.

Vulkan is now the primary API of the player. This is shown in the following image.

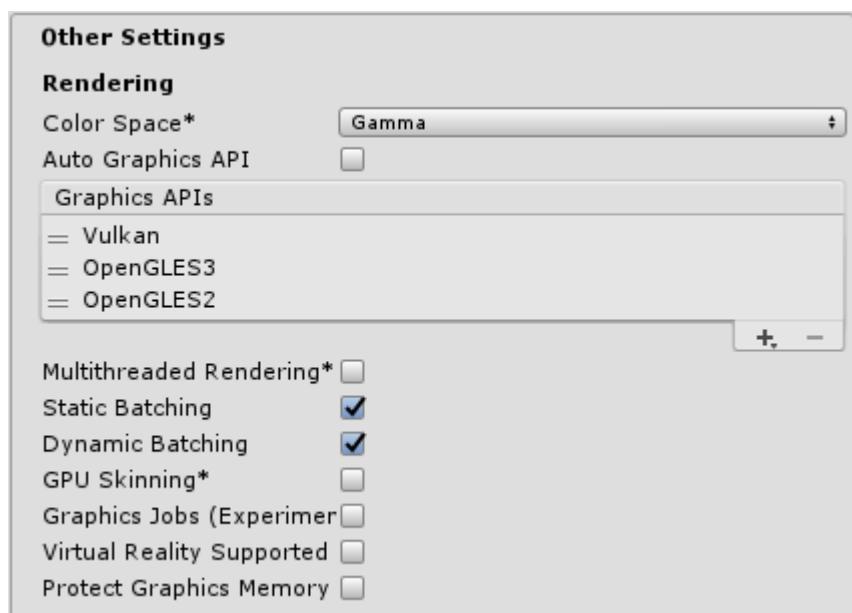


Figure 8-2 Unity API selection

You can leave Vulkan as the only API by removing the other two.

To do this, select an API and press – to remove it.

8.4

Vulkan case study

How game developer Infinite Dreams used Vulkan to improve performance.

Mobile game developers Infinite Dreams develop the popular mobile game Sky Force Reloaded. Sky Force Reloaded features intense action in a rich graphical environment that is highly demanding of both the application processor and the GPU. Sky Force Reloaded was originally built using Unity with OpenGL ES.



Figure 8-3 Sky Force Reloaded

Sky Force Reloaded optimization

The initial development of the game produced the rich graphical experience required. The next stage was optimization.

The game is graphically complex with many things on screen at the same time so the team thought that the biggest area for optimization would be fill rate. Problems with fill rate can usually be eliminated by decreasing the resolution of the framebuffer.

However, when they tried rendering the game in low resolution it still had performance problems. Even on high end devices, it could not always maintain 60 *frames per second* (FPS). The team found the game was making a large number of draw calls. Sometimes, up to 1000 draw calls per frame.

Every draw call has a computation overhead so making large numbers of draw calls is computationally expensive. This causes the OpenGL ES driver to keep the application processor busy for long periods of time preparing data for the GPU. As a result, this can slow down even high end mobile devices.

To optimize, the team could minimize the number of draw calls or modify them so that the game engine could batch them. However, this is not always possible without reducing the quality of the game.

Testing Vulkan

The team heard about Vulkan and they decided to see if it could improve the performance of the game.

In Unity, Vulkan is just a rendering API and Unity does all the hard work. The team were able to use Vulkan by just activating it in Unity.

To test the difference between OpenGL ES and Vulkan, the team analyzed the game. They found one of the slowest parts of the game, where OpenGL could not deliver 60 FPS. To measure the performance difference between OpenGL ES and Vulkan, the development team created a synthetic benchmark based on this part of the game.

The benchmark uses a repeatable scene for the APIs to render, to provide a consistent test for both APIs.

In the following graph, you can see that Vulkan is able to keep 60 FPS most of the time while OpenGL ES struggles. The overall improvement in performance is 15% above OpenGL ES. The team considered this to be a very good result, considering Vulkan is only a graphics API.

The following image shows results of the benchmark comparing OpenGL ES and Vulkan:

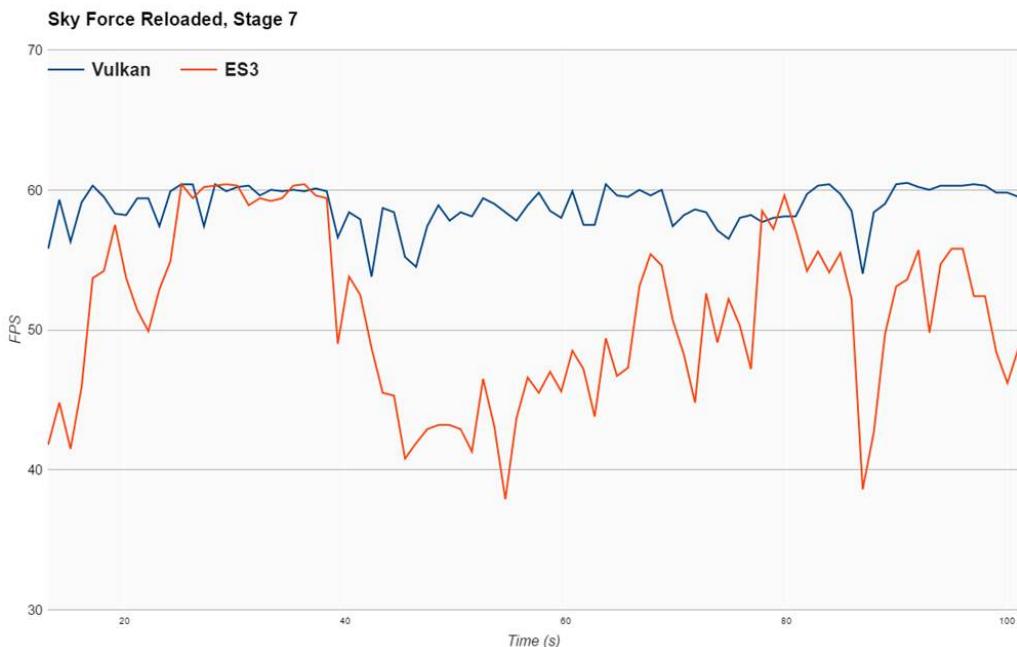


Figure 8-4 Sky Force Reloaded benchmark comparing OpenGL ES and Vulkan

Testing additional performance from Vulkan

With Vulkan able to run at 60 FPS most of the time, the team wondered how much more performance could they achieve with Vulkan.

The developers started adding some additional objects to the benchmark level, to the point where even Vulkan was not able to achieve 60 FPS, causing the gap between OpenGL ES and Vulkan to grow. On average, Vulkan was 32% faster than OpenGL ES.

The following graph shows results of the benchmark comparing OpenGL ES and Vulkan, now with added objects in the scene:

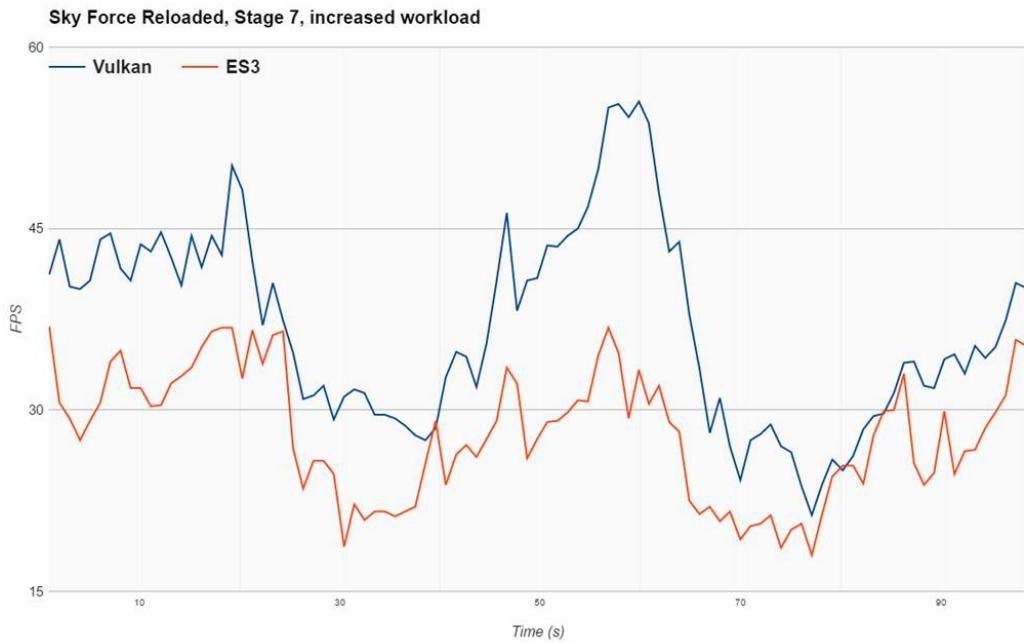


Figure 8-5 Sky Force Reloaded benchmark with added objects

There was definitely some extra performance available, so the team asked themselves how can this extra power be utilized in the game. They added more graphics, particles, objects, and animations. The team found they could make the game look even richer and still maintain 60 FPS. This was all effectively for free, because of Vulkan.

Vulkan video comparison

There is a video comparison available at: <https://www.youtube.com/watch?v=VCSkp-QZ37M>.

The video shows a side-by-side performance comparison of the Vulkan version versus the OpenGL ES version. On the left-hand side, the game uses OpenGL ES. On the right-hand side, the game uses Vulkan. Both versions are running at 60 frames per second.

At this frame rate, Vulkan can render six times more stars and twice as many bullets compared to the OpenGL ES version. With OpenGL ES, rendering the same number of objects causes the frame rate to drop considerably during busy scenes. Vulkan enables you to add more geometry to the screen in comparison to OpenGL ES at the same frame rate.

The following image shows a frame from the comparison video:



Figure 8-6 Sky Force Reloaded frame comparing OpenGL ES and Vulkan

Vulkan power consumption

Sky Force Reloaded is a very application processor and GPU intensive game. Some players have even complained that the game drained their batteries too quickly. The team did not want to take anything away from the console-like quality of game, so they checked to see if Vulkan could help. They tested with Vulkan to see if it would reduce the power consumption and therefore increase the battery life. Enabling Vulkan in Unity resulted in 10-12% lower power consumption, resulting in 10-12% minutes of extra play time.

There is a video comparison available at: <https://www.youtube.com/watch?v=WI7nXq8oozw>

Chapter 9

The Mali Graphics Debugger

This chapter describes the *Mali Graphics Debugger* (MGD).

It contains the following sections:

- [9.1 About Mali Graphics Debugger](#) on page 9-207.
- [9.2 MGD features](#) on page 9-208.
- [9.3 Enabling MGD](#) on page 9-210.
- [9.4 Enabling MGD with Vulkan](#) on page 9-214.

9.1 About Mali Graphics Debugger

The Mali Graphics Debugger is an API trace and debug tool for OpenGL ES, Vulkan, and OpenCL. It helps you to identify possible issues with your applications.

Mali Graphics Debugger provides you with a set of tools that enable you to monitor and analyze the operation of your application.

It can display different kinds of information and statistics that tell you about the data accesses and processing in your application. You can use this information to help you debug applications or to identify performance bottlenecks.

The following image shows Mali Graphics Debugger:

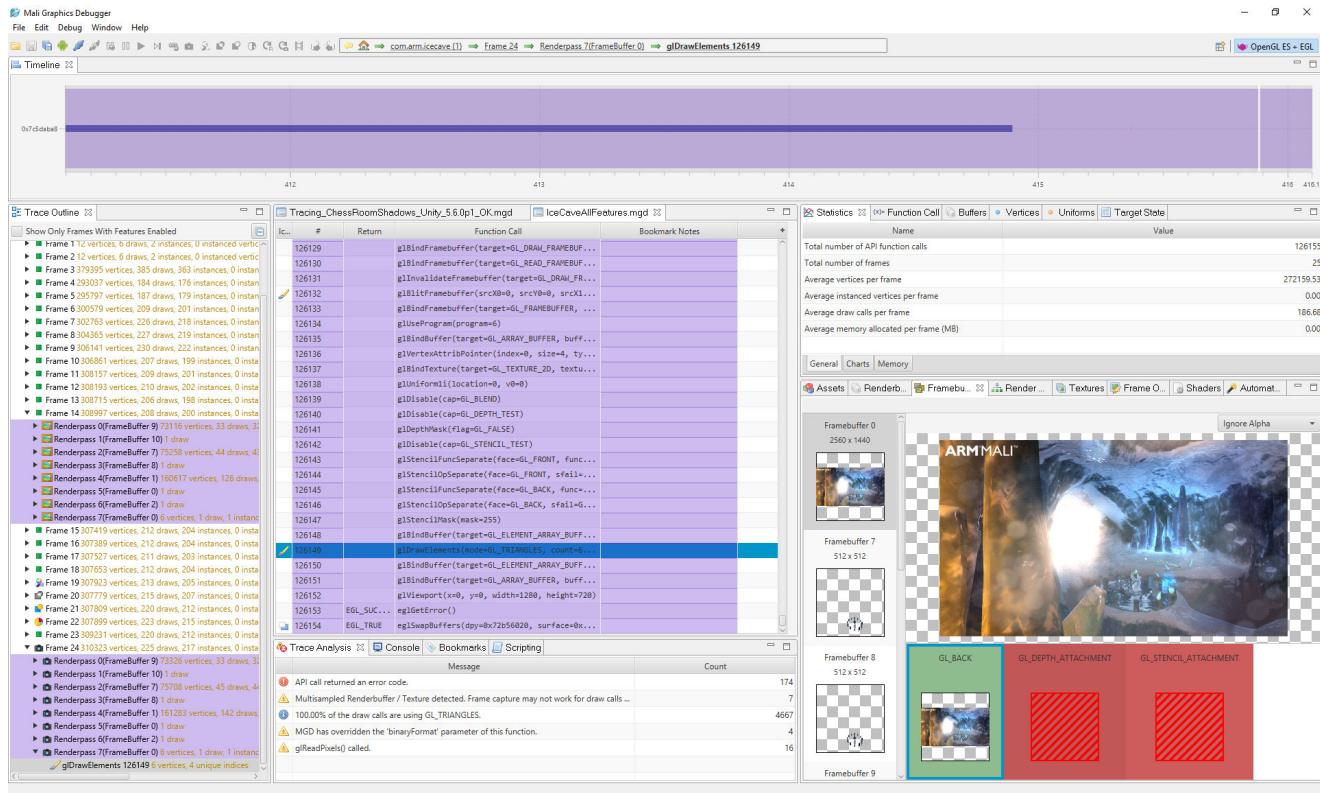


Figure 9-1 Mali Graphics Debugger

For more information, see <https://developer.arm.com/products/software-development-tools/graphics-development-tools/mali-graphics-debugger>.

9.2

MGD features

The Mali Graphics Debugger has a number of features that help you identify possible issues with your applications.

MGD enables you to trace API calls and understand their frame-by-frame effects in your application. It can also highlight common API problems and provide performance improvement advice.

MGD has the following features and benefits:

Trace API calls

Draw-call by draw-call stepping helps you to identify draw-call related issues, redundant draw calls, and other opportunities to optimize.

You can trace OpenGL ES, Vulkan, EGL, and OpenCL API calls with arguments, time, process ID, and thread ID.

There is full trace from the start and it automatically checks for errors. There are search capabilities and it includes multi-process support.

An automated trace view enables you to run a range of standard MGD commands automatically when a specific frame is encountered.

Outline view

The outline view enables you to investigate at the frame level to find what draw calls have a higher geometry impact. It provides you with quick access to important graphics events such as frames, render targets, and draw calls.

Target state debug

Target state debug provides a complete view of the graphics state and state changes, at any point in the trace. You can also use this to discover when and how a state is changed.

Every time a new API call is selected in the trace, the state is updated. This is useful for debugging problems and understanding the causes of performance issues.

You can jump between the outline view and the trace view seamlessly to aid your investigations.

Trace statistics and analysis

You can analyze the trace to find common issues such as performance warnings.

MGD provides per-frame and per draw-call, high-level statistics about the trace.

MGD also shows API misusage by highlighting problematic API calls in the trace view.

Uniform and vertex attribute data views

When you select a draw call, MGD shows all the associated uniform and vertex attribute data. The vertex attribute and uniform views enable you to see vertex and index buffer data. You can analyze these to see their usage in your application.

The Uniform view shows uniform values such as samplers, matrices, and arrays.

The vertex attributes view shows vertex attributes such as vertex names and vertex positions.

Buffer and textures access

Texture view enables you to visualize the texture usage in your applications, helping you to identify opportunities to compress textures or change formats.

Client and server-side buffers are captured every time they change. This enables you to see how each API call affects the buffers.

All textures uploaded are captured at native resolution, enabling you to check the size, format, and type of the textures. You can also access all the large assets including data buffers and textures, for debugging.

Shaders reports and statistics

Shader statistics and cycle counts help you understand what vertex and fragment shaders are the most computationally expensive.

All the shaders used by the application are measured. For each draw call, MGD counts the number of times the shader has been executed, and calculates overall statistics.

Each shader is also compiled with the Mali Offline Compiler. This statically analyzes the shader to display the number of instructions for each GPU pipeline, the number of work registers, and uniform registers.

Frame capture and analysis

MGD can capture frames so you can analyze the effect of each draw call.

A native resolution snapshot of each framebuffer is captured after every draw call. The snapshot capture happens on the target device, so you can investigate target-dependent bugs or precision issues.

You can also export the captured images for separate analysis.

Alternative drawing modes

MGD includes specialist drawing modes that help you analyze your frames. You can use these for both live rendering and frame captures:

Native mode

In Native mode, frames are rendered with the original shaders.

Overdraw mode

The overdraw mode highlights where overdraw is present. You can analyze the impact of overdrawn pixels on each framebuffer by looking at the overdraw captures and histograms.

Shader map mode

In shader map mode, native shaders are replaced with different solid colors. This shows you where the different shaders are used in a frame.

Fragment count mode

In fragment count mode, all the fragments that are processed for each frame are counted. This shows you how many fragments are used in different parts of a frame.

9.3

Enabling MGD

This describes how to enable MGD in Unity.

MGD library loading is integrated in Unity. You can use MGD with Unity to profile OpenGL ES applications.

To enable MGD, you require the following:

- Unity 5.6 or higher.
- Your Android device must have Android 4.2 or above.
- The latest version of MGD from <https://developer.arm.com/graphics>.
- Android *Software Development Kit* (SDK).
- Your PATH must include the path to the adb binary.
- You must have a valid ADB connection to your target device. ADB must be able to return the ID of your device with no permission errors and you must be able to execute applications. For more information, see your Android device documentation.
- Your target device must permit TCP/IP communication on port 5002.

To add MGD support to a Unity application, follow these steps:

1.

Locate the MGD installation folder. For example, on Windows OS the default installation path is C:\Programs Files\ARM\Mali Developer Tools\Mali Graphics Debugger vX.Y.0. Locate the subfolder target\android\arm\ that contains the MGD Android Application MGD.apk.

To install this apk on the Android device where you want to profile your Unity game or application, type this command:

```
adb install -r MGD.apk
```

2.

Locate the libMGD.so library in the MGD installation folder under target\android\arm\unrooted\armeabi-v7a\.

In your Unity project, create the subfolder Assets\Plugins\Android\. Place the libMGD.so library here.

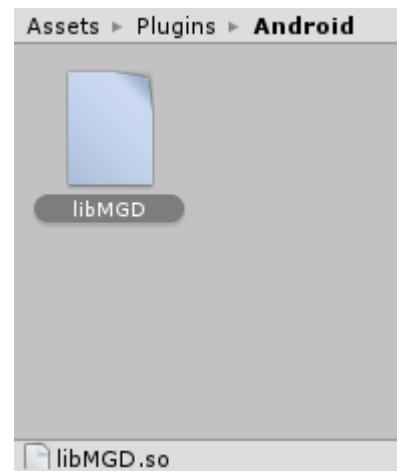


Figure 9-2 Copy libMGD.so to Unity project

3.

In the Unity Build Settings dialog window, check the **Development Build** option. This instructs Unity to package libMGD.so library into the apk and load it at runtime on the host device.

Build your apk.

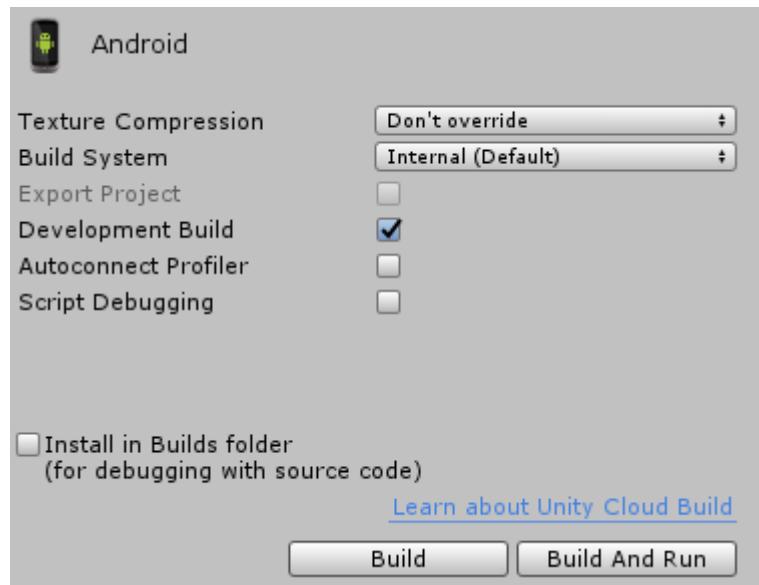


Figure 9-3 Check the Development Build option

4.

To install the Unity application with MGD support on an Android device, type:

```
adb install -r YourApplication.apk
```

5.

To enable MGD to connect to the daemon over USB, run the following command on the host:

```
adb forward tcp:5002 tcp:5002
```

6.

Launch the MGD Android application on the host and enable the MGD Daemon by sliding the button. Your Unity application linked to libMGD.so library is listed.

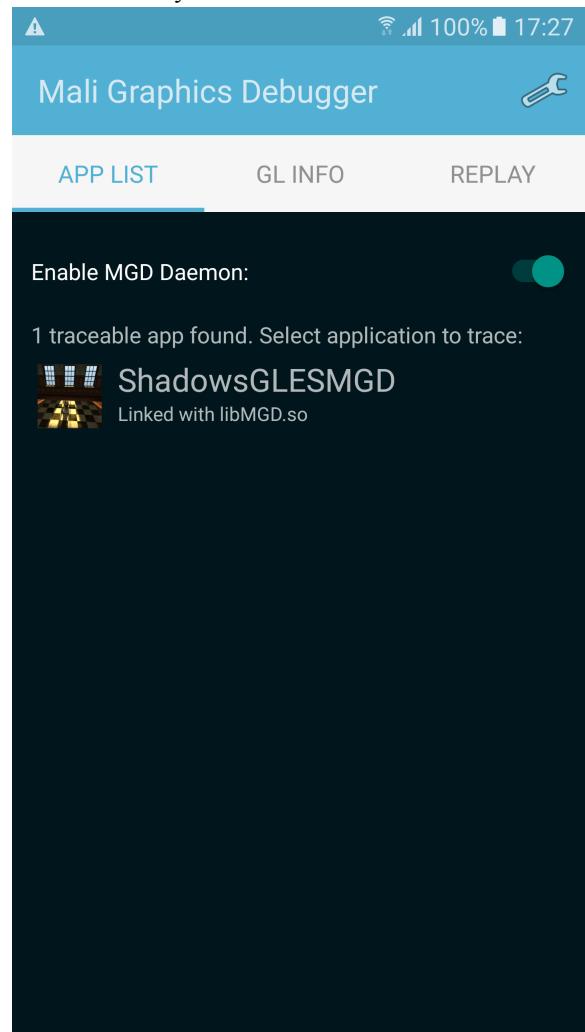


Figure 9-4 Enabling the MGD daemon

7.

Launch MGD on your desktop. Click on the  icon to connect to the target and start tracing. MGD waits for trace input.

8.

Tap on your application where it is listed by the MGD Android application, on your device. The MGD interceptor sends all the OpenGL ES and EGL calls made by your application, to the MGD host application.

The following image shows MGD tracing Unity with OpenGL ES:

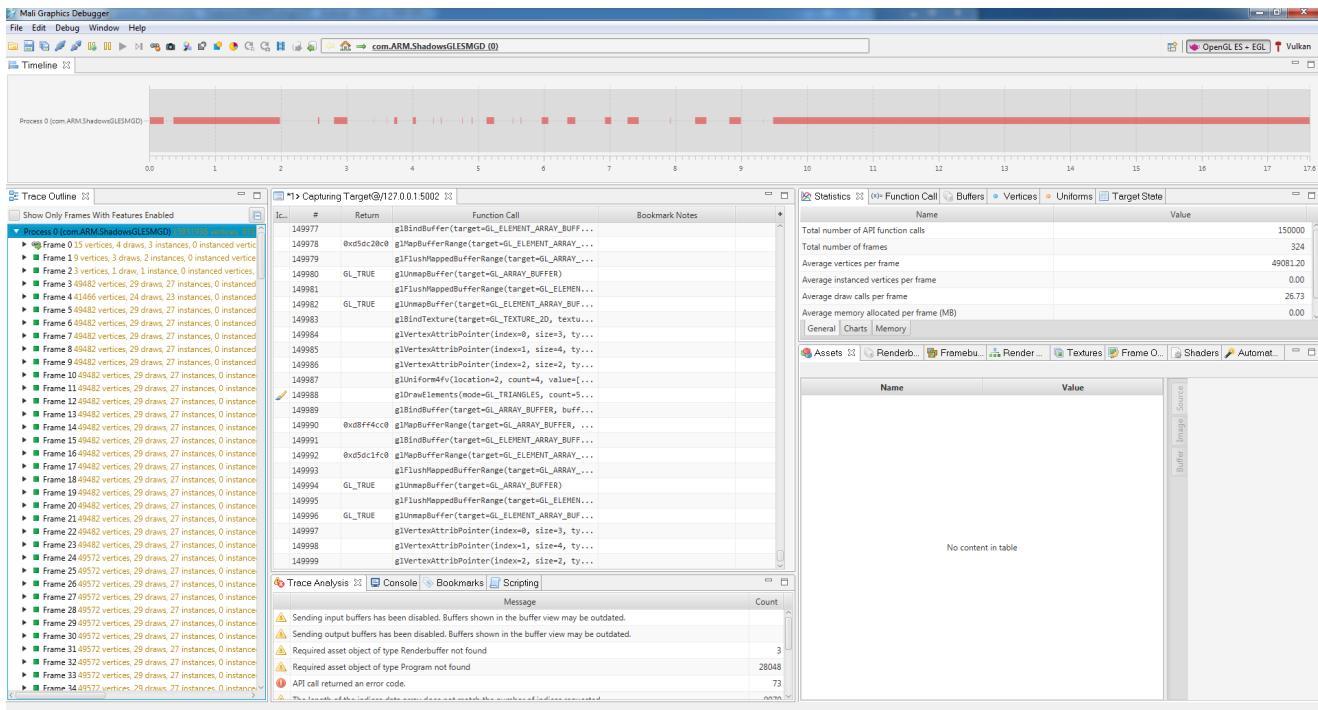


Figure 9-5 MGD tracing Unity with OpenGL ES

9.4 Enabling MGD with Vulkan

This describes how to enable MGD in Unity for Vulkan.

MGD supports debugging of graphics applications developed with the high-performance Vulkan API.

To debug your applications in Vulkan, you use additional layers that you can inject into the system at runtime. MGD includes an interceptor library that you can use as a Vulkan layer.

To add layers to a Unity Android application based on Vulkan, you must copy the layer libraries into the project folder `Assets/Plugins/Android/libs/armeabi-v7a`.

To build a Vulkan application with MGD support, you must also copy the MGD library in the same folder.

To enable MGD for Unity with Vulkan, you require the following:

- Unity 5.6 or higher.
- Your Android device must have Android 4.2 or above.
- The latest version of MGD from <https://developer.arm.com/graphics>.
- Android *Software Development Kit* (SDK).
- Your PATH must include the path to the adb binary.
- You must have a valid ADB connection to your target device. ADB must be able to return the ID of your device with no permission errors and you must be able to execute applications. For more information, see your Android device documentation.
- Your target device must permit TCP/IP communication on port 5002.

To add MGD support to a Unity application developed with Vulkan, follow these steps:

1.

Locate the MGD installation folder. For example, on Windows OS the default installation path is `C:\Programs Files\ARM\Mali Developer Tools\Mali Graphics Debugger vX.Y.0`.
Locate the subfolder `target\android\arm\`.
That contains the MGD Android Application MGD.apk.
To install this apk on the Android device where you want to profile your Unity game or application, type this command:
`adb install -r MGD.apk`

2.

Locate the MGD library libMGD.so in the MGD installation folder, under the folder `target\android\arm\rooted\armeabi-v7a\`.

Make a copy of the library and rename it to `libVkLayerMGD32.so`.

Copy this into the Unity project folder `Assets\Plugins\Android\libs\armeabi-v7a`.

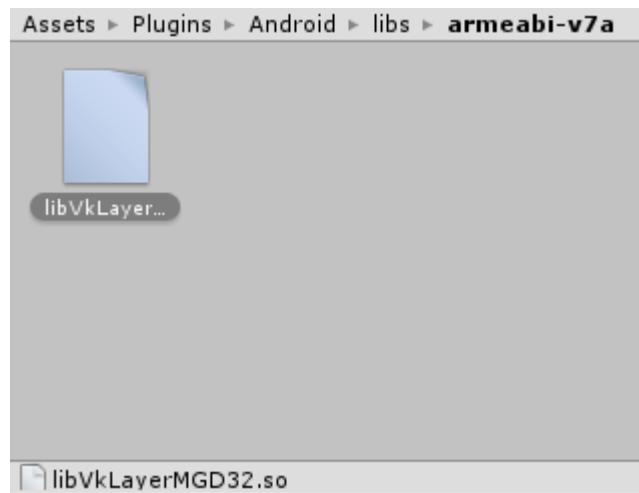


Figure 9-6 Copy libVkLayerMGD32.so to Unity project

3.

In the Unity Build Settings dialog window, check the **Development Build** option. This instructs Unity to enable the Vulkan layer mechanism at runtime. The libVkLayerMGD32.so library is packaged into the apk independently of the state of Development Build option.

Build your apk.

4.

Install your application by running:

```
adb install -r YourApplication.apk
```

5.

To enable MGD to connect to the daemon over USB, run the following command on the host:
`adb forward tcp:5002 tcp:5002`

6.

Launch the MGD Android Application on the host and enable the MGD Daemon by sliding the button.

————— Note ————

If you are using MGD version 4.5 or lower, the MGD Daemon does not list your application.

7.

To enable the MGD Vulkan layer, type the following commands:

```
adb shell
```

```
setprop debug.vulkan.layers VK_LAYER_ARM_MGD
```

This command instructs the Vulkan loader to load the layer with the name `VK_LAYER_ARM_MGD`.

8.

Launch MGD on your desktop. Click on the icon to connect to the target and start tracing. Press the perspective icon at the top right corner of the MGD window, and select Vulkan.

MGD waits for trace input.

9.

Tap on your application on your device. The MGD Vulkan layer sends the Vulkan calls to MGD and MGD displays them.

The following image shows MGD tracing Unity with Vulkan:

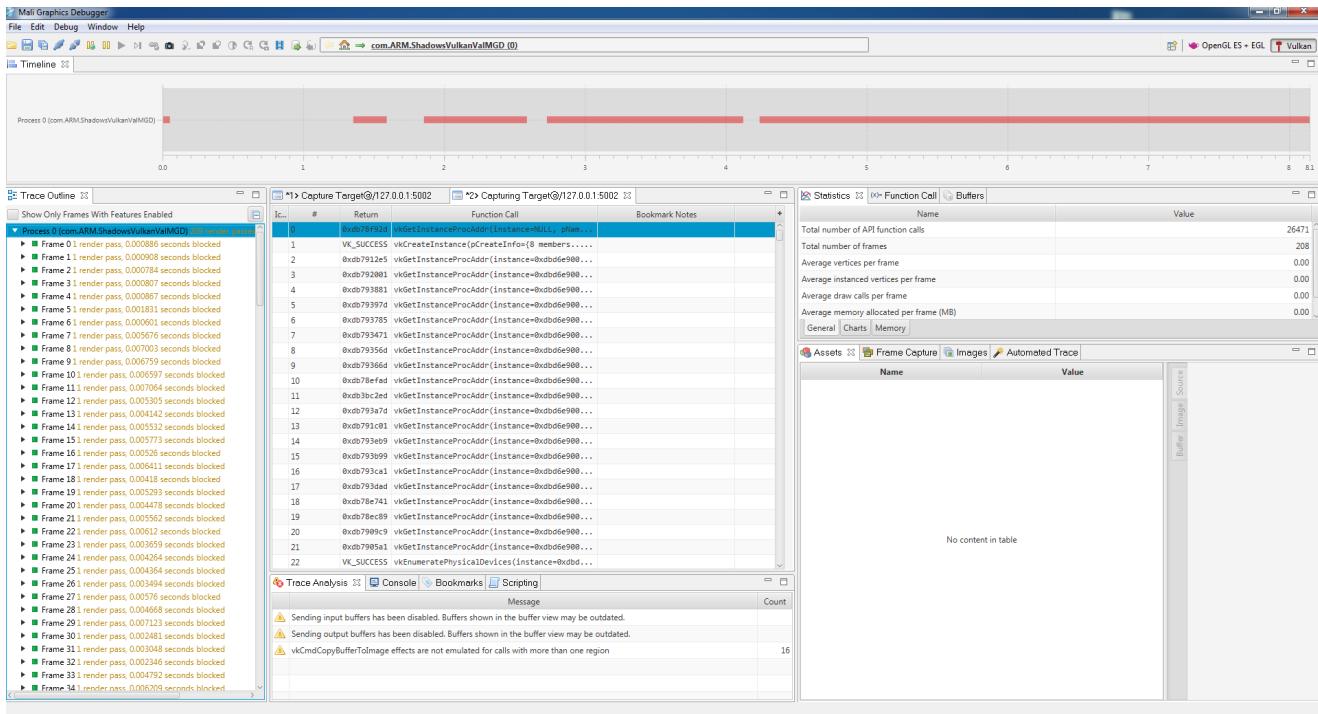


Figure 9-7 MGD tracing Unity with Vulkan

Appendix A

Revisions

This appendix describes the changes between released issues of this book.

It contains the following section:

- *A.1 Revisions* on page Appx-A-218.

A.1 Revisions

This appendix describes the additions to this issue of this book.

Table A-1 Additions in version 3.0

Addition	Location	Affects
Using Enlighten in custom shaders	5.5 Using Enlighten in custom shaders on page 5-80	Version 3.0
Combining reflections	6.3 Combining reflections on page 6-114	Version 3.0
Using Early-z	6.7 Using Early-z on page 6-137	Version 3.0
Dirty lens effect	6.8 Dirty lens effect on page 6-138	Version 3.0
Light shafts	6.9 Light shafts on page 6-141	Version 3.0
Fog effects	6.10 Fog effects on page 6-145	Version 3.0
Bloom	6.11 Bloom on page 6-152	Version 3.0
Icy wall effect	6.12 Icy wall effect on page 6-159	Version 3.0
Procedural skybox	6.13 Procedural skybox on page 6-165	Version 3.0
Fireflies	6.14 Fireflies on page 6-173	Version 3.0
Tangent space to world space conversion tool.	6.15 The tangent space to world space normal conversion tool on page 6-177	Version 3.0

Table A-2 Changes in version 3.0_01

Change	Location	Affects
Removed a list entry from Using Early-z	6.7 Using Early-z on page 6-137	Version 3.0

Table A-3 Changes in version 3.1

Change	Location	Affects
Added a chapter on virtual reality	Chapter 7 Virtual Reality on page 7-184	Version 3.1

Table A-4 Changes in version 3.2

Change	Location	Affects
Updated information and structure of Enlighten chapter	Chapter 5 Global illumination in Unity with Enlighten on page 5-60	Version 3.2 first release

Table A-5 Changes in version 3.3

Change	Location	Affects
Added chapter on Vulkan	Chapter 8 Vulkan on page 8-197	Version 3.3 first release
Added chapter on Mali Graphics Debugger	Chapter 9 The Mali Graphics Debugger on page 9-206	Version 3.3 first release

Table A-6 Changes in version 3.3_01

Change	Location	Affects
Updated the screenshot in section 8.3	8.3 Enabling Vulkan in Unity on page 8-201	Version 3.3 second release