

Chris Dickinson

Unity 2017 Game Optimization

Second Edition

Optimize all aspects of Unity performance



Packt

Unity 2017 Game Optimization

Second Edition

Optimize all aspects of Unity performance

Chris Dickinson

Packt

BIRMINGHAM - MUMBAI

Unity 2017 Game Optimization

Second Edition

Copyright © 2017 Packt Publishing

All rights reserved. No part of this book may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, without the prior written permission of the publisher, except in the case of brief quotations embedded in critical articles or reviews.

Every effort has been made in the preparation of this book to ensure the accuracy of the information presented. However, the information contained in this book is sold without warranty, either express or implied. Neither the author, nor Packt Publishing, and its dealers and distributors will be held liable for any damages caused or alleged to be caused directly or indirectly by this book.

Packt Publishing has endeavored to provide trademark information about all of the companies and products mentioned in this book by the appropriate use of capitals. However, Packt Publishing cannot guarantee the accuracy of this information.

First published: November 2015

Second edition: November 2017

Production reference: 1161117

Published by Packt Publishing Ltd.

Livery Place
35 Livery Street
Birmingham
B3 2PB, UK.

ISBN 978-1-78839-236-5

www.packtpub.com

Credits

Author Chris Dickinson	Copy Editor Dhanya Baburaj
Reviewers Luiz Henrique Bueno Sebastian Koenig	Project Coordinator Ritika Manoj
Commissioning Editor	Proofreader

Kunal Chaudhari

Safis Editing

Acquisition Editor

Shweta Pant

Indexer

Tejal Daruwale Soni

Content Development Editor

Aditi Gour

Graphics

Jason Monteiro

Technical Editor

Shweta Jadhav

Production Coordin

Shantanu Zagade



About the Author

Chris Dickinson grew up in a quiet little corner of England with a strong passion for mathematics, science and, in particular, video games. He loved playing them, dissecting their gameplay, and trying to figure out how they worked. Watching his dad hack the hex code of a PC game to get around the early days of copy protection completely blew his mind! His passion for science won the battle at the time; however, after completing a master's degree in physics with electronics, he flew out to California to work in the field of scientific research in the heart of Silicon Valley. Shortly afterward, he had to admit to himself that research work was an unsuitable career path for his temperament. After firing resumes in all directions, he landed a job that finally set him on the correct course in the field of software engineering (this is not uncommon for physics grads, I hear).

His time working as an automated tools developer for IPBX phone systems fit his temperament much better. Now he was figuring out complex chains of devices, helping its developers fix and improve them, and building tools of his own. Chris learned a lot about how to work with big, complex, real-time, event-based, user-input driven state machines (sounds familiar?). Being mostly self-taught at this point, Chris's passion for video games was flaring up again, pushing him to really figure out how video games were built. Once he felt confident enough, he returned to school for a bachelor's degree in game and simulation programming. By the time he was done, he was already hacking together his own (albeit rudimentary) game engines in C++ and regularly making use of those skills during his day job. However, if you want to build games, you should just build games, and not game engines. So, Chris picked his favorite publically available game engine at the time--an excellent little tool called Unity 3D--and started hammering out some games.

After a brief stint of indie game development, Chris regretfully decided that the demands of that particular career path weren't for him, but the amount of knowledge he had accumulated in just a few short years was impressive by most standards, and he loved to make use of it in ways that enabled other developers with their creations. Since then, Chris has authored a tutorial book

on game physics (*Learning Game Physics with Bullet Physics and OpenGL*, Packt Publishing) and two editions of a Unity performance optimization book (which you are currently reading). He has married the love of his life, Jamie, and works with some of the coolest modern technology as a software development engineer in Test (SDET) at Jaunt Inc. in San Mateo, CA, a Virtual Reality/Augmented Reality startup that focuses on delivering VR and AR experiences, such as 360 videos (and more!).

Outside of work, Chris continues to fight an addiction to board games (particularly *Battlestar: Galactica* and *Blood Rage*), an obsession with Blizzard's *Overwatch* and *Starcraft II*, cater to the ever-growing list of demands from a pair of grumpy yet adorable cats, and gazing forlornly at the latest versions of Unity with a bunch of game ideas floating around on paper. Someday soon, when the time is right (and when he stops slacking off), his plans may come to fruition.

It's been a long road, from my humble beginnings to where I am today. I owe much of it to all of the friends, teachers, tutors, and colleagues I've met along the way. Their instruction, criticism and guidance have made much of what I have accomplished possible. The rest I owe to my family, particularly my wife and best friend Jamie, who have always been nothing but understanding and supportive of my hobbies, passions and aspirations.

About the Reviewers

Luiz Henrique Bueno is a certified ScrumMaster® (CSM) and Unity Certified Developer with over 29 years of experience in software development. He has experimented with the evolution of languages, editors, databases, and frameworks.

In 2002, he wrote the book *Web Applications with Visual Studio .NET, ASP.NET, and C#*, at the launch of Visual Studio .NET. He also participated in the development of a Brazilian magazine called Casa Conectada, about Home Automation.

Based on this magazine's project, he started the development of projects focused on the same subject. He has used technologies such as Crestron, Control4, Marantz, Windows Mobile, and Symbian OS, always implementing touchscreen applications.

Since 2010, he has been developing apps and video games for mobile devices, including VR/AR applications. He has already developed many projects for iPhone, iPad, Apple Watch, Apple TV, and Android using Unity, C#, Xcode, Cocoa Touch, Core Data, SpriteKit, SceneKit, Objective-C, Swift, Git, Photoshop, and Maya.

His motto is "*Do not write code for QA, write code for Production.*" You can reach Luiz Henrique Bueno on his personal website.

Dr. Sebastian Thomas Koenig received his Ph.D. in human interface technology from the University of Canterbury, New Zealand, developing a framework for personalized virtual reality cognitive rehabilitation. He obtained his diploma in psychology from the University of Regensburg, Germany, in the areas of experimental psychology, clinical neuropsychology, and virtual reality rehabilitation.

Dr. Koenig is the founder and CEO of Katana Simulations, where he oversees the design, development, and evaluation of cognitive assessment

and training simulations. His professional experience spans over 10 years of clinical work in cognitive rehabilitation and virtual reality research, development, and human computer interaction. He has been awarded over \$2 million in research funding in the USA, Germany, and Australia as principal investigator and industry partner. He has extensive experience as a speaker at international conferences and as a reviewer of scientific publications in the areas of rehabilitation, cognitive psychology, neuropsychology, software engineering, game development, games user research, and virtual reality.

Dr. Koenig has developed numerous software applications for cognitive assessment and training. For his work on virtual memory tasks, he was awarded the prestigious Laval Virtual Award in 2011 in the Medicine and Health category. Other applications include the Wonderworks Virtual Reality Attention Training in collaboration with the Kessler Foundation, NJ, USA, and the patent-pending Microsoft Kinect-based motor and cognitive training JewelMine/Mystic Isle at the USC Institute for Creative Technologies, CA, USA. Dr. Koenig was awarded the Early Career Investigator Award (2nd place) by the International Society for Virtual Rehabilitation in 2016.

www.PacktPub.com

For support files and downloads related to your book, please visit www.PacktPub.com. Did you know that Packt offers eBook versions of every book published, with PDF and ePub files available? You can upgrade to the eBook version at www.PacktPub.com and as a print book customer, you are entitled to a discount on the eBook copy. Get in touch with us at service@packtpub.com for more details.

At www.PacktPub.com, you can also read a collection of free technical articles, sign up for a range of free newsletters and receive exclusive discounts and offers on Packt books and eBooks.



<https://www.packtpub.com/mapt>

Get the most in-demand software skills with Mapt. Mapt gives you full access to all Packt books and video courses, as well as industry-leading tools to help you plan your personal development and advance your career.

Why subscribe?

- Fully searchable across every book published by Packt
- Copy and paste, print, and bookmark content
- On demand and accessible via a web browser

Customer Feedback

Thanks for purchasing this Packt book. At Packt, quality is at the heart of our editorial process. To help us improve, please leave us an honest review on this book's Amazon page at www.amazon.in/dp/1788392361.

If you'd like to join our team of regular reviewers, you can email us at customerreviews@packtpub.com. We award our regular reviewers with free eBooks and videos in exchange for their valuable feedback. Help us be relentless in improving our products!

Software and Hardware List

Software Specifications:

Chapter Number	Software required (With version)	Free/Proprietary	If proprietary, can code testing be performed using a trial version	If proprietary, then cost of the software	Dow
1-8	Unity 2017.2.0 (although any version of Unity 2017 should work fine)	Free	Yes	N/A	https://unity.com

Detailed installation steps (software-wise)

The following are the steps to install the softwares:

1. Download it
2. Install it

3. Unpack the code files/pull them from git
4. Launch Unity and point it to the folder containing the code files (the folder containing the /Assets folder)

Table of Contents

Preface

- What this book covers
- What you need for this book
- Who this book is for
- Conventions
- Reader feedback
- Customer support
 - Downloading the example code
 - Downloading the color images of this book
- Errata
- Piracy
- Questions

1. Pursuing Performance Problems

- The Unity Profiler
 - Launching the Profiler
 - Editor or standalone instances
 - Connecting to a WebGL instance
 - Remote connection to an iOS device
 - Remote connection to an Android device
 - Editor profiling
 - The Profiler window
 - Profiler controls
 - Add Profiler
 - Record
 - Deep Profile
 - Profile Editor
 - Connected Player
 - Clear
 - Load
 - Save
 - Frame Selection
 - Timeline View
 - Breakdown View Controls
 - Breakdown View
 - The CPU Usage Area
 - The GPU Usage Area
 - The Rendering Area

- The Memory Area
- The Audio Area
- The Physics 3D and Physics 2D Areas
- The Network Messages and Network Operations Areas
- The Video Area
- The UI and UI Details Areas
- The Global Illumination Area
- Best approaches to performance analysis
 - Verifying script presence
 - Verifying script count
 - Verifying the order of events
 - Minimizing ongoing code changes
 - Minimizing internal distractions
 - Minimizing external distractions
 - Targeted profiling of code segments
 - Profiler script control
 - Custom CPU Profiling
- Final thoughts on Profiling and Analysis
 - Understanding the Profiler
 - Reducing noise
 - Focusing on the issue
- Summary

2. Scripting Strategies

- Obtain Components using the fastest method
- Remove empty callback definitions
- Cache Component references
- Share calculation output
- Update, Coroutines, and InvokeRepeating
- Faster GameObject null reference checks
- Avoid retrieving string properties from GameObjects
- Use appropriate data structures
- Avoid re-parenting Transforms at runtime
- Consider caching Transform changes
- Avoid Find() and SendMessage() at runtime
 - Assigning references to pre-existing objects
 - Static Classes
 - Singleton Components
 - A global Messaging System
 - A globally accessible object

- Registration
- Message processing
- Implementing the Messaging System
- Message queuing and processing
- Implementing custom messages
- Message sending
- Message registration
- Message cleanup
- Wrapping up the Messaging System
- Disable unused scripts and objects
 - Disabling objects by visibility
 - Disabling objects by distance
- Consider using distance-squared over distance
- Minimize Deserialization behavior
 - Reduce serialized object size
 - Load serialized objects asynchronously
 - Keep previously loaded serialized objects in memory
 - Move common data into ScriptableObjects
- Load scenes additively and asynchronously
- Create a custom Update() layer
- Summary

3. The Benefits of Batching

- Draw Calls
- Materials and Shaders
- The Frame Debugger
- Dynamic Batching
 - Vertex attributes
 - Mesh scaling
 - Dynamic Batching summary
- Static Batching
 - The Static flag
 - Memory requirements
 - Material references
 - Static Batching caveats
 - Edit Mode debugging of Static Batching
 - Instantiating static meshes at runtime
 - Static Batching summary
- Summary

4. Kickstart Your Art

Audio

- Importing audio files
- Loading audio files
- Encoding formats and quality levels
- Audio performance enhancements**
 - Minimize active Audio Source count
 - Enable Force to Mono for 3D sounds
 - Resample to lower frequencies
 - Consider all compression formats
 - Beware of streaming
 - Apply Filter Effects through Mixer Groups to reduce duplication
 - Use remote content streaming responsibly
 - Consider Audio Module files for background music

Texture files

- Texture compression formats
- Texture performance enhancements**
 - Reduce texture file size
 - Use Mip Maps wisely
 - Manage resolution downscaling externally
 - Adjust Anisotropic Filtering levels
 - Consider Atlasing
 - Adjust compression rates for non-square textures
 - Sparse Textures
 - Procedural Materials
- Asynchronous Texture Uploading

Mesh and animation files

- Reduce polygon count
- Tweak Mesh Compression
- Use Read-Write Enabled appropriately
- Consider baked animations
- Combine meshes

Asset Bundles and Resources

Summary

5. Faster Physics

- Physics Engine internals
- Physics and time
 - Maximum Allowed Timestep
 - Physics updates and runtime changes
- Static Colliders and Dynamic Colliders

- Collision detection
- Collider types
- The Collision Matrix
- Rigidbody active and sleeping states
- Ray and object casting
- Debugging Physics
- Physics performance optimizations
 - Scene setup
 - Scaling
 - Positioning
 - Mass
 - Use Static Colliders appropriately
 - Use Trigger Volumes responsibly
 - Optimize the Collision Matrix
 - Prefer Discrete collision detection
 - Modify the Fixed Update frequency
 - Adjust the Maximum Allowed Timestep
 - Minimize Raycasting and bounding-volume checks
 - Avoid complex Mesh Colliders
 - Use simpler primitives
 - Use simpler Mesh Colliders
 - Avoid complex physics Components
 - Let physics objects sleep
 - Modify the Solver Iteration Count
 - Optimize Ragdolls
 - Reduce Joints and Colliders
 - Avoid inter-Ragdoll collisions
 - Replace, deactivate or remove inactive Ragdolls
 - Know when to use physics
- Summary

6. Dynamic Graphics

- The Rendering Pipeline
 - The GPU Front End
 - The GPU Back End
 - Fill Rate
 - Overdraw
 - Memory Bandwidth
 - Lighting and Shadowing
 - Forward Rendering
 - Deferred Rendering

- Vertex Lit Shading (legacy)
- Global Illumination
- Multithreaded Rendering
- Low-level rendering APIs
- Detecting performance issues
 - Profiling rendering issues
 - Brute-force testing
- Rendering performance enhancements
 - Enable/Disable GPU Skinning
 - Reduce geometric complexity
 - Reduce Tessellation
 - Employ GPU Instancing
 - Use mesh-based Level Of Detail (LOD)
 - Culling Groups
 - Make use of Occlusion Culling
- Optimizing Particle Systems
 - Make use of Particle System Culling
 - Avoid recursive Particle System calls
- optimizing Unity UI
 - Use more Canvases
 - Separate objects between static and dynamic canvases
 - Disable Raycast Target for noninteractive elements
 - Hide UI elements by disabling the parent Canvas Component
 - Avoid Animator Components
 - Explicitly define the Event Camera for World Space Canvases
 - Don't use alpha to hide UI elements
- Optimizing ScrollRects
 - Make sure to use a RectMask2D
 - Disable Pixel Perfect for ScrollRects
 - Manually stop ScrollRect motion
- Use empty UIText elements for full-screen interaction
- Check the Unity UI source code
- Check the documentation

- Shader optimization
- Consider using Shaders intended for mobile platforms
- Use small data types
- Avoid changing precision while swizzling
- Use GPU-optimized helper functions
- Disable unnecessary features
- Remove unnecessary input data

- Expose only necessary variables
- Reduce mathematical complexity
- Reduce texture sampling
- Avoid conditional statements
- Reduce data dependencies
- Surface Shaders
- Use Shader-based LOD
- Use less texture data
- Test different GPU Texture Compression formats
- Minimize texture swapping
- VRAM limits
 - Preload textures with hidden GameObjects
 - Avoid texture thrashing
- Lighting optimization
 - Use real-time Shadows responsibly
 - Use Culling Masks
 - Use baked Lightmaps
- Optimizing rendering performance for mobile devices
 - Avoid Alpha Testing
 - Minimize Draw Calls
 - Minimize Material count
 - Minimize texture size
 - Make textures square and power-of-two
 - Use the lowest possible precision formats in Shaders

Summary

7. Virtual Velocity and Augmented Acceleration

XR Development

- Emulation

- User comfort

Performance enhancements

- The kitchen sink

- Single-Pass versus Multi-Pass Stereo Rendering

- Apply anti-aliasing

- Prefer Forward Rendering

- Image effects in VR

- Backface culling

- Spatialized audio

- Avoid camera physics collisions

- Avoid Euler angles

Exercise restraint
Keep up to date with the latest developments

Summary

8. Masterful Memory Management

- The Mono platform
- Memory Domains
 - Garbage collection
 - Memory Fragmentation
 - Garbage collection at runtime
 - Threaded garbage collection
- Code compilation
 - IL2CPP
- Profiling memory
 - Profiling memory consumption
 - Profiling memory efficiency
- Memory management performance enhancements
 - Garbage collection tactics
 - Manual JIT compilation
 - Value types and Reference types
 - Pass by value and by reference
 - Structs are Value types
 - Arrays are Reference types
 - Strings are immutable Reference types
 - String concatenation
 - StringBuilder
 - String formatting
 - Boxing
 - The importance of data layout
 - Arrays from the Unity API
 - Using InstanceIDs for dictionary keys
 - foreach loops
 - Coroutines
 - Closures
 - The .NET library functions
 - Temporary work buffers
 - Object Pooling
 - Prefab Pooling
 - Poolable Components
 - The Prefab Pooling System
 - Prefab pools

- Object spawning
- Instance prespawning
- Object despawning
- Prefab pool testing
- Prefab Pooling and Scene loading
- Prefab Pooling summary
- IL2CPP optimizations
- WebGL optimizations
- The future of Unity, Mono, and IL2CPP
- The upcoming C# Job System
- Summary

9. Tactical Tips and Tricks

- Editor hotkey tips
 - GameObjects
 - Scene window
 - Arrays
 - Interface
 - In-editor documentation
- Editor UI tips
 - Script Execution Order
 - Editor files
 - The Inspector window
 - The Project window
 - The Hierarchy window
 - The Scene and Game windows
 - Play Mode
- Scripting tips
 - General
 - Attributes
 - Variable attributes
 - Class attributes
 - Logging
 - Useful links
- Custom Editor scripts and menu tips
- External tips
 - Other tips
- Summary

Preface

User experience is a critical component of any game, and this includes not only our game's story and its Gameplay, but also how smoothly the graphics run, how reliably it connects to multiplayer servers, how responsive it is to user input, and even how large the final application file size is due to the prevalence of mobile devices and cloud downloads. The barrier of entry into game development has been lowered considerably thanks to tools such as Unity that offers an enormous array of useful development features while still being accessible to individual developers. However, due to the amount of competition in the gaming industry, the quality level of the final product that our players expect us to provide is increasing with every passing day. We should expect that every facet of our game can and will be scrutinized by players and critics alike.

The goals of performance optimization are deeply entwined with user experience. Poorly optimized games can result in low frame rates, freezes, crashes, input lag, long loading times, inconsistent and jittery runtime behavior, Physics Engine breakdowns, and even excessively high battery power consumption (an often-neglected metric for mobile devices). Having just one of these issues can be a game developer's worst nightmare as reviews will tend to focus on the one thing that we did badly, in spite of all the things that we did well.

One goal of performance optimization is to make the best use of the available resources, which includes CPU resources such as the number of cycles consumed, how much main memory space we're using (known as RAM) as well as Graphics Processing Unit (GPU) resources, which includes its own memory space (known as VRAM), Fill Rate, Memory Bandwidth, and so on. However, the most important goal of performance optimization is to ensure that no single resource causes a bottleneck at an inappropriate time, and that the highest priority tasks get taken care of first. Even small, intermittent hiccups and sluggishness in performance can pull the player out of the experience, breaking immersion and limiting our potential to create the experience we intended. Another consideration is that the more resources we

can save, the more activity we can afford to implement in our games, allowing us to generate more interesting and dynamic gameplay.

It is also important to decide when to take a step back and stop making performance enhancements. In a world with infinite time and resources, there will always be another way to make it better, faster, and more efficient. There must be a point during development where we decide that the product has reached an acceptable level of quality. If not, we risk dooming ourselves to repeatedly implementing changes that result in little or no tangible benefit, while each change also risks the chance that we introduce more bugs.

The best way to decide whether a performance issue is worth fixing is to answer the question "will the user notice it?". If the answer to this questions is "no," then performance optimization will be a wasted effort. There is an old saying in software development:

Premature optimization is the root of all evil.

Premature optimization is the cardinal sin of reworking and refactoring code to enhance performance without any proof that it is necessary. This can mean either making changes without showing that a performance problem even exists, or making changes because we only believe a performance issue might stem from a particular area before it has been proven to be true.

Of course, the original version of this common saying by Donald Knuth goes on to say that we should still write our code to avoid the more straightforward and obvious performance problems. However, the real performance optimization work toward the end of a project can take a lot of time, and we should both plan the time to polish the product properly, while avoiding the desire to implement the more costly and time-consuming changes without verifiable proof. These kinds of mistake have cost software developers, as a collective whole, a depressing number of work hours for nothing.

This book intends to give you the tools, knowledge, and skills you need to both detect and fix performance issues in a Unity application, no matter where they stem from. These bottlenecks can appear within hardware components such as the CPU, GPU, and RAM, or within software subsystems such as Physics, Rendering, and the Unity Engine itself.

Optimizing the performance of our games will give them a much better chance of succeeding and standing out from the crowd in a marketplace that is inundated with new, high-quality games every single day.

What this book covers

[Chapter 1](#), *Pursuing Performance Problems*, provides an exploration of the Unity Profiler and a series of methods to profile our application, detect performance bottlenecks, and perform root cause analysis.

[Chapter 2](#), *Scripting Strategies*, deals with the best practices for our Unity C# Script code, minimizing MonoBehaviour callback overhead, improving inter-object communication, and more.

[Chapter 3](#), *The Benefits of Batching*, explores Unity's Dynamic Batching and Static Batching systems, and how they can be utilized to ease the burden on the Rendering Pipeline.

[Chapter 4](#), *Kickstart Your Art*, helps you understand the underlying technology behind art assets and learn how to avoid common pitfalls with importing, compression, and encoding.

[Chapter 5](#), *Faster Physics*, is about investigating the nuances of Unity's internal Physics Engines for both 3D and 2D games, and how to properly organize our physics objects for improved performance.

[Chapter 6](#), *Dynamic Graphics*, provides an in-depth exploration of the Rendering Pipeline, and how to improve applications that suffer rendering bottlenecks in the GPU, or CPU, how to optimize graphical effects such as lighting, shadows, and Particle Effects, ways in which to optimize Shader code, and some specific techniques for mobile devices.

[Chapter 7](#), *Virtual Velocity and Augmented Acceleration*, focuses on the new entertainment mediums of Virtual Reality (VR) and Augmented Reality (AR), and includes several techniques for optimizing performance that are unique to apps built for these platforms.

[Chapter 8](#), *Masterful Memory Management*, examines the inner workings of the Unity Engine, the Mono Framework, and how memory is managed within these components to protect our application from excessive heap allocations

and runtime garbage collection.

[Chapter 9](#), *Tactical Tips and Tricks*, closes the book with a multitude of useful techniques used by Unity professionals to improve project workflow and scene management.

What you need for this book

The majority of this book will focus on features and enhancements that apply to Unity 2017. Many of the techniques explored within this book can be applied to Unity 5.x projects and older, but feature lists may appear different. These differences will be highlighted where applicable.

Who this book is for

This book is intended for intermediate and advanced Unity developers who have experience with most of Unity's feature set, and those who want to maximize the performance of their game or solve particular bottlenecks. Whether the bottleneck is caused by continuous CPU overload, runtime CPU spiking, slow memory access, memory fragmentation, garbage collection, poor GPU Fill Rate, or Memory Bandwidth, this book will teach you the techniques you need to identify the source of the problem and help explore multiple ways of reducing their impact on your application.

Familiarity with the C# language will be needed for sections involving scripting and memory usage, and a basic understanding of Cg will be needed for areas involving Shader optimization.

Conventions

In this book, you will find a number of text styles that distinguish between different kinds of information. Here are some examples of these styles and an explanation of their meaning.

Code words in text, database table names, folder names, filenames, file extensions, pathnames, dummy URLs, user input, and Twitter handles are shown as follows: "Finally, we will need to implement the `GameLogic` class."

A block of code is set as follows:

```
void Start() {
    GameLogic.Instance.RegisterUpdateableObject(this);
    Initialize();
}

protected virtual void Initialize() {
    // derived classes should override this method for initialization code, and
}
```

New terms and important words are shown in bold. Words that you see on the screen, for example, in menus or dialog boxes, appear in the text like this: "Select System info from the Administration panel."



Warnings or important notes appear in a box like this.



Tips and tricks appear like this.

Reader feedback

Feedback from our readers is always welcome. Let us know what you think about this book—what you liked or disliked. Reader feedback is important for us as it helps us develop titles that you will really get the most out of.

To send us general feedback, simply e-mail feedback@packtpub.com, and mention the book's title in the subject of your message.

If there is a topic that you have expertise in and you are interested in either writing or contributing to a book, see our author guide at www.packtpub.com/authors.

Customer support

Now that you are the proud owner of a Packt book, we have a number of things to help you to get the most from your purchase.

Downloading the example code

You can download the example code files for this book from your account at <http://www.packtpub.com>. If you purchased this book elsewhere, you can visit <http://www.packtpub.com/support> and register to have the files emailed directly to you.

You can download the code files by following these steps:

1. Log in or register to our website using your email address and password.
2. Hover the mouse pointer on the SUPPORT tab at the top.
3. Click on Code Downloads & Errata.
4. Enter the name of the book in the Search box.
5. Select the book for which you're looking to download the code files.
6. Choose from the drop-down menu where you purchased this book from.
7. Click on Code Download.

Once the file is downloaded, please make sure that you unzip or extract the folder using the latest version of:

- WinRAR / 7-Zip for Windows
- Zipeg / iZip / UnRarX for Mac
- 7-Zip / PeaZip for Linux

The code bundle for the book is also hosted on GitHub at <https://github.com/PacktPublishing/Unity-2017-Game-Optimization-Second-Edition>. We also have other code bundles from our rich catalog of books and videos available at <https://github.com/PacktPublishing/>. Check them out!

Downloading the color images of this book

We also provide you with a PDF file that has color images of the screenshots/diagrams used in this book. The color images will help you better understand the changes in the output. You can download this file from https://www.packtpub.com/sites/default/files/downloads/Unity2017GameOptimizationSecondEdition_ColorImages.pdf.

Errata

Although we have taken every care to ensure the accuracy of our content, mistakes do happen. If you find a mistake in one of our books-maybe a mistake in the text or the code-we would be grateful if you could report this to us. By doing so, you can save other readers from frustration and help us improve subsequent versions of this book. If you find any errata, please report them by visiting <http://www.packtpub.com/submit-errata>, selecting your book, clicking on the Errata Submission Form link, and entering the details of your errata. Once your errata are verified, your submission will be accepted and the errata will be uploaded to our website or added to any list of existing errata under the Errata section of that title.

To view the previously submitted errata, go to <https://www.packtpub.com/books/content/support> and enter the name of the book in the search field. The required information will appear under the Errata section.

Piracy

Piracy of copyrighted material on the Internet is an ongoing problem across all media. At Packt, we take the protection of our copyright and licenses very seriously. If you come across any illegal copies of our works in any form on the Internet, please provide us with the location address or website name immediately so that we can pursue a remedy.

Please contact us at copyright@packtpub.com with a link to the suspected pirated material.

We appreciate your help in protecting our authors and our ability to bring you valuable content.

Questions

If you have a problem with any aspect of this book, you can contact us at questions@packtpub.com, and we will do our best to address the problem.

Pursuing Performance Problems

Performance evaluation for most software products is a very scientific process. First, we determine the maximum/minimum supported performance metrics, such as the allowed memory usage, acceptable CPU consumption, the number of concurrent users, and so on. Next, we perform load testing against the application in scenarios with a version of the application built for the target platform, and test it while gathering instrumentation data. Once this data is collected, we analyze and search it for performance bottlenecks. If problems are discovered, we complete a root-cause analysis, make changes in the configuration or application code to fix the issue and repeat it.

Although game development is a very artistic process, it is still exceptionally technical, so there is a good reason to treat it in similarly objective ways. Our game should have a target audience in mind, which can tell us what hardware limitations our game might be operating under and, perhaps, tell us exactly what performance targets we need to meet (particularly in the case of console and mobile games). We can perform runtime testing on our application, gather performance data from multiple subsystems (CPU, GPU, memory, the Physics Engine, the Rendering Pipeline, and so on), and compare them against what we consider to be acceptable. We can use this data to identify bottlenecks in our application, perform additional instrumentation, and determine the root cause of the issue. Finally, depending on the type of problem, we should be capable of applying a number of fixes to improve our application's performance to bring it more in line with the intended behavior.

However, before we spend even a single moment making performance fixes, we will need to prove that a performance problem exists to begin with. It is unwise to spend time rewriting and refactoring code until there is good reason to do so, since pre-optimization is rarely worth the hassle. Once we have proof of a performance issue, the next task is figuring out exactly where the bottleneck is located. It is important to ensure that we understand why the performance issue is happening, otherwise we could waste even more time applying fixes that are little more than educated guesses. Doing so often means that we only fix a symptom of the issue, not its root cause, and so we

risk the chance that it manifests itself in other ways in the future, or in ways we haven't yet detected.

In this chapter, we will explore the following:

- How to gather profiling data using the Unity Profiler
- How to analyze Profiler data for performance bottlenecks
- Techniques to isolate a performance problem and determine its root cause

With a thorough understanding of a given problem, you will then be ready for information presented in the remaining chapters, where you will learn what solutions are available for the issue we've detected.

The Unity Profiler

The Unity Profiler is built into the Unity Editor itself and provides an expedient way of narrowing down our search for performance bottlenecks by generating usage and statistics reports on a multitude of Unity3D subsystems during runtime. The different subsystems it can gather data for are listed as follows:

- CPU consumption (per-major subsystem)
- Basic and detailed rendering and GPU information
- Runtime memory allocations and overall consumption
- Audio source/data usage
- Physics Engine (2D and 3D) usage
- Network messaging and operation usage
- Video playback usage
- Basic and detailed user interface performance (new in Unity 2017)
- Global Illumination statistics (new in Unity 2017)

There are generally two approaches to make use of a profiling tool: **instrumentation** and **benchmarking** (although, admittedly, the two terms are often used interchangeably).

Instrumentation typically means taking a close look into the inner workings of the application by observing the behavior of targeted function calls, where/how much memory is being allocated, and, generally getting an accurate picture of what is happening with the hope of finding the root cause of a problem. However, this is normally not an efficient way of starting to find performance problems because profiling of any application comes with a performance cost of its own.

When a Unity application is compiled in Development Mode (determined by the Development Build flag in the Build Settings menu), additional compiler flags are enabled causing the application to generate special events at runtime, which get logged and stored by the Profiler. Naturally, this will cause additional CPU and memory overhead at runtime due to all of the extra

workload the application takes on. Even worse, if the application is being profiled through the Unity Editor, then even more CPU and memory will be spent, ensuring that the Editor updates its interface, renders additional windows (such as the Scene window), and handles background tasks. This profiling cost is not always negligible. In excessively large projects, it can sometimes cause wildly inconsistent behavior when the Profiler is enabled. In some cases, the inconsistency is significant enough to cause completely unexpected behavior due to changes in event timings and potential race conditions in asynchronous behavior. This is a necessary price we pay for a deep analysis of our code's behavior at runtime, and we should always be aware of its presence.

Before we get ahead of ourselves and start analyzing every line of code in our application, it would be wiser to perform a surface-level measurement of the application. We should gather some rudimentary data and perform test scenarios during a runtime session of our game while it runs on the target hardware; the test case could simply be a few seconds of Gameplay, playback of a cut scene, a partial play through of a level, and so on. The idea of this activity is to get a general feel for what the user might experience and keep watching for moments when performance becomes noticeably worse. Such problems may be severe enough to warrant further analysis.

This activity is commonly known as benchmarking, and the important metrics we're interested in are often the number of **frames per-second (FPS)** being rendered, overall memory consumption, how CPU activity behaves (looking for large spikes in activity), and sometimes CPU/GPU temperature. These are all relatively simple metrics to collect and can be used as a best first approach to performance analysis for one important reason; it will save us an enormous amount of time in the long run, since it ensures that we only spend our time investigating problems that users would notice.

We should dig deeper into instrumentation only after a benchmarking test indicates that further analysis is required. It is also very important to benchmark by simulating actual platform behavior as much as possible if we want a realistic data sample. As such, we should never accept benchmarking data that was generated through Editor Mode as representative of real gameplay, since Editor Mode comes with some additional overhead costs that

might mislead us, or hide potential race conditions in a real application. Instead, we should hook the profiling tool into the application while it is running in a standalone format on the target hardware.



Many Unity developers are surprised to find that the Editor sometimes calculates the results of operations much faster than a standalone application does. This is particularly common when dealing with serialized data like audio files, Prefabs and Scriptable Objects. This is because the Editor will cache previously imported data and is able to access it much faster than a real application would.

Let's cover how to access the Unity Profiler and connect it to the target device so that we can start to make accurate benchmarking tests.



Users who are already familiar with connecting the Unity Profiler to their applications can skip to the section titled The Profiler window.

Launching the Profiler

We will begin with a brief tutorial on how to connect our game to the Unity Profiler within a variety of contexts:

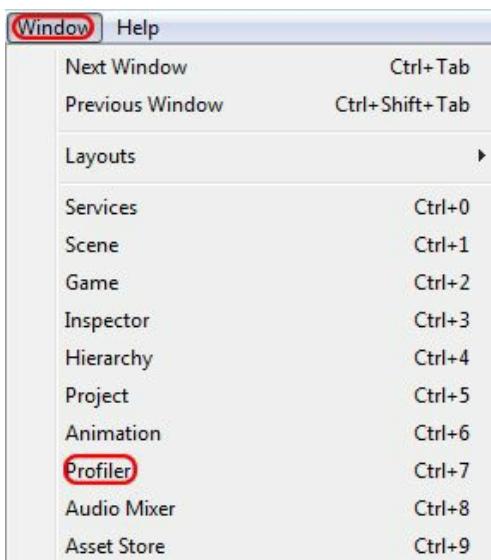
- Local instances of the application, either through the Editor or a standalone instance
- Local instances of a WebGL application running in a browser
- Remote instances of the application on an iOS device (for example, iPhone or iPad)
- Remote instances of the application on an Android device (for example, an Android tablet or phone)
- Profiling the Editor itself

We will briefly cover the requirements for setting up the Profiler in each of these contexts.

Editor or standalone instances

The only way to access the Profiler is to launch it through the Unity Editor and connect it to a running instance of our application. This is the case whether we're executing our game in Play Mode within the Editor, running a standalone application on the local or remote device, or we wish to profile the Editor itself.

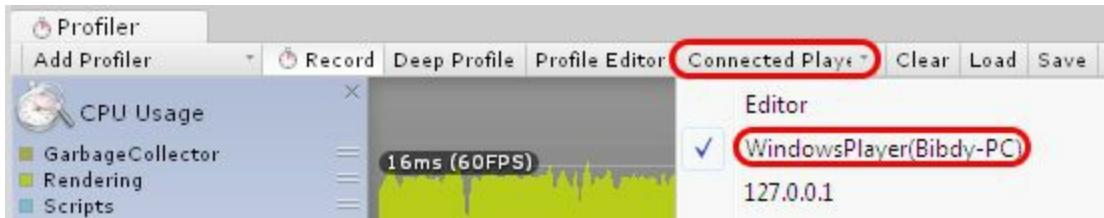
To open the Profiler, navigate to Window | Profiler within the Editor:



If the Editor is already running in Play Mode, then we should see reporting data actively gathering in the Profiler window.

To profile standalone projects, ensure that the Development Build and Autoconnect Profiler flags are enabled when the application is built.

Choosing whether to profile an Editor-based instance (through the Editor's Play Mode) or a standalone instance (built and running separately from the Editor) can be achieved through the **Connected Player** option in the Profiler window:



Note that switching back to the Unity Editor while profiling a separate standalone project will halt all data collection since the application will not be updated while it is in the background.



Note that the Development Build option is named Use Development Mode and the Connected Player option is named Active Profiler in Unity 5.

Connecting to a WebGL instance

The Profiler can also be connected to an instance of the Unity WebGL Player. This can be achieved by ensuring that the Development Build and Autoconnect Profiler flags are enabled when the WebGL application is built and run from the Editor. The application will then be launched through the Operating System's default browser. This enables us to profile our web-based application in a more real-world scenario through the target browser and test multiple browser types for inconsistencies in behavior (although this requires us to keep changing the default browser).

Unfortunately, the Profiler connection can only be established when the application is first launched from the Editor. It currently (at least in early builds of Unity 2017) cannot be connected to a standalone WebGL instance already running in a browser. This limits the accuracy of benchmarking WebGL applications since there will be some Editor-based overhead, but it's the only option we have available for the moment.

Remote connection to an iOS device

The Profiler can also be connected to an active instance of the application running remotely on an iOS device, such as an iPad or iPhone. This can be achieved through a shared Wi-Fi connection.



Note that remote connection to an iOS device is only possible when Unity (and hence the Profiler) is running on an Apple Mac device.

Follow the given steps to connect the Profiler to an iOS device:

1. Ensure that the Development Build and Autoconnect Profiler flags are enabled when the application is built.
2. Connect both the iOS device and Mac device to a local Wi-Fi network, or to an ad hoc Wi-Fi network.
3. Attach the iOS device to the Mac via the USB or Lightning cable.
4. Begin building the application with the Build & Run option as usual.
5. Open the Profiler window in the Unity Editor and select the device under Connected Player.

You should now see the iOS device's profiling data gathering in the Profiler window.



The Profiler uses ports from 54998 to 55511 to broadcast profiling data. Ensure that these ports are available for outbound traffic if there is a firewall on the network.

For troubleshooting problems with building iOS applications and connecting the Profiler to them, consult the following documentation page: <https://docs.unity3d.com/Manual/TroubleShootingiPhone.html>.

Remote connection to an Android device

There are two different methods for connecting an Android device to the Unity Profiler: either through a Wi-Fi connection or using the **Android Debug Bridge (ADB)** tool. Either of these approaches will work from an Apple Mac, or a Windows PC.

Perform the following steps to connect an Android device over a Wi-Fi connection:

1. Ensure that the Development Build and Autoconnect Profiler flags are enabled when the application is built.
2. Connect both the Android and desktop devices to a local Wi-Fi network.
3. Attach the Android device to the desktop device via the USB cable.
4. Begin building the application with the Build & Run option as usual.
5. Open the Profiler window in the Unity Editor and select the device under Connected Player.

The application should then be built and pushed to the Android device through the USB connection, and the Profiler should connect through the Wi-Fi connection. You should then see the Android device's profiling data gathering in the Profiler window.

The second option is to use ADB. This is a suite of debugging tools that comes bundled with the Android **Software Development Kit (SDK)**. For ADB profiling, follow these steps:

1. Ensure that the Android SDK is installed by following Unity's guide for Android SDK/NDK setup: <https://docs.unity3d.com/Manual/android-sdksetup.html>.
2. Connect the Android device to your desktop machine via the USB cable.
3. Ensure that the Development Build and Autoconnect Profiler flags are enabled when the application is built.

4. Begin building the application with the Build & Run option as usual.
5. Open the Profiler window in the Unity Editor and select the device under Connected Player.

You should now see the Android device's profiling data gathering in the Profiler window.

For troubleshooting problems with building Android applications and connecting the Profiler to them, consult the following documentation page: <https://docs.unity3d.com/Manual/TroubleShootingAndroid.html>.

Editor profiling

We can profile the Editor itself. This is normally used when trying to profile the performance of custom Editor Scripts. This can be achieved by enabling the Profile Editor option in the Profiler window and configuring the Connected Player option to Editor, as shown in the following screenshot:



Note that both options must be configured this way if we want to profile the Editor. Setting Connected Player to Editor without enabling the Profile Editor button is the default case, where the Profiler is collecting data for our application while it is running in Play Mode.

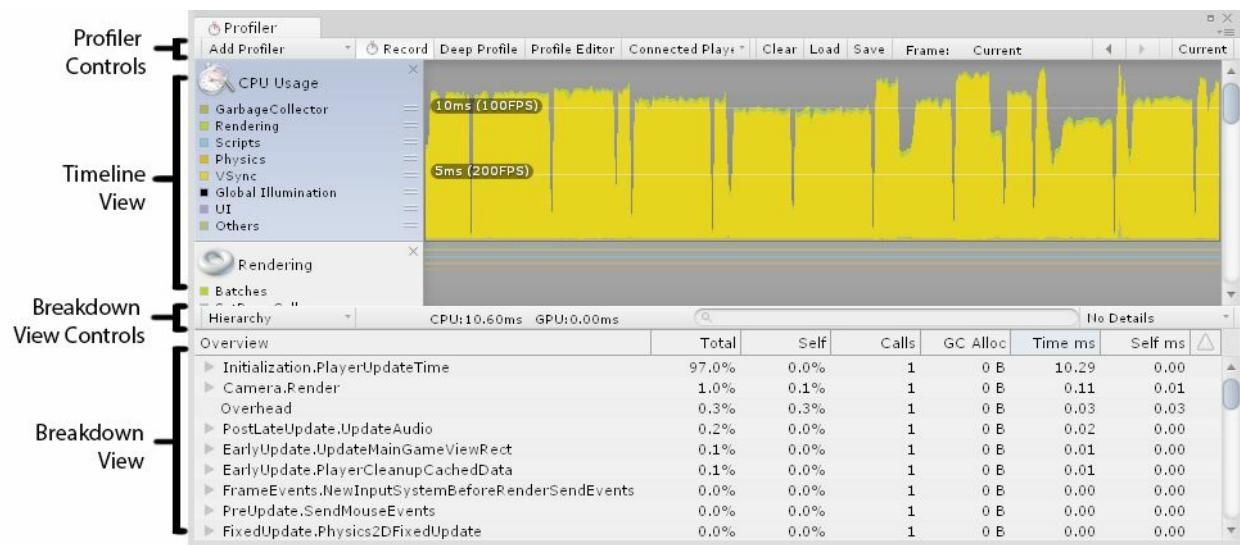
The Profiler window

We will now cover the essential features of the Profiler as they can be found within the interface.

The Profiler window is split into four main sections:

- **Profiler Controls**
- **Timeline View**
- **Breakdown View Controls**
- **Breakdown View**

These sections are shown in the following screenshot:



We'll cover each of these sections in detail.

Profiler controls

The top bar in the previous screenshot contains multiple drop-down and toggle buttons we can use to affect what is being profiled and how deeply in the subsystem that data is gathered from. They are covered in the next subsections.

Add Profiler

By default, the Profiler will collect data for several different subsystems that cover the majority of the Unity's Engine subsystems in the Timeline View. These subsystems are organized into various Areas containing relevant data. The Add Profiler option can be used to add additional Areas or restore them if they were removed. Refer to the *Timeline View* section for a complete list of subsystems we can profile.

Record

Enabling the Record option makes the Profiler record profiling data. This will happen continuously while this option is enabled. Note that runtime data can only be recorded if the application is actively running. For an app running in the Editor, this means that Play Mode must be enabled and it should not be paused; alternatively, for a standalone app, it must be the active window. If Profile Editor is enabled, then the data that appears will be gathered for the Editor itself.

Deep Profile

Ordinary profiling will only record the time and memory allocations made by the common Unity callback methods, such as `Awake()`, `Start()`, `Update()`, and `FixedUpdate()`. Enabling the Deep Profile option re-compiles our scripts with much deeper level of instrumentation, allowing it to measure each and every invoked method. This causes a significantly greater instrumentation cost during runtime than normal, and uses substantially more memory since data is being collected for the entire callstack at runtime. As a consequence, Deep Profiling may not even be possible in large projects, as Unity may run out of memory before testing even begins or the application may run so slowly as to make the test pointless.



Note that toggling Deep Profile requires the entire project to be completely re-compiled before profiling can begin again, so it is best to avoid toggling the option back and forth between tests.

Since this option blindly measures the entire callstack, it would be unwise to keep it enabled during most of our profiling tests. This option is best reserved for when default profiling is not providing enough detail to figure out the root cause, or if we're testing performance of a small test Scene, which we're using to isolate certain activities.

If Deep Profiling is required for larger projects and scenes, but the Deep Profile option is too much of a hindrance during runtime, then there are alternative approaches that can be used to perform more detailed profiling in the upcoming section titled *Targeted profiling of code segments*.

Profile Editor

The Profile Editor option enables Editor profiling, that is, gathering profiling data for the Unity Editor itself. This is useful in order to profile any custom Editor scripts we have developed.



Remember that Connected Player must also be set to the Editor option for Editor profiling to occur.

Connected Player

The Connected Player drop-down offers choices to select the target instance of Unity we want to profile. This can be the current Editor application, a local standalone instance of our application, or an instance of our application running on a remote device.

Clear

The Clear button clears all profiling data from the Timeline View.

Load

The Load button will open up a dialog window to load in any previously-saved Profiler data (from using the Save option).

Save

The Save button saves any Profiler data currently presented in the Timeline View to a file. Only 300 frames of data can be saved in this fashion at a time, and a new file must be manually created for any more data. This is typically sufficient for most situations, since when a performance spike occurs we then have about five to ten seconds to pause the application and save the data for future analysis (such as attaching it to a bug report) before it gets pushed off the left side of the Timeline View. Any saved Profiler data can be loaded into the Profiler for future examination using the Load option.

Frame Selection

The Frame Counter shows how many frames have been profiled and which frame is currently selected in the Timeline View. There are two buttons to move the currently selected frame forward or backward by one frame and a third button (the Current button) that resets the selected frame to the most recent frame and keeps that position. This will cause the Breakdown View to always show the profiling data for the current frame during runtime profiling and will display the word Current.

Timeline View

The Timeline View reveals profiling data that has been collected during runtime, organized into a series of Areas. Each Area focuses on profiling data for a different subsystem of the Unity Engine and each is split into two sections: a graphical representation of profiling data on the right, and a series of checkboxes to enable/disable different activities/data types on the left. These colored boxes can be toggled, which changes the visibility of the corresponding data types within the graphical section of the Timeline View.

When an Area is selected in the Timeline View, more detailed information for that subsystem will be revealed in the Breakdown View (beneath the Timeline View) for the currently selected frame. The kinds of information displayed in the Breakdown View varies depending on which Area is currently selected in the Timeline View.

Areas can be removed from the Timeline View by clicking on the X at the top-right corner of an Area. Recall that Areas can be restored to the Timeline View through the Add Profiler option in the Controls bar.

At any time, we can click at a location in the graphical part of the Timeline View to reveal information about a given frame. A large vertical white bar will appear (usually with some additional information on either side coinciding with the line graphs), showing us which frame is selected.

Depending on which Area is currently selected (determined by which Area is currently highlighted in blue), different information will be available in the Breakdown View, and different options will be available in the Breakdown View Controls. Changing the Area that is selected is as simple as clicking on the relevant box on the left-hand side of the Timeline View or on the graphical side, although clicking inside the graphical Area might also change which frame has been selected, so be careful clicking in the graphical Area if you wish to see Breakdown View information for the same frame.

Breakdown View Controls

Different drop-downs and toggle button options will appear within the Breakdown View Controls, depending on which Area is currently selected in the Timeline View. Different Areas offer different controls, and these options dictate what information is available, and how that information is presented in the Breakdown View.

Breakdown View

The information revealed in the Breakdown View will vary enormously based on which Area is currently selected and which Breakdown View Controls options are selected. For instance, some Areas offer different modes in a drop-down within the Breakdown View Controls, which can provide a simpler or detailed view of the information or even a graphical layout of the same information so that it can be parsed more easily.

Let's cover each Area and the different kinds of information and options available in the Breakdown View.

The CPU Usage Area

This Area shows data for all CPU usage and statistics. This Area is perhaps the most complex and useful since it covers a large number of Unity subsystems, such as `MonoBehaviour` Components, cameras, some rendering and physics processes, user interface (including the Editor's interface, if we're running through the Editor), audio processing, the Profiler itself, and more.

There are three different modes of displaying CPU usage data in the Breakdown View:

- Hierarchy mode
- Raw Hierarchy mode
- Timeline mode

Hierarchy mode reveals most callstack invocations, while grouping similar data elements and global Unity function calls together for convenience. For instance, rendering delimiters, such as `BeginGUI()` and `EndGUI()` calls, are combined together in this mode. Hierarchy mode is helpful as an initial first step to determine which function calls cost the most CPU time to execute.

Raw Hierarchy mode is similar to Hierarchy mode, except it will separate global Unity function calls into separate entries rather than being combined into one bulk entry. This will tend to make the Breakdown View more difficult to read, but may be helpful if we're trying to count how many times a particular global method is invoked or determining whether one of these calls is costing more CPU/memory than expected. For example, each `BeginGUI()` and `EndGUI()` calls will be separated into different entries, making it more clear how many times each is being called compared to the Hierarchy mode.

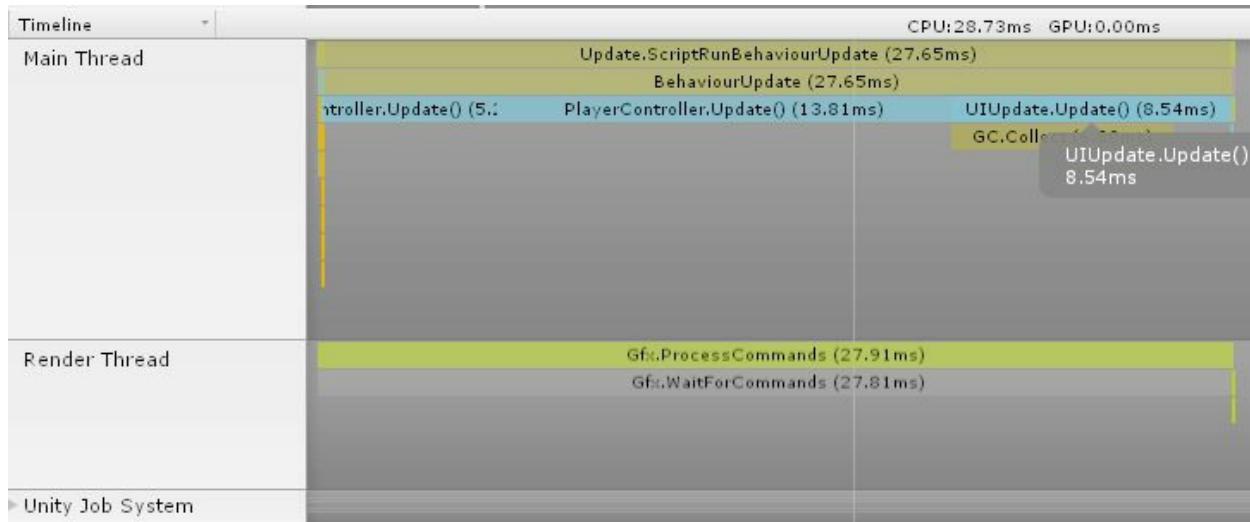
Perhaps, the most useful mode for the CPU Usage Area is the Timeline mode option (not to be confused with the main Timeline View). This mode organizes CPU usage during the current frame by how the call stack expanded and contracted during processing.

Timeline mode organizes the Breakdown View vertically into different sections that represent different threads at runtime, such as Main Thread, Render Thread, and various background job threads called Unity Job System, used for loading activity such as scenes and other assets. The horizontal axis represents time, so wider blocks are consuming more CPU time than narrower blocks. The horizontal size also represents relative time, making it easy to compare how much time one function call took compared to another. The vertical axis represents the callstack, so deeper chains represent more calls in the callstack at that time.

Under Timeline mode, blocks at the top of the Breakdown View are functions (or technically, callbacks) called by the Unity Engine at runtime (such as `Start()`, `Awake()`, or `Update()`), whereas blocks underneath them are functions that those functions had called into, which can include functions on other Components or regular C# objects.

The Timeline mode offers a very clean and organized way to determine which particular method in the callstack consumes the most time and how that processing time measures up against other methods being called during the same frame. This allows us to gauge the method that is the biggest cause of performance problems with minimal effort.

For example, let's assume that we are looking at a performance problem in the following screenshot. We can tell, with a quick glance, that there are three methods that are causing a problem, and they each consume similar amounts of processing time, due to their similar widths:



In the previous screenshot, we have exceeded our 16.667 millisecond budget with calls to three different `MonoBehaviour` Components. The good news is that we have three possible methods through which we can find performance improvements, which means lots of opportunities to find code that can be improved. The bad news is that increasing the performance of one method will only improve about one-third of the total processing for that frame. Hence, all three methods may need to be examined and optimized in order get back under budget.



It's a good idea to collapse the Unity Job System list when using Timeline mode, as it tends to obstruct the visibility of items shown in the Main Thread block, which is probably what we're most interested in.

In general, the CPU Usage Area will be most useful for detecting issues that can be solved by solutions that will be explored in [chapter 2, Scripting Strategies](#).

The GPU Usage Area

The GPU Usage Area is similar to the CPU Usage Area, except that it shows method calls and processing time as it occurs on the GPU. Relevant Unity method calls in this Area will relate to cameras, drawing, opaque and transparent geometry, lighting and shadows, and so on.

The GPU Usage Area offers hierarchical information similar to the CPU Usage Area and estimates time spent calling into various rendering functions such as `Camera.Render()` (provided rendering actually occurs during the frame currently selected in the Timeline View).

The GPU Usage Area will be a useful tool to refer to when you go through [Chapter 6, Dynamic Graphics](#).

The Rendering Area

The Rendering Area provides some generic rendering statistics that tend to focus on activities related to preparing the GPU for rendering, which is a set of activities that occur on the CPU (as opposed to the act of rendering, which is activity handled within the GPU and is detailed in the GPU Usage Area). The Breakdown View offers useful information, such as the number of SetPass calls (otherwise known as Draw Calls), the total number of batches used to render the Scene, the number of batches saved from Dynamic Batching and Static Batching and how they are being generated, as well as memory consumed for textures.

The Rendering Area also offers a button to open the Frame Debugger, which will be explored more in [Chapter 3, *The Benefits of Batching*](#). The rest of this Area's information will prove useful when you go through [Chapter 3, *The Benefits of Batching*](#), and [Chapter 6, *Dynamic Graphics*](#).

The Memory Area

The Memory Area allows us to inspect memory usage of the application in the Breakdown View in the following two modes:

- Simple mode
- Detailed mode

Simple mode provides only a high-level overview of memory consumption of subsystems. This include Unity's low-level Engine, the Mono framework (total heap size that is being watched by the Garbage Collector), graphical assets, audio assets and buffers, and even memory used to store data collected by the Profiler.

Detailed mode shows memory consumption of individual `GameObjects` and `MonoBehaviours` for both their *Native* and *Managed* representations. It also has a column explaining the reason why an object may be consuming memory and when it might be deallocated.



The Garbage Collector is a common feature provided by the various languages Unity supports, which automatically releases any memory we have allocated to store data, but if it is handled poorly it has the potential to stall our application for brief moments. This topic, and many more related topics such as Native and Managed memory spaces, will be explored in [Chapter 8, Masterful Memory Management](#).

Note that information only appears in Detailed mode through manual sampling by clicking on the Take Sample: <TargetName> button. This is the only way to gather information when using Detailed mode, since performing this kind of analysis automatically for each update would be prohibitively expensive.

The Breakdown View also provides a button labelled Gather Object References, which can gather deeper memory information about some

objects.

The Memory Area will be a useful tool to use when we dive into the complexities of memory management, Native versus Managed memory, and the Garbage Collector in [Chapter 8, Masterful Memory Management](#).

The Audio Area

The Audio Area grants an overview of audio statistics and can be used both to measure CPU usage from the audio system and total memory consumed by Audio Sources (both for those that are playing or paused) and Audio Clips.

The Breakdown View provides lots of useful insight into how the Audio System is operating and how various audio channels and groups are being used.

The Audio Area may come in handy as we explore art assets in [Chapter 4](#), *Kickstart Your Art*.



*Audio is often overlooked when it comes to performance optimization, but audio can become a surprisingly large source of bottlenecks if it is not managed properly due to the potential amount of hard disk access and CPU processing required.
Don't neglect it!*

The Physics 3D and Physics 2D Areas

There are two different Physics Areas, one for 3D physics (Nvidia's PhysX) and another for the 2D physics system (Box2D). This Area provides various physics statistics, such as Rigidbody, Collider, and Contact counts.

The Breakdown View for each Physics Area provides some rudimentary insight into the subsystem's inner workings, but we can get further insight by exploring the Physics Debugger, which we will introduce in [Chapter 5, Faster Physics](#).

The Network Messages and Network Operations Areas

These two Areas provide information about Unity's Networking System, which was introduced during the Unity 5 release cycle. The information present will depend on whether the application is using the **High-Level API (HLAPI)** or **Transport Layer API (TLAPI)** provided by Unity. The HLAPI is a more easy-to-use system for managing Player and `GameObject` network synchronization automatically, whereas the TLAPI is a thin layer that operates just above the socket level, allowing Unity developers to conjure up their own networking system.

Optimizing network traffic is a subject that fills an entire book all by itself, where the right solution is typically very dependent on the particular needs of the application. This will not be a Unity-specific problem, and as such, the topic of network traffic optimization will not be explored in this book.

The Video Area

If our application happens to make use of Unity's `VideoPlayer` API, then we might find this Area useful for profiling video playback behavior.

Optimization of media playback is also a complex, non-Unity-specific topic and will not be explored in this book.

The UI and UI Details Areas

These Areas are new in Unity 2017 and provide insight into applications making use of Unity's built-in User Interface System. If we're using a custom-built or 3rd-party User Interface System (such as NGUI), then these Areas will probably provide little benefit.

Poorly optimized user interface can often affect one or both of the CPU and GPU, so we will investigate some code optimization strategies for UI in [Chapter 2, Scripting Strategies](#), and graphics-related approaches in [Chapter 6, Dynamic Graphics](#).

The Global Illumination Area

The Global Illumination Area is another new Area in Unity 2017, and gives us a fantastic amount of detail into Unity's **Global Illumination (GI)** system. If our application makes use of GI, then we should refer to this Area to verify that it is performing properly.

This Area may become useful as we explore lighting and shadowing in [chapter 6, Dynamic Graphics](#).

Best approaches to performance analysis

Good coding practices and project asset management often make finding the root cause of a performance issue relatively simple, at which point the only real problem is figuring out how to improve the code. For instance, if the method only processes a single gigantic `for` loop, then it will be a pretty safe assumption that the problem is either with how many iterations the loop is performing, whether or not the loop is causing cache misses by reading memory in a non-sequential fashion, how much work is done in each iteration, or how much work it takes to prepare for the next iteration.

Of course, whether we're working individually or in a group setting, a lot of our code is not always written in the cleanest way possible, and we should expect to have to profile some poor coding work from time to time.

Sometimes we are forced to implement a *hacky* solution for the sake of speed, and we don't always have the time to go back and refactor everything to keep up with our best coding practices. In fact, many code changes made in the name of performance optimization tend to appear very strange or arcane, often making our codebase more difficult to read. The common goal of software development is to make code that is clean, feature-rich, and fast. Achieving one of these is relatively easy, but the reality is that achieving two will cost significantly more time and effort, while achieving all three is a near-impossibility.

At its most basic level, performance optimization is just another form of problem-solving, and overlooking the obvious while problem-solving can be an expensive mistake. Our goal is to use benchmarking to observe our application looking for instances of problematic behavior, then use instrumentation to hunt through the code for clues about where the problem originates. Unfortunately, it's often very easy to get distracted by invalid data or jump to conclusions because we're being too impatient or missed a subtle detail. Many of us have run into occasions during software debugging where we could have found the root cause of the problem much faster if we had

simply challenged and verified our earlier assumptions. Hunting down performance issues is no different.

A checklist of tasks would be helpful to keep us focused on the issue, and not waste time chasing so-called *ghosts*. Of course, every project is different, with its own unique challenges to overcome, but the following checklist is general enough that it should be able to apply to any Unity project:

- Verifying that the target script is present in the Scene
- Verifying that the script appears in the Scene the correct number of times
- Verifying the correct order of events
- Minimizing ongoing code changes
- Minimizing internal distractions
- Minimizing external distractions

Verifying script presence

Sometimes, there are things we expect to see, but don't. These are usually easy to spot because the human brain is very good at pattern recognition and spotting differences we didn't expect. Meanwhile, there are times where we assume that something has been happening, but it didn't. These are generally more difficult to notice, because we're often scanning for the first kind of problem, and we're assuming that the things we don't see are working as intended. In the context of Unity, one problem that manifests itself this way is verifying that the scripts we expect to be operating are actually present in the Scene.

Script presence can be quickly verified by typing the following into the Hierarchy window textbox:

```
| t:<monobehaviour name>;
```

For example, typing `t:mytestmonobehaviour` (note that it is not case-sensitive) into the Hierarchy textbox will show a shortlist of all `GameObjects` that currently have at least one `MyTestMonoBehaviour` script attached as a Component.



Note that this shortlist feature also includes any `GameObjects` with Components that derive from the given script name.

We should also double-check that the `GameObjects` they are attached to are still enabled, since we may have disabled them during earlier testing since someone or something may have accidentally deactivated the object.

Verifying script count

If we're looking at our Profiler data and note that a certain `MonoBehaviour` method is being executed more times than expected, or is taking longer than expected, we might want to double-check that it only occurs as many times in the Scene as we expect it to. It's entirely possible that someone created the object more times than expected in the Scene file, or that we accidentally instantiated the object more than the expected number of times from code. If so, the problem could be due to conflicting or duplicated method invocations generating a performance bottleneck. We can verify the count using the same shortlist method used in the *Best approaches to performance analysis* section.

If we expected a specific number of Components to appear in the Scene, but the shortlist revealed more (or less!) than this, then it might be wise to write some initialization code that prevents this from ever happening again. We could also write some custom Editor helpers to display warnings to any level designers who might be making this mistake.

Preventing casual mistakes such as this is essential for good productivity, since experience tells us that if we don't explicitly disallow something, then someone, somewhere, at some point, for whatever reason, will do it anyway. This is likely to cost us a frustrating afternoon hunting down a problem that eventually turned out to be caused by human-error.

Verifying the order of events

Unity applications mostly operate as a series of callbacks from *Native code* to *Managed code*. This concept will be explained in more detail in more detail in [chapter 8, Masterful Memory Management](#), but for the sake of a brief summary, Unity's main thread doesn't operate like a simple console application would. In such applications, code would be executed with some obvious starting point (usually a `main()` function), and we would then have direct control of the game engine, where we will initialize major subsystems, then the game runs in a big `while`-loop (often called the *Game Loop*) that checks for user input, updates the game, renders the current Scene, and repeats. This loop only exits once the player chooses to quit the game.

Instead, Unity handles the Game Loop for us, and we expect callbacks such as `Awake()`, `Start()`, `Update()`, and `FixedUpdate()` to be called at specific moments. The big difference is that we don't have fine-grained control over the order in which events of the same type are called. When a new Scene is loaded (whether it's the first Scene of the game, or a later Scene), every `MonoBehaviour` Component's `Awake()` callback gets called, but there's no way of telling which order this will happen in.

So, if we one set of objects that configure some data in their `Awake()` callback, and then another set of objects do something with that configured data in their own `Awake()` callback, some reorganization or recreation of Scene objects or a random change in the codebase or compilation process (it's unclear what exactly causes it) may cause the order of these `Awake()` calls to change, and then the dependent objects will probably try to do things with data that wasn't initialized how we expected. The same goes for all other callbacks provided by `MonoBehaviour` Components, such as `start()` and `update()`.

There's no way of telling the order in which the same type of callback gets called among a group of `MonoBehaviour` Components, so we should be very careful not to assume that object callbacks are happening in a specific order. In fact, its essential practice to never write code in a way that assumes these callbacks will need to be called in a certain order because it could break

at any time.

A better place to handle late-stage initialization is in a `MonoBehaviour` Component's `Start()` callback, which is always called after every object's `Awake()` is called and just before its first `Update()`. Late-stage updates can also be done in the `LateUpdate()` callback.

If we're having trouble determining the actual order of events, then this is best handled by either step-through debugging with an IDE (MonoDevelop, Visual Studio, and so on) or by printing simple logging statements with `Debug.Log()`.



Be warned that Unity's logger is notoriously expensive. Logging is unlikely to change the order of the callbacks, but it can cause some unwanted spikes in performance if used too aggressively. Be smart and do targeted logging only on the most relevant parts of the codebase.

Coroutines are typically used to script some sequence of events, and when they're triggered will depend on what `yield` types are being used. The most difficult and unpredictable type to debug is perhaps the `waitForSeconds` `yield` type. The Unity Engine is non-deterministic, meaning that you'll get a slightly different behavior from one session and the next, even on the same hardware. For example, you might get 60 updates called during the first second of application runtime during one session, 59 in the next, and 62 in the one after that. In another session, you might get 61 updates in the first second, then 60, followed by 59.

A variable number of `Update()` callbacks will be called between when the Coroutine starts and when it ends, and so if the Coroutine depends on something's `Update()` being called a specific number of times, we will run into problems. It's best to keep Coroutine behavior dead-simple and dependency-free of other behavior once it begins. Breaking this rule may be tempting, but it's essentially guaranteed that some future change is going to interact with the Coroutine in an unexpected way leading to a long, painful debugging session of a game-breaking bug that's very hard to reproduce.

Minimizing ongoing code changes

Making code changes to the application in order to hunt down performance issues is best done carefully, as the changes are easy to forget as time wears on. Adding debug logging statements to our code can be tempting, but remember that it costs us time to introduce these calls, recompile our code, and remove these calls once our analysis is complete. In addition, if we forget to remove them, then they can cost unnecessary runtime overhead in the final build since Unity's debug Console window logging can be prohibitively expensive in both CPU and memory.

A good way to combat this problem is to add a flag or comment everywhere we made a change with our name so that it's easy to find and remove it later. Hopefully we're also wise enough to use a source-control tool for our codebase making it easy to differentiate the contents of any modified files and revert them to their original state. This is an excellent way to ensure that unnecessary changes don't make it into the final version. Of course, this is by no means a guaranteed solution if we also applied a fix at the same time and didn't double-check all of our modified files before committing the change.

Making use of breakpoints during runtime debugging is the preferred approach, as we can trace the full callstack, variable data, and conditional code paths (for example, `if-else` blocks), without risking any code changes or wasting time on recompilation. Of course, this is not always an option if, for example, we're trying to figure out what causes something strange to happen in one out of a thousand frames. In this case, it's better to determine a threshold value to look for and add an `if` statement with a breakpoint inside, which will be triggered when the value has exceeded the threshold.

Minimizing internal distractions

The Unity Editor has its own little quirks and nuances, which can sometimes make it confusing to debug some kinds of problems.

Firstly, if a single frame takes a long time to process, such that our game noticeably freezes, then the Profiler may not be capable of picking up the results and recording them in the Profiler window. This can be especially annoying if we wish to catch data during application/Scene initialization. The upcoming section, *Custom CPU profiling*, will offer some alternatives to explore to solve this problem.

One common mistake (that I have admittedly fallen victim to multiple times during the writing of this book) is that if we are trying to initiate a test with a keystroke and have the Profiler open, we should not forget to click back into the Editor's Game window before triggering the keystroke. If the Profiler is the most recently clicked window, then the Editor will send keystroke events to that, instead of the runtime application, and hence no `GameObject` will catch the event for that keystroke. This can also apply to the GameView for rendering tasks and even Coroutines using the `WaitForEndOfFrame` `yield` type. If the Game window is not visible and active in the Editor, then nothing is being rendered to that view, and therefore, no events that rely on the Game window rendering will be triggered. Be warned!

Vertical Sync (otherwise known as *VSync*) is used to match the application's frame rate to the frame rate of the device it is being displayed to, for example, a monitor may run at 60 Hertz (60 cycles per-second), and if a rendering loop in our game is running faster than this then it will sit and wait until that time has elapsed before outputting the rendered frame. This feature reduces *screen-tearing* which occurs when a new image is pushed to the monitor before the previous image was finished, and for a brief moment part of the new image overlaps the old image.

Executing the Profiler with VSync enabled will probably generate a lot of noisy spikes in the CPU Usage Area under the `WaitForTargetFPS` heading, as

the application intentionally slows itself down to match the frame rate of the display. These spikes often appear very large in Editor Mode since the Editor is typically rendering to a very small window, which doesn't take a lot of CPU or GPU work to render.

This will generate unnecessary clutter, making it harder to spot the real issue(s). We should ensure that we disable the VSync checkbox under the CPU Usage Area when we're on the lookout for CPU spikes during performance tests. We can disable the VSync feature entirely by navigating to Edit | Project Settings | Quality and then to the sub-page for the currently selected platform.

We should also ensure that a drop in performance isn't a direct result of a massive number of exceptions and error messages appearing in the Editor Console window. Unity's `Debug.Log()` and similar methods, such as `Debug.LogError()` and `Debug.LogWarning()` are notoriously expensive in terms of CPU usage and heap memory consumption, which can then cause garbage collection to occur and even more lost CPU cycles (refer to [Chapter 8, Masterful Memory Management](#), for more information on these topics).

This overhead is usually unnoticeable to a human being looking at the project in Editor Mode, where most errors come from the compiler or misconfigured objects. However, they can be problematic when used during any kind of runtime process, especially during profiling, where we wish to observe how the game runs in the absence of external disruptions. For example, if we are missing an object reference that we were supposed to assign through the Editor and it is being used in an `Update()` callback, then a single `MonoBehaviour` could throw new exceptions every single update. This adds lots of unnecessary noise to our profiling data.

Note that we can hide different log level types with the buttons shown in the next screenshot. The extra logging still costs CPU and memory to execute, even though they are not being rendered, but it does allow us to filter out the junk we don't want. Although, it is often a good practice to keep all of these options enabled to verify that we're not missing anything important.



Minimizing external distractions

This one is simple but absolutely necessary. We should double-check that there are no background processes eating away CPU cycles or consuming vast swathes of memory. Being low on available memory will generally interfere with our testing, as it can cause more cache misses, hard-drive access for virtual memory page-file swapping, and generally slow responsiveness of the application. If our application is suddenly behaving significantly worse than what we expected, double-check the system's task manager (or equivalent) for any CPU/memory/hard disk activity, which might be causing problems.

Targeted profiling of code segments

If our performance problem isn't resolved by the checklist mentioned previously, then we probably have a real issue on our hands that demands further analysis. The Profiler window is effective at showing us a broad overview of performance; it can help us find specific frames to investigate and can quickly inform us which `MonoBehaviour` and/or method may be causing issues. We would then need to figure out whether the problem is reproducible, under what circumstances a performance bottleneck arises, and where exactly within the problematic code block the issue is originating from.

To accomplish these, we will need to perform some profiling of targeted sections of our code, and there are a handful of useful techniques we can employ for this task. For Unity projects, they essentially fit into two categories:

- Controlling the Profiler from script code
- Custom timing and logging methods



Note that the next section focuses on how to investigate Scripting bottlenecks through C# code. Detecting the source of bottlenecks in other engine subsystems will be discussed in their related chapters.

Profiler script control

The Profiler can be controlled in script code through the `Profiler` class. There are several useful methods in this class that we can explore within the Unity documentation, but the most important methods are the delimiter methods that activate and deactivate profiling at runtime. These can be accessed through the `UnityEngine.Profiling.Profiler` class through its `BeginSample()` and `EndSample()` methods.



Note that the delimiter methods, `BeginSample()` and `EndSample()`, are only compiled in development builds, and as such, they will not be compiled or executed in release builds where Development Mode is unchecked. This is commonly known as non-operation or no-op code.

The `BeginSample()` method has an overload that allows a custom name for the sample to appear in the CPU Usage Area's Hierarchy mode. For example, the following code will profile invocations of this method and make the data appear in the Breakdown View under a custom heading, as follows:

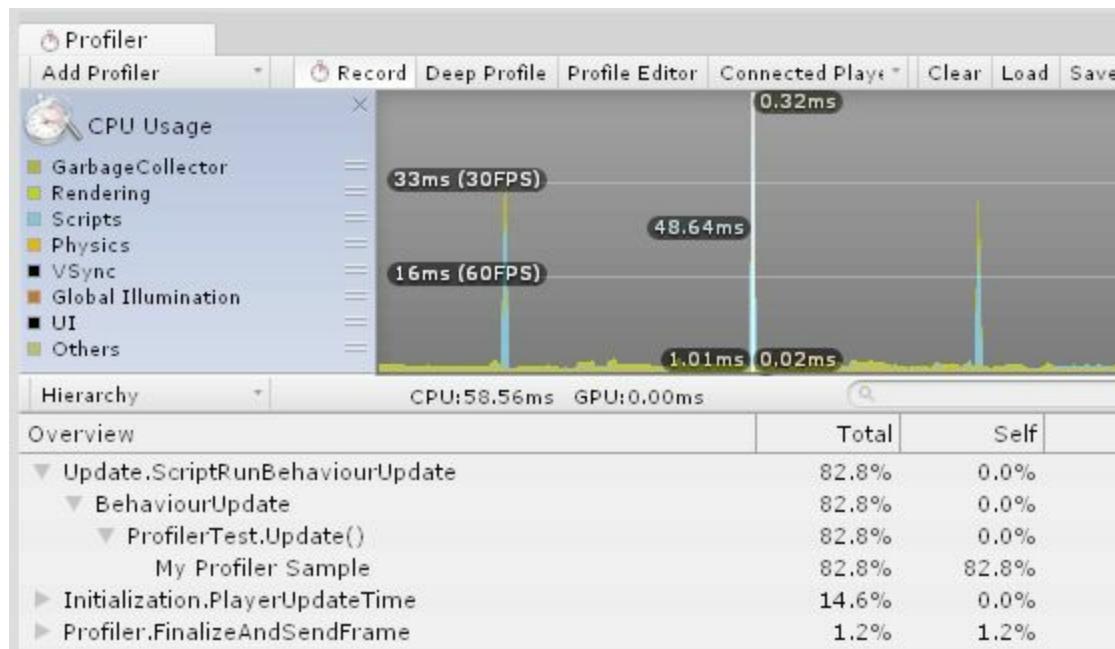
```
void DoSomethingCompletelyStupid() {
    Profiler.BeginSample("My Profiler Sample");
    List<int> listOfInts = new List<int>();
    for(int i = 0; i < 1000000; ++i) {
        listOfInts.Add(i);
    }
    Profiler.EndSample();
}
```



You can download the example code files from your account at <http://www.packtpub.com> for all the Packt Publishing books you have purchased. If you purchased this book elsewhere, you can visit <http://www.packtpub.com/support> and register to have the files emailed directly to you.

We should expect that invoking this poorly designed method (which generates a `List` containing a million integers, and then does absolutely nothing with it) will cause a huge spike in CPU usage, chew up several

megabytes of memory, and appear in the Profiler Breakdown View under the heading My Profiler Sample, as the following screenshot shows:



Custom CPU Profiling

The Profiler is just one tool at our disposal. Sometimes, we may want to perform customized profiling and logging of our code. Maybe we're not confident that the Unity Profiler is giving us the right answer, maybe we consider its overhead cost too great, or maybe we just like having complete control of every single aspect of our application. Whatever our motivations, knowing some techniques to perform an independent analysis of our code is a useful skill to have. It's unlikely we'll only be working with Unity for the entirety of our game development careers, after all.

Profiling tools are generally very complex, so it's unlikely we would be able to generate a comparable solution on our own within a reasonable time frame. When it comes to testing CPU usage, all we should really need is an accurate timing system, a fast, low-cost way of logging that information, and some piece of code to test against. It just so happens that the .NET library (or, technically, the Mono framework) comes with a `Stopwatch` class under the `System.Diagnostics` namespace. We can stop and start a `Stopwatch` object at any time, and we can easily acquire a measure of how much time has passed since the `Stopwatch` was started.

Unfortunately, this class is not perfectly accurate; it is accurate only to milliseconds, or tenths of a millisecond, at best. Counting high-precision real time with a CPU clock can be a surprisingly difficult task when we start to get into it; so, in order to avoid a detailed discussion of the topic, we should try to find a way for the `Stopwatch` class to satisfy our needs.

If precision is important, then one effective way to increase it is by running the same test multiple times. Assuming that the test code block is both easily repeatable and not exceptionally long, we should be able to run thousands, or even millions of tests within a reasonable time frame and then divide the total elapsed time by the number of tests we just performed to get a more accurate time for a single test.

Before we get obsessed with the topic of high precision, we should first ask

ourselves if we even need it. Most games expect to run at 30 FPS or 60 FPS, which means that they only have around 33 milliseconds or 16 milliseconds, respectively, to compute everything for the entire frame. So, hypothetically, if we need to bring only the performance of a particular code block under ten milliseconds, then repeating the test thousands of times to get microsecond precision is too many orders of magnitude away from the target to be worthwhile.

The following is a class definition for a custom timer that uses a stopwatch to count time for a given number of tests:

```
using System;
using System.Diagnostics;

public class CustomTimer : IDisposable {
    private string _timerName;
    private int _numTests;
    private Stopwatch _watch;

    // give the timer a name, and a count of the
    // number of tests we're running
    public CustomTimer(string timerName, int numTests) {
        _timerName = timerName;
        _numTests = numTests;
        if (_numTests <= 0) {
            _numTests = 1;
        }
        _watch = Stopwatch.StartNew();
    }

    // automatically called when the 'using()' block ends
    public void Dispose() {
        _watch.Stop();
        float ms = _watch.ElapsedMilliseconds;
        UnityEngine.Debug.Log(string.Format("{0} finished: {1:0.00} " +
            "milliseconds total, {2:0.000000} milliseconds per-test " +
            "for {3} tests", _timerName, ms, ms / _numTests, _numTests));
    }
}
```



Adding an underscore before member variable names is a common and useful way of distinguishing a class' member variables (also known as fields) from a method's arguments and local variables.

The following is an example of the `CustomTimer` class usage:

```
const int numTests = 1000;
using (new CustomTimer("My Test", numTests)) {
```

```
|   for(int i = 0; i < numTests; ++i) {  
|     TestFunction();  
|   } // the timer's Dispose() method is automatically called here
```

There are three things to note when using this approach. Firstly, we are only making an average of multiple method invocations. If processing time varies enormously between invocations, then that will not be well represented in the final average.

Secondly, if memory access is common, then repeatedly requesting the same blocks of memory will result in an artificially higher cache hit rate (where the CPU can find data in memory very quickly because it's accessed the same region recently), which will bring the average time down when compared to a typical invocation.

Thirdly, the effects of **Just-In-Time (JIT)** compilation will be effectively hidden for similarly artificial reasons, as it only affects the first invocation of the method. JIT compilation is a .NET feature that will be covered in more detail in [Chapter 8, Masterful Memory Management](#).

The `using` block is typically used to safely ensure that unmanaged resources are properly destroyed when they go out of scope. When the `using` block ends, it will automatically invoke the object's `Dispose()` method to handle any cleanup operations. In order to achieve this, the object must implement the `IDisposable` interface, which forces it to define the `Dispose()` method.

However, the same language feature can be used to create a distinct code block, which creates a short-term object, which then automatically processes something useful when the code block ends, which is how it is being used in the preceding code block.



Note that the `using` block should not be confused with the `using` statement, which is used at the start of a script file to pull in additional namespaces. It's extremely ironic that the keyword for managing namespaces in C# has a naming conflict with another keyword.

As a result, the `using` block and the `customTimer` class give us a clean way of

wrapping our test code in a way that makes it obvious when and where it is being used.

Another concern to worry about is application warm-up time. Unity has a significant startup cost when a Scene begins, given the amount of data that needs to be loaded from disk, the initialization of complex subsystems, such as the Physics and Rendering Systems, and the number of calls to various `Awake()` and `Start()` callbacks that need to be resolved before anything else can happen. This early overhead might only last a second, but that can have a significant effect on the results of our testing if the code is also executed during this early initialization period. This makes it crucial that if we want an accurate test, then any runtime testing should begin only after the application has reached a steady state.

Ideally, we would be able to execute the target code block in its own Scene after its initialization has completed. This is not always possible, so as a backup plan, we could wrap the target code block in an `Input.GetKeyDown()` check in order to assume control over when it is invoked. For example, the following code will execute our test method only when the spacebar is pressed:

```
if (Input.GetKeyDown(KeyCode.Space)) {
    const int numTests = 1000;
    using (new CustomTimer("Controlled Test", numTests)) {
        for(int i = 0; i < numTests; ++i) {
            TestFunction();
        }
    }
}
```

As mentioned previously, Unity's Console window logging mechanism is prohibitively expensive. As a result, we should try not to use these logging methods in the middle of a profiling test (or during gameplay, for that matter). If we find ourselves absolutely needing detailed profiling data that prints out lots of individual messages (such as performing a timing test on a loop to figure out which iteration is costing more time than the rest), then it would be wiser to cache the logging data and print it all out at the end, as the `CustomTimer` class does. This will reduce runtime overhead, at the cost of some memory consumption. The alternative is that many milliseconds are lost to printing each `Debug.Log()` message in the middle of the test, which pollutes the

results.

The `CustomTimer` class also makes use of `string.Format()`. This will be covered in more detail in [Chapter 8, *Masterful Memory Management*](#), but the short explanation is that this method is used because generating custom `string` object using the `+` operator (for example, code such as `Debug.Log("Test: " + output);`) can result in a surprisingly large amount of memory allocations, which attracts the attention of the Garbage Collector. Doing otherwise would conflict with our goal of achieving accurate timing and analysis and should be avoided.

Final thoughts on Profiling and Analysis

One way of thinking about performance optimization is *the act of stripping away unnecessary tasks that spend valuable resources*. We can do the same and maximize our own productivity through minimizing any wasted effort. Effective use of the tools we have at our disposal is of paramount importance. It would serve us well to optimize our own workflow by keeping aware of some best practices and techniques.

Most, if not all, advice for using any kind of data-gathering tool properly can be summarized into three different strategies:

- Understanding the tool
- Reducing noise
- Focusing on the issue

Understanding the Profiler

The Profiler is an arguably well-designed and intuitive tool, so understanding the majority of its feature set can be gained by simply spending an hour or two exploring its options with a test project and reading its documentation. The more we know about a tool in terms of its benefits, pitfalls, features, and limitations, the more sense we can make of the information it is giving us, so it is worth spending the time to use it in a playground setting. We don't want to be two weeks away from release, with a hundred performance defects to fix, with no idea how to do performance analysis efficiently.

For example, always remain aware of the relative nature of the Timeline View's graphical display. The Timeline View does not provide values on its vertical axis and automatically readjusts this axis based on the content of the last 300 frames; it can make small spikes appear to be a bigger problem than they really are because of the relative change. So, just because a spike or resting state in the timeline seems large and threatening does not necessarily mean there is a performance issue.

Several Areas in the Timeline View provide helpful benchmark bars, which appear as horizontal lines with a timing and FPS value associated with them. These should be used to determine the magnitude of the problem. Don't let the Profiler trick us into thinking that big spikes are always bad. As always, it's only important if the user will notice it.

As an example, if a large CPU usage spike does not exceed the 60 FPS or 30 FPS benchmark bars (depending on the application's target frame rate), then it would be wise to ignore it and search elsewhere for CPU performance issues, as no matter how much we improve the offending piece of code, it will probably never be noticed by the end user, and therefore isn't a critical issue that affects user experience.

Reducing noise

The classical definition of noise (at least in the realm of computer science) is meaningless data, and a batch of profiling data that was blindly captured with no specific target in mind is always full of data that won't interest us. More sources of data takes more time to mentally process and filter, which can be very distracting. One of the best methods to avoid this is to simply reduce the amount of data we need to process by stripping away any data deemed non-vital to the current situation.

Reducing clutter in the Profiler's graphical interface will make it easier to determine which subsystems are causing a spike in resource usage. Remember to use the colored checkboxes in each Timeline View Area to narrow the search.



Be warned that these settings are auto-saved in the Editor, so ensure that you re-enable them for the next profiling session, as this might cause us to miss something important next time.

Also, `GameObjects` can be deactivated to prevent them from generating profiling data, which will also help reduce clutter in our profiling data. This will naturally cause a slight performance boost for each object we deactivate. However, if we're gradually deactivating objects and performance suddenly becomes significantly more acceptable when a specific object is deactivated, then clearly that object is related to the root cause of the problem.

Focusing on the issue

This category may seem redundant, given that we've already covered reducing noise. All we should have left is the issue at hand, right? Not exactly. Focus is the skill of not letting ourselves become distracted by inconsequential tasks and wild goose chases.

Recall that profiling with the Unity Profiler comes with a minor performance cost. This cost is even more severe when using the Deep Profiling option. We might even introduce more minor performance costs into our application with additional logging. It's easy to forget when and where we introduced profiling code if the hunt continues for several hours.

We are effectively changing the result by measuring it. Any changes we implement during data sampling can sometimes lead us to chase after non-existent bugs in the application, when we could have saved ourselves a lot of time by attempting to replicate the scenario without additional profiling instrumentation. If the bottleneck is reproducible and noticeable without profiling, then it's a candidate to begin an investigation. However, if new bottlenecks keep appearing in the middle of an existing investigation, then keep in mind that they could be bottlenecks we introduced with our test code and not an existing problem that's been newly exposed.

Finally, when we have finished profiling, completed our fixes, and are now ready to move on to the next investigation, we should make sure to profile the application one last time to verify that the changes have had the intended effect.

Summary

You learned a great deal throughout this chapter on how to detect and analyze performance issues within your application. You learned about many of the Profiler's features and secrets, explored a variety of tactics to investigate performance issues with a more hands-on approach, and have been introduced to a variety of different tips and strategies to follow. You can use these to improve your productivity immensely, so long as you appreciate the wisdom behind them and remember to exploit them when the situation makes it possible.

This chapter has introduced us to the tips, tactics, and strategies we will need to find a performance problem that needs improvement. In the remaining chapters, we will explore methods on how to fix issues and improve performance whenever possible. So, give yourself a pat on the back for getting through the boring part first, and let's move on to learning some approaches to optimize our C# scripts.

Scripting Strategies

Since scripting will consume a great deal of our development time, it will be enormously beneficial to learn some best practices. Scripting is a very broad term, so we will try to limit our exposure in this chapter to situations that are very Unity specific, focusing on problems surrounding `MonoBehaviours`, `GameObjects`, and related functionality.



We will discuss the nuances and advanced topics of the C# language, .NET library, and Mono Framework in [Chapter 8](#), Masterful Memory Management.

In this chapter, we will explore ways of applying performance enhancements to the following areas:

- Accessing Components
- Component callbacks (`Update()`, `Awake()`, and so on)
- Coroutines
- `GameObject` and `Transform` usage
- Interobject communication
- Mathematical calculations
- Deserialization such as Scene and Prefab loading

Whether you have some specific problems in mind that you wish to solve or you just want to learn some techniques for future reference, this chapter will introduce you to a wide array of methods that you can use to improve your scripting efforts now and in the future. In each case, we will explore how and why the performance issue arises, an example situation in which the problem is occurring, and one or more solutions to combat the issue.

Obtain Components using the fastest method

There are several variations of the `GetComponent()` method, and they each have a different performance cost, so it becomes prudent to call the fastest possible version of this method. The three overloads available are `GetComponent(string)`, `GetComponent<T>()`, and `GetComponent(typeof(T))`. It turns out that the fastest version depends on which version of Unity we are running since several optimizations have been made to these methods through the years; however, in all later versions of Unity 5 and the initial release of Unity 2017, it is best to use the `GetComponent<T>()` variant.

Let's prove this with some simple testing:

```
int numTests = 1000000;
TestComponent test;
using (new CustomTimer("GetComponent(string)", numTests)) {
    for (var i = 0; i < numTests; ++i) {
        test = (TestComponent)GetComponent("TestComponent");
    }
}

using (new CustomTimer("GetComponent<ComponentName>", numTests)) {
    for (var i = 0; i < numTests; ++i) {
        test = GetComponent<TestComponent>();
    }
}

using (new CustomTimer("GetComponent(typeof(ComponentName))", numTests)) {
    for (var i = 0; i < numTests; ++i) {
        test = (TestComponent)GetComponent(typeof(TestComponent));
    }
}
```

The preceding code tests each of the `GetComponent()` overloads a million times. This is far more tests than would be sensible for a typical project, but it helps make the relative costs clear.

Here is the result we get when the tests complete:

```
! GetComponent(string) finished: 6413.00ms total, 0.006413ms per test for 1000000 tests
UnityEngine.Debug:Log(Object)
! GetComponent<ComponentName> finished: 89.00ms total, 0.000089ms per test for 1000000 tests
UnityEngine.Debug:Log(Object)
! GetComponent(typeof(ComponentName)) finished: 95.00ms total, 0.000095ms per test for 1000000 tests
UnityEngine.Debug:Log(Object)
```

As you can see, the `GetComponent<T>()` method is only a tiny fraction faster than `GetComponent(typeof(T))`, whereas `GetComponent(string)` is significantly slower than the alternatives. Therefore, it is pretty safe to use either of the type-based versions of `GetComponent()` because of the small performance difference. However, we should ensure that we never use `GetComponent(string)` since the outcome is identical, and there is no benefit for the costs incurred. There are some very rare exceptions, like if we were writing a custom debug console for Unity, which might parse a user-inputted string to acquire a Component. In any case, these would be used for debugging and diagnostics situations only, where performance isn't too important. For a production-level application, the use of `GetComponent(string)` is just a needless waste.

Remove empty callback definitions

The primary means of scripting in Unity is to write callback functions in classes derived from `MonoBehaviour`, which we know Unity will call when necessary. Perhaps the four most commonly used callbacks are `Awake()`, `Start()`, `Update()`, and `FixedUpdate()`.

`Awake()` is called the moment a `MonoBehaviour` is first created, whether this occurs during Scene initialization or when a new `GameObject` containing the `MonoBehaviour` is instantiated at runtime from a Prefab. `Start()` will be called shortly after `Awake()` but before its first `Update()`. During Scene initialization, every `MonoBehaviour` Component's `Awake()` callback will be called before any of their `Start()` callbacks are.

After this, `Update()` will be called repeatedly, each time the Rendering Pipeline presents a new image. `Update()` will continue to be called provided the `MonoBehaviour` is still present in the Scene, it is still enabled, and its parent `GameObject` is active.

Finally, `FixedUpdate()` is called just prior to when the Physics Engine updates. Fixed Updates are used whenever we want activity similar in behavior to `Update()` but that isn't tied directly to the render frame rate and is called more consistently over time.

Refer to the following page in the Unity documentation for an accurate picture of when various Unity callbacks are called: <https://docs.unity3d.com/Manual/ExecutionOrder.html>.

Whenever a `MonoBehaviour` is first instantiated in our Scene, Unity will add any defined callbacks to a list of function pointers, which it will call at key moments. However, it is important to realize that Unity will hook into these callbacks even if the function body is empty. The core Unity Engine has no awareness that these function bodies may be empty and only knows that the

method has been defined and, therefore, that it must acquire it and then call it when necessary. Consequently, if we leave empty definitions of these callbacks scattered throughout the codebase, then they will waste a small amount of CPU due to the overhead cost of the engine invoking them.

This can be a problem since any time we create a new `MonoBehaviour` script file in Unity, it will automatically generate two boilerplate callback stubs for us for `start()` and `update()`:

```
// Use this for initialization
void Start () {
}

// Update is called once per-frame
void Update () {
}
```

It can be easy to accidentally leave these empty definitions on scripts that don't actually need them. An empty `start()` definition is liable to cause any object to initialize a little slower, for no good reason. This effect may not be particularly noticeable for a handful of `MonoBehaviours`, but as development on the project continues and we populate our scenes with thousands of custom `MonoBehaviours` with lots of empty `start()` definitions, it could start to become a problem, causing slow Scene initialization and wasting CPU time whenever a new Prefab is created via `GameObject.Instantiate()`.

Such calls typically happen during key gameplay events; for instance, when two objects collide, we might spawn a Particle Effect, create some floating damage text, play a sound effect, and so on. This can be a critical moment for performance because we've suddenly requested that the CPU makes a lot of complicated changes, but with only a finite amount of time to complete them before the current frame ends. If this process takes too long, then we would experience a frame drop as the Rendering Pipeline isn't allowed to present a new frame until all of the `update()` callbacks (counted across all `MonoBehaviours` in the Scene!) have finished. Ergo, a bunch of empty `start()` definitions being called at this time is a needless waste and could potentially cut into our tight time-budget at a critical moment.

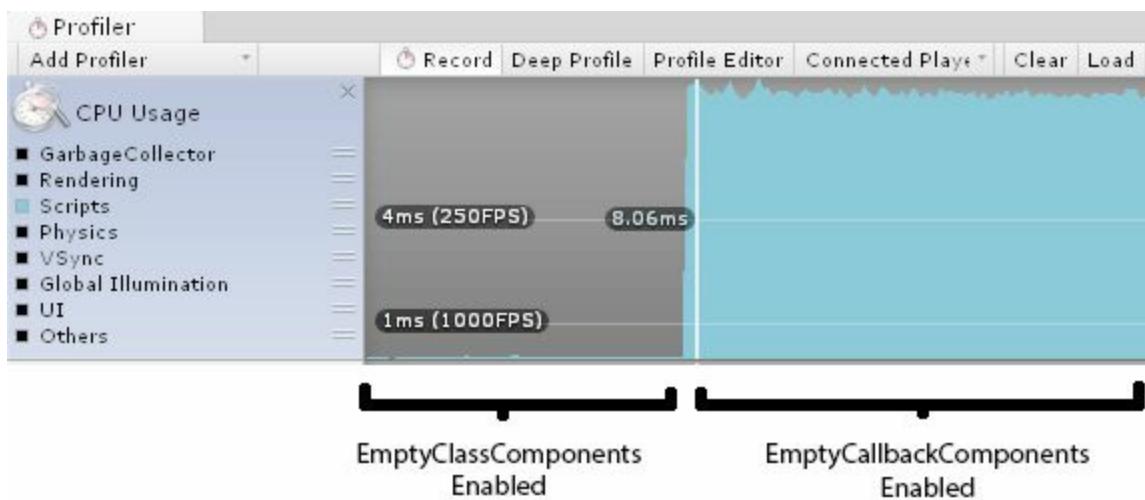
Meanwhile, if our Scene contains thousands of `MonoBehaviours` with these

empty `Update()` definitions, then we would be wasting a lot of CPU cycles every frame, potentially causing havoc on our frame rate.

Let's prove all of this with a simple test. Our test Scene should have `GameObjects` with two types of Component, `EmptyClassComponent` with no methods defined at all and `EmptyCallbackComponent` with an empty `Update()` callback defined:

```
public class EmptyClassComponent : MonoBehaviour {  
}  
  
public class EmptyCallbackComponent : MonoBehaviour {  
    void Update () {}  
}
```

The following are the test results for 30,000 Components of each type. If we enable all `GameObjects` with attached `EmptyClassComponents` during runtime, then nothing interesting happens under the CPU Usage Area of the Profiler. There will be a small amount of background activity, but none of this activity will be caused by the `EmptyClassComponents`. However, as soon as we enable all objects with `EmptyCallbackComponent`, we will observe a huge increase in CPU usage:



It's hard to imagine a Scene with more than 30,000 objects in it, but keep in mind that `MonoBehaviours` contain the `update()` callback, not `GameObjects`. A single `GameObject` can contain multiple `MonoBehaviours` at once, and each of their children can contain even more `MonoBehaviours`, and so on. A few thousand or

even a hundred empty `update()` callbacks will inflict a noticeable impact on frame rate budget, for zero potential gain. This is particularly common with Unity UI Components which tend to attach a lot of different Components in a very deep hierarchy.

The fix for this is simple: delete the empty callback definitions. Unity will have nothing to hook into, and nothing will be called. Finding such empty definitions in an expansive codebase may be difficult, but if we use some basic regular expressions (known as *regex*), we should be able to find what we're looking for relatively easily.



All common code-editing tools for Unity, such as MonoDevelop, Visual Studio, and even Notepad++, provide a way to perform a regex-based search on the entire codebase. Check out the tool's documentation for more information, since the method can vary greatly depending on the tool and its version.

The following regex search should find any empty `update()` definitions in our code:

```
| void\s*Update\s*?\(\s*?\)\s*?\n*?\{\n*?\s*?\}
```

This regex checks for a standard method definition of the `update()` callback, while including any surplus whitespace and newline characters that can be distributed throughout the method definition.

Naturally, all of the above is also true for the non-boilerplate Unity callbacks, such as `onGUI()`, `OnEnable()`, `OnDestroy()`, and `LateUpdate()`. The only difference is that only `start()` and `update()` are defined automatically in a new script. Check out the `MonoBehaviour` Unity Documentation page for a complete list of these callbacks at <http://docs.unity3d.com/ScriptReference/MonoBehaviour.html>.

It might also seem unlikely that someone generated so many empty versions of these callbacks in our codebase, but never say never. For example, if we use a common base class `MonoBehaviour` throughout all of our custom Components, then a single empty callback definition in that base class will permeate the entire game, which can cost us dearly. Be particularly careful of the `onGUI()` method, as it can be invoked multiple times within the same frame

or UI event.

Perhaps the most common source of performance problems in Unity scripting is to misuse the `update()` callback by doing one or more of the following things:

- Repeatedly recalculating a value that rarely or never changes
- Having too many Components perform work for a result which could be shared
- Performing work far more often than is necessary

It's worth getting into the habit of remembering that literally every single line of code we write in an `update()` callback, and in functions called by those callbacks, will eat into our frame rate budget. To hit 60 fps, we have 16.667 milliseconds to complete all of the work in all of our `update()` callbacks, every frame. This seems like plenty of time when we start prototyping, but somewhere in the middle of development, we will probably start noticing things getting slower and less responsive because we've gradually been eating away at that budget, due to an unchecked desire to cram more stuff into our project.

Let's cover some tips that directly address these problems.

Cache Component references

Repeatedly recalculating a value is a common mistake when scripting in Unity, and particularly when it comes to the `GetComponent()` method. For example, the following script code is trying to check a creature's health value, and if its health goes below 0, it will disable a series of Components to prepare it for a death animation:

```
void TakeDamage() {  
  
    Rigidbody rigidbody = GetComponent<Rigidbody>();  
    Collider collider = GetComponent<Collider>();  
    AIControllerComponent ai = GetComponent<AIControllerComponent>();  
    Animator anim = GetComponent<Animator>();  
  
    if (GetComponent<HealthComponent>().health < 0) {  
        rigidbody.enabled = false;  
        collider.enabled = false;  
        ai.enabled = false;  
        anim.SetTrigger("death");  
    }  
}
```

Each time this poorly optimized method executes, it will reacquire five different Component references. This is not very friendly on CPU usage. This is particularly problematic if the main method was called during `Update()`. Even if it is not, it still might coincide with other important events, such as creating particle effects, replacing an object with a Ragdoll (thus invoking various activity in the Physics Engine), and so on. This coding style can seem harmless, but it can cause a lot of long-term problems and runtime work for very little benefit.

It costs us a small amount of memory space (only 32 or 64 bits each time; Unity version, platform, and fragmentation permitting) to cache these references for future use. So, unless you're extremely bottlenecked on memory, a better approach would be to acquire the references during initialization and keep them until they are needed:

```
private HealthComponent _healthComponent;  
private Rigidbody _rigidbody;  
private Collider _collider;
```

```
private AIControllerComponent _ai;
private Animator _anim;

void Awake() {
    _healthComponent = GetComponent<HealthComponent>();
    _rigidbody = GetComponent<Rigidbody>();
    _collider = GetComponent<Collider>();
    _ai = GetComponent<AIControllerComponent>();
    _anim = GetComponent<Animator>();
}

void TakeDamage() {
    if (_healthComponent.health < 0) {
        _rigidbody.detectCollisions = false;
        _collider.enabled = false;
        _ai.enabled = false;
        _anim.SetTrigger("death");
    }
}
```

Caching Component references in this way spares us from reacquiring them each time they're needed, saving us some CPU overhead each time. The cost is a small amount of additional memory consumption, which is very often worth the price.

The same tip applies to literally any piece of data we decide to calculate at runtime. There's no need to ask the CPU to keep recalculating the same value every `Update()` when we can just store it in memory for future reference.

Share calculation output

Performance can be saved by having multiple objects share the result of some calculation; of course, this only works if all of them would generate the same result. Such situations are often easy to spot, but can be tricky to refactor, and so exploiting this would be very implementation dependent.

Some examples might include finding an object in a Scene, reading data from a file, parsing data (such as XML or JSON), finding something in a big list or deep dictionary of information, calculating pathing for a group of **Artificial Intelligence (AI)** objects, complex mathematics-like trajectories, raycasting, and so on.

Think about each time an expensive operation is undertaken, and consider whether it is being called from multiple locations but always results in the same output. If this is the case, then it would be wise to restructure things so that the result is calculated once and then distributed to every object that needs it in order to minimize the amount of recalculation. The biggest cost is typically just a small loss in code simplicity, although some extra overhead may be inflicted by passing the value around.

Note that it's often easy to get into the habit of hiding some big complex function in a base class, and then we define derived classes that make use of that function, completely forgetting how costly that function is because we rarely glance at that code again. It's best to use the Unity Profiler to tell us how many times that expensive function may be called, and as always, don't preoptimize those functions unless it's been proven to be a performance issue. No matter how expensive it may be, if it doesn't cause us to exceed performance restrictions (such as frame rate and memory consumption), then it's not really a performance problem.

Update, Coroutines, and InvokeRepeating

Another habit that's easy to fall into is to call something repeatedly in an `Update()` callback way more often than is needed. For example, we may start with a situation like this:

```
void Update() {
    ProcessAI();
}
```

In this case, we're calling some custom `ProcessAI()` subroutine every single frame. This may be a complex task, requiring the AI system to check some grid system to figure out where it's meant to move or determine some fleet maneuvers for a group of spaceships or whatever our game needs for its AI.

If this activity is eating into our frame rate budget too much, and the task can be completed less frequently than every frame with no significant drawbacks, then a good trick to improve performance is to simply reduce the frequency that `ProcessAI()` gets called:

```
private float _aiProcessDelay = 0.2f;
private float _timer = 0.0f;

void Update() {
    _timer += Time.deltaTime;
    if (_timer > _aiProcessDelay) {
        ProcessAI();
        _timer -= _aiProcessDelay;
    }
}
```

In this case, we've reduced the `Update()` callback's overall cost by only invoking `ProcessAI()` about five times every second, which is an improvement over the previous situation, at the expense of code that can take a bit of time to understand at first glance, and a little extra memory to store some floating-point data. Although, at the end of the day we're still having Unity call an empty callback function more often than not.

This function is a perfect example of a function we could convert into a Coroutine to make use of their delayed invocation properties. As mentioned previously, Coroutines are typically used to script a short sequence of events, either as a one-time, or repeated, action. They should not be confused with threads, which would run on a completely different CPU core in a concurrent manner, and multiple threads can be running simultaneously. Instead, Coroutines run on the main thread in a sequential manner such that only one Coroutine is handled at any given moment, and each Coroutine decides when to pause and resume via `yield` statements. The following code is an example of how we might rewrite the above `update()` callback in the form of a Coroutine:

```
void Start() {
    StartCoroutine(ProcessAICoroutine ());
}

IEnumerator ProcessAICoroutine () {
    while (true) {
        ProcessAI();
        yield return new WaitForSeconds(_aiProcessDelay);
    }
}
```

The preceding code demonstrates a Coroutine that calls `ProcessAI()`, then pause at the `yield` statement for the given number of seconds (the value of `_aiProcessDelay`) before the main thread resumes the Coroutine again. At which point, it will return to the start of the loop, call `ProcessAI()`, pause on the `yield` statement again, and repeat forever (via the `while(true)` statement) until asked to stop.

The main benefit of this approach is that this function will only be called as often as dictated by the value of `_aiProcessDelay`, and it will sit idle until that time, reducing the performance hit inflicted in most of our frames. However, this approach has its drawbacks.

For one, starting a Coroutine comes with an additional overhead cost relative to a standard function call (around three times as slow), as well as some memory allocations to store the current state in memory until it is invoked the next time. This additional overhead is also not a one-time cost because Coroutines often constantly call `yield`, which inflicts the same overhead cost

again and again, so we need to ensure that the benefits of a reduced frequency outweighs this cost.



In a test of 1,000 objects with empty `update()` callbacks, it took 1.1 milliseconds to process, whereas 1,000 Coroutines yielding on `waitForEndOfFrame` (which has identical frequency to `Update()` callbacks) took 2.9 milliseconds. So, the relative cost is almost three times as much.

Secondly, once initialized, Coroutines run independent of the triggering `MonoBehaviour` Component's `Update()` callback and will continue to be invoked regardless of whether the Component is disabled or not, which can make them unwieldy if we're performing a lot of `GameObject` construction and destruction.

Thirdly, the Coroutine will automatically stop the moment the containing `GameObject` is made inactive for whatever reason (whether it was set inactive or one of its parents was) and will not automatically restart if the `GameObject` is set active again.

Finally, by converting a method into a Coroutine, we may have reduced the performance hit inflicted during most of our frames, but if a single invocation of the method body causes us to break our frame rate budget, then it will still be exceeded no matter how rarely we call the method. Therefore, this approach is best used for situations where we are only breaking our frame rate budget because of the sheer number of times the method is called in a given frame, not because the method is too expensive on its own. In those cases, we have no option but to either dig into and improve the performance of the method itself or reduce the cost of other tasks to free up the time it needs to complete its work.

There are several `yield` types available to us when generating Coroutines. `waitForSeconds` is fairly self-explanatory; the Coroutine will pause at the `yield` statement for a given number of seconds. It is not exact an exact timer, however, so expect a small amount of variation when this `yield` type actually resumes.

`waitForSecondsRealTime` is another option and is different to `waitForSeconds` only

in that it uses unscaled time. `WaitForSeconds` compares against scaled time, which is affected by the global `Time.timeScale` property while `WaitForSecondsRealTime` is not, so be careful about which `yield` type you use if you're tweaking the time scale value (for example, for slow motion effects).

There is also `WaitForEndOfFrame`, which would continue at the end of the next `Update()`, and then there's `WaitForFixedUpdate`, which would continue at the end of the next `FixedUpdate()`. Lastly, Unity 5.3 introduced `WaitUntil` and `WaitWhile`, where we provide a delegate function, and the Coroutine will pause until the given delegate returns `true` or `false`, respectively. Note that the delegates provided to these `yield` types will be executed for each `Update()` until they return the Boolean value needed to stop them, which makes them very similar to a Coroutine using `WaitForEndOfFrame` in a `while`-loop that ends on a certain condition. Of course, it is also important that the delegate function we provide is not too expensive to execute.



Delegate functions are incredibly useful constructs in C# that allows us to pass local methods around as arguments to other methods and are typically used for callbacks. Check out the MSDN C# Programming Guide for more information on delegates at <https://docs.microsoft.com/en-us/dotnet/csharp/programming-guide/delegates/>.

The way that some `Update()` callbacks are written could probably be condensed down into simple Coroutines that always call `yield` on one of these types, but we should be aware of the drawbacks mentioned previously. Coroutines can be tricky to debug since they don't follow normal execution flow; there's no caller in the callstack we can directly blame as to why a Coroutine triggered at a given time, and if Coroutines perform complex tasks and interact with other subsystems, then they can result in some impossibly difficult bugs because they happened to be triggered at a moment that some other code didn't expect, which also tend to be the kinds of bugs that are painstakingly difficult to reproduce. If you do wish to make use of Coroutines, the best advice is to keep them simple and independent of other complex subsystems.

Indeed, if our Coroutine is simple enough that it can be boiled down to a `while`-loop that always calls `yield` on `WaitForSeconds`, or `WaitForSecondsRealtime`

as in the above example, then we can usually replace it with an `InvokeRepeating()` call, which is even simpler to set up and has a slightly smaller overhead cost. The following code is functionally equivalent to the previous implementation that used a Coroutine to regularly invoke a `ProcessAI()` method:

```
| void Start() {  
|     InvokeRepeating("ProcessAI", 0f, _aiProcessDelay);  
| }
```

An important difference between `InvokeRepeating()` and Coroutines is that `InvokeRepeating()` is completely independent of both the `MonoBehaviour` and `GameObject`'s state. The only two ways to stop an `InvokeRepeating()` call is to either call `CancelInvoke()` which stops all `InvokeRepeating()` callbacks initiated by the given `MonoBehaviour` (note that they cannot be canceled individually) or to destroy the associated `MonoBehaviour` or its parent `GameObject`. Disabling either the `MonoBehaviour` or `GameObject` does not stop `InvokeRepeating()`.



A test of 1,000 `InvokeRepeating()` calls was processed in about 2.6 milliseconds; slightly faster than 1,000 equivalent Coroutine `yield` calls, which took 2.9 milliseconds.

That covers most of the useful information related to the `update()` callback. Let's look into some other useful scripting tips.

Faster GameObject null reference checks

It turns out that performing a `null` reference check against a `GameObject` will result in some unnecessary performance overhead. `GameObjects` and `MonoBehaviours` are special objects compared to a typical C# object in that they have two representations in memory: one exists within the memory managed by the same system managing the C# code we write (Managed code), whereas the other exists in a different memory space which is handled separately (Native code). Data can move between these two memory spaces, but each time this happens will result in some additional CPU overhead and possibly an extra memory allocation.

This effect is commonly referred to as crossing the *Native-Managed Bridge*. If this happens, it is likely to generate an additional memory allocation for an object's data to get copied across the Bridge, which will require the Garbage Collector to eventually perform some automatic cleanup of memory for us. This subject will be explored in much more detail in [Chapter 8, Masterful Memory Management](#), but for the time-being just consider that there are many subtle ways to accidentally trigger this extra overhead, and a simple `null` reference check against a `GameObject` is one of them:

```
| if (gameObject != null) {  
|   // do stuff with gameObject  
| }
```

An alternative that generates a functionally equivalent output which operates around twice as quickly (although it does obfuscate the purpose of the code a little) is `System.Object.ReferenceEquals()`:

```
| if (!System.Object.ReferenceEquals(gameObject, null)) {  
|   // do stuff with gameObject  
| }
```

This applies to both `GameObjects` and `MonoBehaviours`, as well as other Unity objects, which have both Native and Managed representations like the `WWW`

class. However, some rudimentary testing reveals that either `null` reference check approach still consumes mere nanoseconds on an Intel Core i5 3570K processor. So, unless you are performing massive amounts of `null` reference checks, the gains might be marginal at best. However, this is a warning worth keeping in mind for the future, as it will come up a lot.

Avoid retrieving string properties from GameObjects

Ordinarily, retrieving a `string` property from an object is the same as retrieving any other reference type property in C#; it should be acquired with no additional memory cost. However, retrieving `string` properties from `GameObject`s is another subtle way of accidentally crossing over the Native-Managed Bridge.

The two properties of `GameObject` affected by this behavior are `tag` and `name`. Therefore, it is unwise to use either properties during gameplay, and you should only use them in performance-inconsequential areas, such as Editor Scripts. However, the Tag System is commonly used for runtime identification of objects, which can make this a significant problem for some teams.

For example, the following code would cause an additional memory allocation during every iteration of the loop:

```
for (int i = 0; i < listOfObjects.Count; ++i) {
    if (listOfObjects[i].tag == "Player") {
        // do something with this object
    }
}
```

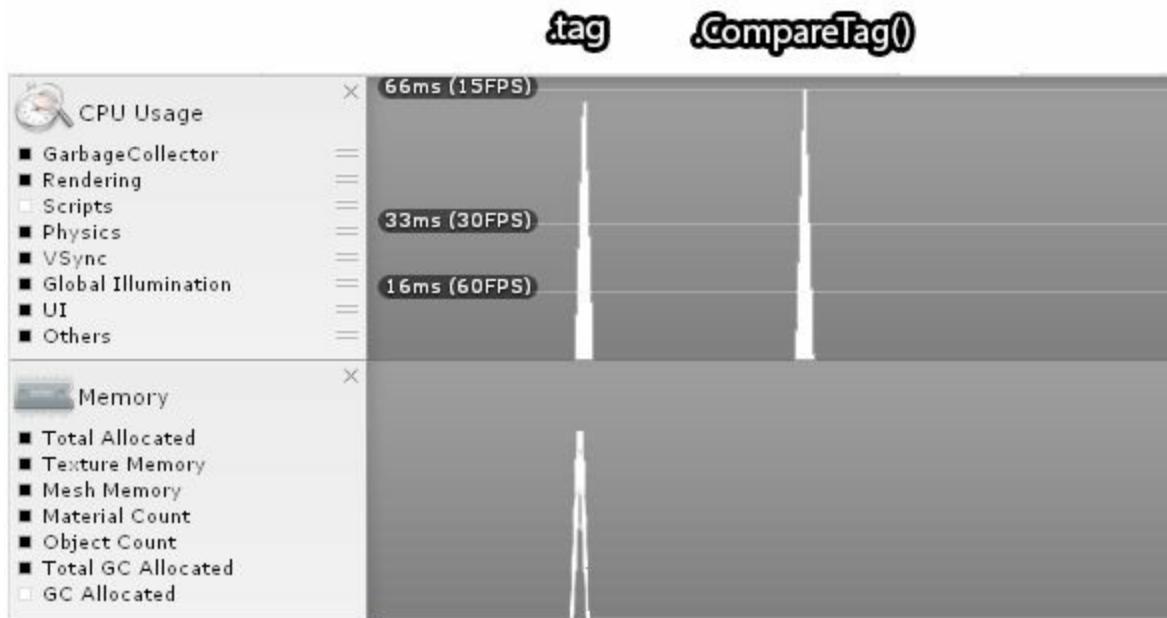
It is often a better practice to identify objects by their Components and class types and to identify values that do not involve `string` objects, but sometimes we're forced into a corner. Maybe we didn't know any better when we started, we inherited someone else's codebase, or we're using it as a workaround for something. Let's assume that, for whatever reason, we're stuck with the Tag system, and we want to avoid the Native-Managed Bridge overhead cost.

Fortunately, the `tag` property is most often used in comparison situations, and `GameObject` provide the `CompareTag()` method which is an alternative way to compare `tag` properties, which avoids the Native-Managed Bridge entirely.

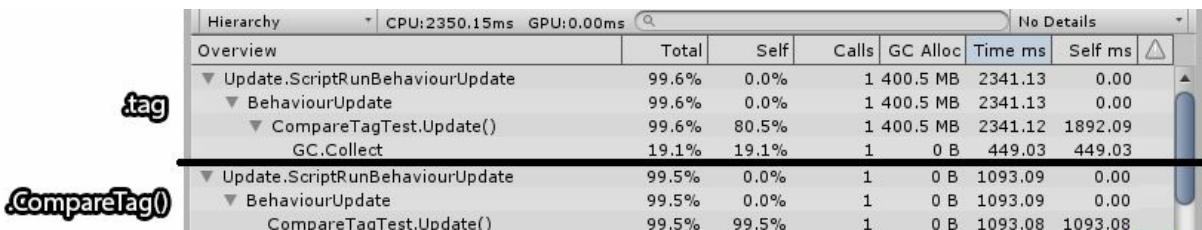
Let's perform a simple test to prove how this simple change can make all the difference in the world:

```
void Update() {  
    int numTests = 10000000;  
  
    if (Input.GetKeyDown(KeyCode.Alpha1)) {  
        for(int i = 0; i < numTests; ++i) {  
            if (gameObject.tag == "Player") {  
                // do stuff  
            }  
        }  
    }  
  
    if (Input.GetKeyDown(KeyCode.Alpha2)) {  
        for(int i = 0; i < numTests; ++i) {  
            if (gameObject.CompareTag ("Player")) {  
                // do stuff  
            }  
        }  
    }  
}
```

We can execute these tests by pressing the 1 and 2 keys to trigger the respective for loop. Here are the results:



Looking at the Breakdown view for each spike, we can see two completely different outcomes:



It's worth noting how the two spikes in the Timeline View appear relatively the same height, and yet one operation took twice as long as the other. The Profiler doesn't have the vertical resolution necessary to generate relatively accurate peaks when we go beyond the 15FPS marker. Both would result in a poor gameplay experience anyway, so the accuracy doesn't really matter.

Retrieving the `tag` property 10 million times (way more than makes sense in reality, but this is useful for comparison) results in about 400 Megabytes of memory being allocated just for `string` objects alone. We can see this memory allocation happening in the spike within the GC Allocated element in the Memory Area of the Timeline View. In addition, this process takes around 2,000 milliseconds to process, where another 400 milliseconds are spent on garbage collection once the `string` objects are no longer needed.

Meanwhile, using `CompareTag()` 10 million times costs around 1,000 milliseconds to process and causes no memory allocations, and hence no garbage collection. This is made apparent from the lack of spike in the GC Allocated element in the Memory Area. This should make it abundantly clear that we must avoid accessing the `name` and `tag` properties whenever possible. If Tag comparison becomes necessary, then we should make use of `CompareTag()`. Unfortunately, there is no equivalent for the `name` property, so we should stick to using Tags where possible.



Note that passing in a string literal, such as "Player", into `CompareTag()` does not result in a runtime memory allocation, since the application allocates hardcoded strings like this during initialization and merely references them at runtime.

Use appropriate data structures

C# offers many different data structures in the `System.Collections` namespace and we shouldn't become too accustomed with using the same ones over and over again. A common performance problem in software development is making use of an inappropriate data structure for the problem we're trying to solve simply because it's convenient. The two most commonly used are perhaps lists (`List<T>`), and dictionaries (`Dictionary<K, V>`).

If we want to iterate through a set of objects then a list is preferred, since it is effectively a dynamic array where the objects and/or references reside next to one another in memory, and therefore iteration causes minimal cache misses. Dictionaries are best used if two objects are associated with one-another and we wish to acquire, insert, or remove these associations quickly. For example, we might associate a level number with a particular Scene file, or an `enum` representing different body parts on a character, with `Collider` Components for those body parts.

However, it's fairly common where we want a data structure that handles both scenarios; we want to quickly figure out which object maps to another, while also being able to iterate through the group. Typically, the developer of this system will use a dictionary, and then iterate over it. However, this process is unfortunately very slow compared to iterating over a list since it must check every potential hash in the dictionary in order to iterate over it fully.

In these cases, it's often better to store data in both a list and a dictionary to better support this behaviour. This will cost additional memory overhead to maintain multiple data structures, and insertion and deletion will require adding and removing objects from both data structures each time, but the benefits on iteration on the list (which tend to happen way more often) will be a stark contrast compared to iterating over a dictionary.

Avoid re-parenting Transforms at runtime

In earlier versions of Unity (version 5.3 and older), the references to `Transform` Components would be laid out in memory in a generally random order. This meant that iteration over multiple `Transforms` was fairly slow due to the likelihood of cache-misses. The upside was that re-parenting a `GameObject` to another one wouldn't really cause a significant performance hit since the `Transforms` operated a lot like a Heap data structure which tend to be relatively fast at insertion and deletion. This behaviour wasn't something we could control, and so we simply lived with it.

However, since Unity 5.4 and beyond, the `Transform` Component's memory layout has changed significantly. Since then, a `Transform` Component's parent-child relationships have operated more like dynamic arrays, whereby Unity attempts to store all `Transforms` that share the same parent sequentially in memory inside a pre-allocated memory buffer and are sorted by their depth in the Hierarchy window beneath the parent. This data structure allows for much, much faster iteration across the entire group, which is particularly beneficial to multiple subsystems like physics and animation. The downside of this change is that if we re-parent a `GameObject` to another one, the parent must fit the new child within its pre-allocated memory buffer as well as sort all of these `Transforms` based on the new depth. Also, if the parent has not pre-allocated enough space to fit the new child, then it must expand its buffer to be able to fit the new child, and all of it's children, in depth-first order. This could take some time to complete for deep and complex `GameObject` structures.

When we instantiate a new `GameObject` through `GameObject.Instantiate()`, one of its arguments is the `Transform` we wish to parent the `GameObject` to, which is `null` by default, and which would place the `Transform` at the root of the Hierarchy window. All `Transforms` at the root of the Hierarchy window need to allocate a buffer to store its current children as well as those that might be added later (child `Transforms` do not need to do this). But, if we then re-parent the `Transform` to another one immediately after instantiation, then it

discards the buffer we just allocated! To avoid this, we should provide the parent `Transform` argument into the `GameObject.Instantiate()` call, which skips this buffer allocation step.

Another way to reduce the costs of this process is to make root `Transform` pre-allocate a larger buffer before we need it so that we don't need to both expand and re-parent another `GameObject` into the buffer in the same frame. This can be accomplished by modifying a `Transform` Component's `hierarchyCapacity` property. If we can estimate the number of child `Transforms` the parent will contain, then we can save a lot of unnecessary memory allocations.

Consider caching Transform changes

The `Transform` Component stores data only relative to its own parent. This means that accessing and modifying a `Transform` Component's position, rotation, and/or scale properties can potentially result in a lot of unanticipated matrix multiplication calculations to generate the correct `Transform` representation for the object through its parent `Transforms`. The deeper the object is in the Hierarchy window, the more calculations are needed to determine the final result.

However, this also means that using `localPosition`, `localRotation`, and `localScale` have a relatively trivial cost associated with them, since these are the values stored directly in the given `Transform` and can be retrieved without any additional matrix multiplication. Therefore, these local property values should be used whenever possible.

Unfortunately, changing our mathematical calculations from world-space to local-space can over-complicate what were originally simple (and solved) problems, so making such changes risks breaking our implementation and introducing a lot of unexpected bugs. Sometimes, it's worth absorbing a minor performance hit in order to solve a complex 3D mathematical problem more easily.

Another problem with constantly changing a `Transform` Component's properties is that it also sends internal notifications to Components like `Collider`, `Rigidbody`, `Light`, and `Camera` which must also be processed since the Physics and Rendering Systems both need to know the new value of the `Transform` and update accordingly.



The speed of these internal notifications were improved dramatically in Unity 5.4 due to the reorganization of `Transforms` in memory in Unity 5.4, as mentioned previously, but we still need to be aware of their costs.

It is not uncommon, during a complex event chain, that we replace a `Transform` Component's properties multiple times in the same frame (although this is probably a warning sign of over-engineered design). This would cause the internal messages to fire each and every time this happens, even if they occur during the same frame, or even the same function call. Ergo, we should consider minimizing the number of times we modify `Transform` properties by caching them in a member variable and committing them only at the end of the frame, as follows:

```
private bool _positionChanged;
private Vector3 _newPosition;

public void SetPosition(Vector3 position) {
    _newPosition = position;
    _positionChanged = true;
}

void FixedUpdate() {
    if (_positionChanged) {
        transform.position = _newPosition;
        _positionChanged = false;
    }
}
```

This code will only commit changes to `position` in the next `FixedUpdate()` method.

Note that changing the `Transform` in this manner does not result in strange-looking behavior or teleporting objects during gameplay. The whole purpose of those internal events are to make sure the Physics and Rendering Systems are always synchronized with the current `Transform` state. Hence, Unity doesn't skip a beat and fires the internal events every time changes come through the `Transform` Component just to be sure nothing gets missed.

Avoid Find() and SendMessage() at runtime

The `SendMessage()` method and family of `GameObject.Find()` methods are notoriously expensive and should be avoided at all costs. The `SendMessage()` method is about 2,000 times slower than a simple function call, and the cost of the `Find()` method scales very poorly with Scene complexity since it must iterate through every `GameObject` in the Scene. It is sometimes forgivable to call `Find()` during initialization of a Scene, such as during an `Awake()` or `Start()` callback. Even in this case, it should only be used to acquire objects that we know for certain already exist in the Scene and for scenes that have only a handful of `GameObjects` in them. Regardless, using either of these methods for interobject communication at runtime is likely to generate a very noticeable overhead and potentially dropped frames.

Relying on `Find()` and `sendMessage()` is typically symptomatic of poor design, inexperience in programming with C# and Unity, or just plain laziness during prototyping. Their usage has become something of an epidemic among beginner-level and intermediate-level projects, so much so that Unity Technologies feels the need to keep reminding users to avoid using them in a real game over and over again in their documentation and at their conferences. They only exist as a less *programmer-y* way to introduce new users to interobject communication, and for some special cases where they can be used in a responsible way (which are few and far between). In other words, they're so ridiculously expensive that they break the rule of not preoptimizing our code, and it's worth going out of our way to avoid using them if our project is going beyond the prototyping stage (which is a distinct possibility since you're reading this book).

To be fair, Unity targets a wide demographic of users, from hobbyists, students, professionals, to those with delusions of grandeur, and also in team sizes from individual developers to hundreds of people on the same team. This results in an incredibly wide range of software development ability. When you're starting out with Unity, it can be difficult to figure out on your

own what you should be doing differently, especially given how the Unity engine does not adhere to the design paradigms of many other game engines we might be familiar with. It has some foreign and quirky concepts surrounding scenes and Prefabs, does not have a built-in *God Class* entry point, nor any obvious raw data storage systems to work with.



A God Class is a fancy name for the first object we might create in our application and whose job would be to create everything else we need based on the current context (what level to load, which subsystems to activate, and so on). These can be particularly useful if we want a single centralized location that controls the order of events as they happen during the entire lifecycle of our application.

This is an important topic not just for performance, but also for any real-time event-driven system design (including, but not limited to games), so it is worth exploring the subject in some detail, evaluating some alternative methods for interobject communication.

Let's start by examining a worst case example, which uses both `Find()` and `SendMessage()` to communicate between objects, and then look into ways to improve upon it.

The following is a class definition for a simple `EnemyManagerComponent` that tracks a list of `GameObjects` representing enemies in our game, and provides a `KillAll()` method to destroy them all when needed:

```
using UnityEngine;
using System.Collections.Generic;

class EnemyManagerComponent : MonoBehaviour {
    List<GameObject> _enemies = new List<GameObject>();

    public void AddEnemy(GameObject enemy) {
        if (!_enemies.Contains(enemy)) {
            _enemies.Add(enemy);
        }
    }

    public void KillAll() {
        for (int i = 0; i < _enemies.Count; ++i) {
            GameObject.Destroy(_enemies[i]);
        }
        _enemies.Clear();
    }
}
```

```
| } }
```

We would then place a `GameObject` in our Scene containing this Component, and name it *EnemyManager*.

The following example method attempts to instantiate a number of enemies from a given Prefab, and then notifies the *EnemyManager* object of their existence:

```
public void CreateEnemies(int numEnemies) {
    for(int i = 0; i < numEnemies; ++i) {
        GameObject enemy = (GameObject)GameObject.Instantiate(_enemyPrefab,
            5.0f * Random.insideUnitSphere,
            Quaternion.identity);
        string[] names = { "Tom", "Dick", "Harry" };
        enemy.name = names[Random.Range(0, names.Length)];
        GameObject enemyManagerObj = GameObject.Find("EnemyManager");
        enemyManagerObj.SendMessage("AddEnemy",
            enemy,
            SendMessageOptions.DontRequireReceiver);
    }
}
```

Initializing data and putting method calls inside any kind of loop, which always output to the same result, is a big red flag for poor performance, and when we're dealing with expensive methods, such as `Find()`, we should always look for ways to call them as few times as possible. Ergo, one improvement we can make is to move the `Find()` call outside of the `for`-loop and cache the result in a local variable so that we don't need to keep reacquiring the *EnemyManager* object over and over again.



Moving the initialization of the `names` variable outside of the `for`-loop is not necessarily critical, since the Compiler is often smart enough to realize it doesn't need to keep reinitializing data that isn't being changed elsewhere. However, it does often make code easier to read.

Another big improvement we can implement, is to optimize our usage of the `SendMessage()` method by replacing it with a `GetComponent()` call. This replaces a very costly method with an equivalent and much cheaper alternative.

This gives us the following result:

```

public void CreateEnemies(int numEnemies) {
    GameObject enemyManagerObj = GameObject.Find("EnemyManager");
    EnemyManagerComponent enemyMgr = enemyManagerObj.GetComponent<EnemyManagerC
    string[] names = { "Tom", "Dick", "Harry" };

    for(int i = 0; i < numEnemies; ++i) {
        GameObject enemy = (GameObject)GameObject.Instantiate(_enemyPrefab,
            5.0f * Random.insideUnitSphere,
            Quaternion.identity);
        enemy.name = names[Random.Range(0, names.Length)];
        enemyMgr.AddEnemy(enemy);
    }
}

```

If this method is called during the initialization of the Scene, and we're not overly concerned with loading time, then we can probably consider ourselves finished with our optimization work.

However, we will often need new objects that are instantiated at runtime to find an existing object to communicate with. In this example, we want new enemy objects to register with our `EnemyManagerComponent` so that it can do whatever it needs to do to track and control the enemy objects in our Scene. We would also like the *EnemyManager* to handle all enemy-related behavior, so that objects calling its functions don't need to perform work on its behalf. This will improve the *Coupling* (how well our codebase separates related behavior) and *Encapsulation* (how well our classes prevent outside changes to the data they manage) of our application. The ultimate aim is to find a reliable and fast way for new objects to find existing objects in the Scene without unnecessary usage of the `Find()` method, so that we can minimize complexity and performance costs.

There are multiple approaches we can take to solving this problem, each with their own benefits and pitfalls:

- Assign references to preexisting objects
- Static Classes
- Singleton Components
- A global Messaging System

Assigning references to pre-existing objects

A simple approach to the problem of interobject communication is to use Unity's built-in Serialization System. Software design purists tend to get a little combative about this feature, since it breaks *Encapsulation*; it makes any field (the C# term for a member variable) marked `private` act in a way that treats them like a `public` field. However, it is a very effective tool for improving development workflow. This is particularly true when artists, designers, and programmers are all tinkering with the same product, where each has wildly varying levels of computer science and software programming knowledge, and with some of whom would prefer to stay away from modifying code files. Sometimes, it's worth bending a few rules in the name of productivity.

Whenever we create a `public` field in a `MonoBehaviour`, Unity automatically serializes and exposes the value in the Inspector window when the Component is selected. However, `public` fields are always dangerous from a software design perspective. These variables can be changed through code at anytime from anywhere, making it hard to keep track of the variable and is liable to introduce a lot of unexpected bugs.

A better solution, is to take any `private` or `protected` member variable of a class and expose it to the Inspector window with the `[SerializeField]` attribute. The value will now behave like a `public` field with respect to the Inspector window, allowing us to change it through the Editor interface for convenience, but keeps the data safely encapsulated from other parts of our codebase.

For example, the following class exposes three `private` fields to the Inspector window:

```
using UnityEngine;  
public class EnemyCreatorComponent : MonoBehaviour {
```

```

[SerializeField] private int _numEnemies;
[SerializeField] private GameObject _enemyPrefab;
[SerializeField] private EnemyManagerComponent _enemyManager;

void Start() {
    for (int i = 0; i < _numEnemies; ++i) {
        CreateEnemy();
    }
}

public void CreateEnemy() {
    _enemyManager.CreateEnemy(_enemyPrefab);
}

```



Note that the `private` access specifiers shown in the preceding code are redundant keywords in C#, since fields and methods default to `private` unless specified otherwise. However, it is often best practice to be explicit about the intended access level.

Looking at this Component in the Inspector window reveals three values, initially given default values of `0`, or `null`, which can be modified through the Editor interface:



We can drag-and-drop a Prefab reference from the Project window into the Enemy Prefab field revealed in the Inspector window.



Note how Unity automatically takes a camel-cased field name and creates a convenient Inspector window name for it. `_numEnemies` becomes `Num Enemies`, `_enemyPrefab` becomes `Enemy Prefab` and so on.

Meanwhile, the `_enemyManager` field is interesting because it is a reference to a specific `MonoBehaviour` class type. If a `GameObject` is dragged-and-dropped into this field, then it will refer to the Component on the given object as opposed to the `GameObject` itself. Note that if the `GameObject` does not contain the expected `MonoBehaviour`, then nothing will be assigned to the field.



A common usage of this Component reference technique is to obtain references to other Components attached to the very same GameObject a Component is attached to. This is an alternative means of caching Components with zero cost, as discussed in the section entitled Cache Component references earlier in this chapter.

There is some danger to using this method. Much of our code would assume that a Prefab is assigned to a field that is used like a Prefab and a `GameObject` is assigned to a field that refers to an instance of a `GameObject`. However, since Prefabs are essentially `GameObject`s, any Prefab or `GameObject` can be assigned to a serialized `GameObject` reference field, which means we could assign the wrong type by accident.

If we do assign the wrong type then we could accidentally instantiate a new `GameObject` from an existing `GameObject` that was previously modified, or we could make changes to a Prefab, which would then change the state of all `GameObject`s instantiated from it. To make matters worse, any accidental changes to a Prefab become permanent since Prefabs occupy the same memory space whether Play Mode is active or not. This is the case even if the Prefab is only modified during Play Mode.

Therefore, this approach is a very team-friendly way of solving the problem of interobject communication, but it is not ideal due to all of the risks involved with team members accidentally leaving `null` references in place, assigning Prefabs to references that expect an instance of a `GameObject` from the Scene, or vice versa.

It is also important to note that not all objects can be serialized and revealed in the Inspector window. Unity can serialize all primitive data types (`int`, `float`, `string`, and `bool`), various built-in types (`vector3`, `Quaternion`, and so on); `enum`, `class`, `struct`, and various data structures containing other serializable types such as `List`. However, it is unable to serialize `static` fields, `readonly` fields, properties, and dictionaries.



Some Unity developers like to implement pseudo-serialization of dictionaries via two separate lists for keys and values, along with a Custom Editor script, or via a single list of struct objects,

TIP which contain both keys and values. Both of these solutions are a little clumsy, and rarely as performant as a proper dictionary, but they can still be useful.

Another solution to the problem of interobject communication is to try and make use of globally accessible objects in order to minimize the number of custom assignments we need to make.

Static Classes

This approach involves creating a class that is globally accessible to the entire codebase at any time. Any kind of global *manager* class is often frowned upon in software engineering circles, partly since the name *manager* is vague and doesn't say much about what it's meant to do, but mostly because problems can be difficult to debug. Changes can occur from anywhere and at any point during runtime, and such classes tend to maintain state information that other systems rely upon. It is also perhaps the most difficult approach to change or replace, since many of our classes might contain direct function calls into it, requiring each to be modified at a future date if it were to be replaced. Despite all of these drawbacks, it is by far the easiest solution to understand and implement.

The Singleton design pattern is a common way of ensuring only one instance of a certain object type ever exists in memory. This design pattern is implemented by giving the class a `private` constructor, a `static` variable is maintained to keep track of the object instance, and the class can only be accessed through a `static` property it provides. Singletons can be useful for managing shared resources or *heavy data traffic*, such as file access, downloads, data parsing, and messaging. A Singleton ensures that we have a single entry point for such activities, rather than having tons of different subsystems competing for shared resources and potentially bottlenecking one another.

Singletons don't necessarily have to be globally accessible objects--their most important feature is that only a single instance of the object exists at a time. However, the way that Singletons are primarily used in most projects is to be a global access point to some shared functionality, and are designed to be created once during application initialization, persist through the entire lifecycle of the application, and only destroyed during application shutdown. As such, a simpler way of implementing this kind of behavior in C# is to use a Static Class. In other words, implementing the typical Singleton design pattern in C# just provides the same behavior as a Static Class, but takes more time and code to implement.

A Static Class that functions in much the same way as the `EnemyManagerComponent` as demonstrated in the previous example, can be defined as follows:

```
using System.Collections.Generic;
using UnityEngine;

public static class StaticEnemyManager {
    private static List<Enemy> _enemies;

    public static void CreateEnemy(GameObject prefab) {
        string[] names = { "Tom", "Dick", "Harry" };
        GameObject enemy = GameObject.Instantiate(prefab, 5.0f * Random.insideUnitSphere, Quaternion.identity);
        Enemy enemyComp = enemy.GetComponent<Enemy>();
        enemy.gameObject.name = names[Random.Range(0, names.Length)];
        _enemies.Add(enemyComp);
    }

    public static void KillAll() {
        for (int i = 0; i < _enemies.Count; ++i) {
            _enemies[i].Die();
            GameObject.Destroy(_enemies[i].gameObject);
        }
        _enemies.Clear();
    }
}
```

Note that every method, property, and field in a Static Class must have the `static` keyword attached, which implies that only one instance of this object will ever reside in memory. This also means that its `public` methods and fields are accessible from anywhere. Static Classes, by definition, do not allow any non-static fields to be defined.

If a Static Class' fields need to be initialized (such as the `_enemies` field, which is initially set to `null`), then Static Class fields can be initialized inline like so:

```
| private static List<Enemy> _enemies = new List<Enemy>();
```

However, if object construction needs to be more complicated than this, then Static Classes can be given a `static` constructor, instead. The Static Class' constructor is automatically called the moment the class is first accessed through any it's fields, properties, or methods, and can be defined like so:

```
static StaticEnemyManager() {
    _enemies = new List<Enemy>();
    // more complicated initialization activity goes here
}
```

This time we have implemented the `CreateEnemy()` method so that it handles much of the activity for creating an enemy object. However, the Static Class must still be given a reference to a Prefab from which it can instantiate an enemy object from. A Static Class can only contain static member variables, and therefore cannot easily interface with the Inspector window in the same way that `MonoBehaviours` can, therefore requiring the caller to provide some implementation-specific information to it. To solve this problem, we could implement a companion-Component for our Static Class to keep our code properly *Decoupled*. The following code demonstrates what this class might look like:

```
using UnityEngine;

public class EnemyCreatorCompanionComponent : MonoBehaviour {
    [SerializeField] private GameObject _enemyPrefab;

    public void CreateEnemy() {
        StaticEnemyManager.CreateEnemy(_enemyPrefab);
    }
}
```

Despite these drawbacks, the `StaticEnemyManager` class illustrates a simple example of how a Static Class might be used to provide information or communication between external objects, providing a better alternative than using `Find()` or `SendMessage()`.

Singleton Components

As mentioned previously, Static Classes have difficulty interfacing with Unity-related functionality, and cannot directly make use of `MonoBehaviour` features, such as event callbacks, Coroutines, hierarchical design, and Prefabs. Also, since there's no object to select in the Inspector window, we lose the ability to inspect the a Static Class' data at runtime through the Inspector window, which can make debugging difficult. These are features that we may wish to make use of in our global classes.

A common solution to this problem is to implement a Component that acts like a Singleton--it provides `static` methods to grant global access, and only one instance of the `MonoBehaviour` is ever allowed to exist at any given time.

The following is the definition for a `singletonComponent` class:

```
using UnityEngine;

public class SingletonComponent<T> : MonoBehaviour where T : SingletonComponent<T>
{
    private static T __Instance;

    protected static SingletonComponent<T> _Instance {
        get {
            if(!__Instance) {
                T [] managers = GameObject.FindObjectsOfType(typeof(T)) as T[];
                if (managers != null) {
                    if (managers.Length == 1) {
                        __Instance = managers[0];
                        return __Instance;
                    } else if (managers.Length > 1) {
                        Debug.LogError("You have more than one " +
                            typeof(T).Name +
                            " in the Scene. You only need " +
                            "one - it's a singleton!");
                        for(int i = 0; i < managers.Length; ++i) {
                            T manager = managers[i];
                            Destroy(manager.gameObject);
                        }
                    }
                }
                GameObject go = new GameObject(typeof(T).Name, typeof(T));
                __Instance = go.GetComponent<T>();
                DontDestroyOnLoad(__Instance.gameObject);
            }
        return __Instance;
    }
}
```

```

        }
        set {
            _Instance = value as T;
        }
    }
}

```

This class works by creating a `GameObject` containing a Component of itself the first time it is accessed. Since we wish this to be a global and persistent object, we will need to call `DontDestroyOnLoad()` shortly after the `GameObject` is created. This is a special function that tells Unity that we want the object to persist between scenes for as long as the application is running. From that point onward, when a new Scene is loaded, the object will not be destroyed and will retain all of its data.

This class definition assumes two things. Firstly, because it is using *generics* to define its behavior, it must be derived from in order to create a concrete class. Secondly, a method must be defined to assign the `_Instance` property (which in turns sets the private `_instance` field) and cast it to/from the correct class type.

For example, the following is the minimum amount of code that is needed to successfully generate a new `SingletonComponent` derived class called `EnemyManagerSingletonComponent`:

```

public class EnemyManagerSingletonComponent : SingletonComponent< EnemyManager
{
    public static EnemyManagerSingletonComponent Instance {
        get { return ((EnemyManagerSingletonComponent)_Instance); }
        set { _Instance = value; }
    }

    public void CreateEnemy(GameObject prefab) {
        // same as StaticEnemyManager
    }

    public void KillAll() {
        // same as StaticEnemyManager
    }
}

```

This class can be used at runtime by having any other object access the `Instance` property at any time. If the Component does not already exist in our Scene, then the `SingletonComponent` base class will instantiate its own `GameObject` and attach an instance of the derived class to it as a Component. From that point forward, access through the `Instance` property will reference the

Component that was created, and only one instance of this Component will exist at a time.

Note that this means we don't need to implement `static` methods in a Singleton Component class definition. For example, we could simply call `EnemyManagerSingletonComponent.Instance.KillAll()` to access the `KillAll()` method.

Note that it is possible to place an instance of a `SingletonComponent` in a Hierarchy window since it derives from `MonoBehaviour`. Although, be warned that the `DontDestroyOnLoad()` method would never be called, which would prevent the Singleton Component's `GameObject` from persisting when the next Scene is loaded. We will perhaps need to call `DontDestroyOnLoad()` in the derived class' `Awake()` callback to make this work, unless, of course, we actually want destructible Singletons. Sometimes it makes sense to allow such Singletons to be destroyed between scenes so that it can start fresh each time; it all depends on our particular use cases.

In either case, shutdown of a Singleton Component can be a little convoluted because of how Unity tears down scenes. An object's `OnDestroy()` callback is called whenever it is destroyed during runtime. The same method is called during application shutdown, whereby every Component on each `GameObject` has its `OnDestroy()` callback called by Unity. The same activities take place when we end Play Mode in the Editor, thus returning to Edit Mode. However, destruction of objects occurs in a random order, and we cannot assume that the `SingletonComponent` object will be the last object destroyed.

Consequently, if any object attempts to do anything with the Singleton Component in the middle of their `OnDestroy()` callback, then they may be calling the `SingletonComponent` object's `Instance` property. However, if the Singleton Component has already been destroyed prior to this moment, then a new instance of the `SingletonComponent` will be created in the middle of application shutdown. This can corrupt our Scene files, as instances of our Singleton Components will be left behind in the Scene. If this happens, then Unity will throw the following error message:

"Some objects were not cleaned up when closing the scene. (Did you spawn new GameObjects from OnDestroy?)"

The obvious workaround is to simply never call into a `SingletonComponent` object during any `MonoBehaviour` Component's `OnDestroy()` callback. However, there are some legitimate reasons we may wish to do so: most notable is that Singletons are often designed to make use of the Observer design pattern. This design pattern allows other objects to register/deregister with them for certain tasks, similar to how Unity latches onto callback methods, such as `start()` and `update()`, but in a more strict fashion.

With the Observer design pattern, objects will typically register with the system when they are created, will make use of it during runtime, and then either deregister from it during runtime when they are finished using it or deregister during their own shutdown for the sake of cleanup. We will see an example of this design pattern in the upcoming section, *A global Messaging System*, but if we imagine a `MonoBehaviour` making use of such a system, then the most convenient place to perform shutdown deregistration would be within an `OnDestroy()` callback. Consequently, such objects are likely to run into the aforementioned problem, where a new `GameObject` for the `SingletonComponent` is accidentally created during application shutdown.

To solve this problem, we will need to make three changes. Firstly, we need to add an additional flag to the `SingletonComponent`, which keeps track of its active state and disable it at the appropriate times. This includes the Singleton's own destruction, as well as application shutdown (`OnApplicationQuit()` is another useful Unity callback for `MonoBehaviours`, which is called during this time):

```
| private bool _alive = true;
| void OnDestroy() { _alive = false; }
| void OnApplicationQuit() { _alive = false; }
```

Secondly, we should implement a way for external objects to verify the Singleton's current state:

```
| public static bool IsAlive {
|     get {
|         if (_Instance == null)
|             return false;
|         return _Instance._alive;
|     }
| }
```

Finally, any object that attempts to call into the Singleton during its own `OnDestroy()` method must first verify the state using the `IsAlive` property before calling instance, as follows:

```
public class SomeComponent : MonoBehaviour {
    void OnDestroy() {
        if (MySingletonComponent.IsAlive) {
            MySingletonComponent.Instance.SomeMethod();
        }
    }
}
```

This will ensure that nobody attempts to access the Singleton instance during destruction. If we don't follow this rule, then we will run into problems where instances of our Singleton object will be left behind in the Scene after returning to Edit Mode.

The irony of the `SingletonComponent` approach is that we are using a `Find()` call to determine whether or not one of these `SingletonComponent` objects already exists in the Scene before we attempt to assign the `_instance` reference variable. Fortunately, this will only happen when the Singleton Component is first accessed, which is usually not a problem if there aren't too many `GameObjects` in the Scene, but it's possible that the initialization of the Singleton Component may not necessarily occur during Scene initialization and can, therefore, cost us a performance spike at a bad moment during gameplay when an instance is first acquired and `Find()` gets called. The workaround for this is to have some God Class confirm that the important Singletons are instantiated during Scene initialization by simply accessing the `Instance` property on each one.

Another downside to this approach is that if we later decide that we want more than one of these Singletons executing at once or we wish to separate out its behavior to be more modular, then there would be a lot of code that needs to change.

The final approach we will explore will attempt to solve many of the problems revealed by the previous solutions and provide a way to gain all of their benefits, by combining ease of implementation, ease of extension, and strict usage that also reduces the likelihood of human-error during configuration.

A global Messaging System

The final suggested approach to solve the problem of interobject communication is to implement a global Messaging System that any object can access and send messages through to any object that may be interested in listening to that specific type of message. Objects can send messages or listen for them (sometimes both!), and the responsibility is on the listener to decide what messages they are interested in. The message sender can broadcast the message without caring at all who is listening, and a message can be sent through the system regardless of the specific contents of the message. This approach is by far the most complex and may require some effort to implement and maintain, but it is an excellent long-term solution to keep our object communication modular, decoupled, and fast as our application gets more and more complex.

The kinds of message we wish to send can take many forms, including data values, references, instructions for listeners, and more, but they should all have a common, basic definition that our Messaging System can use to determine what the message is and who it is intended for.

The following is a simple class definition for a `Message` object:

```
public class Message {  
    public string type;  
    public Message() { type = this.GetType().Name; }  
}
```

The `Message` class' constructor caches the message's type in a local `string` property to be used later for cataloging and distribution purposes. Caching this value is important, as each call to `GetType().Name` will result in a new string being allocated, and we've previously learned that we want to minimize this activity as much as possible.

Any custom messages can contain whatever superfluous data they wish so long as they derive from this base class, which will allow it to be sent through our Messaging System. Take note that despite acquiring the `type` from the object during its base class constructor, the `name` property will still contain the

name of the derived class, not the base class.

Moving on to our `MessagingSystem` class, we should define its features by what kind of requirements we need it to fulfill:

- It should be globally accessible
- Any object (`MonoBehaviour` or not) should be able to register/deregister as listeners to receive specific message types (that is, the Observer design pattern)
- Registering objects should provide a method to call when the given message is broadcasted from elsewhere
- The system should send the message to all listeners within a reasonable time frame, but not choke on too many requests at once

A globally accessible object

The first requirement makes the Messaging System an excellent candidate for a Singleton object, since we would only ever need one instance of the system. Although, it is wise to think long and hard if this is truly the case before committing to implementing a Singleton.

If we later decide that we want multiple instances of this object to exist, wish to allow the systems to be created/destroyed during runtime, or even wish to create test cases that allow us to fake or create/destroy them in the middle of a test, then it can be difficult task to refactor a Singleton out of our codebase. This is due to all of the dependencies we will gradually introduce to our code as we use the system more and more.

If we wish to avoid Singletons due to the above drawbacks, then it may be easier to create a single instance of the Messaging System during initialization and then pass it around from subsystem to subsystem as needed, or we might wish to go further and explore the concept of Dependency Injection, which attempts to solve problems like these. However, for the sake of simplicity, we will assume that a Singleton fits our needs and design our `MessagingSystem` class accordingly.

Registration

The second and third requirements can be achieved by offering some public methods that allow registration with the Messaging System. If we force the listening object to provide us a delegate function to call when the message is broadcast, then this allows listeners to customize which method is called for which message. We can make our codebase very easy to understand if we name the delegate after the message it is intended to process.

In some cases, we might wish to broadcast a general notification message and have all listeners do something in response, such as an *Enemy Created* message. Other times, we might be sending a message that specifically targets a single listener among a group. For example, we might want to send an *Enemy Health Value Changed* message that is intended for a specific *Health Bar* object that is attached to the enemy that was damaged. However, we may have many health bar objects in the Scene, all of which are interested in this message type, but each is only interested in hearing health update messages for the enemy they're providing health information for. So, if we implement a way for the system to stop checking after it has been handled, then we can probably save a good number of CPU cycles when there are many listeners waiting for the same message type.

The delegate we define should, therefore, provide a way to retrieve the message via an argument and return a response that determines whether or not processing for the message should stop if and when the listener is done with it. The decision on whether to stop processing or not can be achieved by returning a simple Boolean, where `true` implies that this listener has handled the message and the processing for the message must stop, and `false` implies that this listener has not handled the message and the Messaging System should try the next listener.

Here is the definition for the delegate:

```
| public delegate bool MessageHandlerDelegate(Message message);
```

Listeners must define a method of this form and pass a delegate reference to

the Messaging System during registration, thus providing a means for the Messaging System to tell the listening object when the message is being broadcast.

Message processing

The final requirement for our Messaging System is that this object should have some kind of timing-based mechanism built in to prevent it from choking on too many messages at once. This means that, somewhere in the codebase, we will need to make use of `MonoBehaviour` event callbacks to tell our Messaging System to perform work during Unity's `Update()`, essentially enabling it to count time.

This could be achieved with the Static Class Singleton (which we defined earlier), which would require some other `MonoBehaviour`-based God Class to call into it, informing it that the Scene has been updated. Alternatively, we can use the Singleton Component to achieve the same thing, which has its own means of determining when `Update()` is called, and hence handle its workload independently of any God Class. The most notable difference between the two approaches is whether or not the system is dependent on the control of other objects and the various pros and cons of managing a Singleton Component (such that it won't get destroyed between scenes; we don't want to accidentally recreate it during shutdown).

The Singleton Component approach is probably the best, since there aren't too many occasions where we wouldn't want this system acting independently, even if much of our game logic depends upon it. For example, even if the game was paused, we wouldn't want the game logic to pause our Messaging System. We still want the Messaging System to continue receiving and processing messages so that we can, for example, keep UI-related Components communicating with one another while the gameplay is in a paused state.

Implementing the Messaging System

Let's define our Messaging System by deriving from the `SingletonComponent` class and provide a method for objects to register with it:

```
using System.Collections.Generic;
using UnityEngine;

public class MessagingSystem : SingletonComponent<MessagingSystem> {
    public static MessagingSystem Instance {
        get { return ((MessagingSystem)_Instance); }
        set { _Instance = value; }
    }

    private Dictionary<string, List<MessageHandlerDelegate>> _listenerDict = new Dictionary<string, List<MessageHandlerDelegate>>();

    public bool AttachListener(System.Type type, MessageHandlerDelegate handler) {
        if (type == null) {
            Debug.Log("MessagingSystem: AttachListener failed due to having no " +
                      "message type specified");
            return false;
        }

        string msgType = type.Name;
        if (!_listenerDict.ContainsKey(msgType)) {
            _listenerDict.Add(msgType, new List<MessageHandlerDelegate>());
        }

        List<MessageHandlerDelegate> listenerList = _listenerDict[msgType];
        if (listenerList.Contains(handler)) {
            return false; // listener already in list
        }

        listenerList.Add(handler);
        return true;
    }
}
```

The `_listenerDict` field is a dictionary of strings mapped to lists containing `MessageHandlerDelegate`. This dictionary organizes our listener delegates into lists by which message type they wish to listen to. Thus, if we know what message type is being sent, then we can quickly retrieve a list of all delegates that have been registered for that message type. We can then iterate through the list, querying each listener to check whether one of them

wants to handle it.

The `AttachListener()` method requires two parameters: a message type in the form of its `System.Type` and a `MessageHandlerDelegate` to send the message to when the given message type comes through the system.

Message queuing and processing

In order to process messages, our Messaging System should maintain a queue of incoming message objects so that we can process them in the order they were broadcast:

```
private Queue<Message> _messageQueue = new Queue<Message>();  
  
public bool QueueMessage(Message msg) {  
    if (!listenerDict.ContainsKey(msg.type)) {  
        return false;  
    }  
    _messageQueue.Enqueue(msg);  
    return true;  
}
```

The `QueueMessage()` method simply checks whether the given message type is present in our dictionary before adding it to the queue. This effectively tests whether or not an object actually cares to listen to the message before we queue it to be processed later. We have introduced a new `private` field, `_messageQueue`, for this purpose.

Next, we'll add a definition for `Update()`. This callback will be called regularly by the Unity Engine. Its purpose is to iterate through the current contents of the message queue, one message a time, verify whether or not too much time has passed since we began processing, and if not, pass them along to the next stage in the process:

```
private const int _maxQueueProcessingTime = 16667;  
private System.Diagnostics.Stopwatch timer = new System.Diagnostics.Stopwatch  
  
void Update() {  
    timer.Start();  
    while (_messageQueue.Count > 0) {  
        if (_maxQueueProcessingTime > 0.0f) {  
            if (timer.ElapsedMilliseconds > _maxQueueProcessingTime) {  
                timer.Stop();  
                return;  
            }  
        }  
    }  
}
```

```

        Message msg = _messageQueue.Dequeue();
        if (!TriggerMessage(msg)) {
            Debug.Log("Error when processing message: " + msg.type);
        }
    }
}

```

The time-based safeguard is in place to make sure that it does not exceed a processing time limit threshold. This prevents the Messaging System from freezing our game if too many messages get pushed through the system too quickly. If the total time limit is exceeded, then all message processing will stop, leaving any remaining messages to be processed during the next frame.



Note that we use the full namespace when creating the stopwatch object. We could have added using System.Diagnostics, but this would lead to a namespace conflict between

System.Diagnostics.Debug and UnityEngine.Debug. Omitting it allows us to continue to call Unity's debug logger with Debug.Log(), without having to explicitly call UnityEngine.Debug.Log() each time.

Lastly, we will need to define the `TriggerMessage()` method, which distributes messages to listeners:

```

public bool TriggerMessage(Message msg) {
    string msgType = msg.type;
    if (!_listenerDict.ContainsKey(msgType)) {
        Debug.Log("MessagingSystem: Message '" + msgType + "' has no listeners!");
        return false; // no listeners for message so ignore it
    }

    List<MessageHandlerDelegate> listenerList = _listenerDict[msgType];

    for(int i = 0; i < listenerList.Count; ++i) {
        if (listenerList[i](msg))
            return true; // message consumed by the delegate
    }
    return true;
}

```

The preceding method is the main workhorse behind the Messaging System. The `TriggerEvent()` method's purpose is to obtain the list of listeners for the given message type and give each of them an opportunity to process it. If one of the delegates returns `true`, then processing of the current message ceases

and the method exits, allowing the `Update()` method to process the next message.

Normally, we would want to use `QueueEvent()` to broadcast messages, but we also provide direct access to `TriggerEvent()` as an alternative. Using `TriggerEvent()` directly allows message senders to force their messages to be processed immediately without waiting for the next `Update()` event. This bypasses the throttling mechanism, which might be necessary for messages that need to be sent during critical moments of gameplay, where waiting an additional frame might result in a strange-looking behavior.

For example, if we intend for two objects to be destroyed and create a Particle Effect the moment they collide with one another, and this work is handled by another subsystem (hence an event needs to be sent for it), then we would want to send the message via `TriggerEvent()` to prevent the objects from continuing to exist for one frame before the event is handled.

Conversely, if we wanted to do something less frame-critical, such as create a pop-up message when the player walks into a new area, we could safely use a `QueueEvent()` call to handle it.

Try to avoid habitually using `TriggerEvent()` for all events, as we could end up handling too many calls simultaneously in the same frame, causing a sudden drop in frame rate. Decide which events are frame-critical, and which are not, and use the `QueueEvent()` and `TriggerEvent()` methods appropriately.

Implementing custom messages

We've created the Messaging System, but an example of how to use it would help us wrap our heads around the concept. Let's start by defining a pair of simple classes that derive from `Message`, which we can use to create a new enemy, as well as notify other parts of our codebase that an enemy was created:

```
public class CreateEnemyMessage : Message {}

public class EnemyCreatedMessage : Message {
    public readonly GameObject enemyObject;
    public readonly string enemyName;

    public EnemyCreatedMessage(GameObject enemyObject, string enemyName) {
        this.enemyObject = enemyObject;
        this.enemyName = enemyName;
    }
}
```

`CreateEnemyMessage` is the simplest form of message that contains no special data, while `EnemyCreatedMessage` will contain a reference to the enemy's `GameObject` as well as its name. Good practice for message objects is to make their member variables not only `public`, but also `readonly`. This ensures that the data is easily accessible but cannot be changed after the object's construction. This safeguards the content of our messages against being altered, as they're passed between one listener and another.

Message sending

To send one of these message objects, we simply need to call either `QueueEvent()` or `TriggerEvent()` and pass it an instance of the message we wish to send. The following code demonstrates how we would broadcast a `CreateEnemyMessage` object when the *Space Bar* is pressed:

```
public class EnemyCreatorComponent : MonoBehaviour {
    void Update() {
        if (Input.GetKeyDown(KeyCode.Space)) {
            MessagingSystem.Instance.QueueMessage(new CreateEnemyMessage());
        }
    }
}
```

If we were to test this code right now, nothing would happen, because even though we are sending a message through the Messaging System, there are no listeners for this message type. Let's cover how to register listeners with the Messaging System.

Message registration

The following code contains a pair of simple classes that register with the Messaging System, each requesting to have one of their methods called whenever certain types of messages have been broadcast from anywhere in our codebase:

```
public class EnemyManagerWithMessagesComponent : MonoBehaviour {
    private List<GameObject> _enemies = new List<GameObject>();
    [SerializeField] private GameObject _enemyPrefab;

    void Start() {
        MessagingSystem.Instance.AttachListener(typeof(CreateEnemyMessage),
                                                this.HandleCreateEnemy);
    }

    bool HandleCreateEnemy(Message msg) {
        CreateEnemyMessage castMsg = msg as CreateEnemyMessage;
        string[] names = { "Tom", "Dick", "Harry" };
        GameObject enemy = GameObject.Instantiate(_enemyPrefab,
                                                   5.0f * Random.insideUnitSphere,
                                                   Quaternion.identity);
        string enemyName = names[Random.Range(0, names.Length)];
        enemy.gameObject.name = enemyName;
        _enemies.Add(enemy);
        MessagingSystem.Instance.QueueMessage(new EnemyCreatedMessage(enemy,
                                                                     enemyName));
        return true;
    }
}

public class EnemyCreatedListenerComponent : MonoBehaviour {
    void Start () {
        MessagingSystem.Instance.AttachListener(typeof(EnemyCreatedMessage),
                                                HandleEnemyCreated);
    }

    bool HandleEnemyCreated(Message msg) {
        EnemyCreatedMessage castMsg = msg as EnemyCreatedMessage;
        Debug.Log(string.Format("A new enemy was created! {0}",
                               castMsg.enemyName));
        return true;
    }
}
```

During initialization, the `EnemyManagerWithMessagesComponent` class registers to receive messages of the type `CreateEnemyMessage`, and will process the message

through its `HandleCreateEnemy()` delegate. During this method, it can typecast the message into the appropriate derived message type and resolves the message in its own unique way. Other classes can register for the same message and resolve it differently through its own custom delegate method (assuming that an earlier listener didn't return `true` from its own delegate).

We know what type of message will be provided by the `msg` argument of the `HandleCreateEnemy()` method, because we defined it during registration through the `AttachListener()` call. Due to this, we can be certain that our typecasting is safe, and we can save time by not having to do a `null` reference check, although, technically, there is nothing stopping us using the same delegate to handle multiple message types. In these cases, though, we will need to implement a way to determine which message object is being passed and treat it accordingly. However, the best approach is to define a unique method for each message type in order to keep things appropriately decoupled. There really is little benefit in trying to use one monolithic method to handle all message types.

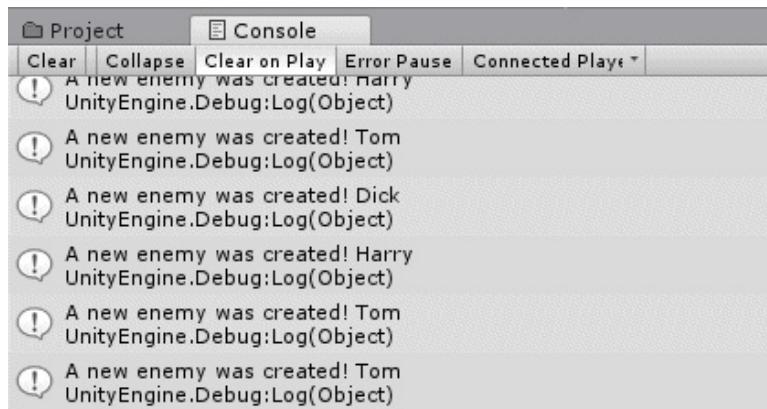
Note how the `HandleEnemyCreated()` method definition matches the function signature of `MessageHandlerDelegate` (that is, it has the same return type and argument list), and that it is being referenced in the `AttachListener()` call. This is how we tell the Messaging System what method to call when the given message type is broadcast, and how delegates ensure type-safety. If the function signature had a different return value or a different list of arguments, then it would be an invalid delegate for the `AttachListener()` method, and we would get compiler errors. Also, note that `HandleEnemyCreated()` is also a `private` method, and yet our `MessagingSystem` class is able to call it. This is a useful feature of delegates in that we can allow only systems we give permission to call this message handler. Exposing the method publicly might lead to some confusion in our code's API, and developers may think that they're meant to call the method directly, which is not its intended use.

The beautiful part is that we're free to give the delegate method whatever name we want. The most sensible approach is to name the method after the message which it handles. This makes it clear to anyone reading our code what the method is used for and what message object type must be sent in order to call it. This makes future parsing and debugging of our code much

more straight-forward since we can follow the chain of events by the matching names of the messages and their handler delegates.

During the `HandleCreateEnemy()` method, we also queue another event, which broadcasts an `EnemyCreatedMessage` instead. The second class, `EnemyCreatedListenerComponent` registers to receive these messages, and then prints out a message containing that information. This is how we would implement a way for subsystems to notify other subsystems of changes. In a real application, we might register a UI system to listen for these types of messages, and update a counter on the screen to show how many enemies are now active. In this case, the enemy management and UI systems are appropriately *decoupled* such that neither needs to know any specific information about how the other operates in order to do their assigned tasks.

If we now add the `EnemyManagerWithMessagesComponent`, `EnemyCreatorComponent` and `EnemyCreatedListenerComponent` to our Scene, and press the *Space Bar* several times, we should see log messages appear in the Console window, informing us of a successful test:



Note that a `MessagingSystem` Singleton object will be created during Scene initialization, when either the `EnemyManagerWithMessagesComponent` object's or `EnemyCreatedListenerComponent` object's `start()` methods are called (whichever happens first), since that is when they register their delegates with the Messaging System, which accesses the `Instance` property, and hence creates the necessary `GameObject` containing the Singleton Component. No additional effort is required on our part to create the `MessagingSystem` object.

Message cleanup

Since message objects are classes, they will be created dynamically in memory and will be disposed of shortly afterward when the message has been processed and distributed among all listeners. However, as you will learn in [chapter 8, Masterful Memory Management](#), this will eventually result in a garbage collection, as memory accumulates over time. If our application runs for long enough, it will eventually result in the occasional garbage collection, which is the most common cause of unexpected and sudden CPU performance spikes in Unity applications. Therefore, it is wise to use the Messaging System sparingly and avoid spamming messages too frequently on every update.

The more important clean-up operation to consider is deregistration of delegates if an object needs to be destroyed. If we don't handle this properly, then the Messaging System will hang on to delegate references that prevent objects from being fully destroyed and freed from memory.

Essentially, we will need to pair every `AttachListener()` call with an appropriate `DetachListener()` call when the object is destroyed, disabled, or we otherwise decide that we no longer need it to be queried when messages are being sent.

The following method definition in the `MessagingSystem` class will detach a listener for a specific event:

```
public bool DetachListener(System.Type type, MessageHandlerDelegate handler)
{
    if (type == null) {
        Debug.Log("MessagingSystem: DetachListener failed due to having no " +
                  "message type specified");
        return false;
    }

    string msgType = type.Name;

    if (!listenerDict.ContainsKey(type.Name)) {
        return false;
    }

    List<MessageHandlerDelegate> listenerList = _listenerDict[msgType];
```

```
| if (!listenerList.Contains (handler)) {  
|     return false;  
| }  
| listenerList.Remove(handler);  
| return true;  
| }
```

Here is an example usage of the `DetachListener()` method added to our `EnemyManagerWithMessagesComponent` class:

```
| void OnDestroy() {  
|     if (MessagingSystem.IsAlive) {  
|         MessagingSystem.Instance.DetachListener(typeof(EnemyCreatedMessage),  
|                                         this.HandleCreateEnemy);  
|     }  
| }
```

Note how this definition makes use of the `IsAlive` property declared in the `SingletonComponent` class. This safeguards us against the aforementioned problem of accidentally creating a new `MessagingSystem` during application shutdown, since we can never guarantee that the Singleton gets destroyed last.

Wrapping up the Messaging System

Congratulations are in order, as we have finally built a fully functional global Messaging System that any and all objects can interface with, and use it to send messages between one another. A useful feature of this approach is that it is type-agnostic, meaning that the message senders and listeners do not even need to derive from any particular class in order to interface with the Messaging System; it just needs to be a class that provides a message type and a delegate function of the matching function signature, which makes it accessible to both ordinary classes and `MonoBehaviours`.

As far as benchmarking the `MessagingSystem` class goes, we will find that it is capable of processing hundreds, if not thousands, of messages in a single frame with minimal CPU overhead (depending on the CPU, of course). The CPU usage is essentially the same, whether one message is being distributed to 100 different listeners or 100 messages are distributed to just one listener. Either way, it costs about the same.

Even if we're predominantly sending messages during UI or gameplay events, this probably has far more power than we need. So, if it does seem to be causing performance problems, then it's far more likely to be caused by what the listener delegates are doing with the message than the Messaging System's ability to process those messages.

There are many ways to enhance the Messaging System to provide more useful features we may need in the future, as follows:

- Allow message senders to suggest a delay (in time or frame count) before a message is delivered to its listeners
- Allow message listeners to define a priority for how urgently it should receive messages compared to other listeners waiting for the same message type. This is a means for listeners to skip to the front of the queue if it was registered later than other listeners

- Implement some safety checks to handle situations where a listener gets added to the list of message listeners for a particular message while a message of that type is still being processed. Currently, C# will throw an `EnumerationException` at us since the delegate list will be changed by `AttachListener()`, while it is still being iterated through in `TriggerEvent()`

At this point, we probably explored the Messaging System enough, so these tasks will be left as an academic exercise for you to undertake, if you become comfortable using this solution in your games. Let's continue to explore more ways to improve performance through script code.

Disable unused scripts and objects

Scenes can get pretty busy sometimes, especially when we're building large, open worlds. The more objects invoking code in an `update()` callback, the worse things it will scale and the slower the game becomes. However, much of what is being processed may be completely unnecessary if it is outside of the player's view or simply too far away to matter. This may not be a possibility in large city-building simulation games, where the entire simulation must be processed at all times, but it is often possible in first-person and racing games since the player is wandering around a large expansive area, where non-visible objects can be temporarily disabled without having any noticeable effect on gameplay.

Disabling objects by visibility

Sometimes, we may want Components or `GameObject`s to be disabled when they're not visible. Unity comes with built-in rendering features to avoid rendering objects that are not visible to the player's Camera view (through a technique known as *Frustum Culling*, which is an automatic process) and to avoid rendering objects that are hidden behind other objects (occlusion culling, which will be discussed in [Chapter 6, Dynamic Graphics](#)), but these are only rendering optimizations. It does not affect Components that perform tasks on the CPU, such as AI scripts, User Interface, and gameplay logic. We must control that behavior ourselves.

A good solution to this problem is using the `OnBecameVisible()` and `OnBecameInvisible()` callbacks. As the names imply, these callback methods are invoked when a renderable object has become visible or invisible with respect to any Cameras in our Scene. In addition, when there are multiple Cameras in a Scene (for example, a local multiplayer game), the callbacks are only invoked if the object becomes visible to any one Camera and becomes invisible to all Cameras. This means that the aforementioned callbacks will be called at exactly the right times we expect; if nobody can see it, `OnBecameInvisible()` gets called, and if at least one player can see it, `OnBecameVisible()` gets called.

Since the visibility callbacks must communicate with the Rendering Pipeline, the `GameObject` must have a renderable Component attached, such as a `MeshRenderer` or `SkinnedMeshRenderer`. We must ensure that the Components we want to receive the visibility callbacks from are also attached to the same `GameObject` as the renderable object and are not a parent or child `GameObject`, otherwise they won't be invoked.



Note that Unity also counts the hidden Camera of the Scene window toward the `OnBecameVisible()` and `OnBecameInvisible()` callbacks. If we find that these methods are not being invoked properly during Play Mode testing, ensure that you turn the Scene window Camera away from everything or disable the

Scene window entirely.

To enable/disable individual Components with the visibility callbacks, we can add the following methods:

```
| void OnBecameVisible() { enabled = true; }  
| void OnBecameInvisible() { enabled = false; }
```

Also, to enable/disable the entire `GameObject` the Component is attached to, we can implement the methods this way instead:

```
| void OnBecameVisible() { gameObject.SetActive(true); }  
| void OnBecameInvisible() { gameObject.SetActive(false); }
```

Although, be warned that disabling the `GameObject` containing the renderable object, or one of its parents, will make it impossible for `OnBecameVisible()` to be called since there's now no graphical representation for the Camera to see and trigger the callback with. We should place the Component on a child `GameObject`, and have the script disable that instead, leaving the renderable object always visible (or find another way to re-enable it later).

Disabling objects by distance

In other situations, we may want Components or `GameObjects` to be disabled after they are far enough away from the player such that they may be barely visible, but too far away to matter. A good candidate for this type of activity is roaming AI creatures that we want to see at a distance, but where we don't need it to process anything, and it can sit idle until we get closer.

The following code is a simple Coroutine that periodically checks the total distance from a given target object and disables itself if it strays too far away from it:

```
[SerializeField] GameObject _target;
[SerializeField] float _maxDistance;
[SerializeField] int _coroutineFrameDelay;

void Start() {
    StartCoroutine(DisableAtADistance());
}

IEnumerator DisableAtADistance() {
    while(true) {
        float distSqrD = (transform.position - _target.transform.position).sqrMagnitude;
        if (distSqrD < _maxDistance * _maxDistance) {
            enabled = true;
        } else {
            enabled = false;
        }

        for (int i = 0; i < _coroutineFrameDelay; ++i) {
            yield return new WaitForEndOfFrame();
        }
    }
}
```

We should assign the player's character object (or whatever object we want it to compare with) to the `_target` field in the Inspector window, define the maximum distance in `_maxDistance`, and modify the frequency with which the Coroutine is invoked using the `_coroutineFrameDelay` field. Any time the object goes further than `_maxDistance` distance away from the object assigned to `_target`, it will be disabled. It will be re-enabled if it returns within that distance.

A subtle performance-enhancing feature of this implementation is comparing against distance-squared instead of the raw distance. This leads us conveniently to our next section.

Consider using distance-squared over distance

It is safe to say that CPUs are relatively good at multiplying floating-point numbers together, but relatively dreadful at calculating square-roots from them. Every time we ask a `Vector3` to calculate a distance with the `magnitude` property or with the `Distance()` method, we're asking it to perform a square-root calculation (as per Pythagorean theorem), which can cost a lot of CPU overhead compared to many other types of vector math calculations.

However, the `Vector3` class also offers a `sqrMagnitude` property, which provides the same result as `distance`, only the value is squared. This means that if we also square the value we wish to compare `distance` against, then we can perform essentially the same comparison without the cost of an expensive square-root calculation.

For example, consider the following code:

```
float distance = (transform.position - other.transform.position).Distance();
if (distance < targetDistance) {
    // do stuff
}
```

This can be replaced with the following and achieve a nearly identical result:

```
float distanceSqrD = (transform.position - other.transform.position).sqrMagnitude;
if (distanceSqrD < (targetDistance * targetDistance)) {
    // do stuff
}
```

The reason the result is nearly identical is because of the floating-point precision. We're likely to lose some of the precision that we would have had from using the square-root values, since the value will be adjusted to an area with a different density of representable numbers; it can land exactly on, or closer to, a more accurate representable number, or, more likely, it will land on a number with less accuracy. As a result, the comparison is not exactly the same, but, in most cases, it is close enough to be unnoticeable, and the

performance gain can be quite significant for each instruction we replace in this manner.

If this minor precision loss is not important, then this performance trick should be considered. However, if precision is very important (such as running an accurate large-scale galactic space simulation), then you may want to give this tip a pass.

Note that this technique can be used for any square-root calculations, not just for distance. This is simply the most common example you might run across, and it brings to light the important `sqrMagnitude` property of the `Vector3` class. This is a property which Unity Technologies intentionally exposed for us to make use of in this manner.

Minimize Deserialization behavior

Unity's Serialization system is mainly used for Scenes, Prefabs, ScriptableObjects and various Asset types (which tend to derive from ScriptableObject). When one of these object types is saved to disk, it is converted into a text file using the **Yet Another Markup Language (YAML)** format, which can be deserialized back into the original object type at a later time. All `GameObjects` and their properties get serialized when a Prefab or Scene is serialized, including `private` and `protected` fields, all of their Components as well as its child `GameObjects` and their Components, and so on.

When our application is built, this serialized data is bundled together in large binary data files internally called Serialized Files in Unity. Reading and deserializing this data from disk at runtime is an incredibly slow process (relatively speaking) and so all deserialization activity comes with a significant performance cost.

This kind of deserialization takes place any time we call `Resources.Load()` for a file path found under a folder named *Resources*. Once the data has been loaded from disk into memory, then reloading the same reference later is much faster, but disk activity is always required the first time it is accessed. Naturally, the larger the data set we need to deserialize, the longer this process takes. Since every Component of a Prefab gets serialized, then the deeper the hierarchy is, the more data needs to be serialized. This can be a problem for Prefabs with very deep hierarchies, Prefabs with many empty `GameObjects` (since every `GameObject` always contains at least a `Transform` Component), and particularly problematic for User Interface (UI) Prefabs, since they tend to house many more Components than a typical Prefab.

Loading in large serialized data sets like these could cause a significant spike in CPU the first time they are loaded, which tend to increase loading time if they're needed immediately at the start of the Scene. More importantly, they can cause frame drops if they are loaded at runtime. There are a couple of

approaches we can use to minimize the costs of deserialization.

Reduce serialized object size

We should aim to make our serialized objects as small as possible, or partition them into smaller data pieces we combine together piece-by-piece so that they can be loaded one piece at a time over time. This can be tricky to manage for Prefabs since Unity does not inherently support nested Prefabs, and so we would be implementing such a system ourselves, which is a notoriously difficult problem to solve in Unity. UI Prefabs are good candidates for separating into smaller pieces, since we don't normally need the entire UI at any given moment, and so we can usually afford to load them in one piece at a time.

Load serialized objects asynchronously

Prefabs and other serialized content can be loaded in asynchronously via `Resources.LoadAsync()`, which will offload reading from disk onto a worker thread that eases the burden on the main thread. It will take some time for the serialized object to become available, which can be checked by calling the `isDone` property on the `ResourceRequest` object returned by the previous method call.

This is not ideal for Prefabs we need immediately at the start of the game, but all future Prefabs are good candidates for asynchronous loading if we're willing to create systems that manage this behaviour.

Keep previously loaded serialized objects in memory

As previously mentioned, once a serialized object has been loaded into memory, then it remains there and can be copied if we need it again later, such as instantiating more copies of a Prefab. We can free this data later with explicit calls to `Resources.Unload()`, which will release the memory space to be reused later. But, if we have a lot of surplus memory in the application's budget, then we could choose to keep this data in memory, which would reduce the need to reload it again from disk later. This naturally consumes a lot of memory with more and more serialized data, making it a risky strategy for memory management, and so we should only do this when necessary.

Move common data into ScriptableObject

If we have a lot of different Prefabs with Components that contain a lot of properties that tend to share data, such as game design values like hit points, strength, speed, and so on, then all of this data will be serialized into every Prefab that uses them. A better approach is to serialize this common data in a `ScriptableObject` which they load and use instead. This reduces the amount of data stored within the Serialized File for the Prefab and could significantly reduce the loading time of our scenes by avoiding too much repetitive work.

Load scenes additively and asynchronously

Scenes can be loaded either to replace the current Scene or can be loaded additively to add its contents to the current Scene without unloading the preceding one. This can be toggled via the `LoadSceneMode` argument of the `SceneManager.LoadScene()` family of functions.

Another mode of Scene loading is to complete it either synchronously or asynchronously, and there are good reasons to use both. Synchronous loading is the typical means of loading a Scene by calling `SceneManager.LoadScene()` where the main thread will block until the given Scene completes loading. This normally results in a poor user experience, as the game appears to freeze as the contents are loaded in (whether as a replacement or additively). This is best used if we want to get the player into the action as soon as possible, or we have no time to wait for Scene objects to appear. This would normally be used if we're loading into the first level of the game or returning to the main menu.

For future Scene loading, however, we may wish to reduce the performance impact so that we can continue to keep the player in action. Loading a Scene can take a lot of work, and the larger the Scene, the longer it will take. However, the option of asynchronous additive loading offers a huge benefit; we can let the Scene gradually load in the background without causing significant impact on the user experience. This can be accomplished with `SceneManager.LoadSceneAsync()` combined with passing in `LoadSceneMode.Additive` for the loading mode argument.

It's important to realize that scenes do not strictly follow the concept of a game level. In most games, players are normally trapped in one level at a time, but Unity can support multiple scenes being loaded simultaneously through additive loading, allowing each Scene to represent a small chunk of a level. Ergo, we could initialize the first Scene for the level (*Scene-1-1a*), and as the player nears the next section, asynchronously and additively load in the

next (*Scene-1-1b*), and repeat this continuously as the player travels through the level.



Unity Technologies restructured the Scene system in a mid-release of Unity 5 into a global SceneManager class, although this was more of a naming convention change to make it clear that levels and scenes are not the same thing. The exact same features are available in older releases of Unity, just with a slightly different API through the Application class.

Exploiting this feature would require a system that either constantly checks the player's position in the level until they get close or uses Trigger Volumes to broadcast a message that the player is nearing the next section and begin asynchronous loading at the appropriate time. Another important consideration is that the Scene's contents won't appear immediately since asynchronous loading effectively spreads the loading out over a handful of frames in order to cause as little visible impact as possible. We need to make sure that we trigger asynchronous Scene loading with more than enough time to spare so that the player won't see objects popping into the game.

Scenes can also be unloaded to clear them out of memory. This will save some memory or runtime performance in the form of removing any Components making use of `update()` that we no longer need. Again, this can be done both synchronously and asynchronously with `SceneManager.UnloadScene()` and `SceneManager.UnloadSceneAsync()`. This can be an enormous performance benefit because we're only using what we need due to the player's location in the level, but note that it is not possible to unload small chunks of a monolithic Scene. If the original Scene file was enormous, then unloading it would unload everything. The original Scene would have to be broken up into smaller scenes and then loaded and unloaded as needed. Similarly, we should only begin to unload a Scene if we're certain the player can no longer see its constituent objects, otherwise they would witness objects disappearing out of nowhere. One last consideration is that Scene unloading would trigger the destruction of many objects, which is likely to free up a lot of memory and trigger the garbage collector. An efficient use of memory is also important when making use of this tip.

This approach would require a significant amount of Scene redesign work,

script writing, testing, and debugging, which is not to be underestimated, but the benefits of improving user experience are exceptional. Having seamless transitions between areas in a game is a benefit that is often praised by players and critics, because it doesn't take the player out of the action. If we use it appropriately, it can save a significant amount of runtime performance, improving the user experience further.

Create a custom Update() layer

Earlier in this chapter, in the "*Update, Coroutines and InvokeRepeating*" section, we discussed the relative pros and cons of using these Unity Engine features as a means of avoiding excessive CPU workload during most of our frames. Regardless of which of these approaches we might adopt, there is an additional risk of having lots of `MonoBehaviours` written to periodically call some function, which is having too many methods triggering in the same frame simultaneously.

Imagine thousands of `MonoBehaviours` that initialized together at the start of a Scene, each starting a Coroutine at the same time that will process their AI tasks every 500 milliseconds. It is highly likely that they would all trigger within the same frame, causing a huge spike in its CPU usage for a moment, which settles down temporarily and then spikes again a few moments later when the next round of AI processing is due. Ideally, we would want to spread these invocations out over time.

The following are the possible solutions to this problem:

- Generating a random time to wait each time the timer expires or Coroutine triggers
- Spread out Coroutine initialization so that only a handful of them are started at each frame
- Pass the responsibility of calling updates to some God Class that places a limit on the number of invocations that occur each frame

The first two options are appealing since they're relatively simple and we know that Coroutines can potentially save us a lot of unnecessary overhead. However, as we discussed, there are many dangers and unexpected side effects associated with such drastic design changes.

A potentially better approach to optimize updates is to not use `update()` at all, or more accurately, to use it only once. When Unity calls `update()`, and in fact, any of its callbacks, it crosses the aforementioned Native-Managed Bridge,

which can be a costly task. In other words, the processing cost of executing 1,000 separate `update()` callbacks will be more expensive than executing one `update()` callback, which calls into 1,000 regular functions. As we witnessed in the "*Remove empty callback definitions*" section, calling `update()` thousands of times is not a trivial amount of work for the CPU to undertake, primarily because of the Bridge. We can, therefore, minimize how often Unity needs to cross the Bridge by having a God Class `MonoBehaviour` use its own `Update()` callback to call our own custom update-style system used by our custom Components.

In fact, many Unity developers prefer implementing this design right from the start of their projects, as it gives them finer control over when and how updates propagate throughout the system; this can be used for things such as menu pausing, cool time manipulation effects, or prioritizing important tasks and/or suspending low priority tasks if we detect that we're about to reach our CPU budget for the current frame.

All objects wanting to integrate with such a system must have a common entry point. We can achieve this through an Interface Class with the `interface` keyword. Interface Classes essentially set up a contract whereby any class that implements the Interface Class Class must provide a specific series of methods. In other words, if we know the object implements an Interface Class, then we can be certain about what methods are available. In C#, classes can only derive from a single base class, but they can implement any number of Interface Classes (this avoids the *deadly diamond of death* problem that C++ programmers will be familiar with).

The following Interface Class definition will suffice, which only requires the implementing class to define a single method called `onUpdate()`:

```
| public interface IUpdateable {  
|     void OnUpdate(float dt);  
| }
```

It's common practice to start an Interface Class definition with a capital 'I' to make it clear that it is an Interface Class we're dealing with. The beauty of Interface Classes is that they improve the decoupling of our codebase, allowing huge subsystems to be replaced, and as long as the Interface Class is



adhered to, we will have greater confidence that it will continue to function as intended.

Next, we'll define a custom `MonoBehaviour` type which implements this Interface Class:

```
| public class UpdateableComponent : MonoBehaviour, IUpdateable {  
|     public virtual void OnUpdate(float dt) {}  
| }
```

Note that we're naming the method `OnUpdate()` rather than `Update()`. We're defining a custom version of the same concept, but we want to avoid name collisions with the built-in `Update()` callback.

The `OnUpdate()` method of the `UpdateableComponent` class retrieves the current delta time (`dt`), which spares us from a bunch of unnecessary `Time.deltaTime` calls, which are commonly used in `Update()` callbacks. We've also made the function `virtual` to allow derived classes to customize it.

This function will never be called as it's currently being written. Unity automatically grabs and invokes methods defined with the `Update()` name, but has no concept of our `OnUpdate()` function, so we will need to implement something that will call this method when the time is appropriate. For example, some kind of `GameLogic` God Class could be used for this purpose.

During the initialization of this Component, we should do something to notify our `GameLogic` object of both its existence and its destruction so that it knows when to start and stop calling its `OnUpdate()` function.

In the following example, we will assume that our `GameLogic` class is a `SingletonComponent`, as defined earlier in the "*Singleton Components*" section, and has appropriate `static` functions defined for registration and deregistration. Bear in mind that it could just as easily use the aforementioned `MessagingSystem` to notify the `GameLogic` of its creation/destruction.

For `MonoBehaviours` to hook into this system, the most appropriate place is within their `Start()` and `OnDestroy()` callbacks:

```
| void Start() {
```

```

        GameLogic.Instance.RegisterUpdateableObject(this);
    }

    void OnDestroy() {
        if (GameLogic.Instance.IsAlive) {
            GameLogic.Instance.DeregisterUpdateableObject(this);
        }
    }
}

```

It is best to use the `Start()` method for the task of registration, since using `Start()` means that we can be certain all other pre-existing Components will have at least had their `Awake()` methods called prior to this moment. This way, any critical initialization work will have already been done on the object before we start invoking updates on it.

Note that because we're using `Start()` in a `MonoBehaviour` base class, if we define a `start()` method in a derived class, it will effectively override the base class definition, and Unity will grab the derived `start()` method as a callback instead. It would, therefore, be wise to implement a `virtual Initialize()` method so that derived classes can override it to customize initialization behavior without interfering with the base class's task of notifying the `GameLogic` object of our Component's existence.

The following code provides an example of how we might implement a `virtual Initialize()` method.

```

void Start() {
    GameLogic.Instance.RegisterUpdateableObject(this);
    Initialize();
}

protected virtual void Initialize() {
    // derived classes should override this method for initialization code, and
}

```

Finally, we will need to implement the `GameLogic` class. The implementation is effectively the same whether it is a `SingletonComponent` or a `MonoBehaviour`, and whether or not it uses the `MessagingSystem`. Either way, our `UpdateableComponent` class must register and deregister as `IUpdateable` objects, and the `GameLogic` class must use its own `Update()` callback to iterate through every registered object and call their `OnUpdate()` function.

Here is the definition for our `GameLogic` class:

```

public class GameLogicSingletonComponent : SingletonComponent<GameLogicSingle
public static GameLogicSingletonComponent Instance {
    get { return ((GameLogicSingletonComponent)_Instance); }
    set { _Instance = value; }
}

List<IUpdateable> _updateableObjects = new List<IUpdateable>();

public void RegisterUpdateableObject(IUpdateable obj) {
    if (!_updateableObjects.Contains(obj)) {
        _updateableObjects.Add(obj);
    }
}

public void DeregisterUpdateableObject(IUpdateable obj) {
    if (_updateableObjects.Contains(obj)) {
        _updateableObjects.Remove(obj);
    }
}

void Update()
{
    float dt = Time.deltaTime;
    for (int i = 0; i < _updateableObjects.Count; ++i) {
        _updateableObjects[i].OnUpdate(dt);
    }
}
}

```

If we make sure that all of our custom Components inherit from the `UpdateableComponent` class, then we've effectively replaced " N " invocations of the `update()` callback with just one `update()` callback, plus " N " virtual function calls. This can save us a large amount of performance overhead because even though we're calling virtual functions (which cost a small overhead more than a non-virtual function call because it needs to redirect the call to the correct place), we're still keeping the overwhelming majority of update behavior inside our Managed code and avoiding the Native-Managed Bridge as much as possible. This class can even be expanded to provide priority systems, to skip low-priority tasks if it detects that the current frame has taken too long, and many other possibilities.

Depending on how deep you already are into your current project, such changes can be incredibly daunting, time-consuming, and likely to introduce a lot of bugs as subsystems are updated to make use of a completely different set of dependencies. However, the benefits can outweigh the risks if time is on your side. It would be wise to do some testing on a group of objects in a Scene that is similarly designed to your current Scene files to verify that the

benefits outweigh the costs.

Summary

This chapter introduced you to many methods that improve your scripting practices in the Unity Engine, with the aim of improving performance if (and only if) you have already proven them to be the cause of a performance problem. Some of these techniques demand some forethought and profiling investigation before being implemented, since they often come with introducing additional risks or obfuscating our codebase for new developers. Workflow is often just as important as performance and design, so before you make any performance changes to the code, you should consider whether or not you're sacrificing too much on the altar of performance optimization.

We will investigate more advanced scripting improvement techniques later, in [Chapter 8, *Masterful Memory Management*](#), but let's take a break from staring at code and explore some ways to improve graphics performance using a pair of built-in Unity features known as Dynamic Batching and Static Batching.

The Benefits of Batching

In 3D graphics and games, batching is a very general term describing the process of grouping a large number of wayward pieces of data together and processing them as a single, large block of data. This situation is ideal for CPUs, and particularly GPUs, which can handle simultaneous processing of multiple tasks with their multiple cores. Having a single core switching back and forth between different locations in memory takes time, so the less this needs to be done, the better.

In some cases, the act of batching refers to large sets of meshes, vertices, edges, UV coordinates, and other different data types that are used to represent a 3D object. However, the term could just as easily refer to the act of batching audio files, sprites, texture files, and other large datasets.

So, just to clear up any confusion, when the topic of batching is mentioned in Unity, it is usually referring to the two primary mechanisms it offers for batching mesh data: Dynamic Batching and Static Batching. These methods are essentially two different forms of geometry merging, where we combine mesh data of multiple objects together and render them all in a single instruction, as opposed to preparing and drawing each one separately.

The process of batching together multiple meshes into a single mesh is possible because there is no reason a mesh object must fill a contiguous volume of 3D space. The Rendering Pipeline is perfectly happy with accepting a collection of vertices that are not attached together with edges, and so we can take multiple separate meshes that might have resulted in multiple render instructions and combine them together into a single mesh, thus rendering it out using a single instruction.

There has been a lot of confusion over the years surrounding the conditions under which the Dynamic Batching and Static Batching systems activate and where we might even see a performance improvement. After all, in some cases, batching can actually degrade performance if it is not used wisely. A proper understanding of these systems will give us the knowledge we need to

improve the graphics performance of our application in significant ways.

This chapter intends to dispel much of the misinformation floating around about these systems. We will observe, via explanation, exploration, and examples, just how these two batching methods operate. This will enable us to make informed decisions, using most of them to improve our application's performance.

We will cover the following topics in this chapter:

- A brief introduction to the Rendering Pipeline and the concept of Draw Calls
- How Unity's Materials and Shaders work together to render our objects
- Using the Frame Debugger to visualize rendering behavior
- How Dynamic Batching works, and how to optimize it
- How Static Batching works, and how to optimize it

Draw Calls

Before we discuss Dynamic Batching and Static Batching independently, let's first understand the problems that they are both trying to solve within the Rendering Pipeline. We will try to keep fairly light on the technicalities as we will explore this topic in greater detail in [Chapter 6, Dynamic Graphics](#).

The primary goal of these batching methods is to reduce the number of Draw Calls required to render all objects in the current view. At its most basic form, a Draw Call is a request sent from the CPU to the GPU asking it to draw an object.



Draw Call is the common industry vernacular for this process, although they are sometimes referred to as SetPass Calls in Unity, since some low-level methods are named as such. Think of it as configuring options before initiating the current rendering pass. We will refer to them as Draw Calls throughout the remainder of this book.

Before a Draw Call can be requested, several tasks need to be completed. Firstly, mesh and texture data must be pushed from the CPU memory (RAM) into GPU memory (VRAM), which typically takes place during initialization of the Scene, but only for textures and meshes the Scene file knows about. If we dynamically instantiate objects at runtime using texture and mesh data that hasn't appeared in the Scene yet, then they must be loaded at the time they are instantiated. The Scene cannot know ahead of time which Prefabs we're planning to instantiate at runtime, as many of them are hidden behind conditional statements and much of our application's behavior depends upon user input. Next, the CPU must prepare the GPU by configuring the options and rendering features that are needed to process the object that is the target of the Draw Call.

These communication tasks between the CPU and GPU take place through the underlying Graphics API, which could be DirectX, OpenGL, OpenGLES, Metal, WebGL, or Vulkan, depending on the platform we're targeting and

certain graphical settings. These API calls go through a library, called a *driver*, which maintains a long series of complex and interrelated settings, state variables, and datasets that can be configured and executed from our application (although drivers are designed to service multiple applications simultaneously, as well as render calls coming from multiple threads). The available features change enormously based on the graphics card we're using and the version of the Graphics API we're targeting; more advanced graphics cards support more advanced features, which would need to be supported by newer versions of the API, so updated drivers would be needed to enable them. The sheer number of settings, supported features, and compatibility levels between one version and another that have been created over the years (particularly for the older APIs such as DirectX and OpenGL) can be nothing short of mind-boggling. Thankfully, at a certain level of abstraction, all of these APIs tend to operate in a similar fashion; hence Unity is able to support many different Graphics APIs through a common interface.

This utterly massive array of settings that must be configured to prepare the Rendering Pipeline just prior to rendering an object is often condensed into a single term known as the *Render State*. Until these Render State options are changed, the GPU will maintain the same Render State for all incoming objects and render them in a similar fashion.

Changing the Render State can be a time-consuming process. So, for example, if we were to set the Render State to use a blue texture file and then ask it to render one gigantic mesh, then it would be rendered very rapidly with the whole mesh appearing blue. We could then render 9 more, completely different meshes, and they would all be rendered blue, since we haven't changed which texture is being used. If, however, we wanted to render 10 meshes using 10 different textures, then this will take longer. This is because we will need to prepare the Render State with the new texture just prior to sending the Draw Call instruction for each mesh.

The texture being used to render the current object is effectively a global variable in the Graphics API, and changing a global variable within a parallel system is much easier said than done. In a massively parallel system such as a GPU, we must effectively wait until all of the current jobs have reached the same synchronization point (in other words, the fastest cores need to stop and

wait for the slowest ones to catch up, wasting processing time that they could be using on other tasks) before we can make a Render State change, at which point we will need to spin up all of the parallel jobs again. This can waste a lot of time, so the less we need to ask the Render State to change, the faster the Graphics API will be able to process our requests.

Things that can trigger Render State synchronization include--but are not limited to--an immediate push of a new texture to the GPU and changing a Shader, lighting information, shadows, transparency, and pretty much any graphical setting we can think of.

Once the Render State is configured, the CPU must decide what mesh to draw, what textures and Shader it should use, and where to draw the object based on its position, rotation, and scale (all represented within a 4x4 matrix known as a *transform*, which is where the `Transform` Component gets its name from) and then send an instruction to the GPU to draw it. In order to keep the communication between CPU and GPU very dynamic, new instructions are pushed into a queue known as the *Command Buffer*. This queue contains instructions that the CPU has created and that the GPU pulls from each time it finishes the preceding command.

The trick to how batching improves the performance of this process is that a new Draw Call does not necessarily mean that a new Render State must be configured. If two objects share the exact same Render State information, then the GPU can immediately begin rendering the new object since the same Render State is maintained after the last object is finished. This eliminates the time wasted due to a Render State synchronization. It also serves to reduce the number of instructions that need to be pushed into the Command Buffer, reducing the workload on both the CPU and GPU.

Materials and Shaders

Render State in Unity is essentially exposed to us via Materials. Materials are a container around Shaders, short programs that define how the GPU should render incoming vertex and texture data. A Shader on its own does not have the necessary knowledge of state to accomplish anything of value. A Shader requires inputs such as diffuse textures, Normal maps, and lighting information and effectively dictates what Render State variables need to be set in order to render the incoming data.



Shaders are named this way because their original implementation many years ago was to only handle lighting and shading of an object (applying shadows, where originally there were none). Their purpose has grown enormously since then, and now they have a much more generic purpose of being a programmable access point to many different kinds of parallel tasks, but the old name still remains.

Every Shader needs a Material, and every Material must have a Shader. Even newly imported meshes introduced into the Scene without an assigned Material are automatically assigned a default (hidden) Material, which gives them a basic diffuse Shader and a white coloration. So, there is no way of getting around this relationship.



Note that a single Material can only support a single Shader. The use of multiple Shaders on the same mesh requires separate Materials to be assigned to different parts of the same mesh.

Therefore, if we want to minimize how often the Render State changes, then we can do so by reducing the number of Materials we use during a Scene. This would result in two performance improvements simultaneously; the CPU will spend less time generating and transmitting instructions to the GPU each frame and the GPU won't need to stop and re-synchronize state changes as often.

Let's begin with a simple Scene in order to visualize the behavior of Materials and Batching. However, before we start, we should disable a few rendering options as they will contribute some extra Draw Calls, which might be distracting:

1. Navigate to Edit | Project Settings | Quality and set Shadows to Disable Shadows (or select the default Fastest quality level)
2. Navigate to Edit | Project Settings | Player, open the Other Settings tab, and disable Static Batching and Dynamic Batching if they are enabled

Next, we'll create a Scene that contains a single Directional Light with four cubes and four spheres, where each object has its own unique Material, position, rotation, and scale:



In the preceding screenshot, in the Batching value in the Game window's Stats popup, we can see 9 total batches. This value closely represents the number of Draw Calls used to render the Scene. The current view will consume one of these batches rendering the background of the Scene, which could be set to Skybox or Solid Color. This is determined by the Camera

object's Clear Flags settings.

The remaining 8 batches are used to draw our 8 objects. In each case, the Draw Call involves preparing the Rendering Pipeline using the Material's properties and asking the GPU to render the given mesh at its current transform. We have ensured that each Material is unique by giving them each a unique texture file to render. Ergo, each mesh requires a different Render State, and, therefore, each of our 8 meshes requires a unique Draw Call.

As previously mentioned, we can theoretically minimize the number of Draw Calls by reducing how often we cause the system to change Render State information. So, part of the goal is to reduce the amount of Materials we use. However, if we set all objects to use the same Material, we still won't see any benefit and the number of batches remains at 9:



This is because we're not actually reducing the number of Render State

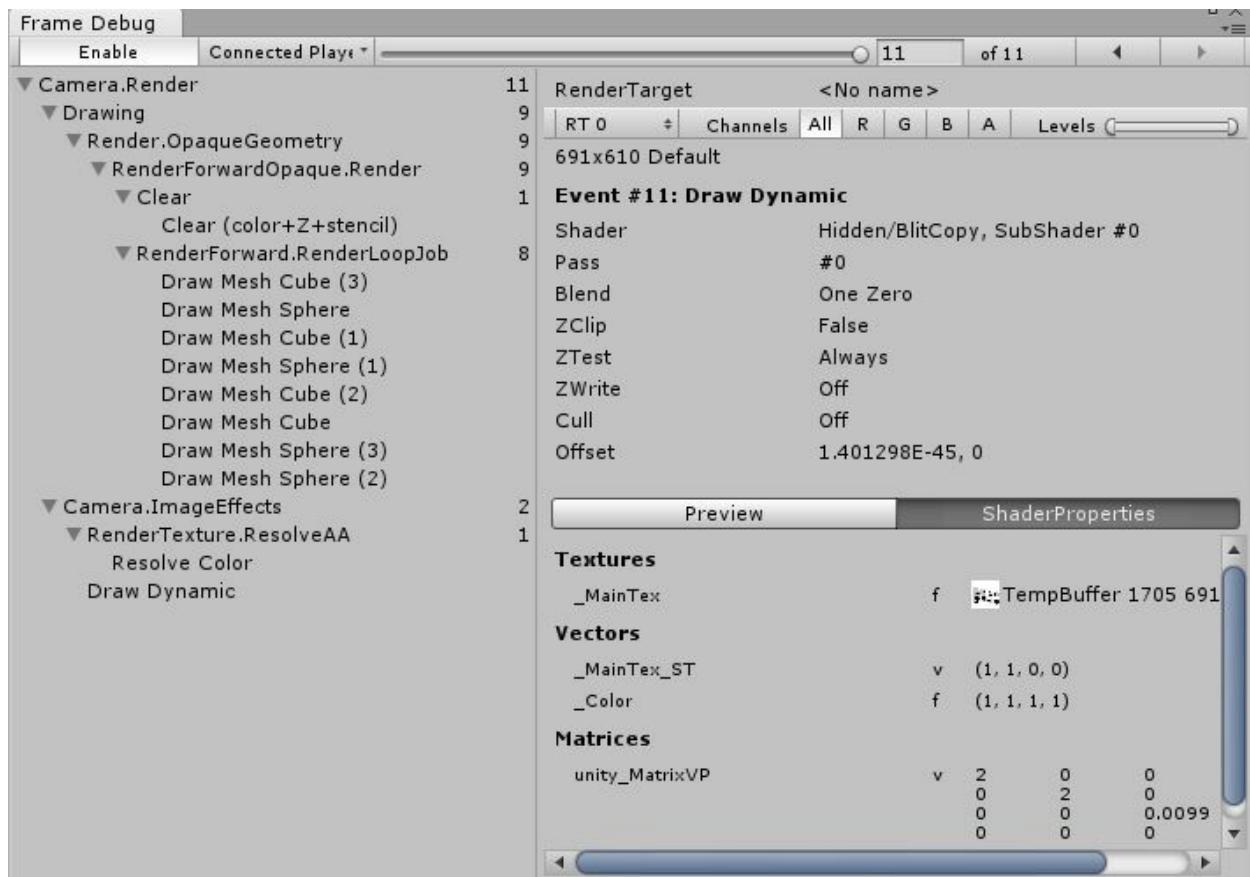
changes nor efficiently grouping mesh information. Unfortunately, the Rendering Pipeline is not smart enough to realize we're overwriting the exact same Render State values, and then asking it to render the same meshes, over and over again.

The Frame Debugger

Before we dive into how batching can save us Draw Calls, let's explore a useful tool, which can help us determine how batching is affecting our Scene—the Frame Debugger.

We can open the Frame Debugger by selecting Window | Frame Debugger from the main window or clicking on the Frame Debugger button in the Breakdown View Options in the Rendering Area of the Profiler. Either approach will open the Frame Debug window.

Clicking on the Enable button in the Frame Debug window will allow us to observe how our Scene is being constructed, one Draw Call at a time. The following screenshot shows the user interface of the Frame Debugger, with a list of GPU instructions in the left-hand panel, and more detailed information in the right-hand panel:



There is a lot of information in this window, which can provide useful information if we want to debug the behavior of a single Draw Call, but the most useful area to look at is the Drawing section in the left-hand panel, which lists all of the Draw Calls in our Scene.

Each item under this section represents a unique Draw Call and what was rendered by it. An amazingly useful feature of this tool is the ability to click on any one of these items and immediately observe only the Draw Calls needed to render the Scene up to that point in the Game window. This lets us perform a quick, visual *diff* of two sequential Draw Calls. This can make it easy to spot exactly which object(s) were rendered by a given Draw Call. This can help determine whether or not a set of objects were batched together by looking at how many of them appear during that Draw Call.

A weird bug with the Frame Debugger (which still exists in early builds of Unity 2017) is that if we are observing a Scene making use of a Skybox and click on various items under the Drawing section, only the final Scene presentation can be





observed in the Game window. We would need to temporarily disable the Skybox via the Camera's Clear Flags setting to take a look at how Draw Call progression appears in the Game window by setting it to Solid Color instead.

As we can see in the preceding Frame Debugger screenshot, one Draw Call is being consumed to clear the screen (the item labelled Clear), and then our 8 meshes are being rendered in 8 separate Draw Calls (the item labelled RenderForward.RenderLoopJob).



Note that the number next to each item in the left-hand panel actually represents a Graphics API call, of which a Draw Call is but one type of API call. These can be seen in the Camera.Render, Camera.ImageEffects, and RenderTexture.ResolveAA items. Any API call can be just as costly as a Draw Call, but the overwhelming majority of API calls we will make in a complex Scene is in the form of Draw Calls, so it is often best to focus on minimizing Draw Calls before worrying about the API communication overhead of things such as post-processing effects.

Dynamic Batching

Dynamic Batching has the following three important qualities:

- Batches are generated at runtime (batches are dynamically generated)
- The objects that are contained within a batch can vary from one frame to the next, depending on what meshes are currently visible to the Main Camera view (batch contents are dynamic)
- Even objects that can move around the Scene can be batched (it works on dynamic objects)

Hence, these attributes lead us to the name *Dynamic* Batching.

If we return to the Player Settings page and enable Dynamic Batching, we should see that the number of batches drops from 9 down to 6. Dynamic Batching automatically recognizes that our objects share Material and mesh information and is, therefore, combining some of them into a larger batch for processing. We should also see a different list of items in the Frame Debugger, demonstrating that meshes are now being dynamically batched.



As we can see from the Frame Debugger, our four boxes have been combined into a single Draw Call named Dynamic Batch, but our four spheres are still being rendered with four separate Draw Calls. This is because the four spheres do not fit the requirements of Dynamic Batching. Despite the fact that they all use the same Material, there are many more requirements we must fulfill.

The list of the requirements needed to successfully dynamically batch a mesh can be found in the Unity documentation at <http://docs.unity3d.com/Manual/DrawCallBatching.html>.



At the time of publication of this book, the previous page is a little out of date, and a better breakdown can be found in this blog post: <https://blogs.unity3d.com/2017/04/03/how-to-see-why-your-draw-calls-are-not-batched-in-5-6/>.

The following list covers the requirements to enable Dynamic Batching for a given mesh:

- All mesh instances must use the same Material reference
- Only `ParticleSystem` and `MeshRenderer` Components are dynamically batched. `SkinnedMeshRenderer` Components (for animated characters) and all other renderable Component types cannot be batched
- There is a limit of 300 vertices per mesh
- The total number of vertex attributes used by the Shader must be no greater than 900
- Either all mesh instances should use a uniform scale or all meshes should use a nonuniform scale, but not a mixture of the two
- Mesh instances should refer to the same Lightmap file
- The Material's Shader should not depend on multiple passes
- Mesh instances must not receive real-time shadows
- There is an upper limit on the total number of mesh indices in the entire batch, which varies per-Graphics API and platform used, which is around 32k-64k indices (check out the documentation/aforementioned blog post for specifics)

It is important to note the term *Material references*, because if we happen to use two different Materials with identical settings, the Rendering Pipeline is

not smart enough to realize that, and they will be treated as different Materials and, therefore, will be disqualified from Dynamic Batching. Most of the rest of these requirements have either already been explained; however, a couple of these requirements are not completely intuitive or clear from the description, which merits some additional explanation.

Vertex attributes

A vertex attribute is simply a piece of information contained within a mesh file on a per-vertex basis, and each is normally represented as a group of multiple floating-point values. This includes, but is not limited to, a vertex's position (relative to the root of the mesh), a normal vector (a vector pointing away from the object's surface, most often used in lighting calculations), one or more sets of texture UV coordinates (used to define how one or more textures wrap around the mesh), and possibly even color information per-vertex (normally used in custom lighting or for a flat-shaded, low-poly style object). Only meshes with less than 900 total vertex attributes used by the Shader can be included in Dynamic Batching.



Note that looking into a mesh's raw data file may contain less vertex attribute information than what Unity loads into memory because of how the engine converts mesh data from one of several raw data formats into an internal format. So, don't assume that the number of attributes our 3D modeling tool tells us the mesh uses will be the final count. The best way to verify the attribute count is to either drill down into the mesh object in the Project window until you find the `MeshFilter` Component and look at the `verts` value that appears in the Preview subsection of the Inspector window.

Using more attribute data per vertex within the accompanying Shader will consume more from our 900-attribute budget and hence reduce the number of vertices the mesh is allowed to have before it can no longer be used in Dynamic Batching. For example, a simple diffuse Shader might only use 3 attributes per-vertex: position, normal, and a single set of UV coordinates. Dynamic Batching would therefore be able to support meshes using this Shader, which have a combined total of 300 vertices. However, a more complex Shader, requiring 5 attributes per-vertex, would only be able support Dynamic Batching with meshes using no more than 180 vertices. Also, note that even if we are using less than 3 vertex attributes per vertex in our Shader, Dynamic Batching still only supports meshes with a maximum of 300

vertices, so only relatively simple objects are candidates for Dynamic Batching.

These restrictions are why our Scene saves only 3 Draw Calls with Dynamic Batching enabled, despite having all objects share the same Material reference. The cube mesh that is autogenerated by Unity contains a mere 8 vertices, each with position, normal, and UV data, for 24 attributes in total. This is far less than the 300-vertex limit and 900-vertex attribute limit. However, an autogenerated sphere mesh contains 515 vertices and hence 1,545 total vertex attributes. These meshes clearly exceed both the 300-vertex and 900-vertex attribute limits and, therefore, cannot be dynamically batched.

If we click on one of the Draw Call items in the Frame Debugger, there will be a section labeled "Why this draw call can't be batched with the previous one". Most of the time, the explanation text beneath tells us which requirement we failed and (or at least the first one it detected) which can be useful for debugging batching behavior.



Be warned that the explanation text that appears is currently a little buggy in an early build of Unity 2017 and sometimes the actual reason appears when you click on the DynamicBatch item instead. For instance, in the preceding Scene example, the explanation given for Draw Mesh Sphere (3) was "The material doesn't have GPU instancing enabled", which is a bit confusing. However, looking at the DynamicBatch item gave the real reason, that is, "a submesh we are trying to dynamic-batch has more than 300 vertices".

Mesh scaling

The documentation suggests that objects should either share a uniform scale or each should have a unique nonuniform scale in order to be included in dynamic batching. A uniform scale means that all three components of the scale vector (x , y , and z) are identical to each other (but not necessarily across meshes), and a nonuniform scale means that at least one of these values is different from the others, and meshes that fall within these two groups will be placed in two different batches.

Let's explain this via an example. Consider that we have the following four objects: A scaled at $(1, 1, 1)$, B scaled at $(2, 1, 1)$, C scaled at $(2, 2, 1)$, and D scaled at $(2, 2, 2)$. Objects A and D have uniform scales, because all three components have the same value. Even though both meshes don't have exactly the same scale, their scales still count as uniform, and so they will be placed together in one dynamic batch. Meanwhile, B and C both have nonuniform scales, because at least one component is different than the rest. They're both scaled very differently (C being scaled on the y -axis twice as much as B), but they both still count as nonuniform, so they will be combined together in a separate batch.

However, be warned that using negative scaling has a strange effect on Dynamic Batching. Negative scaling is often a quick way to mirror a mesh in our Scene, which can save us from having to create and import a completely different mesh for something that's only flipped on one axis. This trick is commonly used for pairs of doors, or just to make a Scene look more varied. However, if we only negatively scale the mesh on one or three axes, then it will be placed into a different Dynamic Batch than meshes negatively scaled on zero or two axes. It does not matter which of the three values (x , y , or z) are negative, only whether the total number of negative values is odd or even.

Another strange by-product of this batch-splitting behavior is how the rendering order of the objects can determine what gets batched together. If the previous object would have appeared in a different batch group than the current one, then it cannot be batched. Again, this is best explained by the

following example. Consider that we have four objects again: V scaled at $(1, 1, 1)$, W scaled at $(-1, 1, 1)$, X scaled at $(-1, -1, 1)$, Y scaled at $(-1, -1, -1)$, and finally Z scaled similarly to V at $(1, 1, 1)$. Objects V and Z share a common uniform scale, so we might expect them to be batched together. However, if all of these objects were rendered to the Scene in the above order, then object V will be rendered, and Unity will test to check whether objects W and V could share a batch. They cannot do so, because of object W's odd negative scaling, so no batching takes place. Unity will then compare object X with object W to check whether they can be batched, which they cannot because W has an odd negative scaling and X has even negative scaling. The next comparisons between objects W-Y and Y-Z fail for the same reason. The end result is that all five objects will be rendered with five separate Draw Calls, and there was no opportunity given to combine objects V and Z. Note that this weird effect only comes into play when negative scaling is used. If all scaling differences are uniform versus nonuniform, then Unity should be capable of batching the objects appropriately.

Presumably, this is all a by-product of the algorithm used to detect valid batchable groups, since mirroring a mesh in two dimensions is mathematically equivalent to rotating the mesh about both of the same axes 180 degrees, while there is no rotational equivalent to mirroring a mesh on one or three axes. Thus, the behavior we observe is perhaps just the Dynamic Batching system automatically transforming the object for us, although this isn't completely clear. Regardless, hopefully this prepares us for many of the weird situations we might run in to when it comes to generating dynamic batches.

Dynamic Batching summary

Dynamic Batching is a very useful tool when we want to render very large groups of simple meshes. The design of the system makes it ideal to use when we're making use of large numbers of simple meshes, which are nearly identical in appearance. Possible situations to apply Dynamic Batching could be as follows:

- A large forest filled with rocks, trees, and bushes
- A building, factory, or space station with many simple, common elements (computers, corridor pieces, pipes, and so on)
- A game featuring many dynamic, non-animated objects with simple geometry and particle effects (a game such as Geometry Wars springs to mind)

If the only requirement preventing two objects from being Dynamically Batched together is the fact that they use different texture files, be aware that it only takes a bit of development time and effort to combine textures, and regenerate mesh UVs so that they can be Dynamically Batched together (commonly known as *Atlasing*). This may cost us in texture quality or the overall size of a texture file (which can have drawbacks we will understand once we dive into the topic of GPU Memory Bandwidth in [Chapter 6, Dynamic Graphics](#)), but it is worth considering.

Perhaps the only situation where Dynamic Batching may be a detriment on performance is if we were to set up a Scene with hundreds of simple objects, where only a few objects are put into each batch. In these cases, the overhead cost of detecting and generating so many small batches might cost more time than we'd save by just making a separate Draw Call for each mesh. Even still, this is unlikely.

If anything, we're far more likely to inflict performance losses on our application by simply assuming that Dynamic Batching is taking place, when we've actually forgotten one of the essential requirements. We can accidentally break the vertex limit by pushing a new version of a mesh, and

in the process of Unity converting a raw Object (with the *.obj* extension) file into its own internal format, it generates more vertex attributes than we expected. We could also exceed it by tweaking some Shader code or adding additional passes without realizing it would disqualify it from Dynamic Batching. We might even set up the object to enable shadows or Light Probes, which breaks another requirement.

There will be no warning when these accidents occur, save for the number of Draw Calls increasing after changes are made and causing performance to degrade further and further. Maintaining a healthy amount of Dynamic Batching in our Scenes requires constant vigilance in checking our Draw Call count and looking at Frame Debugger data to make sure that we didn't accidentally disqualify objects from Dynamic Batching during our latest changes. However, as always, we only need to worry about our Draw Call performance if we've already proven that it's causing a performance bottleneck.

Ultimately, every situation is unique, so it is worth experimenting with our mesh data, Materials, and Shaders to determine what can and cannot be dynamically batched, and performing some testing in our Scene from time to time to ensure that the number of Draw Calls we're using remains reasonable.

Static Batching

Unity offers a second batching mechanism through Static Batching. This batching feature is similar to Dynamic Batching in a couple of ways, such as which objects are batched is determined at runtime based on what's visible to the Camera, and the contents of these batches will vary from frame to frame. However, there is one very important difference: it only works on objects that are marked Static, hence the name *Static* Batching.

The Static Batching system has its own set of requirements:

- As the name implies, the meshes must be flagged as Static (specifically, Batching Static)
- Additional memory must be set aside for each mesh being statically batched
- There is an upper limit on the number of vertices that can be combined in a static batch that varies per Graphics API and platform, which is around 32k-64k vertices (check out the documentation/aforementioned blog post for specifics)
- The mesh instances can come from any source mesh, but they must share the same Material reference

Let's cover some of these requirements in more detail.

The Static flag

Static Batching can only be applied to objects with the static flag enabled or, more specifically, the Batching Static subflag (these subflags are known as `StaticEditorFlags`). Clicking on the small down-pointing arrow next to the Static option for a `GameObject` will reveal a dropdown of the `StaticEditorFlags`, which can alter the object's behavior for various Static processes.

An obvious side effect of this is that the object's transform cannot be changed, and, hence, any object wishing to make use of Static Batching cannot be moved, rotated, or scaled in any way.

Memory requirements

The additional memory requirement for Static Batching will vary, depending on the amount of replication occurring within the batched meshes. Static Batching works by copying the data for all flagged and visible meshes into a single, large mesh data buffer, and passing it into the Rendering Pipeline through a single Draw Call, while ignoring the original mesh. If all of the meshes being statically batched are unique, then this would cost us no additional memory usage compared to rendering the objects normally, as the same amount of memory space is required to store the meshes.

However, since the data is effectively copied, these statically batched duplicates cost us additional memory equal to the number of meshes, multiplied by the size of the original mesh. Ordinarily, rendering one, ten, or a million clones of the same object costs us the same amount of memory, because they're all referencing the same mesh data. The only difference between objects in this case is the transform of each object. However, because Static Batching needs to copy the data into a large buffer, this referencing is lost, since each duplicate of the original mesh is copied into the buffer with a unique set of data with a hardcoded transform baked into the vertex positions.

Therefore, using Static Batching to render 1,000 identical tree objects will cost us 1,000 times more memory than rendering the same trees without Static Batching. This causes some significant memory consumption and performance issues if Static Batching is not used wisely.

Material references

We are already aware that sharing Material references is a means of reducing Render State changes, so this requirement is fairly obvious. In addition, sometimes, we statically batch meshes that require multiple Materials. In this case, all meshes using a different Material will be grouped together in their own static batch and for each unique Material being used.

The downside to this requirement is that, at best, Static Batching can only render all of the static meshes using a number of Draw Calls equal to the number of Materials they need.

Static Batching caveats

The Static Batching system has some additional drawbacks. Owing to how it approaches the batching solution, by combining meshes into a single greater mesh, the Static Batching system has a few caveats that we need to be aware of. These concerns range from minor inconveniences to major drawbacks, depending on the Scene:

- Draw Call savings are not immediately visible from the Stats window until runtime
- Objects marked Batching Static introduced in the Scene at runtime will not be automatically included in Static Batching

Let's explore these problems in a little more detail.

Edit Mode debugging of Static Batching

Trying to determine the overall effect that Static Batching will have on our Scene can be a little tricky since nothing is being Statically Batched while in Edit Mode. All of the magic happens during runtime, which can make it difficult to determine what benefits Static Batching would provide without manual testing. We should use the Frame Debugger to verify that our static batches are being properly generated, and that they contain the expected objects.

This can be especially problematic if we leave implementing this feature until late in the project lifecycle, where we can spend a lot of time launching, tweaking, and relaunching our Scene to ensure that we're getting the Draw Call savings we're expecting. Consequently, it is best to start working on Static Batching optimization early in the process of building a new Scene.

It goes without saying that static batch creation work is not completely trivial, and it may also massively inflate Scene initialization time if there are many batches to create and/or many large objects to batch.

Instantiating static meshes at runtime

Any new objects we add into the Scene at runtime will not be automatically combined into any existing batch by the Static Batching system, even if they were marked as Batching Static. To do so would cause an enormous runtime overhead between recalculating the mesh and synchronizing with the Rendering Pipeline, so Unity does not even attempt to do it automatically.

For the most part, we should try to keep any meshes we want to be statically batched present in the original Scene file. However, if dynamic instantiation is necessary, or we are making use of additive Scene loading, then we can control static batch eligibility with the `StaticBatchUtility.Combine()` method. This utility method has two overloads: either we provide a root `GameObject`, in which case all child `GameObjects` with meshes will be turned into new static batch groups (multiple could be created if they share multiple Materials) or we provide a list of `GameObjects` and a root `GameObject`, and it will automatically attach them as children to the root and generate new static batch groups in the same manner.

We should profile our usage of this function, as it can be quite an expensive operation if there are many vertices to combine. It will also not combine the given meshes with any preexisting statically batched groups, even if they share the same Material. So we will not be able to save Draw Calls by instantiating or additively loading Static meshes that use the same Material as other statically batched groups already present in the Scene (it can only combine with meshes it was grouped with in the `Combine()` call).

Note that if any of the `GameObjects` we batch with the `StaticBatchUtility.Combine()` method are not marked as `Static` before batching, the `GameObjects` will remain non-`Static`, but the mesh itself will be `Static`. This means that we could accidentally move the `GameObject`, its `Collider` Component, and any other important objects, but the mesh will remain in the same



location. Be careful about accidentally mixing Static and non-Static states in statically batched objects.

Static Batching summary

Static Batching is a powerful, but dangerous tool. If we don't use it wisely, we can very easily inflict enormous performance losses via memory consumption (potentially leading to application crashes) and rendering costs on our application. It also takes a good amount of manual tweaking and configuration to ensure that batches are being properly generated, and that we aren't accidentally introducing any unintended side effects of using various Static flags. However, it does have a significant advantage in that it can be used on meshes of different shapes and enormous sizes, which Dynamic Batching cannot provide.

Summary

It is clear that the Dynamic Batching and Static Batching systems are not a silver bullet. We cannot blindly apply them to any given Scene and expect improvements. If our application and Scene happen to fit a particular set of parameters, then these methods are very effective at reducing CPU load and rendering bottlenecks. However, if not, then some additional work is required to prepare our Scene to meet batching feature requirements. Ultimately, only a good understanding of these batching systems and how they function can help us determine where and when this feature can be applied, and, hopefully, this chapter has given us all of the information we need to make informed decisions.

You will learn more about the Rendering Pipeline and performance improvement techniques in [Chapter 6, *Dynamic Graphics*](#). Until then, let's move onto a different topic and look into some of the more subtle performance improvements that we can achieve through managing our art assets in intelligent ways.

Kickstart Your Art

Art is a very subjective area, dominated by personal opinion and preference. It can be difficult to say whether, and why, one piece of art is better than the other, and it's likely that we won't be able to find complete consensus on our opinions. The technical aspects behind art assets that support a game's artistry can also be very subjective. There are multiple workarounds that can be implemented to improve performance, but these tend to result in a loss of quality for the sake of speed. If we're trying to reach peak performance then it's important that we consult with our team members whenever we decide to make any changes to our art assets, as it is primarily a balancing act, which can be an art form in itself.

Whether we're trying to minimize our runtime memory footprint, keeping the smallest possible executable size, maximizing loading speed, or maintaining consistency in frame rate, there are plenty of options to explore. There are some methods that are clearly always ideal, but most require a little more care and forethought before being adopted, as they would result in reduced quality or increase the chances of developing bottlenecks in other subsystems.

In this chapter, we will explore how to improve performance for the following asset types:

- Audio files
- Texture files
- Mesh and animation files
- Asset Bundles and Resources

In each case, we will investigate how Unity stores, loads, and manipulates these assets both during application build time and runtime. We will also examine our options in the event of performance issues, and what we can do to avoid behavior that might generate performance bottlenecks.

Audio

Unity, as a framework, can be used to build anything from small applications that require only a handful of sound effects and a single background track to huge role playing games that need millions of lines of spoken dialog, music tracks, and ambient sound effects. Regardless of the actual scope of the application, audio files are often a large contributor to the application size after it is built (sometimes called its *disk footprint*). Moreover, many developers are surprised to find that runtime audio processing can turn into a significant source of CPU and memory consumption.

Audio is often neglected on both sides of the gaming industry; developers tend not to commit many resources to it until the last minute, whereas users rarely draw their attention to it. Nobody notices when audio is handled well, but we all know what bad audio sounds like--it's instantly recognizable, jarring, and guaranteed to draw unwanted attention. This makes it crucial not to sacrifice too much audio clarity in the name of performance.

Audio bottlenecks can come from a variety of sources. Excessive compression, too much audio manipulation, too many active Audio Components, inefficient memory storage methods, and access speeds are all ways to invite poor memory and CPU performance. However, with a little effort and understanding, all it takes is a few tweaks here and there to save us from a user experience disaster.

Importing audio files

When we select an imported audio file in the Project window, the Inspector window will reveal multiple import settings. These settings dictate everything from loading behavior, compression behavior, quality, sample rate, and (in later versions of Unity) whether to support ambisonic audio (multichannel audio, which combines tracks via spherical harmonics in order to create more realistic audio experiences).



Many of the audio import options can be configured on a per-platform basis, allowing us to customize behavior between different target platforms.

Loading audio files

The following are the three settings that dictate the manner in which an audio file is loaded:

- Preload Audio Data
- Load In Background
- Load Type

Our audio files are initially packaged as binary data files that are bundled with our application, which reside on the hard disk of the device (although in some cases they downloaded from somewhere on the internet). *Loading* audio data simply means pulling it into main memory (RAM) so that it can be later processed by audio decoders, which then convert the data into audio signals to our headphones, or speakers. However, how loading happens will vary enormously based on the previous three settings. The first setting, Preload Audio Data determines whether audio data will be automatically loaded during Scene initialization, or loaded at a later time. When loading of audio data does occur, the second setting, Load In Background, determines whether this activity blocks the main thread until it is finished, or loads it asynchronously in the background. Finally, the Load Type setting defines what kind of data gets pulled into memory and how much data gets pulled at a time. All three of these settings can have a dramatically negative effect on performance if they are not used wisely.

The typical use case of an audio file is to assign it to the `audioclip` property of an `object, which will wrap it in an AudioClip object. We can then trigger playback via AudioSource.Play() or AudioSource.PlayOneShot(). Each Audio Clip assigned in this way would be loaded into memory during Scene initialization since the Scene contains immediate references to these files, which it must resolve before they are needed. This is the default case when Preload Audio Data is enabled.`

Disabling Preload Audio Data tells the Unity Engine to skip audio file asset loading during Scene initialization, which defers loading activity to the first

moment it is needed, in other words, when `Play()` or `PlayOneShot()` are called. Disabling this option will speed up Scene initialization, but it also means that the first time we play the file, the CPU will need to immediately access the disk, retrieve the file, load it into memory, decompress it, and play it. This is a synchronous operation and will block the main thread until it is completed. We can prove this with a simple test:

```
public class PreloadAudioDataTest : MonoBehaviour {
    [SerializeField] AudioSource _source;

    void Update() {
        if (Input.GetKeyDown(KeyCode.Space)) {
            using (new CustomTimer("Time to play audio file", 1)) {
                _source.Play();
            }
        }
    }
}
```

If we add an `to our Scene, assign a large audio file to it, and assign it to the _source field of the PreloadAudioDataTest Component, we can press the spacebar and take a look at a printout of how long the Play() function took to complete. A simple test of this code against a 10 MB audio file with Preload Audio Data enabled will reveal that the call was practically instantaneous. However, disabling Preload Audio Data, applying the changes to the file, and repeating the test reveals that it takes significantly longer (around 700 ms on a desktop PC with an Intel i5 3570K). This completely blows past our budget for a single frame, so in order to use this toggle responsibly we will need to load the majority of our audio assets into memory ahead of time.`

This can be achieved by calling `AudioClip.LoadAudioData()` (which can be acquired through an `Component's clip property). However, this activity will still block the main thread for the same amount of time it takes to load it in the previous example, and so loading our audio file will still cause frame drops, regardless of whether or not we choose to load it ahead of time. Data can also be unloaded through AudioClip.UnloadAudioData().`

This is where the Load In Background option comes in. This changes audio loading into an asynchronous task; therefore, loading will not block the main thread. With this option enabled, the actual call to `AudioClip.LoadAudioData()` would complete instantly, but keep in mind that

the file won't be ready to play until loading completes on a separate thread. We can double-check an `AudioClip` Component's current loading state through the `AudioClip.loadState` property. If `Load In Background` is enabled, and we call `without loading the data first, Unity will still require the file to be loaded into memory before it can be played, and so there will be a delay between when we called and when the audio file actually begins playback. This risks introducing jarring behavior if we try to access a sound file before it is fully loaded, causing it to be out of sync with other tasks, such as animations.`

Modern games typically implement convenient stopping points in levels to perform tasks such as loading or unloading audio data--for example, an elevator between floors or long corridors, where very little action is taking place. Solutions involving custom loading and unloading of audio data via these methods would need to be tailor-made to the particular game, depending on when audio files are needed, how long they're needed for, how Scenes are put together, and how players traverse them.

This can require a significant number of special case changes, testing, and asset management tweaks. So, it is recommended that you save this approach as a *Nuclear Option* to be used late in production, in the event that all other techniques have not succeeded as well as we hoped.

Finally, there is the Load Type option, which dictates how audio data loads when it occurs. There are three options available:

- Decompress On Load
- Compressed in Memory
- Streaming

These three options are explained in detail:

- **Decompress On Load:** This setting compresses the file on disk to save space and decompresses it into memory when it is first loaded. This is the standard method of loading an audio file and should be used in most cases. It takes some time to decompress the file, which leads to a little extra overhead during loading, but reduces the amount of work required when the audio file is played.

- Compressed In Memory: This setting simply copies the compressed file straight from disk into memory when it is loaded. It will only decompress the audio file during runtime when it is being played. This will sacrifice runtime CPU when the Audio Clip is played, but improves loading speed and reduces runtime memory consumption while the Audio Clip remains dormant. Hence, this option is best used for very large audio files that are used fairly frequently, or if we're incredibly bottlenecked on memory consumption and are willing to sacrifice some CPU cycles to play the Audio Clip.
- Streaming: Finally, this setting (also known as *Buffered*) will load, decode, and play files on the fly at runtime by gradually pushing the file through a small buffer where only one small piece of the overall file is present in memory at a time. This method uses the least amount of memory for a particular Audio Clip, but the largest amount of runtime CPU. Since each instance of playback of the file will need to generate its own buffer, this setting comes with the unfortunate drawback where referencing the Audio Clip more than once leads to multiple copies of the same Audio Clip in memory that must all be processed separately, resulting in a runtime CPU cost if used recklessly. Consequently, this option is best reserved for single-instance Audio Clips that play regularly and never need to overlap with other instances of itself or even with other streamed Audio Clips. For example, this setting is best used with background music and ambient sound effects that need to be played during the majority of a Scene's lifetime.

So, let's recap. The default case, with Preload Audio Data enabled, Load In Background disabled, and a Load Type of Decompress On Load, causes a long Scene loading time, but ensures that every Audio Clip we reference in the Scene is ready immediately when we need it. There will be no loading delays when the Audio Clip is needed, and the Audio Clip will playback the moment we call `Play()`. A good compromise to improve Scene loading time is to enable Load In Background for Audio Clips we won't need until later, but this should not be used for Audio Clips we need shortly after Scene initialization. We then control when our audio data is loaded manually through `AudioClip.LoadAudioData()` and `AudioClip.UnloadAudioData()`. We should be willing to use all of these methods in a single Scene in order to reach optimal performance.

Encoding formats and quality levels

Unity supports three general case encoding formats for Audio Clips, which is determined by the Compression Format option when we view an Audio Clip's properties in the Inspector window:

- Compressed (the actual text for this option can appear differently, depending on the platform)
- PCM
- ADPCM

The audio files we import into the Unity Engine can be one of many popular audio file formats, such as Ogg Vorbis, MPEG-3 (MP3), and Wave, but the actual encoding that is bundled into the executable will be converted into a different format.

The compression algorithm used with the Compressed setting will depend on the platform being targeted. Stand-alone applications and other non-mobile platforms will convert the file into the Ogg-Vorbis format, whereas mobile platforms use MP3.



There are a few platforms that always use a specific type of compression, such as HEVAG for the PS Vita, XMA for XBox One, and AAC for WebGL.

Statistics are provided in the Inspector window for the currently selected format in the area following the Compression Format option, providing an idea of how much disk space is being saved by the compression. Note that the first value displays the original file size and the second displays the size cost on disk. How much memory the audio file will consume at runtime once loaded will be determined by how efficient the chosen compression format is. For example, Ogg Vorbis compression will generally decompress to about 10 times its compressed size, whereas ADPCM will decompress to about four

times the compressed size.



The cost savings displayed in the Inspector window for an audio file only apply for the currently selected platform and most recently applied settings. Ensure that the Editor is switched to the correct platform in File | Build Settings and that you click on Apply after making changes in order to see the actual cost savings (or cost inflation) for the current configuration. This is particularly important for WebGL applications, since the AAC format generally leads to very inflated audio file sizes.

The encoding/compression format used can have a dramatic effect on the quality, file size, and memory consumption of the audio file during runtime, and only the Compressed setting gives us the ability to alter the quality without affecting the sampling rate of the file. Meanwhile, the PCM and ADPCM settings do not provide this luxury, and we're stuck with whatever file size those compression formats decide to give us--that is, unless we're willing to reduce audio quality for the sake of file size by reducing the sampling rate.

The PCM format is a lossless and uncompressed audio format, providing a close approximation of analog audio. It trades large file sizes for higher audio quality and is best used for very short sound effects that require a lot of clarity where any compression would otherwise distort the experience.

Meanwhile, the ADPCM format is far more efficient in both size and CPU consumption than PCM, but compression results in a fair amount of noise. This noise can be hidden if it is reserved for short sound effects with a lot of chaos, such as explosions, collisions, and impact sounds where we might not be aware of any generated artifacts.

Finally, the Compressed format will result in small files that have lower quality than PCM, but significantly better quality than ADPCM at the expense of additional runtime CPU usage. This format should be used in most cases. This option allows us to customize the resultant quality level of the compression algorithm to tweak quality against file size. Best practices with the Quality slider are to search for a quality level that is as small as possible, but unnoticeable to users. Some user testing may be required to find

the *sweet spot* for each file.



Do not forget that any additional audio effects applied to the file at runtime will not play through the Editor in Edit Mode, so any changes should be fully tested through the application in Play Mode.

Audio performance enhancements

Now that we have a better understanding of audio file formats, loading methods, and compression modes; let's explore some approaches that we can make to improve performance through tweaking audio behavior.

Minimize active Audio Source count

Since each actively playing Audio Source consumes a particular amount of CPU, it stands to reason that we can save CPU cycles by disabling redundant Audio Sources in our Scene. One approach is to limit how many instances of an Audio Clip can be played simultaneously. This involves sending audio playback requests through an intermediary that controls our Audio Sources in such a way that puts a hard cap on how many instances of an Audio Clip can be played simultaneously.

Almost every Audio Management Asset available in the Unity Asset Store implements an audio-throttling feature of some kind (often known as *Audio Pooling*), and for good reason; it's the best trade-off in minimizing excessive audio playback with the least cost in quality. For example, having 20 footstep sounds playing simultaneously won't sound too much different to playing 10 of them simultaneously and is less likely to become distracting by being too loud. For this reason, and because these tools often provide many more subtle performance-enhancing features, it is recommended that you use a preexisting solution rather than rolling out your own, as there is a lot of complexity to consider from audio file types, stereo/3D audio, layering, compression, filters, cross-platform capability, efficient memory management, and so on.

When it comes to ambient sound effects, they still need to be placed at specific locations in the Scene to make use of the logarithmic volume effect, which gives it the pseudo-3D effect, so an Audio Pooling system would probably not be an ideal solution. Limiting playback on ambient sound effects is best achieved by reducing the total number of Audio Sources. The best approach is to either remove some of them or reduce them down to one larger, louder Audio Source. Naturally, this approach affects the quality of the user experience since it would appear that the sound is coming from a single source and not multiple sources, therefore it should be used with care.

Enable Force to Mono for 3D sounds

Enabling the Force to Mono setting on a stereo audio file will mix together the data from both audio channels into a single channel, saving 50 percent of the file's total disk and memory space usage effectively. Enabling this option is generally not a good idea for some 2D sound effects where the stereo effect is often used to create a specific audio experience. However, we can enable this option for some good space savings on 3D positional Audio Clips, where the two channels are effectively identical. These Audio Source types will let the direction between the Audio Source and the player determine how the audio file gets played into the left/right ear, and playing a stereo effect in this case is generally meaningless. Forcing 2D sounds (sounds that play into the player's ears at full volume, regardless of distance/direction to the Audio Source) to mono might also make sense if there is no need for a stereo effect.

Resample to lower frequencies

Resampling imported audio files to lower frequencies will reduce the file size and runtime memory footprint. This can be achieved by setting an audio file's Sample Rate Setting to Override Sample Rate, at which point we can configure the sample rate through the Sample Rate option. Some files require high sample rates to sound reasonable, such as files with high pitches and most music files. However, lower settings can reduce the file's size without noticeable quality degradation in most cases. 22,050 Hertz is a common sampling rate for sources that involve human speech and classical music. Some sound effects may be able to get away with even lower frequency values. However, each sound effect will be affected by this setting in a unique way, so it would be wise to spend some time running a few tests before we finalize our decision on the sampling rate.

Consider all compression formats

Each of the Compressed, PCM and ADPCM compression formats have their own benefits and drawbacks as explained previously. It's possible to make some compromises in memory footprint, disk footprint, CPU usage, and audio quality using different encoding formats for different files where appropriate. We should be willing to use all of them in the same application and come up with a system that works for the kinds of audio files we're using so that we don't need to treat each file individually. Otherwise, we would need to do a prohibitive amount of testing to ensure that audio quality hasn't been degraded for each file.

Beware of streaming

The upside of the Streaming loading type is a low runtime memory cost since a small buffer is allocated and the file is continuously pushed through it like a data queue. This can seem quite appealing, but streaming files from disk should be restricted to large, single-instance files only, as it requires runtime hard disk access; which is one of the slowest forms of data access available to us (second only to pulling a file through a network). Layered or transitioning music clips may run into major hiccups using the Streaming option, at which point, it would be wise to consider using a different Load Type and control loading/unloading manually. We should also avoid streaming more than one file at a time, as it's likely to inflict a lot of cache misses on the disk that interrupt gameplay. This is why this option is primarily used for background music/ambient sound effects since we only need one at a time.

Apply Filter Effects through Mixer Groups to reduce duplication

Filter Effects can be used to modify the sound effect playing through an Audio Source and can be accomplished through `FilterEffect` Components. Each individual Filter Effect will cost some amount of both memory and CPU and can be a good way to achieve disk space savings while maintaining a lot of variety in audio playback, since one file could be tweaked by a different set of filters to generate completely different sound effects.

Due to the additional overhead, overusing Filter Effects in our Scene can result in dire consequences in performance. A better approach is to make use of Unity's Audio Mixer utility (Window | Audio Mixer) to generate common Filter Effect templates that multiple Audio Sources can reference to minimize the amount of memory overhead.

The official tutorial on Audio Mixers covers the topic in excellent detail:

<https://unity3d.com/learn/tutorials/modules/beginner/5-pre-order-beta/audiomixer-and-audiomixer-groups>

Use remote content streaming responsibly

It is possible to dynamically load game content via the Web through Unity, which can be an effective means of reducing an application's disk footprint since less data files need to be bundled into the executable, and also provides a means to present dynamic content using web services to determine what is presented to the user at runtime. Asset streaming can either be accomplished through the `WWW` class in Unity 5 or the `UnityWebRequest` class in Unity 2017.

The `WWW` class provides the `audioClip` property, which is used to access an `AudioClip` object if it was an audio file we downloaded via a `WWW` object. However, be warned that accessing this property will allocate a whole new `AudioClip` resource each time it is invoked, and similarly with other `WWW` resource-acquiring methods. This resource must be freed with the `Resources.UnloadAsset()` method once it is no longer required.

Unlike managed resources, discarding the reference (setting it to `null`) will not automatically free the resource, so it will continue to consume memory. Ergo, we should only obtain the `AudioClip` through the `audioClip` property once, and only use that reference from that point forward, releasing it when it is no longer required.

Meanwhile, in Unity 2017, the `WWW` class has been effectively replaced by the `UnityWebRequest` class, which makes use of the new HLAPI and LLAPI networking layers. This class provides various utilities to download and access what are primarily text files. Multimedia-based requests should go through the `UnityWebRequestMultimedia` helper class. So, if an `AudioClip` is requested, we should call `UnityWebRequestMultimedia.GetAudioClip()` to create the request, and `DownloadHandlerAudioClip.GetContent()` to retrieve it once the download is complete.

This new version of the API is designed to be more efficient at storing and providing the data we requested, and so reacquiring an `AudioClip` multiple

times through `DownloadHandlerAudioClip.GetContent()` will not lead to additional allocations. Instead, it will merely return a reference to the originally downloaded `AudioClip`.

Consider Audio Module files for background music

Audio Module files, also known as **Tracker Modules**, are an excellent means of saving a significant amount of space without any noticeable quality loss. Supported file extensions in Unity are `.it`, `.s3m`, `.xm`, and `.mod`. Unlike the common audio formats, which are read as streams of bits that must be decoded at runtime to generate a specific sound, Tracker Modules contain lots of small, high-quality samples and organize the entire track similar to a music sheet; defining when, where, how loud, with what pitch, and with what special effects each sample should be played with. This can provide significant size savings while maintaining high-quality sampling. So, if the opportunity is available to us to make use of Tracker Module versions of our music files, then it is worth exploring.

Texture files

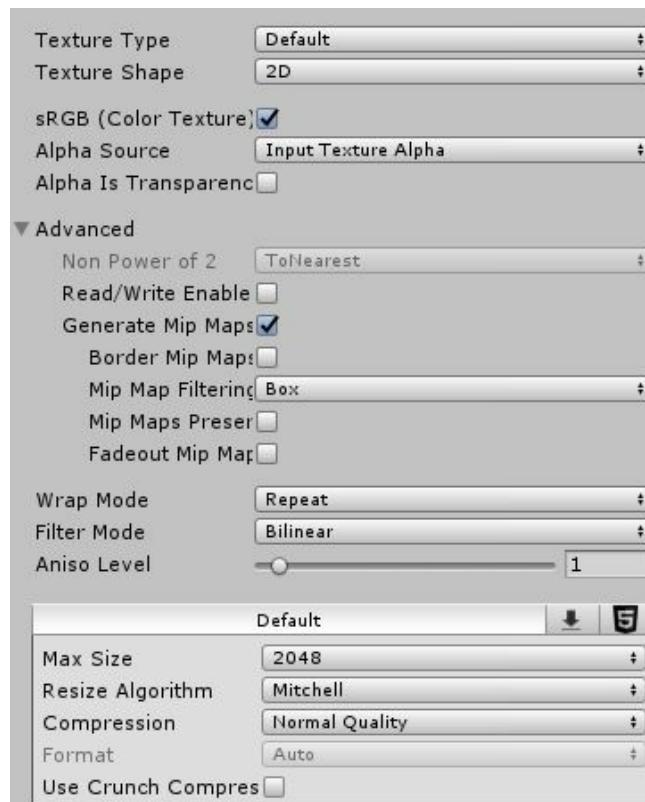
The terms *texture* and *sprite* often get confused in game development, so it's worth making the distinction--a texture is simply an image file, a big list of color data telling the interpreting program what color each pixel of the image should be, whereas a sprite is the 2D equivalent of a mesh, which is often just a single *quad* (a pair of triangles combined to make a rectangular mesh) that renders flat against the current Camera. There are also things called Sprite Sheets, which are large collections of individual images contained within a larger texture file, commonly used to contain the animations of a 2D character. These files can be split apart by tools, such as Unity's Sprite Atlas tool, to form individual textures for the character's animated frames.

Both meshes and sprites use textures to render an image onto its surface. Texture image files are typically generated in tools such as Adobe Photoshop or Gimp and then imported into our project in much the same way as audio files. At runtime, these files are loaded into memory, pushed to the GPU's VRAM, and rendered by a Shader over the target sprite or mesh during a given Draw Call.

Texture compression formats

Much like audio files, Unity will import texture files with a default list of settings that tend to keep things simple and work okay in the general case, but there are many import settings available, allowing us to improve a texture's quality and performance with some custom tweaking. Of course, making changes is just as likely to reduce quality and performance if we blindly make changes without fully understanding the internal processes going on.

The first option is the file's Texture Type. This setting will determine what other options are available, particularly under the Advanced dropdown. Not all importing options are available to all types, so it is best to configure this option for the texture's intended purpose, whether it is set to Normal Map, Sprite, Lightmap, and so on, as this will reveal the options appropriate for that type:



Similar to audio files, we can import texture files in multiple common

formats (such as `.jpg`, and `.png`), but the actual compression format built into the application could be one of many different texture compression formats ideally suited for GPUs of the given platform. These formats represent different ways of organizing the texture's color information; this includes different numbers of bits used to represent each channel (the more bits used, the more colors that can be represented), different numbers of bits per-channel (for example, the red channel may use more bits than the green channel), different total number of bits used for all channels (more bits naturally mean larger textures and more disk and memory consumption), whether or not an alpha channel is included, and perhaps the most important, different ways of packing the data together, which can allow for efficient memory access for the GPU (or incredibly inefficient access if the wrong packing type is chosen!).

The simple way of altering compression is to use the Compression texture import option to select one among the following:

- None
- Low Quality
- Normal Quality
- High Quality

Selecting None means that no compression will be applied. In this case, the final texture will still change format from the file type we imported, but it will select a format that makes no attempt at compression, and so we should see little or no quality loss at the expense of large texture files. The other three settings will pick a compression format, which, again, will vary depending on the platform, and Unity will try to pick a compression format that matches the option. For instance, selecting Low Quality will mean Unity picks a compression format that greatly reduces the texture size, but will generate some compression artifacts, while selecting High Quality will consume more memory with much larger texture sizes and minimal artifacting. Again, this is an automatic selection made by Unity.



The exact formats Unity picks for each platform for each of these Compression settings can be found at <https://docs.unity3d.com/Manual/class-TextureImporterOverride.html>.

The exact compression format Unity chooses can be overridden, although the available options vary per platform, since practically every platform has its own custom formats that work best for it. If we click on one of the platform-specific tabs beside the Default tab (just above the Max Size option), we will expose settings for a specific platform and can choose the exact compression format we want Unity to use.



There is also the Crunch Compression setting, which will apply an additional level of lossy compression on top of the DXT compression format. This option is only revealed if the other compression settings resulted in a DXT level of compression. This setting can save even more space at the cost of potentially glaring compression artifacts, depending on the Compressor Quality setting.

Several of a texture's import settings are fairly mundane, such as determining whether the file contains an alpha channel, how to wrap the texture at its extents, filtering method, and maximum possible resolution of the file (a global limit so that we don't accidentally overscale the texture beyond its original size on certain platforms). However, there are several other interesting options in these import settings, which we will cover in other sections where appropriate.

Texture performance enhancements

Let's explore some changes we can make to our texture files, which might help improve performance, depending on the situation and the content of the files we're importing. In each case, we'll explore the changes that need to be made and the overall effect they have, whether this results in a positive or negative impact on memory or CPU, an increase or decrease in the texture quality, and under what conditions we can expect to make use of these techniques.

Reduce texture file size

The larger a given texture file, the more GPU Memory Bandwidth will be consumed, pushing the texture when it is needed. If the total memory pushed per second exceeds the graphics card's total Memory Bandwidth, then we will have a bottleneck, as the GPU must wait for all textures to be uploaded before the next rendering pass can begin. Smaller textures are naturally easier to push through the pipeline than larger textures, so we will need to find a good middle ground between high quality and performance.

A simple test to find out if we're bottlenecked in Memory Bandwidth is to reduce the resolution of our games most abundant and largest texture files and relaunch the Scene. If the frame rate suddenly improves, then the application was most likely bound by texture throughput. If the frame rate does not improve or improves very little, then either we still have some Memory Bandwidth to make use of or there are bottlenecks elsewhere in the Rendering Pipeline, preventing us from seeing any further improvement.

Use Mip Maps wisely

There would be no point rendering small, distant objects, such as rocks and trees, with a high-detail texture if there's no way the player would ever be able to see that detail. Of course, they may see some slight improvement, but the performance cost may not be worth the minor detail increase. Mip Maps were invented as a way to solve this problem (as well as helping eliminate aliasing problems that were plaguing video games at around the same time), by pregenerating lower-resolution alternatives of the same texture and keeping them together in the same memory space. At runtime, the GPU picks the appropriate Mip Map level based on how large the surface appears within the perspective view (essentially based on the texel-to-pixel ratio when the object is rendered).

By enabling the Generate Mip Maps setting, Unity automatically handles the generation of these lower-resolution copies of the texture. These alternatives are generated using high-quality resampling and filtering methods within the Editor, rather than during runtime. There are several other options available for Mip Map generation, which can affect the quality of the generated levels, so some tweaking may be required to get a high quality set of Mip Maps. We will need to decide whether the time spent tweaking these values is worth it, since the whole purpose of Mip Maps is to intentionally reduce quality to save performance in the first place.

The following image shows how a 1024 x 1024 image that has been Mip Mapped into multiple lower-resolution duplicates:



These images will be packed together to save space, essentially creating a final texture file that will be about 33 percent larger than the original image. This will cost some disk space and GPU Memory Bandwidth to upload.

It's possible to observe which Mip Map levels are being used by our application at certain points by changing the Draw Mode setting of the Scene window to Mipmaps. This will highlight textures in red if they are larger than they should be, given the player's current view (the extra detail is wasted), while being highlighted blue means that they are too small (the player is observing a low quality texture with a poor texel-to-pixel ratio).

Remember that Mip Mapping is only useful if we have textures that need to be rendered at varying distances from the Camera. If we have textures that always render at a common distance from the Main Camera, such that the Mip Mapped alternatives are never used, then enabling Mip Maps is just a waste of space. Similarly, if we happen to have a texture that always resolves to the same Mip Map level since the player's Camera never gets too close/far away to switch levels, then it would be wiser to simply downscale the original texture.

Good examples of this would be any texture file used in a 2D game, textures used by UI systems, or those used in a Skybox or distant background since, by design, these textures will always be about the same distance from the Camera, so Mip Mapping would be essentially pointless. Other good examples include objects that only appear near the player, such as player-centric Particle Effects, characters, objects that only appear near the player, and objects that only the player can hold/carry.

Manage resolution downscaling externally

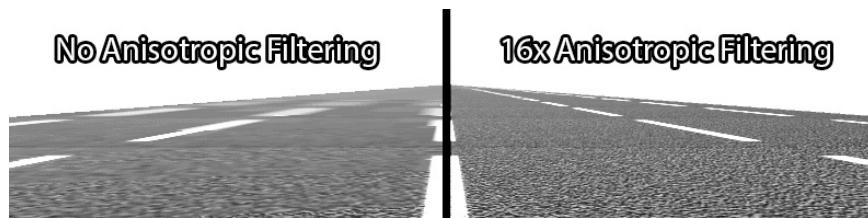
Unity puts a lot of effort into making things as easy to use as possible and provides us with the ability to place the project files from external tools to our project workspace, such as `.PSD` and `.TIFF` files, which are often large and split into multiple layered images. Unity automatically generates a texture file from the file's contents for the rest of the Engine to make use of, which can be very convenient, as we only need to maintain a single copy of the file through Source Control, and the Unity copy is automatically updated when an artist makes changes.

The problem is that the aliasing introduced by Unity's auto-texture generation and compression techniques from these files may not be as good as what the texture-editing tools we use could generate for us. Unity is very feature-rich and, first and foremost, focuses on being a game development platform, which means that it can have difficulty competing in areas that other software developers work on full time. Unity may be introducing artifacts through aliasing as a result of downscaling the image for us, and so we might find ourselves working around it by importing image files with a higher resolution than necessary just to keep the intended quality level. However, had we downsampled the image through the external application first, we might have suffered much less aliasing. In these cases, we may achieve an acceptable level of quality with a lower resolution, while consuming less overall disk and memory space.

We can either avoid using `.PSD` and `.TIFF` files within our Unity project as a matter of habit (storing them elsewhere and importing the downsampled version into Unity) or just perform some occasional testing to ensure that we're not wasting file size, memory, and GPU Memory Bandwidth using larger resolution files than necessary. This costs us some convenience in project file management, but can provide some significant savings for some textures, if we're willing to spend the time comparing the different downsampled versions.

Adjust Anisotropic Filtering levels

Anisotropic Filtering is a feature that improves the image quality of textures when they are viewed at very oblique (shallow) angles. The following screenshot shows the classic example of painted lines on a road with and without Anisotropic Filtering applied:



In either case, the painted lines close to the Camera appear fairly clear, but things change as they get further away from the Camera. Without Anisotropic Filtering, the distant painted lines get more and more blurry and distorted, whereas these lines remain crisp and clear with Anisotropic Filtering applied.

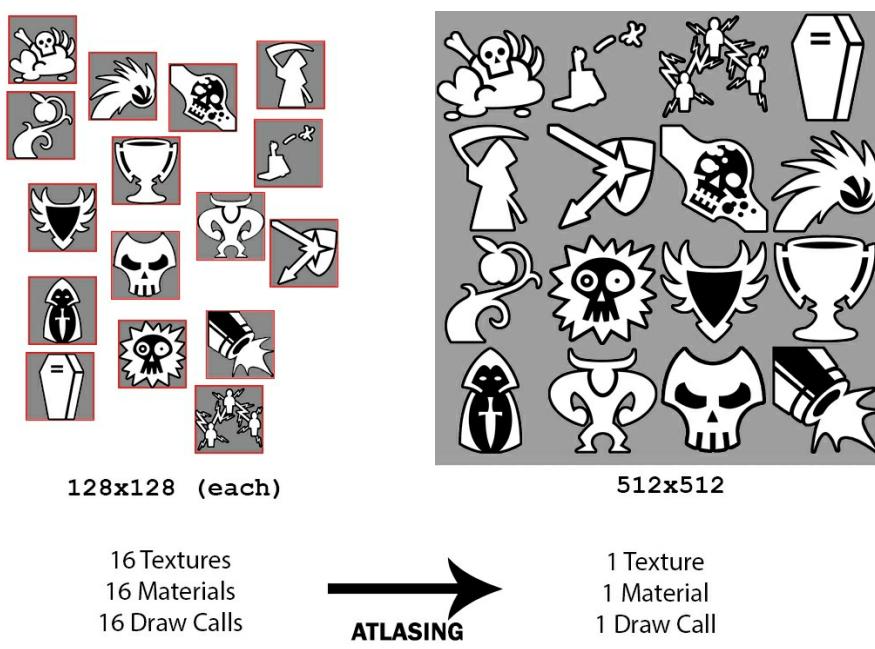
The strength of Anisotropic Filtering applied to the texture can be hand-modified on a per-texture basis with the Aniso Level setting as well as globally enabled/disabled using the Anisotropic Textures option within the Edit | Project | Quality settings.

Much like Mip Mapping, this effect can be costly and, sometimes, unnecessary. If there are textures in our Scene that we are certain will never be viewed at an oblique angle (such as distant background objects, UI elements, and billboard Particle Effect textures), then we can safely disable Anisotropic Filtering for them to save runtime overhead. We can also consider adjusting the strength of the Anisotropic Filtering effect on a per-texture basis to find the magic spot between quality and performance.

Consider Atlasing

Atlasing is the technique of combining lots of smaller, isolated textures together into a single, large texture file in order to minimize the number of Materials, and hence Draw Calls, we need to use. This is effectively a means to exploit Dynamic Batching. Conceptually, this technique is very similar to the approaches of minimizing Material usage you learned in [Chapter 3, The Benefits of Batching](#).

Each unique Material will require an additional Draw Call, but each Material only supports a single primary texture. Of course, they can also support multiple secondary textures, such as Normal Maps and Emission Maps. However, by combining multiple primary textures into a single large texture file, we can minimize the number of Draw Calls used to render objects that share this texture:



Extra work is required to modify the UV coordinates used by the mesh or sprite object to only sample the portion of the larger texture file that it needs, but the benefits are clear; reducing Draw Calls results in reduction of CPU workload and improvement in the frame rate if our application is

bottlenecked on CPU. Assuming that the merged texture file's resolution is equivalent to that of all of the combined images, there will be no loss of quality, and memory consumption will be essentially identical. Note that Atlasing does not result in reduced Memory Bandwidth consumption since the amount of data being pushed to the GPU would also be identical. It just happens to be bundled together in one bigger texture file.



Atlasing is only an option when all of the given textures require the same Shader. If some of the textures need unique graphical effects applied through Shaders, then they must be isolated into their own Materials and Atlased in separate groups.

Atlasing is a common tactic applied to UI elements and in games that feature a lot of 2D graphics. Atlasing becomes practically essential when developing mobile games with Unity since Draw Calls tend to be the most common bottleneck on those platforms. However, we would not want to generate these Atlas files manually. Life would be much simpler if we could continue to edit our textures individually and automate the task of combining them into a larger file.

Many GUI-related tools in the Unity Asset Store provide an automated texture-Atlasing feature. There are some stand-alone programs scattered across the internet, which can handle this work, and Unity can generate Atlases for Sprites in the form of Assets. These can be created via Asset | Create | Sprite Atlas.

Check out the Unity documentation to discover more about this useful feature at <https://docs.unity3d.com/Manual/SpriteAtlas.html>.



Note that the Sprite Atlas feature effectively supplants the Sprite Packer tool from older versions of Unity.

Atlasing does not need to be applied to 2D graphics and UI elements, either. We can apply this technique to 3D meshes if we happen to be creating a lot of low-resolution textures. 3D games that feature simple texture resolutions, or a flat-shaded low-poly art style, are ideal candidates for Atlasing in this way.

However, because Dynamic Batching affects only non-animated meshes (that is, `MeshRenderer`, but not `SkinnedMeshRenderer`), there is no reason to combine texture files for animated characters into an Atlas. Since they are animated, the GPU needs to multiply each object's bones by the transform of the current animation state. This means a unique calculation is needed for each character, and they will result in an extra Draw Call regardless of any attempts we make to have them share Materials.

As a result, combining textures for animated characters should be done only as a matter of convenience and space saving; for example, in a flat-shaded low-poly art style game, where everything happens to use a common color palette, we can make some space savings using a single texture for the entire game world, objects, and characters.

The disadvantages of Atlasing are mostly in terms of development time and workflow costs. It requires a lot of effort to overhaul an existing project to make use of Atlasing, which can be a lot of work just to figure out whether it is worth the effort or not. In addition, we need to be aware of generating texture files, which are too large for the target platform.

Some devices (specifically mobile devices) have a relatively low limit on the size of textures that can be pulled into the lowest memory cache of the GPU. If the Atlased texture file is too large, then it must be broken up into smaller textures in order to fit the target memory space. If the device's GPU happens to need textures from different pieces of the Atlas every other Draw Call, then we not only will inflict a lot of cache misses, but also might find that we choke the Memory Bandwidth, as textures are constantly pulled from VRAM and the lower-level cache.

We would probably not have this problem if the Atlas was left as individual textures. The same texture swapping will occur, but will result in much smaller files being swapped at the cost of additional Draw Calls. Our best options at this stage would be to lower the Atlas resolution or generate multiple smaller Atlases to have better control over how they will be dynamically batched.

Atlasing is clearly not a perfect solution, and if it is not clear whether it would result in a performance benefit, then we should be careful not to waste

too much time on its implementation. Speaking very generally, mobile games with a very simplistic 2D art style probably won't need to make use of Atlasing. However, mobile games attempting to compete with high quality assets or use any kind of 3D graphics should probably start integrating Atlasing from the very beginning of development, since it is likely that the project will reach texture throughput limits very quickly. They may even need to apply many per-platform and per-device optimizations in order to reach a wide audience.

Meanwhile, we should consider applying Atlasing to high-quality desktop games only if our Draw Call count exceeds reasonable hardware expectations since we will want many of our textures to maintain high resolutions for maximum quality. Low-quality Desktop games can probably afford to avoid Atlasing since Draw Calls are unlikely to be the biggest bottleneck.

Of course, no matter what the product is, if we're ever limited in CPU by too many Draw Calls and have already exhausted many of the alternative techniques, then Atlasing is a very effective performance enhancement in most cases.

Adjust compression rates for non-square textures

Texture files are normally stored in a square, *power-of-two* format, meaning that their height and width are equal in length, and its size is a power of two. For example, typical sizes are 256 x 256 pixels, 512 x 512, 1024 x 1024, and so on.

It is possible to provide rectangular power-of-two textures (such as 256 x 512) or those with a non-power-of-two format (such as 192 x 192), but creating textures such as these is not recommended. Some GPUs require square texture formats, so Unity will compensate by automatically expanding the texture to include additional empty space in order to fit the form factor that the GPU expects, which will result in additional Memory Bandwidth costs, pushing what is essentially unused and useless data to the GPU. Other GPUs may support non-power-of-two textures, but it is likely to result in slower sampling than a square texture.

So, the first recommendation is to avoid non-square and/or non-power-of-two textures altogether. If the image can be placed within a square, power-of-two texture, and does not result in too much quality degradation due to squeezing/stretching, then we should apply those changes just to keep the CPU and GPU happy. As a second option, we can customize this scaling behavior in Unity through the texture file's `Non Power of 2` import setting, though because this is an automated process, it might not give us the graphical quality we expect.

Sparse Textures

Sparse Textures, also known as **Mega-Textures** or **Tiled-Textures**, provide a way of effectively streaming texture data from disk at runtime. Relatively speaking, if the CPU performs operations in the order of seconds, then the disk would operate in the order of days. So, the common advice is that hard disk access during gameplay should be avoided as much as possible since any such technique risks inflicting more disk access than available, causing our application to grind to a halt.

However, Sparse Texturing offers some interesting performance-saving techniques if we're smart about starting data transfer for portions of the texture before we need them. Sparse Texturing is prepared by combining many textures into an enormous texture file, which would be far too large to load into graphics memory as a single texture file. This is similar to the concept of Atlasing, except the file containing the textures is incredibly large—for example, 32,768 x 32,768 pixels—and would contain considerable color detail, such as 32 bits per-pixel (this would result in a texture file that consumes 4 GBs of disk space). The idea is to save large amounts of runtime memory and Memory Bandwidth by hand-picking small subsections of the texture to load from disk dynamically, pulling them from the disk moments before they are needed in the game. The main cost of this technique is the file size requirement and the potentially continuous disk access. Other costs for this technique can be overcome, but normally take a great deal of Scene preparation work.

The game world needs to be created in such a way that it minimizes the amount of texture swapping taking place. In order to avoid very noticeable *texture popping* problems, texture subsections must be pulled from a disk into RAM with just enough time to spare that the GPU does not need to wait before the transfer to VRAM can begin (in much the same way that it normally doesn't need to wait for ordinary texture files that are preloaded into RAM). This takes place in the design of the texture file itself by keeping common elements for a given Scene in the same general area of the texture, and the design of the Scene, by triggering new texture subsection loading at

key moments during gameplay and making sure that disk access of the new tile is quickly located by the disk without extreme cache misses. If it is handled with care, then Sparse Texturing can result in impressive benefits in both Scene quality and memory savings.

It is a highly specialized technique in the gaming industry and has not yet been widely adopted partly because it requires specialized hardware and platform support and partly because it is difficult to pull it off well. The Unity documentation on Sparse Texturing has improved somewhat over time and provides an example Scene showing the effect at work, which can be found at <http://docs.unity3d.com/Manual/SparseTextures.html>.

For Unity developers who consider themselves advanced enough to experiment with Sparse Texturing, it might be worth taking the time to perform some research to check whether Sparse Texturing is right for their project since it promises some significant performance savings.

Procedural Materials

Procedural Materials, also known as **Substances**, are a means of procedurally generating textures at runtime by combining small, high-quality texture samples with custom mathematical formulas. The goal of Procedural Materials is to greatly minimize the application disk footprint at the cost of additional runtime memory and CPU processing during initialization to generate the texture via mathematical operations rather than static color data.

Texture files are, sometimes, the biggest disk space consumer of a game project, and it's fairly common knowledge that download times have a tremendous negative impact on the completed download rate and getting people to try our game (even if it's free). Procedural Materials offer us the ability to sacrifice some initialization and runtime processing power for much faster downloads. This is very important for mobile games that are trying to compete via graphical fidelity.

The Unity documentation on Procedural Materials is fairly extensive, so it is recommended to work through the documentation for a clearer picture of how Substances work and how they can provide us with performance benefits. Check out the following Unity documentation page for more information on Procedural Materials: <http://docs.unity3d.com/Manual/ProceduralMaterials.html>.

Asynchronous Texture Uploading

The last texture import option we haven't covered is the Read/Write Enabled option. By default, this option is disabled, which is good, because this allows textures to make use of the Asynchronous Texture Uploading feature, which allows two benefits: the texture will be uploaded asynchronously from disk to RAM, and when the texture data is needed by the GPU, the transfer happens on the render thread, not the main thread. Textures will be pushed into a circular buffer, which pushes data to the GPU continuously so long as the buffer contains new data. If not, then it early-exits the process and waits until new texture data is requested.

Ultimately, this reduces the time spent preparing the Render States for each frame and allows more CPU resources to be spent on gameplay logic, the Physics Engine, and so on. Of course, some time is still spent on the main thread preparing the Render State, but moving the texture uploading task to a separate thread saves a significant chunk of CPU time on the main thread.

However, enabling read/write access to the texture essentially tells Unity we might be reading and editing this texture at any time. This implies that the GPU will need fresh access to it every time, so it will disable Asynchronous Texture Uploading for that texture; all uploading must occur on the main thread. We might want to enable this option for things such as simulating painting colors onto a canvas or writing image data from the internet into a premade texture, but the downside is that the GPU must always wait for any changes to be made to the texture before it can be uploaded since it cannot predict when those changes will happen.

In addition, Asynchronous Texture Uploading only works for textures we explicitly imported into the project and that were present during build time, since the feature only works if the texture was packed together into special streamable assets. Therefore, any textures generated via `LoadImage(byte[])`, texture assets imported/downloaded from external locations, or loaded from a

Resources folder via `Resources.Load()` (which all implicitly call `LoadImage(byte[])` themselves) will not be converted into streamable content, and thus will be unable to make use of Asynchronous Texture Uploading.

It is possible to tweak both the upper limit on the maximum allowed time to spend on Asynchronous Texture Uploads and the total circular buffer size Unity should use to push the textures we want to upload. These settings can be tweaked under *Edit | Project Settings | Quality | Other* and are named *Async Upload Time Slice* and *Async Upload Buffer Size*, respectively. We should set the *Async Upload Time Slice* value to the maximum number of milliseconds we want Unity to spend on Asynchronous Texture Uploads on the render thread. It might be wise to set the *Async Upload Buffer Size* value to the largest texture file we might need to use, plus a little extra buffer if multiple fresh textures are needed in the same frame. The circular buffer that texture data is copied into will expand as needed, but this is often costly. Since we probably already know ahead of time how large we need that circular buffer to be, we might as well set it to the maximum expected size to avoid potential frame drops when it needs to resize the buffer.

Mesh and animation files

Finally, let's cover mesh and animation files. These file types are essentially large arrays of vertex and skinned bone data, and there are a variety of techniques we can apply to minimize file size, while keeping a similar, if not identical, appearance. There are also ways to lower the cost of rendering large groups of these objects through batching techniques. Let's take a look at a series of performance-enhancing techniques we can apply to such files.

Reduce polygon count

This is the most obvious way to gain performance and should always be considered. In fact, since we cannot batch objects using Skinned Mesh Renderers, it's one of the good ways of reducing CPU and GPU runtime overhead for animated objects.

Reducing the polygon count is simple, straightforward, and provides both CPU and memory cost savings for the time required for artists to clean up the mesh. Much of an object's detail is provided almost entirely by detailed texturing and complex shading in this day and age, so we can often get away with stripping away a lot of vertices on modern meshes and most users would be unable to tell the difference.

Tweak Mesh Compression

Unity offers four different Mesh Compression settings for imported mesh files: Off, Low, Medium, and High. Increasing this setting will convert floating-point data into fixed values, reducing the accuracy in vertex position/Normal direction, simplifying vertex color information, and so on. This can have a noticeable effect on meshes that contain lots of small parts near one another, such as a fence or grate. If we're generating meshes procedurally, we can achieve the same type of compression by calling the `Optimize()` method of a `MeshRenderer` Component (of course, this will take some time to complete).

There are also two global settings found in `Edit | Project Settings | Player | Other Settings`, which can affect how mesh data is imported. These are named Vertex Compression and Optimize Mesh Data.

We can use the Vertex Compression option to configure the type of data that will be optimized when we import a mesh file with Mesh Compression enabled, so if we want accurate Normal data (for lighting), but have less worry over positional data, then we can configure it here. Unfortunately, this is a global setting and will affect all imported meshes (although it can be configured on a per-platform basis since it is a Player setting).

Enabling Optimize Mesh Data will strip away any data from the mesh that isn't required by the Material(s) assigned to it. So, if the mesh contains tangent information, but the Shader never requires it, then Unity will ignore it during build time.

In each case, the benefits are reducing the application's disk footprint at the cost of extra time loading the mesh since extra time must be spent decompressing the data before it's needed.

3D mesh building/animation tools often provide their own built-in ways of automated mesh optimization in the form of estimating the overall shape and stripping the mesh down to



fewer total polygons. This can cause significant loss of quality and should be tested vigorously if used.

Use Read-Write Enabled appropriately

The Read-Write Enabled flag allows changes to be made to the mesh at runtime either via Scripting or automatically by Unity during runtime, similar to how it is used for texture files. Internally, this means that it will keep the original mesh data in memory until we want to duplicate it and make changes dynamically. Disabling this option will allow Unity to discard the original mesh data from memory once it has determined the final mesh to use, since it knows it will never change.

If we use only a uniformly scaled version of a mesh throughout the entire game, then disabling this option will save runtime memory since we will no longer need the original mesh data to make further rescaled duplicates of the mesh (incidentally, this is how Unity organizes objects by scale factor when it comes to Dynamic Batching). Unity can, therefore, discard this unwanted data early since we will never need it again until the next time the application is launched.

However, if the mesh often reappears at runtime with different scales, then Unity needs to keep this data in memory so that it can recalculate a new mesh more quickly, hence it would be wise to enable the Read-Write Enabled flag. Disabling it will require Unity to not only reload the mesh data each time the mesh is reintroduced, but also make the rescaled duplicate at the same time, causing a potential performance hiccup.

Unity tries to detect the correct behavior for this setting at initialization time, but when meshes are instantiated and scaled in a dynamic fashion at runtime, we must force the issue by enabling this setting. This will improve instantiation speed of the objects, but cost some memory overhead since the original mesh data is kept around until it's needed.



Note that this potential overhead cost also applies when using the Generate Colliders option.

—

Consider baked animations

This tip will require changes in the asset through the 3D rigging and animation tool that we are using since Unity does not provide such tools itself. Animations are normally stored as key frame information where it keeps track of specific mesh positions and interpolates between them at runtime using skinning data (bone shapes, assignments, animation curves, and so on). Meanwhile, baking animations means effectively sampling and hardcoding each position of each vertex per-frame into the mesh/animation file without the need for interpolation and skinning data.

Using baked animations can, sometimes, result in much smaller file sizes and memory overhead than blended/skinned animations for some objects since skinning data can take up a surprisingly large amount of space to store. This is most likely to be the case for relatively simple objects, or objects with short animations, since we would effectively be replacing procedural data with a hardcoded series of vertex positions. So, if the mesh's polygon count is low enough where storing lots of vertex information is cheaper than skinning data, then we may see some significant savings through this simple change.

In addition, how often the baked sample is taken can usually be customized by the exporting application. Different sample rates should be tested to find a good value where the key moments of the animation still shine-through what is essentially a simplified estimate.

Combine meshes

Forcefully combining meshes into a large, single mesh can be convenient to reduce Draw Calls, particularly if the meshes are too large for Dynamic Batching and don't play well with other statically batched groups. This is essentially the equivalent of Static Batching, but performed manually, so sometimes it's wasted effort if Static Batching could take care of the process for us.

Beware that if any single vertex of the mesh is visible in the Scene, then the entire object will be rendered together as one whole. This can lead to a lot of wasted processing if the mesh is only partially visible most of the time. This technique also comes with the drawback that it generates a whole new mesh asset file that we must deposit into our Scene, which means any changes we make to the original meshes will not be reflected in the combined one. This results in a lot of tedious workflow effort every time changes need to be made, so if Static Batching is an option, it should be used instead.

There are several tools available online, which can combine mesh files together for us in Unity. They are only an Asset Store or Google search away.

Asset Bundles and Resources

We touched upon the topic of Resources and Serialization in [Chapter 2](#), *Scripting Strategies*, and it should be fairly clear that the Resource System can be a great benefit during prototyping, as well as during the early stages of our project, and can be used relatively effectively in games of limited scope.

However, professional Unity projects should instead favor the Asset Bundle System. There are a number of reasons for this. Firstly, the Resource System is not very scalable when it comes to builds. All Resources are merged together into a single massive Serialized File binary data blob with an index list of where various assets can be found within it. This can be hard to manage, and take a long time to build, as we add more data to the list.

Secondly, the Resource System's ability to acquire data from the Serialized File scales in an $N \log(N)$ fashion, which should make us very wary of increasing the value of N . Thirdly, the Resource System makes it unwieldy for our application to provide different asset data on a per-device basis, whereas Asset Bundles tend to make this matter trivial. Finally, Asset Bundles can be used to provide small, periodic custom content updates to the application, while the Resource System would require updates that completely replace the entire application to achieve the same effect.

Asset Bundles share a lot of common functionality with Resources such as loading from files, loading data asynchronously, and unloading data we no longer need. However, they also offer much more functionality such as content streaming, content updates, and content generation and sharing. These can all be used to improve the performance of our application to great effect. We can deliver applications with much smaller disk footprints and have the user download additional content before or during gameplay, stream assets at runtime to minimize the initial loading time of the application, and provide more optimized assets to the application on a per-platform basis without the need to push a complete application overwrite to the user.

Of course, there are downsides to Asset Bundles. They are much more complicated to set up and maintain than Resources, they're more complicated to understand since they use a much more sophisticated system for accessing asset data than the Resources System, and making full use of their functionality (such as streaming and content updates) would require a lot of additional QA testing to make sure that server is delivering content properly, and that the game is reading and updating its content to match. Ergo, Asset Bundles are best used only when our team size is able to support the extra workload they require.

A tutorial on the Asset Bundle system is beyond the scope of this book, but there are dozens of useful guides online and in the Unity Documentation.

Check out the following Unity tutorial to find out more about the Asset Bundle system:

<https://unity3d.com/learn/tutorials/topics/best-practices/guide-assetbundles-and-resources>.

If you require further convincing, then the following Unity blog post from April 2017 should help reveal how the Asset Bundle system can use memory more efficiently during runtime in ways that the Resources System cannot provide through memory pooling: <https://blogs.unity3d.com/2017/04/12/asset-bundles-vs-resources-a-memory-showdown/>.

Summary

There are many different opportunities that we can explore to achieve performance gains for our application just by tinkering with our imported assets. Alternatively, from another perspective, there are plenty of ways to ruin our application's performance through asset mismanagement. Almost every single import configuration opportunity is a trade-off between one performance metric or workflow task and another. Typically, this means saving disk footprint via compression at the expense of CPU at runtime to decompress the data, or faster access while reducing the quality level of the final presentation. So, we must remain vigilant and only pick the right techniques for the right assets for the right reasons.

This concludes our exploration of improving performance through art asset manipulation. In the next chapter, we will be investigating how to improve our usage of Unity's Physics Engine.

Faster Physics

Each of the performance-enhancing suggestions we've explored thus far have been primarily centered on reducing resource costs and avoiding frame rate issues. However, at its most fundamental level, seeking peak performance means improving the user experience. This is because every frame rate hiccup, every crash, and every system requirement that is too costly for a given market ultimately detracts from the quality of the product. Physics Engines are in a unique category of subsystems whose behavior and consistency contributes a major factor toward product quality, and spending the time to improve its behavior is often worth the cost.

If important collision events get missed, the game freezes while it calculates a complex physics event, or the player falls through the floor, then these have an obvious and significant negative impact on the quality of gameplay. A few glitches are often bearable, but continuous problems will get in the way of gameplay. This often results in pulling the player out of the experience, and it's a coin-toss whether the user finds it inconvenient, obnoxious, or hilarious. Unless our game is specifically targeting the Comedy Physics genre (games such as QWOP or Goat Simulator), these are situations we should strive to avoid.

Some games may not use physics at all, whereas others require the Physics Engine to handle a huge number of tasks during gameplay, such as collision detection between hundreds of objects, Trigger Volumes to initiate cutscenes, Raycasting for player attacks and UI behavior, gathering lists of objects in a given region, or even just using physics as eye candy with lots of physical particles flying around. Its importance also varies depending on the type of game being created. For example, it is essential in platformer and action games to tune the physics properly--how the player character reacts to input and how the world reacts to the player character are two of the most important aspects that make the game feel responsive and fun. Whereas, accurate physics may be somewhat less important in **Massively Multiplayer Online (MMO)** games, which tend to have limited physics interaction.

Therefore, in this chapter, we will cover ways to reduce CPU spikes, overhead, and memory consumption through Unity's Physics Engine, but also include ways to alter physics behavior with the aim of improving, or at least maintaining, gameplay quality while optimizing performance. In this chapter, we will cover the following areas:

- Understanding how Unity's Physics Engine works:
 - Timesteps and FixedUpdatees
 - Collider types
 - Collisions
 - Raycasting
 - Rigidbody active states
- Physics performance optimizations:
 - How to structure Scenes for optimal physics behavior
 - Using the most appropriate types of Collider
 - Optimizing the Collision Matrix
 - Improving physics consistency and avoiding error-prone behavior
 - Ragdolls and other Joint-based objects

Physics Engine internals

Unity technically features two different Physics Engines: Nvidia's PhysX for 3D physics, and the Open Source project Box2D for 2D physics. However, their implementations are highly abstracted, and from the perspective of the higher-level Unity API that we configure through the main Unity Engine, both Physics Engine solutions operate in a functionally identical fashion.

In either case, the more we understand about Unity's Physics Engines, the more sense we can make of possible performance enhancements. So, first we'll cover some theory about how Unity implements these systems.

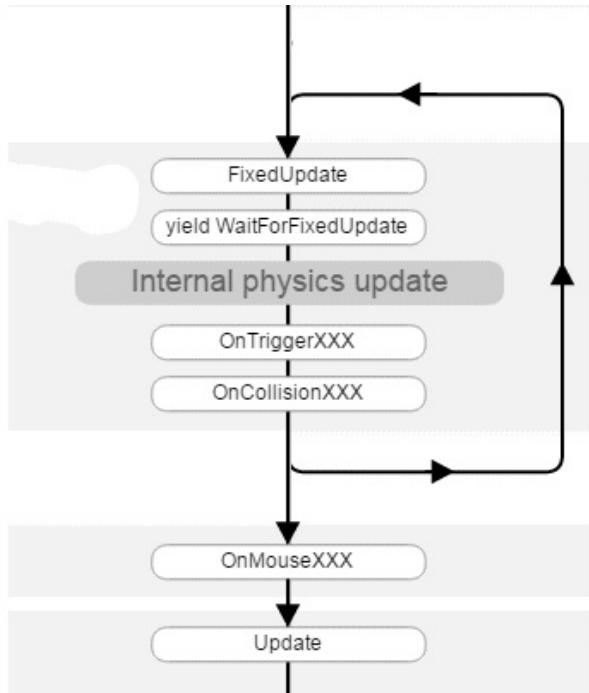
Physics and time

Physics Engines generally operate under the assumption that time advances by fixed values, and both of Unity's Physics Engines operate in this manner. Each of these iterations is known as the **Timestep**. The Physics Engine will only resolve each Timestep using very specific values of time, which is independent of how much time it took to render the previous frame. This is known in Unity as the Fixed Update Timestep, and it is set to a value of 20 milliseconds by default (50 updates per second).



It can be very difficult to generate consistent results for collisions and forces between two different computers if a Physics Engine uses a variable timestep due to differences in architecture (in how floating-point values are represented) as well as latency between clients. Such Physics Engines tend to generate very inconsistent results between multiplayer clients or during recorded replays.

The following diagram shows an important snippet of the Unity Order of Execution diagram:



The full Order of Execution diagram can be found at <http://docs.unity3d.com/Manual/ExecutionOrder.html>.

As we can see in the preceding diagram, Fixed Updates are processed just before the Physics Engine performs its own update and the two are inextricably linked. The process begins with determining whether enough time has passed to begin the next Fixed Update. Once this is determined, the outcome of resolving it will vary, depending on how much time has passed since the last Fixed Update.

If enough time has passed, then the Fixed Update processes will invoke all `FixedUpdate()` callbacks defined across all active MonoBehaviours in the Scene, followed by any Coroutines tied to Fixed Updates (specifically those that `yield` to `waitForFixedUpdate`). Note that there is no guarantee of execution order for methods invoked during either of these processes, so we should never write code under this assumption. Once these tasks are done, the Physics Engine can begin to process the current timestep and invoke any necessary Trigger and Collider callbacks.

Conversely, if too little time has passed since the last Fixed Update (that is, less than 20 milliseconds), then the current Fixed Update is skipped, and all of the tasks listed previously don't happen during the current iteration. At this

point, input, gameplay logic, and rendering will be allowed to happen as normal. Once this activity is complete, Unity checks whether the next Fixed Update is required.

At high frame rates, rendering updates are likely to complete multiple times before the Physics Engine gets a chance to update itself. This process repeats itself during runtime and effectively gives Fixed Updates and the Physics Engine a higher priority over rendering while also forcing the physics simulation into a fixed frame rate.



In order to ensure that objects move smoothly between Fixed Updates, Physics Engines (including Unity's) interpolate the visible location of each object between where it was during the previous state and where it should be after resolving the current state based on how much time remains until the next Fixed Update. This interpolation ensures that objects appear to move smoothly despite the fact that their physical positions, velocities, and so on are being updated less frequently than the render frame rate.

The `FixedUpdate()` callback is a useful place to define any gameplay behavior that we want to be frame-rate independent. AI calculations are commonly resolved in Fixed Updates since they tend to be easier to work with if we assume a fixed update frequency.

Maximum Allowed Timestep

It is important to note that if a lot of time has passed since the last Fixed Update (for example, the game froze momentarily), then Fixed Updates will continue to be calculated within the same Fixed Update loop until the Physics Engine has *caught up* with the current time. For example, if the previous frame took 100 ms to render (for example, a sudden CPU spike caused the main thread to block for a long time), then the Physics Engine will need to be updated five times. The `FixedUpdate()` method will, therefore, be called five times before `Update()` can be called again due to the default Fixed Update Timestep of 20 milliseconds. Of course, if there is a lot of physics activity to process during these five Fixed Updates, such that it takes more than 20 milliseconds to resolve them all, then the Physics Engine will need to invoke a sixth update.

Consequently, it's possible during moments of heavy physics activity that the Physics Engine takes more time to process a Fixed Update than the amount of time it is simulating. For example, if it took 30 ms to process a Fixed Update simulating 20 ms of Gameplay, then it has fallen behind, requiring it to process more Timesteps to try and keep up, but this could cause it to fall behind even further, requiring it to process even more Timesteps, and so on. In these situations the Physics Engine is never able to escape the Fixed Update loop and allow another frame to render. This problem is often known as the **spiral of death**. However, to prevent the Physics Engine from locking up our game during these moments, there is a maximum amount of time that the Physics Engine is allowed to process each Fixed Update loop. This threshold is called the **Maximum Allowed Timestep**, and if the current batch of Fixed Updates takes too long to process, then it will simply stop and forgo further processing until the next render update completes. This design allows the Rendering Pipeline to at least render the current state and allow for user input and gameplay logic to make some decisions during rare moments where the Physics Engine has gone ballistic (pun intended).

This setting can be accessed through Edit | Project Settings | Time | Maximum Allowed Timestep.

Physics updates and runtime changes

When the Physics Engine processes a given timestep, it must move any active Rigidbody objects (`GameObjects` with a `Rigidbody` Component), detect any new collisions, and invoke the collision callbacks on the corresponding objects. The Unity documentation makes an explicit note that changes to Rigidbody objects should be handled within `FixedUpdate()` and other physics callbacks for exactly this reason. These methods are tightly coupled with the update frequency of the Physics Engine as opposed to other parts of the Game Loop, such as `Update()`.

This means that callbacks such as `FixedUpdate()` and `OnTriggerEnter()` are safe places to make Rigidbody changes, whereas methods such as `Update()` and Coroutines yielding on `WaitForSeconds` or `WaitForSecondsFrame` are not. Ignoring this advice could cause unexpected physics behavior, as multiple changes may be made to the same object before the Physics Engine is given a chance to catch and process all of them.

It's particularly dangerous to apply forces or impulses to objects in `Update()` callbacks without taking into account the frequency of those calls. For instance, applying a 10-Newton force each Update while the player holds down a key would result in completely different resultant velocity between two different devices than if we did the same thing in Fixed Update since we can't rely on the number of `Update()` calls being consistent. However, doing so in a `FixedUpdate()` callback will be much more consistent. Therefore, we must ensure that all physics-related behavior is handled in the appropriate callbacks or we will risk introducing some especially confusing gameplay bugs that are very hard to reproduce.

It logically follows that the more time we spend in any given Fixed Update iteration, the less time we have for the next gameplay and rendering pass. Most of the time this results in minor, unnoticeable background processing tasks, since the Physics Engine barely has any work to do, and

the `FixedUpdate()` callbacks have a lot of time to complete their work. However, in some games, the Physics Engine could be performing a lot of calculations during each and every Fixed Update. This kind of bottlenecking in physics processing time will affect our frame rate, causing it to plummet as the Physics Engine is tasked with greater and greater workloads. Essentially, the Rendering Pipeline will try to proceed as normal, but whenever it's time for a Fixed Update, in which the Physics Engine takes a long time to process, the Rendering Pipeline would have very little time to generate the current display before the frame is due, causing a sudden stutter. This is in addition to the visual effect of the Physics Engine stopping early because it hit the Maximum Allowed Timestep. All of this together would generate a very poor user experience.

Hence, in order to keep a smooth and consistent frame rate, we will need to free up as much time as we can for rendering by minimizing the amount of time the Physics Engine takes to process any given timestep. This applies in both the best-case scenario (nothing moving) and worst-case scenario (everything smashing into everything else at once). There are a number of time-related features and values we can tweak within the Physics Engine to avoid performance pitfalls such as these.

Static Colliders and Dynamic Colliders

There is a rather extreme namespace conflict with the terms *static* and *dynamic* in Unity. When *static* is used, it normally means that the object or process under discussion is not moving, remains unchanged, or exists in only one location, whereas *dynamic* simply means the opposite--objects or processes that tend to move or change. However, it's important to remember that each of these are separate topics and usage of the terms *static* and *dynamic* mean something different in each case. We have already introduced the Static subflags for `GameObjects`, the Dynamic Batching and Static Batching systems, and the concepts of Static classes, static variables, and static functions in the C# language. So, just to be extra-confusing, Unity also has the concept of Static and Dynamic Colliders.

Dynamic Colliders simply mean `GameObjects` that contain both a `Collider` Component (which could be one of several types) and a `Rigidbody` Component. By attaching a `Rigidbody` to the same object as a `Collider`, the Physics engine will treat that `Collider` as the bounding volume of a physical object that must react to outside forces (such as gravity) and collisions with other `Rigidbodies`. If we collide one Dynamic Collider into another, they will both react based on Newton's Laws of Motion (or at least as best as a computer using floating-point arithmetic is capable of).

We can also have Colliders that do not have a `Rigidbody` Component attached, and these are called Static Colliders. These effectively work as invisible barriers that Dynamic Colliders can collide into, but the Static Collider will not react in response. To think of it another way, imagine objects without a `Rigidbody` Component as having infinite mass. No matter how hard you throw a rock into an object of infinite mass, it will never move, but you can still expect the rock to react like it just hit a solid wall. This makes Static Colliders ideal for world barriers and other obstacles that must not move.

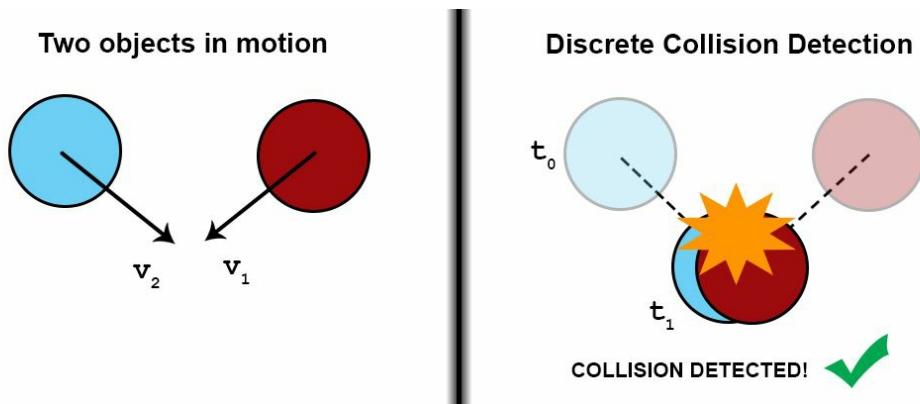
The Physics Engine automatically separates Dynamic Colliders and Static

Colliders into two different data structures, each optimized to handle the types of Collider present. This helps to simplify future processing tasks since, for example, there's no point resolving collisions and impulses between two Static Colliders.

Collision detection

There are three settings for collision detection in Unity, which can be configured in a `Rigidbody` Component's Collision Detection property: Discrete, Continuous, and ContinuousDynamic. The Discrete setting enables Discrete collision detection, which effectively teleports objects a small distance every timestep based on their velocity and how much time has passed. Once all the objects have been moved, it then performs a bounding volume check for any overlaps, treats them as collisions, and resolves them based on their physical properties and how they overlap. This method risks collisions being missed if small objects move too quickly.

The following image shows how Discrete collision detection works to catch two objects as they teleport from one location to the next:

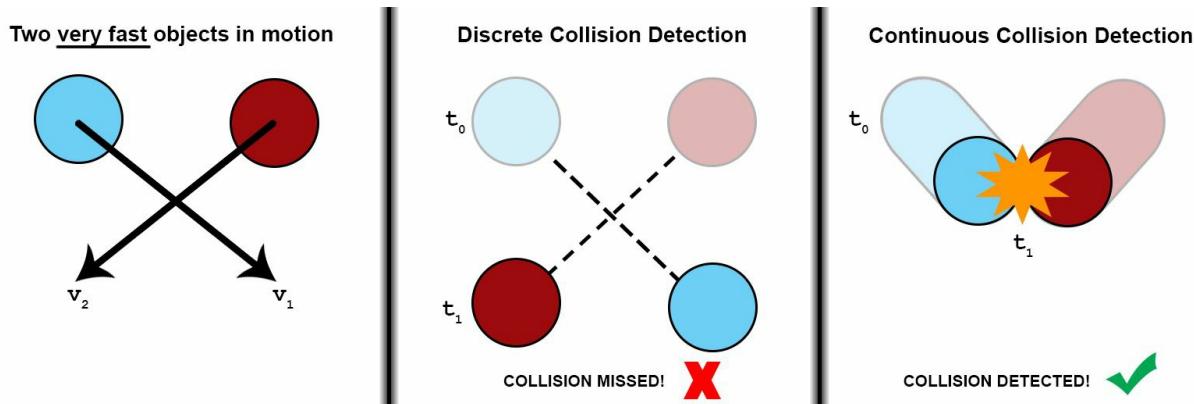


Either of the remaining settings will enable Continuous collision detection, which works by interpolating Colliders from their starting and ending positions for the current timestep and checking for any collisions along the way. This reduces the risk of missed collisions and generates a more accurate simulation at the expense of a significantly greater CPU overhead compared to Discrete collision detection.

The Continuous setting enables Continuous collision detection only between the given Collider and Static Colliders. Collisions between the same Collider and Dynamic Colliders will still make use of Discrete collision detection. Meanwhile, the ContinuousDynamic setting enables Continuous collision

detection between the Collider and all Static and Dynamic Colliders, making it the most expensive in terms of resources.

The following diagram shows how the Discrete and Continuous collision detection methods work for a pair of small, fast-moving objects:



This is an extreme example for the sake of illustration. In the case of Discrete collision detection, we can observe that the objects are teleporting a distance around four times their own size in a single timestep, which would typically only happen with very small objects with very high velocities, and is, hence, very rare if our game is running optimally. In the overwhelming majority of cases, the distances the objects travel in a single 20-milliseconds timestep are much smaller relative to the size of the object, and so the collision is easily caught by Discrete collision detection methods.

Collider types

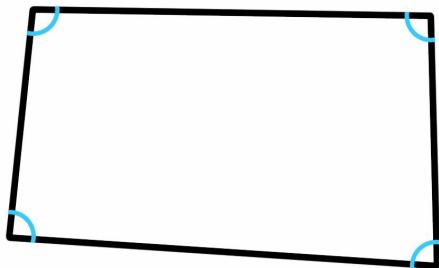
There are four different types of 3D Colliders in Unity. In order of the lowest performance cost to the greatest, they are Sphere, Capsule, Box, and Mesh Colliders. The first three Collider types are often called *primitives* and maintain very specific shapes, although they can generally be scaled in different directions to meet certain needs. Mesh Colliders can, however, be customized to a particular shape depending on the assigned mesh. There are also three types of 2D Colliders--Circle, Box, and Polygon--that are functionally similar to Sphere, Box, and Mesh Colliders, respectively. All of the following information is essentially transferable to the equivalent 2D shape.



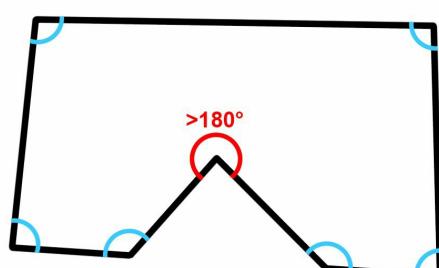
Note that we can also generate cylindrical 3D objects in Unity, but this is only for its graphical representation. Autogenerated Cylinder shapes use Capsule Colliders to represent their physical bounding volume, which may not create the expected physics behavior.

In addition, there are two varieties of Mesh Collider: **Convex** and **Concave**. The difference is that a Concave shape features at least one internal angle (an angle between two inside edges of the shape) of greater-than 180 degrees. To illustrate this, the following diagram shows the difference between Convex and Concave shapes:

Convex



Concave





An easy way to remember the difference between a Convex and Concave shape is that a Concave shape has at least one cave within it.

Both Mesh Collider types use the same Component (a `MeshCollider` Component). The type of Mesh Collider that gets generated is toggled using the Convex checkbox. Enabling this option will allow the object to collide with all primitive shapes (Spheres, Boxes, and so on) as well as other Mesh Colliders with Convex enabled. In addition, if the Convex checkbox is enabled for a Mesh Collider with a Concave shape, then the Physics Engine will automatically simplify it, generating a Collider with the nearest Convex shape it can. In the preceding example, if we imported the Concave mesh on the right and enable the Convex checkbox, it would generate a Collider shape closer to the Convex shape on the left. In either case, the Physics Engine will attempt to generate a Collider that matches the shape of the attached mesh with an upper limit of 255 vertices. If the target mesh has more vertices than this, it will throw an error during mesh generation.

Collider Components also contain the `isTrigger` property, allowing them to be treated as nonphysical objects, but still invoke physics events when other Colliders enter or leave them. These are called Trigger Volumes. Normally, a Collider's `onCollisionEnter()`, `onCollisionStay()`, and `onCollisionExit()` callbacks are called when another Collider touches, keeps touching (each timestep), or stops touching it, respectively. However, when the Collider is used as a Trigger Volume, the `onTriggerEnter()`, `OnTriggerStay()`, and `OnTriggerExit()` callbacks will be used instead.



Note that due to the complexity of resolving inter-object collisions, Concave Mesh Colliders cannot also be Dynamic Colliders. Concave shapes can only be used as Static Colliders or Trigger Volumes. If we attempt to add a `Rigidbody` Component to a Concave Mesh Collider, Unity will simply ignore it.

The Collision Matrix

The Physics Engine features a Collision Matrix that defines which objects are allowed to collide with which other objects. Objects that do not fit this matrix are automatically ignored by the Physics Engine when the time comes to resolve bounding volume overlaps and collisions. This saves on physics processing during collision detection stages and also allows the objects to move through one another without any collisions taking place.

The Collision Matrix can be accessed through `Edit | Project Settings | (Physics / Physics2D) | Layer Collision Matrix`.

The Collision Matrix system works through Unity's Layer system. The matrix represents every possible Layer-to-Layer combination that might be possible, and enabling a checkbox means that Colliders in both of those Layers will be checked during the collision detection phase. Note that there's no way to allow only one of the two objects to respond to the collision. If one Layer can collide with another, then they must both respond to the collision. However, Static Colliders are an exception since they aren't allowed to respond physically to collisions (although they still receive the `onCollision...()` callbacks).

Note that we are limited to only 32 total Layers for our entire project (since the Physics Engine uses a 32-bit bitmask to determine inter-Layer collision opportunities), so we must organize our objects into sensible Layers that will extend throughout the entire lifetime of the project. If, for whatever reason, 32 Layers are not enough for our project, then we might need to find cunning ways to reuse Layers or remove Layers that aren't necessary.

Rigidbody active and sleeping states

Every modern Physics Engine shares a common optimization technique, whereby objects that have come to rest have their internal state changed from an active state to a sleeping state. While a Rigidbody is in the sleeping state, little-to-no processor time will be spent during Fixed Updates to update the object until it has been awoken by an external force or collision event.

The value of measurement that is used to determine what *at rest* means tends to vary among different Physics Engines; it could be calculated using linear and rotational speed, kinetic energy, momentum, or some other physical properties of the Rigidbody. Both of Unity's Physics Engines work by evaluating the object's *mass-normalized-kinetic-energy*, which essentially boils down to *the magnitude of its velocity-squared*.

If the object's velocity has not exceeded some threshold value after a short time, then the Physics Engine will assume that the object will no longer need to move again until it has undergone a new collision, or a new force has been applied to it. Until then, the sleeping object will maintain its current position. Setting the threshold value too low would mean objects are much less likely to go to sleep, so we will keep paying a small processing cost within the Physics Engine every Fixed Update, even though it is not doing anything important. Meanwhile, setting the threshold value too high would mean slow-moving objects will appear to jerk to a sudden stop once the Physics Engine decides that they need to go to sleep. The threshold value that controls the sleeping state can be modified under Edit | Project Settings | Physics | Sleep Threshold. We can also get a count of the total number of active Rigidbody objects from the Physics Area of the Profiler window.

Note that sleeping objects are not removed entirely from the simulation. If a moving Rigidbody approaches the sleeping object, then it must still perform checks to see whether nearby objects have collided with it, which would reawaken the sleeping object, reintroducing it to the simulation for

processing.

Ray and object casting

Another common feature of Physics Engines is the ability to *cast a ray* from one point to another and generate collision information with one or more of the objects in its path. This is known as **Raycasting**. It is pretty common to implement several gameplay mechanics through Raycasting, such as firing a gun. This is typically implemented by performing Raycasts from the player to the target location and finding any viable targets in its path (even if it's just a wall).

We can also obtain a list of targets within a finite distance of a fixed point in space using a `Physics.OverlapSphere()` check. This is typically used to implement area-of-effect gameplay features, such as grenade or fireball explosions. We can even cast entire objects forward in space using `Physics.SphereCast()` and `Physics.CapsuleCast()`. These methods are often used to simulate wide laser beams, or if we simply want to see what would be in the path of a moving character.

Debugging Physics

Physics bugs usually fall into two categories: an object pair collided/didn't collide when it shouldn't have/should have or the objects collided but something unexpected happened after the fact. The former case is generally easier to debug; it is often either due to mistakes in the Collision Matrix, incorrect Layers used in Raycasting, or object Colliders being the wrong size or shape. The latter case is often much more difficult to resolve because of three big problems:

- Determining which collided objects caused the issue
- Determining the conditions of the collision just prior to resolution
- Reproducing the collision

Any of these three pieces of information would make resolution much easier, but they can all be difficult to obtain in some circumstances.

The Profiler provides some measure of information in the Physics and Physics (2D) Areas (for 3D and 2D physics, respectively), which can be moderately useful. We can get a measure of how much CPU activity is being spent on all Rigidbodies and groups of Rigidbodies isolated to different types such as Dynamic Colliders, Static Colliders, Kinematic objects, Trigger Volumes, Constraints (used to simulate hinges and other connected physics objects), and Contacts. The Physics 2D Area contains a little more information such as the number of sleeping and active Rigidbodies and how long processing the Timestep took. The Detailed Breakdown View provides even more information in both cases. This information helps us keep an eye on Physics performance, but it doesn't tell us much about what went wrong in the event we find a bug in our physics behavior.

A tool that is better suited to the task of helping us debug physics issues is the Physics Debugger, which can be opened via Window | Physics Debugger. This tool can help us filter out different types of Colliders from the Scene window to give us a better idea of which objects collided with one another. Of course, this does not help too much in determining the conditions of the

problem and reproducing issues.



Note that settings in the Physics Debugger do not affect object visibility in the Game window.

Unfortunately, there isn't much secret advice to be given for the remaining problems. Catching information about collisions before or when they happen normally involves a lot of targeted breakpoints in an `OnCollisionEnter()` or `OnTriggerEnter()` callback to catch the problem in the act and using step-through debugging until the source of the issue becomes clear. As a last resort, we can add `Debug.Log()` statements to log important information just before the problem occurs, although this can be a frustrating exercise because we sometimes don't know what information we need to log, or which objects to log from, and so we end up adding logs to everything.

Another frequent source of headaches is trying to reproduce physics problems. Reproducing collisions is always a challenge due to the nondeterministic nature between user input (normally handled in `Update()`) and physics behavior (handled in `FixedUpdate()`). Even though physics timesteps occur with relative regularity, the simulation will have different timings on each `Update()` between one session and the next, so even if we recorded user input timings and automatically replayed the Scene, trying to apply the recorded inputs at the moments they were applied isn't going to be exactly the same every time, and so we may not get the exact same result.

Moving user input handling to `FixedUpdate()` is possible, and helpful, if user input controls Rigidbody behavior such as applying forces in different directions while the player holds down certain keys. However, this will tend to lead to input latency or *lag*, since it will be anywhere from 0 to 20 milliseconds (based on the Fixed Update timestep frequency) before the Physics Engine can respond to the key being pressed. Instantaneous inputs, such as jumping or activating an ability, are always best handled in `Update()` in order to avoid missing keystrokes. Helper functions such as `Input.GetKeyDown()` would only return `true` for the frame the player presses the given key and will return `false` during the next `Update()`. If we tried to read a *key-down* event during a `FixedUpdate()`, we would never know that the user

pressed the key, unless a physics timestep just happens to occur between these two frames. This can be worked around with an input buffering/tracking system, but this is certainly more trouble than its worth if we're implementing it merely to replicate a physics bug.

Ultimately, experience and persistence are the only good ways to debug most physics problems. The more experience we have with the Physics Engine, the more intuition we will have to find the source of the problem, but unfortunately they almost always take a lot of time to resolve due to their limited reproducibility and sometimes nebulous behavior, and so we should expect physics issues to take longer than most logic bugs to resolve and plan extra time before it can be resolved.

Physics performance optimizations

Now that we have an understanding of the majority of features of the Unity Physics Engine, we can cover several optimization techniques to improve our game's physics performance.

Scene setup

Firstly, there are a number of best practices we can apply to our Scenes to improve consistency of the physics simulation. Note that several of these techniques will not necessarily improve CPU or memory usage, but they will result in a reduced likelihood of instability from the Physics Engine.

Scaling

We should try to keep all physics object scales in the world as close to $(1, 1, 1)$ as we possibly can. By default, Unity assumes that we are trying to simulate gameplay equivalent to being on the surface of the Earth. The force of gravity at the surface of the Earth is 9.81 meters-per-second-squared, and hence the default gravity value is set to -9.81 to match. 1 unit in Unity's world space is equivalent to 1 meter, and the negative sign means that it will pull the object downward. Our object sizes should reflect our effective world scale, since scaling them too large will cause gravity to appear to move the objects much more slowly than we would expect. If all of our objects are scaled five times too big, then gravity will appear to be five times weaker. The converse is also true; scaling objects too small will make them appear to fall too quickly and will not seem realistic.

We can tweak the world's implied scale by modifying the strength of gravity under `Edit | Project Settings | Physics / Physics 2D | Gravity`. However, note that any floating-point arithmetic will be more accurate with values closer to 0, so if we have some objects that have scale values far above $(1, 1, 1)$, even if they match the implied world-scale, then we could still observe erratic physics behavior. So, early in the project, we should import and scale our most common physics objects around a scale value of $(1, 1, 1)$ and then adjust the value of gravity to match. This will give us a reference point to work with as we introduce new objects.

Positioning

Similarly, keeping all objects close to $(0, 0, 0)$ in the world-space position will result in better floating-point accuracy, improving the consistency of the simulation. Space simulator and free-running games try to simulate incredibly large spaces and typically use a trick of either secretly teleporting the player back toward the center of the world or fixing their position there, in which case, either volumes of space are compartmentalized so that physics calculations are always calculated with values close to 0, or everything else is moved to simulate travel, and the player's motion is only an illusion.

Most games are not at risk of introducing floating-point inaccuracy, since most game levels tend to last around 10 to 30 minutes, which doesn't give the player much time to travel absurdly long distances, but if we're working with exceptionally large Scenes or asynchronously loading Scenes throughout the course of the entire game to the point that the player travels tens of thousands of meters, then we may start to notice some strange physics behavior the further they travel.

So, unless we're already far too deep into our project such that changing and retesting everything at a late stage would be too much hassle, we should try to keep all of our physics objects close to $(0, 0, 0)$. Plus, this is good practice for our project workflow, as it makes it much quicker to find objects and tweak things in our game world.

Mass

Mass is stored as a floating-point value under a `Rigidbody` Component's `mass` property, and documentation on its usage has changed a fair amount over the years due to updates in its Physics Engine. In the recent versions of Unity (late Unity 5 and early Unity 2017), we are essentially free to choose whatever we want the value `1.0` to represent and then scale other values appropriately.

Traditionally, a mass value of `1.0` is used to represent a mass of 1 kilogram, but we could decide that a human being has a mass of `1.0` (~130 kilogram), in which case, a car would be given a mass value of `10.0` (~1,300 kilogram), and physics collisions will resolve similarly to what we expect. The most important part is the relative difference in mass, which allows collisions between these objects to look believable without stressing the engine too much. Floating-point precision is also a concern, so we don't want to use large mass values that are too ridiculous.



Note that if we intend to use Wheel Colliders, their design assumes that a mass of `1.0` represents 1 kilogram, so we should assign our mass values appropriately.

Ideally, we would maintain mass values around `1.0` and ensure a maximum relative mass-ratio of around `100`. If two objects collide with a mass-ratio much greater than this, then large momentum differences can turn into sudden, immense velocity changes from the impulse, resulting in some unstable physics and potential loss of floating-point precision. Object pairs that have a significant scale difference should probably be culled with the Collision Matrix to avoid problems (more on this shortly).



Improper mass-ratios are the most common cause for physics instability and erratic behavior in Unity. This is particularly true when using Joints for objects such as Ragdolls.

Note that the force of gravity at the center of the Earth affects all objects

equally, regardless of their mass, so it does not matter if we consider a mass property value of 1.0 to be the mass of a rubber ball or the mass of a warship. There's no need to adjust the force of gravity to compensate. What does matter, however, is the amount of air resistance the given object undergoes while falling (which is why a parachute falls slowly). So, to maintain realistic behavior, we may need to customize the `drag` property for such objects or customize the force of gravity on a per-object basis. For example, we could disable the Use Gravity checkbox and apply our own custom gravitational force during Fixed Updates.

Use Static Colliders appropriately

As mentioned previously, the Physics Engine automatically generates two separate data structures to contain Static Colliders separately from Dynamic Colliders. Unfortunately, if new objects are introduced into the Static Collider data structure at runtime, then it must be regenerated, similar to calling `StaticBatchingUtility.Combine()` for Static Batching. This is likely to cause a significant CPU spike. This makes it vital that we avoid instantiating new Static Colliders during gameplay.

In addition, merely moving, rotating, or scaling Static Colliders also triggers this regeneration process and should be avoided. If we have Colliders that we wish to move around without reacting to other objects colliding with them in a physical way, then we should attach a Rigidbody to make it a Dynamic Collider and enable the Kinematic flag. This flag prevents the object from reacting to external impulses from inter-object collisions, similar to Static Colliders, except the object can still be moved through its `Transform` Component or through forces applied to its `Rigidbody` Component (preferably during Fixed Updates). Since a Kinematic object won't respond to other objects hitting it, it will tend to simply push other Dynamic Colliders out of its way as it moves.



It's for this reason that player character objects are often made into Kinematic Colliders.

Use Trigger Volumes responsibly

As mentioned previously, we can treat our physics objects as normal Colliders or as Trigger Volumes. An important distinction between these two types is that the `OnCollider...()` callbacks provide a `Collision` object as a parameter to the callback, which contains useful information such as the exact location of collision (useful to position a Particle Effect) and the contact normal (useful if we want to manually move the object after the collision). Whereas, the `OnTrigger...()` callbacks do not provide this kind of information.

As a result, we should not try to use Trigger Volumes for collision-reactive behavior since we won't have enough information to make the collision appear accurate. Trigger Volumes are best used for their intended purpose of tracking when an object enters/exits a specific area, such as dealing with damage while a player stays in a lava pit, triggering a cutscene when a player enters a building, and initiating asynchronous loading/unloading of a Scene when the player approaches/moves far enough away from another major area.

If contact information is absolutely needed for a Trigger Volume collision, then common workarounds are to do any of the following:

- Generate a rough estimate for the contact point by halving the distance between the Trigger Volume and colliding objects' center of mass (this assumes that they're of roughly equal size)
- Perform a Raycast upon collision from the center of the Trigger Volume to the center of mass of the colliding object (works best if both of the objects are spherical)
- Create a non-Trigger Volume object, give it an infinitesimally small mass (so that the colliding object is barely effected by its presence), and immediately destroy it upon collision (since collision with such a large mass differential will probably send this small object into orbit)

Of course, each of these approaches have their drawbacks; limited physical

accuracy, extra CPU overhead during collision, and/or additional Scene setup (and rather hacky-looking collision code), but they can be useful in a pinch.

Optimize the Collision Matrix

As we know, the Physics Engine's Collision Matrix defines which objects assigned to certain Layers are allowed to collide with objects assigned to other Layers, or to put it more succinctly, which object collision pairs are even considered viable by the Physics Engine. Every other object-Layer pair is simply ignored by the Physics Engine, which makes this an important avenue for minimizing Physics Engine workload since it reduces the number of bounding volume checks that must be performed each and every Fixed Update and how many collisions would ever need to be processed during the lifecycle of the application (which would save on battery life for a mobile device).



Note that the Collision Matrix can be accessed through Edit | Project Settings | Physics (or Physics2D) | Layer Collision Matrix.

The following screenshot shows a common Collision Matrix for an arcade shooter game:

		Layer Collision Matrix										
		Default	TransparentFX	Ignore Raycast	Water	UI	Player	Enemies	Player Missiles	Enemy Missiles	Powerups	World
Default		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>				
TransparentFX		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>				
Ignore Raycast		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>				
Water		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>				
UI		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>				
Player		<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>				
Enemies		<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Player Missiles		<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Enemy Missiles		<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Powerups		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>				
World		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>				

In the preceding example, we have objects flagged as Player, Enemies, Player

Missiles, Enemy Missiles, and Powerups, and we have minimized the number of possible inter-object collisions for the Physics Engine to check.

Starting with the first row of checkmarks labeled Player, we want the Player object to be able to collide with World objects, pickup Powerups, get hit by Enemy Missiles, and collide with Enemies. However, we do not want them to collide with their own Player Missiles or themselves (although there would probably only be one object in this Layer, anyway). Hence, the enabled checkboxes in the Player row reflect these requirements. We only want Enemies to collide with World objects and Player Missiles, so these are checked in the Enemies row. Note that the Player-to-Enemies collision pair would have been handled already by the previous row; hence there is no need for it to appear in the Enemies row. We also want both Player Missiles and Enemy Missiles to explode when they hit World objects, so these are marked, and finally we don't care about Powerups colliding with anything but the Player, nor do we want World objects to collide with other World objects, so no checkboxes are marked on the final two rows.

At any given moment, we might have a single Player object, 2 Powerups, 7 Player Missiles, 10 Enemies, and 20 Enemy Missiles, which is 780 potential collision pairs (this is calculated as each of 40 different objects could collide with 39 other ones, giving us 1,560 potential collision pairs, but then we divide the total by 2 to ignore duplicate pairs). By merely optimizing this matrix, we have reduced this to less than 100, for an almost 90% reduction in potential collision checks. Of course, the Unity Physics Engine efficiently culls away many of these object pairs if they are too far apart from one another, hence there is little to no chance that they could collide (this is calculated during a hidden process known as **Broadphase Culling**), so the actual savings will probably never be this good, but it will free up some CPU cycles with next to no effort. Another big benefit is that it simplifies our game logic coding; there's no need to figure out what's supposed to happen if Powerups and Enemy Missiles collide if we tell the Physics Engine to ignore collisions between them.

We should perform logical sanity checks like this for all potential Layer combinations in the Collision Matrix to see whether we're wasting precious time checking for inter-object collisions between object pairs that aren't

necessary.

Prefer Discrete collision detection

Discrete collision detection is fairly cheap since teleporting objects once and performing a single overlap check between nearby object pairs is a fairly trivial amount of work to perform in a single timestep. The amount of calculation required to perform Continuous collision detection is significantly greater, since it involves interpolating both objects between their starting and ending positions, while analyzing for any slight bounding volume overlaps between these points as they might occur during the timestep.

Consequently, the Continuous collision detection option is an order of magnitude more expensive than the Discrete detection method, whereas the ContinuousDynamic collision detection setting is an order of magnitude even more expensive than Continuous. Having too many objects configured to use either of the Continuous collision detection types will cause serious performance degradation in complex Scenes. In either case, the costs are multiplied by the number of objects that need to be compared during any given frame and whether or not the comparison Collider is Static or Dynamic.

Ergo, we should favor the Discrete collision detection setting for the overwhelming majority of our objects, while using either of the Continuous collision detection settings only in extreme circumstances. The Continuous setting should be used when important collisions are frequently missed with the more static parts of our game world. For instance, if we wish to be certain that the player characters never fall through the game world or never accidentally teleport through walls if they move too quickly, then we might want to apply Continuous collision detection only for those objects. Finally, the ContinuousDynamic setting should only be used if the same situation applies and we wish to catch collisions between pairs of very fast moving Dynamic Colliders.

Modify the Fixed Update frequency

In some cases, Discrete collision detection might not work well enough on a large scale. Perhaps our entire game revolves around a lot of small physics objects, and Discrete collision detection simply isn't catching enough collisions to maintain product quality. However, applying one of the Continuous collision detection settings to everything would be far too prohibitive on performance. In this case, there is one option we can try: we can customize the physics timestep to give the Discrete collision detection system a better chance of catching such collisions by modifying how frequently the engine checks for Fixed Updates.

As mentioned previously, Fixed Updates and physics timestep processing are strongly coupled; so, by modifying the frequency of Fixed Update checks, we not only change the frequency that the Physics Engine will calculate and resolve the next callback, but we also change how frequently the `FixedUpdate()` callbacks and Coroutines are being invoked. Consequently, changing this value can be risky if we're deep into our project and have a lot of behavior that depends on these callbacks since we will be changing a very important assumption about how often these methods are invoked.

Altering the Fixed Update frequency can be accomplished using the `Edit | Project Settings | Time | Fixed Timestep` property in the Editor or through the `Time.fixedDeltaTime` property in script code.

Reducing this value (increasing the frequency) will force the Physics Engine to process more frequently, giving it a better chance of catching collisions with Discrete collision detection. Naturally, this comes with an added CPU cost since we're invoking more `FixedUpdate()` callbacks and asking the Physics Engine to update more frequently, having it move objects and verify collisions more often.

Conversely, increasing this value (decreasing the frequency) provides more

time for the CPU to complete other tasks before it must handle physics processing again, or looking at it from another perspective, giving the Physics Engine more time to process the last timestep before it begins processing the next one. Unfortunately, lowering the Fixed Update frequency would essentially lower the maximum velocity at which objects can move before the Physics Engine can no longer capture collisions with Discrete collision detection (depending on the objects' sizes). We might also start to see objects changing velocities in strange ways because it is essentially becoming a weaker approximation of real-world physics behavior.

This makes it absolutely vital to perform a significant amount of testing each time the Fixed Timestep value is changed. Even with a complete understanding of how this value works, it is difficult to predict what the overall outcome will look like during gameplay and whether the result is passable for quality purposes. Hence, changes to this value should be made early in the project's lifecycle and then made infrequently in order to get a sufficient amount of testing against as many physics situations as possible.

It might help to create a test Scene that flings some of our high-velocity objects at one another to verify that the results are acceptable and run through this Scene whenever Fixed Timestep changes are made. However, actual gameplay tends to be rather complex, with many background tasks and unanticipated player behavior that causes additional work for the Physics Engine or gives it less time to process the current iteration. Actual gameplay conditions are impossible to replicate in a vacuum. Also, there's no substitute for the real thing, so the more testing we can accomplish against the current value of Fixed Timestep, the more confident we can be that the changes meet acceptable quality standards.

Take it from someone who's in the career of developing software automation tools: automation of software testing is helpful in a lot of situations, but when it comes to real-time event and user input-driven applications, which synchronize with multiple hardware devices and complex subsystems such as Physics Engines, and which tend to change rapidly due to iterations on feedback, the support and maintenance costs of automated testing often becomes more effort than its worth,



making manual testing the most sensible approach.

We always have Continuous collision detection as a last resort to offset some of the resulting instability we're observing. Unfortunately, even if the changes are targeted, it is more likely that this will cause further performance issues than we started with due to the overhead costs of Continuous collision detection. It would be wise to profile our Scene before and after enabling Continuous collision detection to verify that the benefits are outweighing the costs.

Adjust the Maximum Allowed Timestep

If we're regularly exceeding the Maximum Allowed Timestep (which, as a reminder, determines how much time the Physics Engine has to resolve a timestep before it must exit early), then it will result in some pretty bizarre-looking physics behavior. Rigidbodies will appear to slow down or jerk to a stop since the Physics Engine needs to keep exiting timestep calculations early before it has fully resolved its entire time quota. In this case, it is a clear sign that we need to optimize our physics behavior from other angles. However, at the very least, we can be confident that the threshold will prevent the game from completely locking up from a spike in the middle of physics processing.



Reminder: this setting can be accessed through Edit | Project Settings | Time | Maximum Allowed Timestep.

The default setting is to consume a maximum of 0.333 seconds, which would manifest itself as a very noticeable drop in frame rate (a mere 3 FPS) if it were exceeded. If you ever feel the need to change this setting, then you obviously have some big problems with your physics workload, so it is recommended that you only tweak this value if you have exhausted all other approaches.

Minimize Raycasting and bounding-volume checks

All of the Raycasting methods are incredibly useful, but they are relatively expensive, particularly `capsulecast()` and `spherecast()`. We should avoid calling these methods regularly within the `update()` callbacks or Coroutines, saving them only for key events in our script code.

If we're making use of persistent line, ray, or area-of-effect collision areas in our Scene (examples include security lasers, continuously burning fires, beam weapons, and so on) and the object remains relatively stationary, then they would perhaps be better simulated using a simple Trigger Volume.

If such replacements are not possible and we truly need persistent casting checks using these methods, we should minimize the amount of processing each Raycast makes by exploiting `LayerMasks`. This is particularly true if we're making use of `Physics.RaycastAll()`. For example, a poorly optimized usage of this kind of Raycasting would look as follows:

```
void PerformRaycast() {
    RaycastHit[] hits;
    hits = Physics.RaycastAll(transform.position, transform.forward,
    100.0f);
    for (int i = 0; i < hits.Length; ++i) {
        RaycastHit hit = hits[i];
        EnemyComponent e = hit.transform.GetComponent<EnemyComponent>();
        if (e.GetType() == EnemyType.Orc) {
            e.DealDamage(10);
        }
    }
}
```

In the preceding example, we're collecting Raycast collision data for every object in the path of this Ray, but we're only processing its effects on objects which hold a specific `EnemyComponent`. Consequently, we're asking the Physics Engine to complete much more work than is necessary.

A better approach will be to use a different overload of `RaycastAll()`, which

accepts a `LayerMask` value as an argument. This will filter collisions for the Ray in much the same way as the Collision Matrix so that it only tests against objects in the given Layer(s). The following code contains a subtle improvement by providing an additional `LayerMask` property; we would configure the `LayerMask` through the Inspector window for this Component, and it will filter the list much faster and only contain hits for objects matching the mask:

```
[SerializeField] LayerMask _layerMask;

void PerformRaycast() {
    RaycastHit[] hits;
    hits = Physics.RaycastAll(transform.position, transform.forward, 100.0f, _layerMask);
    for (int i = 0; i < hits.Length; ++i) {
        // as before ...
    }
}
```

This optimization doesn't work as well for the `Physics.RaycastHit()` function since that version only provides ray collision information for the first object the Ray collides with, regardless of whether we're using a `LayerMask` or not.



Note that because the `RaycastHit` and `Ray` classes are managed by the Native memory space of the Unity Engine, they don't actually result in memory allocations that draw the attention of the Garbage Collector. We will learn more about such activity in [Chapter 8, Masterful Memory Management](#).

Avoid complex Mesh Colliders

In order of collision detection efficiency, the various Colliders are Spheres, Capsules, Boxes, Convex Mesh Colliders, followed by Concave Mesh Colliders being far and away the most expensive. Collisions always involve pairs of objects, and the amount of work (math) needed to resolve the collision will depend on the complexity of both objects. Detecting collisions between two primitive objects can be reduced down to a relatively simple set of mathematical equations, which are highly optimized. Performing comparisons against a pair of Convex Mesh Colliders is a much more complex equation, making them an order of magnitude more expensive than collisions between two primitives. Then, there are collisions between two Concave Mesh Colliders, which are so complex that they cannot be reduced down to a simple formula and require collision checks to be resolved between each pair of triangles across both meshes, easily making them orders of magnitude more expensive than collisions between other Collider types. The amount of work involved scales similarly when we resolve collisions between shapes of different groups. For example, a collision between a primitive and Concave Mesh Collider would be slower than a collision between two primitives, but faster than a collision between two Concave Mesh Colliders.

There is also the question of whether one, or both, of the objects involved in the collision is moving (one of the objects being a Static Collider is easier to process than both objects being Dynamic Colliders). There is also the matter of how many of these objects are within our Scene since the total processing costs of collision detection will grow rapidly if we're not careful with how many shapes we introduce into the simulation.

A great irony between representing physics and graphics in 3D applications is how difficult it is to handle spherical and cube objects between the two of them. The perfect spherical mesh would require an infinite number of polygons to be generated, making such an object impossible to represent graphically.

However, handling collisions between two spheres in a Physics Engine is perhaps the simplest problem to solve for contact points and collisions (the contact point is always at the edge of either of the sphere's radius, and the contact normal is always the vector between their centers of mass).

Conversely, a cube is one of the simplest objects to represent graphically (as little as 8 vertices and 12 triangles) and yet takes significantly more mathematics and processing power to find contact points and resolve collisions for (and the mathematics to resolve it depends on whether the collision occurred between faces, edges, corners, or a mixed pairing).

Anecdotally, this implies that the most efficient way of creating the largest number of objects would be to populate our world with cube objects, which use spherical Colliders. However, this would make absolutely no sense to a human observer, as they would witness cubes rolling around like balls.

The previous anecdote serves as a reminder that the physical representation of an object does not necessarily need to match its graphical representation. This is beneficial, as a graphical mesh can often be condensed down into a much simpler shape, while still generating very similar physics behavior and simultaneously removing the need to use an overly complex Mesh Collider.

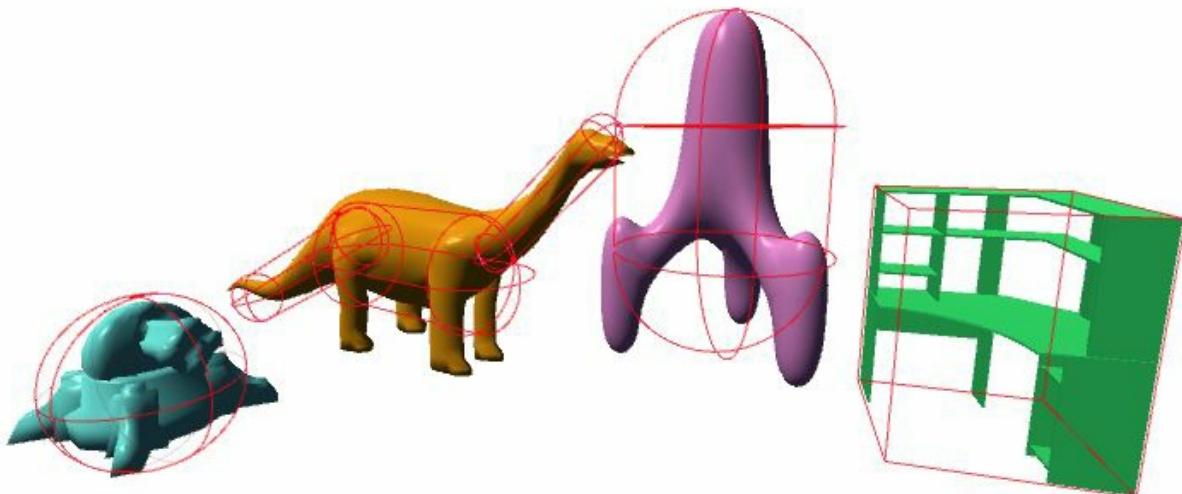
This separation of representations between graphics and physics allows us to optimize the performance of one system without (necessarily) negatively affecting the other. So long as there are no noticeable repercussions on gameplay (or we're willing to make the sacrifice), then we are free to represent complex graphical objects with much simpler physics shapes without players noticing. Also, if the player never notices, then no harm is done.

So, we can solve this problem in one of the two ways: either by approximating the physics behavior of the complex shape using one (or more) of the standard primitives or by using a much simpler Mesh Collider.

Use simpler primitives

Most shapes can be approximated using one of the three primitive Colliders. In fact, we do not need to represent the object using only a single Collider. We are free to use several Colliders if they serve our needs for creating a complex collision shape by attaching additional child `GameObjects` with their own Colliders. This is almost always less expensive than using a single Mesh Collider and should be preferred.

The following image shows a handful of complex graphical objects represented by one or more simpler primitive Collider shapes:



Using a Mesh Collider for any one of these objects would be significantly more expensive than the primitive Colliders shown here due to the number of polygons they contain. It is worth exploring any and all opportunities to simplify our objects down using these primitives as much as we can, as they can provide significant performance gains.

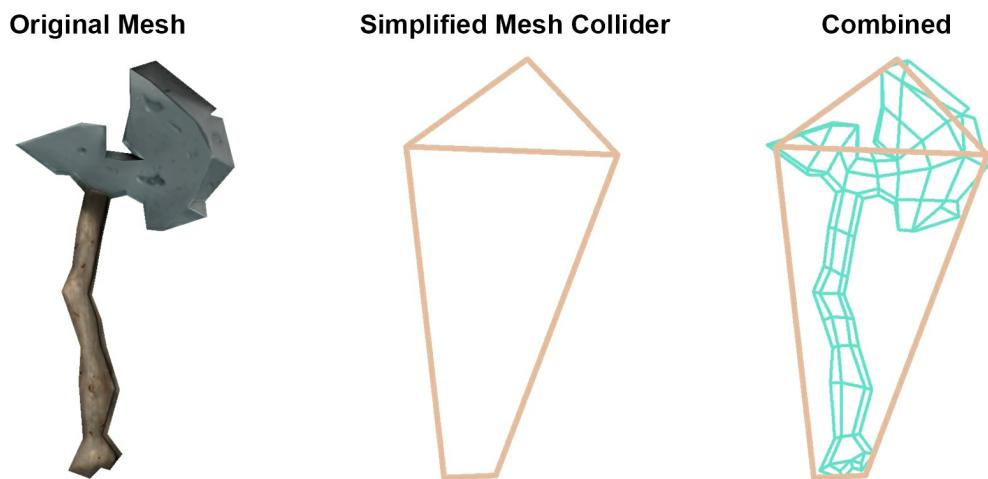
For example, Concave Mesh Colliders are unique in that they can feature gaps or holes that allow other meshes to fall into, or even through them, which introduces opportunities for the objects to fall through the world if

such Colliders are used for world collision areas. It is often better to place Box Colliders in strategic locations for this purpose.

Use simpler Mesh Colliders

Similarly, the mesh assigned to a Mesh Collider does not necessarily need to match the graphical representation of the same object (Unity simply picks it as the default). This gives us an opportunity to assign a simpler mesh to the Mesh Collider's `mesh` property, which is different from the one we use for its graphical representation.

The following image shows an example of a complex graphical mesh that has been given a much more simplified mesh for its Mesh Collider:



Simplifying the rendered mesh into convex shapes with lower polygon counts in this way will greatly reduce the overhead needed to determine bounding volume overlaps with other Colliders. Depending on how well the original object is estimated, there should be minimal noticeable gameplay differences, especially in the case of this axe, which we expect to be moving quickly as creatures swing it during attacks, making it unlikely that players will notice the difference between the two meshes as Colliders. In fact, the simplified mesh is much less likely to be missed by Discrete collision detection and is also preferable for that reason.

Avoid complex physics Components

Certain special physics Collider Components, such as `TerrainCollider`, `Cloth`, and `WheelCollider`, are orders of magnitude more expensive than all primitive Colliders and even Mesh Colliders in some cases. We should simply not include such Components in our Scenes unless they are absolutely necessary. For instance, if we have terrain objects in the distance that the player will never approach, there's little reason to include an attached `TerrainCollider`.

Games featuring `Cloth` Components should consider instantiating different objects without them when running in lower-quality settings, or simply animating cloth behavior (although it is totally understandable if the team has grown attached and fallen in love with how the stuff moves around).

Games using `WheelCollider` Components should simply try to use fewer Wheel Colliders. Large vehicles, with more than four wheels, may be able to simulate similar behavior using only four wheels, while faking the graphical representation of additional wheels.

Let physics objects sleep

The Physics Engine's sleep feature can pose several problems for our game. Firstly, some developers don't realize that many of their Rigidbodies are sleeping during most of the lifetime of their application. This tends to lead developers to assume that they can get away with (for example) doubling the number of Rigidbodies in their game and the overall costs would simply double to match it. This is unlikely. The frequency of collisions and total accumulated time of active objects is more likely to increase in an exponential fashion rather than a linear one. This leads to unexpected performance costs every time new physics objects are introduced into the simulation. We should keep this in mind when we decide to increase the physical complexity of our Scenes.

Secondly, changing any properties on a `Rigidbody` Component at runtime, such as `mass`, `drag`, and `useGravity`, will also reawaken an object. If we're regularly changing these values (such as a game where object sizes and masses change over time), then they will remain active for longer periods of time than usual. This is also the case for applying forces, so if we're using a custom gravity solution (such as suggested in the *Mass* section), we should try to avoid applying the gravitational force every Fixed Update, otherwise the object will be unable to fall asleep. We could check its mass-normalized-kinetic-energy (just take the value of `velocity.sqrMagnitude`) and manually disable our custom gravity when we detect that it is very low.

Thirdly, there is the danger of islands of sleeping physics objects being generated. Islands are created when a large number of Rigidbodies are touching one another and have gradually gone to sleep once the kinetic energy of the system has fallen low enough. However, because they're all still touching one another, as soon as one of these objects is awoken, it will start a chain reaction, awakening all other nearby Rigidbodies. Suddenly, we have a large spike in CPU usage because dozens of objects have re-entered the simulation. Even worse, because the objects are so close together, there will be many potential collision pairs that must keep being resolved until the objects fall asleep again.

Avoiding these situations is best done by reducing the complexity of our Scenes, but if we find ourselves unable to do so, we could look for ways to detect that islands are forming, and then strategically destroy/despawn some of them to prevent too many large islands from being generated. However, performing regular distance comparisons between all of our Rigidbodies is not a cheap task to accomplish and could be costly. The Physics Engine already performs such checks itself, during Broadphase Culling, but, unfortunately, Unity doesn't expose this data through the Physics Engine API. Any workarounds for this problem will be dependent on how the game is designed; for example, a game that requires the player to move lots of physics objects into an area (for example, a game that involves herding sheep into a pen) could choose to remove the sheep's Collider as soon as the player moves it into position, locking the object to its final destination, easing the workload on the Physics Engine and preventing islands from becoming a problem.

Sleeping objects can be a blessing and a curse. They can save us a lot of processing power, but if too many of them reawaken at the same time or our simulation is too busy to allow enough of them to fall asleep, then we could be incurring some unfortunate performance costs during gameplay. We should strive to limit these situations as much as possible by letting our objects enter the sleeping state as much as possible and avoiding grouping them together in large clusters.



Note that the sleep threshold can be modified under Edit | Project Settings | Physics | Sleep Threshold.

Modify the Solver Iteration Count

Using Joints, Springs, and other ways to connect Rigidbodies together are fairly complex simulations in Physics Engines. Owing to the codependent interactivity (internally represented as movement constraints) that occurs due to joining two objects together, the system must often make several attempts at solving the necessary mathematical equations. This multi-iteration approach is required to calculate an accurate result whenever there is a change in velocity to any single part of the object chain.

It, therefore, becomes a balancing act of limiting the maximum number of attempts the *Solver* makes to resolve a particular situation versus how accurate a result we can get away with. We don't want the Solver to spend too much time on a single collision, because there are a lot of other tasks that the Physics Engine has to complete within the same iteration. However, we also don't want to reduce the maximum number of iterations too far, as it will only approximate what the final solution would have been, making its motion look much less believable than if it had been given more time to calculate the result.



The same Solver also gets involved when resolving inter-object collisions and contacts. It can almost always determine the correct result for simple collisions with a single iteration, with the exception of some very rare and complex collision situations with Mesh Colliders. It is mostly when attached objects will be affected through Joints that the Solver requires additional effort to integrate the final result.

The maximum number of iterations the Solver is allowed to attempt is called the *Solver Iteration Count*, which can be modified under *Edit | Project Settings | Physics | Default Solver Iterations*. In most cases, the default value of six iterations is perfectly acceptable. However, games that include very complex Joint systems may wish to increase this count to suppress any erratic

(or downright explosive) characterJoint behaviors, whereas some projects may be able to get away with reducing this count. Testing must be performed after changing this value to check whether the project still maintains the intended levels of quality. Note that this value is the *default* Solver Iteration Count--the value that gets applied to any newly created Rigidbodies. We can change this value at runtime through the `Physics.defaultSolverIterations` property, but this still won't affect preexisting Rigidbodies. If necessary, we can modify their Solver Iteration Count after they are constructed through the `Rigidbody.solverIterations` property.

If we find our game regularly runs into jarring, erratic, and physics-breaking situations with complex Joint-based objects (such as Ragdolls), then we should consider gradually increasing the Solver Iteration Count until the problems are suppressed. These problems typically occur if our Ragdolls absorb too much energy from colliding objects and the Solver is unable to iterate the solution down to something reasonable before it is asked to give up. At this point, one of the Joints goes supernova, dragging the rest of them into orbit along with it. Unity has a separate setting for this problem, which can be found under Edit | Project Settings | Physics | Default Solver Velocity Iterations. Increasing this value will give the Solver more opportunity to calculate a sensible velocity during Joint-based object collisions and help to avoid the above scenario. Again, this is a default value, hence it is only applied to newly created Rigidbodies. The value can be modified at runtime through the `Physics.defaultSolverVelocityIterations` property and can be customized on specific Rigidbodies through the `Rigidbody.solverVelocityIterations` property.

In either case, increasing the number of iterations will consume more CPU resources during every Fixed Update where the Joint objects remain active.



Note that the Physics 2D settings for solver iterations are named Position Iterations and Velocity Iterations.

Optimize Ragdolls

Speaking of Joint-based objects, Ragdolls are incredibly popular features for good reason; they're tons of fun! Ignoring the morbidity of flinging corpses around a game world for the moment, there's something about watching a complex chain of objects flail around and smash into things that hits a lot of psychological fun buttons. This makes it very tempting to allow many Ragdolls to coexist within our Scene at the same time, but as we quickly discover, this risks an enormous performance hit when too many Ragdolls are in motion and/or collide with other objects due to the amount of iterations the Solver would need to resolve them all. So, let's explore some ways to improve the performance of Ragdolls.

Reduce Joints and Colliders

Unity provides a simple Ragdoll-generation tool (the Ragdoll Wizard) under GameObject | 3D Object | Ragdoll.... This tool can be used to create Ragdolls from a given object by selecting the appropriate child GameObjects to attach Joint and Collider Components to for any given body part or limb. This tool always creates 13 different Colliders and associated Joints (pelvis, chest, head, two Colliders per arm, and three Colliders per leg).



Note that a bug causes the Ragdoll Wizard not to complain if nothing is assigned to Left Foot or Right Foot transform Component references like it does for the rest of them, but Unity will throw a NullReferenceException if we try to create the mesh without them assigned. Ensure that all 13 transform Component references have been assigned when we try to create a Ragdoll.

However, it's possible to use only seven Colliders (pelvis, chest, head, and one Collider per limb) to greatly reduce the overhead cost at the expense of Ragdoll realism. This can be achieved by deleting unwanted Colliders and manually reassigning the Character Joint's connectedBody properties to the proper parent joints (connect the arm Colliders to the chest, and connect the leg Colliders to the pelvis).

Note that we assign a mass value during Ragdoll creation using the Ragdoll Wizard. This mass value is spread across the various Joints as appropriate and, therefore, represents the total mass of the object. We should ensure that we don't apply a mass value too high or too low compared to other objects in our game to avoid potential instability.

Avoid inter-Ragdoll collisions

The performance cost of Ragdolls grows exponentially when Ragdolls are allowed to collide with other Ragdolls, since any Joint collision requires the Solver to calculate the resultant velocity applied to all of the Joints connected to it and then each of the Joints connected to them, such that both Ragdolls must be completely resolved multiple times. Also, it gets significantly more complicated if multiple parts of the Ragdolls are likely to collide with one another during the same collision.

This is a tough task for the Solver to handle, so we should avoid it. The best way to do this is to simply use the Collision Matrix. It is wise to assign all Ragdolls to their own Layer and uncheck the corresponding checkbox in the Collision Matrix so that objects in the given Layer cannot collide with objects in the same Layer.

Replace, deactivate or remove inactive Ragdolls

In some games, once a Ragdoll has reached its final *destination*, we no longer need it to remain in the game world as an interactable object. We could then either deactivate, destroy, or replace the Ragdoll with a simpler alternative when they are no longer needed (a good trick is to replace them with the simpler version that uses only seven Joints as suggested earlier). Such simplifications are often implemented as a means of reducing overhead for weaker hardware/lower quality settings or as a compromise to allow more Ragdolls to coexist in our Scene. It could even be used dynamically if a particular number of Ragdolls is already present.

We would need some object to keep track of all of our Ragdolls, being notified any time a Ragdoll is created, keeping track of how many Ragdolls currently exist, watching each of them until they fall asleep through `RigidBody.IsSleeping()` and then do something appropriate with them. The same object could also choose to instantiate simpler Ragdoll variations if the Scene already contains more Ragdolls than is reasonable. This would be another good opportunity to make use of the Messaging System we explored in [Chapter 2, Scripting Strategies](#).

Whichever approach we choose to improve the performance of our Ragdolls will no doubt result in limiting Ragdolls as a gameplay feature, either by instantiating fewer of them, giving them less complexity, or giving them a shorter lifetime, but these are reasonable compromises to make given the performance-saving opportunities.

Know when to use physics

The most obvious method to improve the performance of a feature is to avoid using it as much as possible. For all moveable objects in our game, we should take a moment to ask ourselves if getting the Physics Engine involved is even necessary. If not, we should look for opportunities to replace them with something simpler and less costly.

Perhaps we're using physics to detect whether the player fell into a kill-zone (water, lava, a death-plummet, and so on), but our game is simple enough that we only have kill-zones at a specific height. In this case, we could avoid physics Colliders altogether and get away with only checking whether the player's `y`-position falls below a particular value.

Consider the following example--we're trying to simulate a meteor shower, and our first instinct was to have many falling objects that move via physics Rigidbodies, detect collisions with the ground via Colliders, and then generate an explosion at the point of impact. However, perhaps the ground is consistently flat or we have access to the Terrain's Height Map for some rudimentary collision detection. In this case, object travel could be simplified by manually tweening the objects' `transform.position` over time to simulate the same traveling behavior without requiring any Physics Components. In both cases, we can reduce the physics overhead by simplifying the situation and pushing the work into Script code.



Tweening is a common short-hand term for in-betweening, which is the act of interpolating variables from one value to another gradually over time. There are many useful (and free) tweening libraries available on the Unity Asset Store that can provide a lot of useful functionality. Although, be careful of potentially poor optimization in these libraries.

The reverse is also possible. There might be occasions where we're performing a great deal of calculation through Script code that could be handled through physics relatively simply. For example, we might have

implemented an inventory system with many objects that can be picked up. When the player hits the `Pick up object` key, each of these objects might be compared against the player's position to figure out which object is the closest.

We could consider replacing all of the Script code with a single `Physics.OverlapSphere()` call to get nearby objects when the key is pressed, and then figure out the closest pickup object from the result (or, better yet, just automatically pick up all of them. Why make the player repeatedly click more than necessary?). This could greatly reduce the total number of objects that must be compared each time the key is pressed, although comparisons should be made to ensure that this is the case.

Ensure that you seek to remove unnecessary physics grunt work from your Scenes or use physics to replace behavior that is costly when performed through Script code. The opportunities are as wide and far reaching as your own ingenuity. The ability to recognize opportunities like this takes experience, but is a vital skill that will serve you well when saving performance in current and future game development projects.

Summary

We've covered numerous methods to improve our game's physics simulation both in terms of performance and consistency. The best technique when it comes to costly systems such as Physics Engines is simply avoidance. The less we need to use the system, the less we need to worry about it generating bottlenecks. In the worst-case scenario, we may need to reduce the scope of our game to condense physics activity down to only the essentials, but as we've learned, there are plenty of ways to reduce physics complexity without causing any noticeable gameplay effects.

In the next chapter, we will immerse ourselves in Unity's Rendering Pipeline and discover how to maximize the graphical fidelity of our application, by making use of all of the CPU cycles we've freed up using the performance enhancements from earlier chapters.

Dynamic Graphics

There is no question that the Rendering Pipeline of a modern graphics device is complicated. Even rendering a single triangle to the screen requires a multitude of Graphics API calls, covering tasks such as creating a buffer for the Camera view that hooks into the Operating System (usually via some kind of windowing system), allocating buffers for vertex data, setting up data channels to transfer vertex and texture data from RAM to VRAM, configuring each of these memory spaces to use a specific set of data formats, determining the objects that are visible to the Camera, setting up and initiating a Draw Call for the triangle, waiting for the Rendering Pipeline to complete its task(s), and finally presenting the rendered image to the screen. However, there's a simple reason for this seemingly convoluted and over engineered way of drawing such a simple object--rendering often involves repeating the same tasks over and over again, and all of this initial setup makes future rendering tasks very fast.

CPUs are designed to handle virtually any computational scenario, but can't handle too many tasks simultaneously, whereas GPUs are designed for incredibly large amounts of parallelism, but they are limited in the complexity they can handle without breaking that parallelism. Their parallel nature requires immense amounts of data to be copied around very rapidly. During the setup of the Rendering Pipeline, we configure memory data channels for our graphics data to flow through. So, if these channels are properly configured for the types of data we will be passing, then they will operate more efficiently. However, setting them up poorly will result in the opposite.

Both the CPU and GPU are used during all graphics rendering, making it a high-speed dance of processing and memory management that spans software; hardware; multiple memory spaces, programming languages (each suited to different optimizations), processors, and processor types; and a large number of special-case features that can be thrown into the mix.

To make matters even more complicated, every rendering situation we will

come across is different in its own way. Running the same application against two different GPUs often results in an apples-versus-oranges comparison due to the different capabilities and APIs they support. It can be difficult to determine where a bottleneck resides within such a complex web of hardware and software systems, and it can take a lifetime of industry work in 3D graphics if we want to have a strong, immediate intuition about the source of performance issues in modern Rendering Pipelines.

Thankfully, Profiling comes to the rescue once again, which makes becoming a *Rendering Pipeline wizard* less of a necessity. If we can gather data about each device, use multiple performance metrics for comparison, and tweak our Scenes to observe how different rendering features affect their behavior, then we should have sufficient evidence to find the root cause of an issue and make appropriate changes. So, in this chapter, you will learn how to gather the right data, dig just deep enough into the Rendering Pipeline to find the true source of the problem, and explore various solutions and workarounds for a multitude of potential problems.

There are many topics to be covered when it comes to improving rendering performance. So, in this chapter, we will explore the following topics:

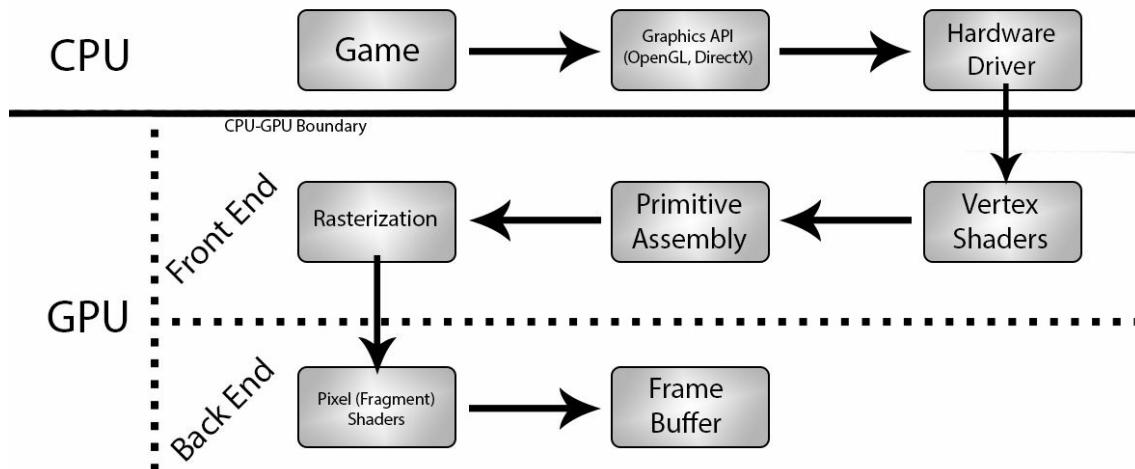
- A brief exploration of the Rendering Pipeline, focusing on the parts where the CPU and GPU come into play
- General techniques on how to determine whether our rendering is limited by the CPU or by GPU
- A series of performance optimization techniques and features, as follows:
 - GPU Instancing
 - Level Of Detail (LOD) and other Culling Groups
 - Occlusion Culling
 - Particle Systems
 - Unity UI
 - Shader optimization
 - Lighting and Shadow optimization
 - Mobile-specific rendering enhancements

The Rendering Pipeline

Poor rendering performance can manifest itself in a number of ways, depending on whether the device is limited by CPU activity (we are CPU bound) or by GPU activity (we are GPU bound). Investigating a CPU-bound application can be relatively simple since all of the CPU work is wrapped up in loading data from disk/memory and calling Graphics API instructions. However, a GPU-bound application can be more difficult to analyze since the root cause could originate from one of a large number of potential places within the Rendering Pipeline. We might find that we need to rely on a little guesswork or *process of elimination* in order to determine the source of a GPU bottleneck. In either case, once the problem is discovered and resolved, we can expect significant improvements since small fixes tend to reap big rewards when it comes to fixing issues in the Rendering Pipeline.

We briefly touched on the Rendering Pipeline in [Chapter 3, The Benefits of Batching](#). To briefly summarize the essential points, we know that the CPU sends rendering instructions through the Graphics API that funnels through the hardware driver to the GPU device, which results in a list of rendering instructions being accumulated in a queue known as the Command Buffer. These commands are processed by the GPU one by one until the Command Buffer is empty. So long as the GPU can keep up with the rate and complexity of instructions before the next frame is due to begin, we will maintain our frame rate. However, if the GPU falls behind, or the CPU spends too much time generating commands, the frame rate will start to drop.

The following is a greatly simplified diagram of a typical Rendering Pipeline on a modern GPU (which can also vary based on device, technology support, and custom optimizations), showing a broad view of the steps that take place:



The top row represents the work that takes place in the CPU, which includes both the act of calling into the Graphics API through the hardware driver and pushing commands into the GPU. The next two rows represent the steps that take place in the GPU. Owing to the GPU's complexity, its internal processes are often split into two different sections--the **Front End** and the **Back End**, which require a little added explanation.

The GPU Front End

The Front End refers to the part of the rendering process where the GPU handles vertex data. It will receive mesh data from the CPU (a big bundle of vertex information), and a Draw Call will be issued. The GPU then gathers all pieces of vertex information from the mesh data and passes them through Vertex Shaders, which are given an opportunity to modify them and output them in a *1-to-1* manner. From this, the GPU now has a list of Primitives to process (triangles--the most primitive shapes in 3D graphics). Next, the Rasterizer takes these Primitives and determines which pixels of the final image will need to be drawn on to create the Primitive based on the positions of its vertices and the current Camera view. The list of pixels generated from this process is known as **fragments**, which will be processed in the Back End.

Vertex Shaders are small C-like programs that determine the input data that they are interested in and the way that they will manipulate it, and then will output a set of information for the Rasterizer to generate fragments with. It is also home to the process of Tessellation, which is handled by Geometry Shaders (sometimes called **Tessellation Shaders**), similar to a Vertex Shader in that they are small scripts uploaded to the GPU, except that they are allowed to output vertices in a *1-to-many* manner, thus generating additional geometry programmatically.



*The term **Shader** is an anachronism from back when these scripts primarily handled Lighting and Shading tasks, before their role was expanded to include all of the tasks they are used for today.*

The GPU Back End

The Back End represents the part of the Rendering Pipeline where fragments are processed. Each fragment will pass through a Fragment Shader (also known as a **Pixel Shader**). These Shaders tend to involve a lot more complex activity compared to Vertex Shaders, such as depth testing, alpha testing, colorization, texture sampling, Lighting, Shadows, and various Post-Processing effects to name a few of the possibilities. This data is then drawn onto the Frame Buffer, which holds the current image that will eventually be sent to the display device (our monitor) once rendering tasks for the current frame are complete.

There are normally two Frame Buffers in use by Graphics APIs by default (although more could be generated for custom rendering scenarios). At any given moment, one of the Frame Buffers contains the data from the frame we just rendered to, and is being presented to the screen, while the other is actively being drawn to by the GPU while it completes commands from the Command Buffer. Once the GPU reaches a `swap buffers` command (the final instruction the CPU asks it to complete for the given frame), the Frame Buffers are flipped around so that the new frame is presented. The GPU will then use the old Frame Buffer to draw the next frame. This process repeats each time a new frame is rendered, hence the GPU only needs two Frame Buffers to handle this task.

This entire process, from making Graphics API calls to swapping Frame Buffers, repeats continuously for every mesh, vertex, fragment, and frame, so long as our application is still rendering.

There are two metrics that tend to be the source of bottlenecks in the Back End--Fill Rate and Memory Bandwidth. Let's explore them a little.

Fill Rate

Fill Rate is a very broad term referring to the speed at which the GPU can draw fragments. However, this only includes fragments that have survived all of the various conditional tests we might have enabled within the given Fragment Shader. A fragment is merely a *potential pixel*, and if it fails any of the enabled tests, then it is immediately discarded. This can be an enormous performance-saver, as the Rendering Pipeline can skip the costly drawing step and begin working on the next fragment instead.

One such example of a test that might cull a fragment is *Z-testing*, which checks whether the fragment from a closer object has already been drawn to the same fragment location (the Z refers to the depth dimension from the point of view of the Camera). If so, the current fragment is discarded. If not, then the fragment is pushed through the Fragment Shader and drawn over the target pixel, which consumes exactly one *fill* from our Fill Rate. Now, imagine multiplying this process by thousands of overlapping objects, each of which generates hundreds or thousands of possible fragments (higher screen resolutions require more fragments to be processed). This could easily lead to millions of fragments to process each and every frame due to all of the possible overlap from the perspective of the Main Camera. On top of this, we're trying to repeat this process dozens of times every second. This is why performing so much initial setup in the Rendering Pipeline is important, and it should be fairly obvious that skipping as many of these draws as we can will result in big rendering cost savings.

Graphics card manufacturers typically advertise a particular Fill Rate as a feature of the card, usually in the form of Gigapixels per-second, but this is a bit of a misnomer, as it would be more accurate to call it Gigafragments per-second; however, this argument is mostly academic. Either way, larger values tell us that the device can potentially push more fragments through the Rendering Pipeline. So, with a budget of 30 Gigapixels per-second and a target frame rate of 60 Hz, we can afford to process $30,000,000,000/60 = 500$ million *fragments* per-frame before being bottlenecked on Fill Rate. With a resolution of 2560 x 1440, and a best-case scenario where each pixel is drawn

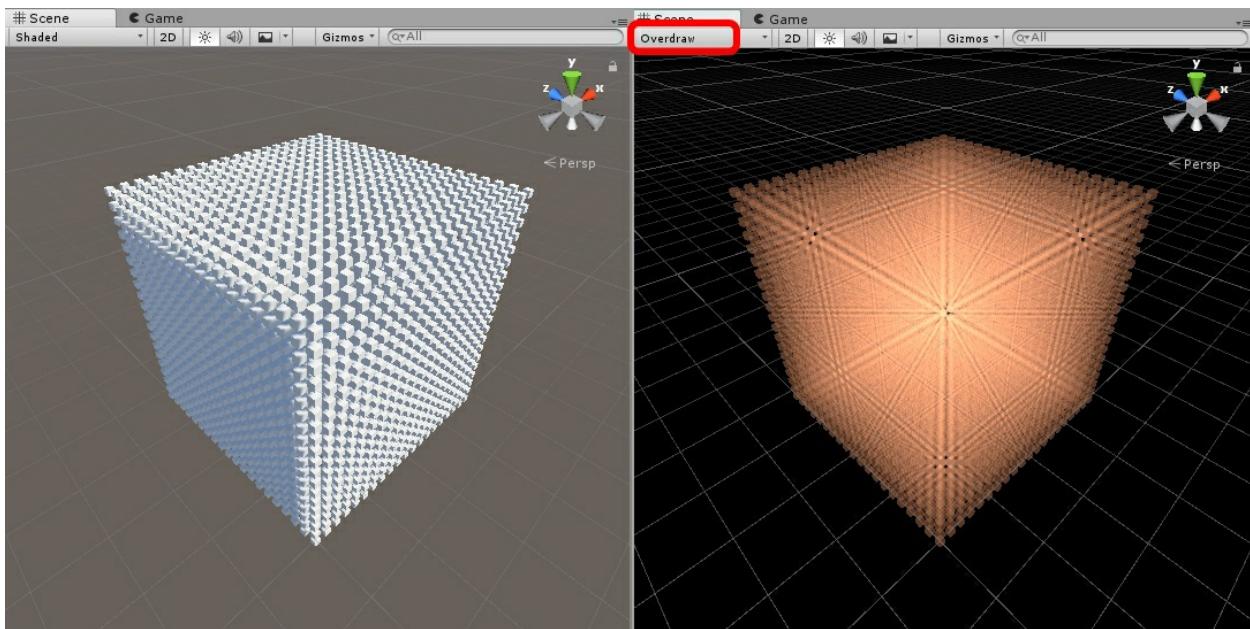
over only once, we could theoretically draw the entire Scene about 125 times without any noticeable problems.

Sadly, this is not a perfect world. Fill Rate is also consumed by other advanced rendering techniques, such as Shadows and Post-Processing effects that needs to take the same fragment data and perform their own *passes* on the Frame Buffer. Even so, we will always end up with some amount of redraw over the same pixels due to the order in which objects are rendered. This is known as **Overdraw**, and it is a useful metric to measure how efficiently we are making use of our Fill Rate.

Overdraw

How much Overdraw we have can be represented visually by rendering all objects with additive alpha blending and a flat coloration. Areas of high Overdraw will show up more brightly, as the same pixel is drawn over with additive blending multiple times. This is precisely how the Scene window's Overdraw Shading mode reveals how much Overdraw our Scene is undergoing.

The following screenshot shows a Scene with several thousand boxes drawn normally (left) versus the Scene window's Overdraw Shading mode (right):



The more Overdraw we have, the more Fill Rate we are wasting by overwriting fragment data. There are several techniques we can apply to reduce Overdraw, which we will explore later.

Note that there are actually several different queues used for rendering, which can be separated into two types: Opaque Queues and Transparent Queues. Objects rendered in one of the Opaque Queues can cull away fragments via Z-testing as explained previously. However, objects rendered in a





Transparent Queue cannot do so since their transparent nature means we can't assume that they won't need to be drawn no matter how many other objects are in the way, which leads to a lot of Overdraw. All Unity UI objects always render in a Transparent Queue, making them a significant source of Overdraw.

Memory Bandwidth

The other potential source of bottlenecks in the Back End comes from Memory Bandwidth. Memory Bandwidth is consumed whenever a texture must be pulled from a section of the GPU's VRAM down into the lower memory levels. This typically happens when a texture is *sampled*, where a Fragment Shader attempts to pick the matching texture pixel (or *texel*) to draw for a given fragment at a given location. The GPU contains multiple cores that each have access to the same area of VRAM, but they also each contain a much smaller, local Texture Cache that stores the texture(s) the GPU has been most recently working with. This is similar in design with the multitude of CPU memory cache levels that allow memory transfer up and down the chain. This is a hardware design workaround for the fact that faster memory will, invariably, be more difficult and expensive to produce. So, rather than having a giant, expensive block of VRAM, we have a large, cheap block of VRAM, but use a smaller, very fast, lower-level Texture Cache to perform sampling with, which gives us the best of both worlds, that is, fast sampling with lower costs.

If a texture is needed that is already within the core's local Texture Cache, then sampling often becomes lightning fast and is barely perceivable. If not, then the texture must be pulled in from VRAM before it can be sampled. This is effectively a cache miss for the Texture Cache since it will now take time to find and pull the required texture from VRAM. This transfer consumes a certain amount of our available Memory Bandwidth, specifically an amount equal to the total size of the texture file stored within VRAM (which may not be the exact size of the original file or the size in RAM, due to GPU-level compression techniques).

In the event that we are bottlenecked in Memory Bandwidth, the GPU will keep fetching the necessary texture files, but the entire process will be throttled, as the Texture Cache keeps waiting for data to appear before it can process a given batch of fragments. The GPU won't be able to push data back to the Frame Buffer in time to be rendered onto the screen, blocking the whole process and culminating in a poor frame rate.

Proper usage of Memory Bandwidth is another budgeting concern. For example, with a Memory Bandwidth of 96 GBs per-second per-core and a target frame rate of 60 frames-per-second, the GPU can afford to pull roughly $1.6 \text{ GBs} (96/60)$ worth of texture data every frame before being bottlenecked in Memory Bandwidth. This is not an exact budget, of course, because of the potential for cache misses, but it does give us a rough value to work with.



Memory Bandwidth is normally listed on a per-core basis, but some GPU manufacturers may try to mislead by multiplying Memory Bandwidth with the number of cores in order to list a bigger, but less practical number. Due to this, research may be necessary to compare apples to apples.

Note that this value is not the maximum limit on the amount of texture data that our game can contain in the project, nor in CPU RAM, nor even in VRAM. It is a metric that essentially limits how much texture swapping can occur during one frame. The same texture could be pulled back and forth multiple times in a single frame, depending on how many Shaders need to use them, the order that the objects are rendered, and how often texture sampling must occur. It's possible for only a handful of objects to consume whole Gigabytes of Memory Bandwidth since there is only a finite amount of Texture Cache space available. Having a Shader that needs a lot of large textures is more likely to cause cache misses thus causing a bottleneck on Memory Bandwidth. This is surprisingly easy to trigger if we consider multiple objects requiring different high quality textures and multiple secondary texture maps (Normal maps, Emission maps, and so on), which are not batched together. In this case, the Texture Cache will be unable to hang on to a single texture file long enough to immediately sample it during the next rendering pass.

Lighting and Shadowing

In modern games, a single object rarely finishes rendering completely in a single step, primarily due to Lighting and Shadowing. These tasks are often handled in multiple *passes* of a Fragment Shader, once for each of the several Light sources, and the final result is combined so that multiple Lights are given a chance to be applied. The result appears much more realistic, or at least, more visually appealing.

Several passes are required to gather Shadowing information. We will first set up our Scene to have Shadow Casters and Shadow Receivers, which will create or receive Shadows, respectively. Then, each time a Shadow Receiver is rendered, the GPU renders any Shadow Caster objects from the point of view of the Light source into a texture with the goal of collecting distance information for each of their fragments. It then does the same for the Shadow Receiver, except now that it knows which fragments the Shadow Casters would overlap from the Light source, it can render those fragments darker since they will be in the Shadow created by the Light source bearing down on the Shadow Caster.

This information then becomes an additional texture known as a *Shadowmap* and is blended with the surface for the Shadow Receiver when it is rendered from the point of view of the Main Camera. This will make its surface appear darker in certain spots where other objects stand between the Light source and the given object. A similar process is used to create *Lightmaps*, which are pregenerated Lighting information for the more static parts of our Scene.

Lighting and Shadowing tends to consume a lot of resources throughout all parts of the Rendering Pipeline. We need each vertex to provide a Normal direction (a vector pointing away from the surface) to determine how Lighting should reflect off that surface, and we might need additional vertex color attributes to apply some extra coloring. This gives the CPU and Front End more information to pass along. Since multiple passes of Fragment Shaders are required to complete the final rendering, the Back End is kept

busy both in terms of Fill Rate (lots and lots of pixels to draw, redraw, and merge) and in terms of Memory Bandwidth (extra textures to pull in or out for Lightmaps and Shadowmaps). This is why real-time Shadows are exceptionally expensive compared to most other rendering features and will inflate Draw Call counts dramatically when enabled.

However, Lighting and Shadowing are perhaps two of the most important parts of game art and design to get right, often making the extra performance requirements worth the cost. Good Lighting and Shadowing can turn a mundane Scene into something spectacular, as there is something magical about professional coloring that makes it visually appealing. Even a low-poly art style (for example, the mobile game Monument Valley) relies heavily on a good Lighting and Shadowing profile in order to allow the player to distinguish one object from another and create a visually pleasing Scene.

Unity offers multiple features that affect Lighting and Shadows, from real-time Lighting and Shadows (of which there are multiple types of each) to static Lighting called *Lightmapping*. There are a lot of options to explore, and of course, a lot of things that can cause performance issues if we're not careful.



The Unity documentation covers all of the various Lighting features in an excellent amount of detail. Start with the following pages and work through them; doing so will be well worth the time since these systems affect the entire Rendering Pipeline; refer to

<https://docs.unity3d.com/Manual/LightingOverview.html>
<https://unity3d.com/learn/tutorials/topics/graphics/introduction-1-lighting-and-rendering>.

There are two different rendering formats, which can greatly affect our Lighting performance, known as **Forward Rendering and Deferred Rendering**. The setting for these Rendering options can be found under Edit | Project Settings | Player | Other Settings | Rendering and configured on a per-platform basis.

Forward Rendering

Forward Rendering is the traditional form of rendering Lights in our Scene, as explored previously. During Forward Rendering, each object will be rendered in multiple passes through the same Shader. How many passes are required will be based on the number, distance, and brightness of Light sources. Unity will try to prioritize the `DirectionalLight` Component that is affecting the object the most and render the object in a *base* pass as a starting point. It will then take several of the most powerful `PointLight` Components nearby and re-render the same object multiple times through the same Fragment Shader. Each of these Point Lights will be processed on a per-vertex basis and all remaining Lights will be condensed into an *average* color by means of a technique called Spherical Harmonics.

Some of this behavior can be simplified by setting a Light's Render Mode to values such as Not Important and changing the value of Edit | Project Settings | Quality | Pixel Light Count. This value limits the number of Lights that will be gathered for Forward Rendering, but is overridden by any Lights with a Render Mode set to Important. It is, therefore, up to us to use this combination of settings responsibly.

As we might imagine, using Forward Rendering can utterly explode our Draw Call count very quickly in Scenes with a lot of Point Lights present due to the number of Render States being configured and Shader passes required.



More information on Forward Rendering can be found in the Unity documentation at <http://docs.unity3d.com/Manual/RenderTech-ForwardRendering.html>.

Deferred Rendering

Deferred Rendering, or *Deferred Shading* as it is sometimes known, is a technique that has been available on GPUs for about a decade or so, but it has not resulted in a complete replacement of the Forward Rendering method due to the caveats involved and somewhat limited support on mobile devices.

Deferred Shading is named as such because actual Shading does not occur until much later in the process, that is, it is deferred until later. It works by creating a geometry buffer (called a *G-Buffer*), where our Scene is initially rendered without any Lighting applied. With this information, the Deferred Shading system can generate a Lighting profile within a single pass.

From a performance perspective, the results are quite impressive as it can generate very good per-pixel Lighting with little Draw Call effort. One disadvantage is that effects such as anti-aliasing, transparency, and applying Shadows to animated characters cannot be managed through Deferred Shading alone. In this case, the Forward Rendering technique is applied as a fallback to cover those tasks, thus requiring extra Draw Calls to complete it. A bigger issue with Deferred Shading is that it often requires more powerful and more expensive hardware and is not available for all platforms, so fewer users will be able to make use of it.



The Unity documentation contains an excellent source of information on the Deferred Shading technique and its benefits and pitfalls, which is found at [http://docs.unity3d.com/Manual/Rend erTech-DeferredShading.html](http://docs.unity3d.com/Manual/RenderTech-DeferredShading.html).

Vertex Lit Shading (legacy)

Technically, there are more than two Lighting methods. The remaining two are *Vertex Lit Shading* and a very primitive, feature-lax version of Deferred Rendering. Vertex Lit Shading is a massive simplification of Lighting, as Lighting will only be considered per-vertex and not per-pixel. In other words, entire faces are colored the same based on the incoming Light color rather than blending Lighting colors across the face through individual pixels.

It is not expected that many, or really any, 3D games will make use of this legacy technique, as a lack of Shadows and proper Lighting make visualizations of depth very difficult. It is mostly used by simple 2D games that don't need to make use of Shadows, Normal maps, and various other Lighting features.

Global Illumination

Global Illumination, or *GI* for short, is an implementation of *baked Lightmapping*. Lightmapping is similar to the Shadowmaps created by Shadowing techniques in that one or more textures are generated for each object that represents extra Lighting information and is later applied to the object during its Lighting pass of a Fragment Shader to simulate static Lighting effects.

The main difference between these Lightmaps and other forms of Lighting is that Lightmaps are pregenerated (or *baked*) in the Editor and packaged into the game build. This ensures that we don't need to keep regenerating this information at runtime, saving numerous Draw Calls and significant GPU activity. Since we can bake this data, we have the luxury of time to generate very high-quality Lightmaps (at the expense of larger generated texture files we need to work with, of course).

Since this information is baked ahead of time, it cannot respond to real-time activity during gameplay, and so by default, any Lightmapping information will only be applied to static objects that were present in the Scene when the Lightmap was generated and at the exact location they were placed.

However, Light Probes can be added to the Scene to generate an additional set of Lightmap textures that can be applied to nearby dynamic objects that move, allowing such objects to benefit from pregenerated Lighting. This won't have pixel-perfect accuracy and will cost disk space for the extra Light Probe maps and Memory Bandwidth at runtime to swap them around, but it does generate a more believable and pleasant Lighting profile.

There have been several techniques of generating Lightmaps developed throughout the years, and Unity has used a couple of different solutions since its initial release. Global Illumination is simply the latest generation of the mathematical techniques behind Lightmapping, which offers very realistic coloring by calculating not only how Lighting affects a given object, but also how light reflects off nearby surfaces, allowing an object to affect the Lighting profile of those around it. This effect is calculated by an internal

system called Enlighten. This tool is used both to create static Lightmaps, as well as create something called *Pre-computed Real-time GI*, which is a hybrid of real-time and static Shading and allows us to simulate effects such as *time-of-day* (where the direction of light from the Sun changes over time) without relying on expensive real-time Lighting effects.

A typical issue with generating Lightmaps is the length of time it can take to generate them and get visual feedback on the current settings, because the Lightmapper is often trying to generate full-detail Lightmaps in a single pass. If the user attempts to modify its configuration, then the entire job must be canceled and started over. To solve this problem, Unity Technologies implemented the Progressive Lightmapper, which performs Lightmapping tasks more gradually over time, but also allows them to be modified while they are being calculated. This makes Lightmaps of the Scene appear to get progressively more detailed as it works in the background while also allowing us to change certain properties when it is still working and without having to restart the entire job. This provides almost immediate feedback and improves the workflow of generating Lightmaps immensely.

Multithreaded Rendering

Multithreaded Rendering is enabled by default on most systems, such as desktop and console platforms whose CPUs provide multiple cores. Other platforms still support many low-end devices to enable this feature by default, so it is a toggleable option for them. For Android, it can be enabled via a checkbox under Edit | Project Settings | Player | Other Settings | Multithreaded Rendering, whereas for iOS, Multithreaded Rendering can be enabled by configuring the application to make use of Apple's Metal API under Edit | Project Settings | Player| Other Settings | Graphics API. At the time of writing this book, WebGL does not support Multithreaded Rendering.

For each object in our Scene, there are three tasks to complete: determine whether the object needs to be rendered (through a technique known as *Frustum Culling*), and if so, generate commands to render the object (since rendering a single object can result in dozens of different commands), and then send the command to the GPU using the relevant Graphics API. Without Multithreaded Rendering, all of these tasks must happen on the main thread of the CPU, thus any activity on the main thread becomes part of the critical path for all rendering. When Multithreaded Rendering is enabled, the task of pushing commands into the GPU are handled by a *render thread*, whereas other tasks such as culling and generating commands get spread across multiple *worker threads*. This setup can save an enormous number of CPU cycles for the main thread, which is where the overwhelming majority of other CPU tasks take place, such as physics and script code.

Enabling this feature will affect what it means to be CPU bound. Without Multithreaded Rendering, the main thread is performing all of the work necessary to generate instructions for the Command Buffer, meaning that any performance we can save elsewhere will free up more time for the CPU to generate commands. However, when Multithreaded Rendering is taking place, a good portion of the workload is pushed onto separate threads, meaning that improvements to the main thread will have less of an impact on rendering performance via the CPU.



Note that being GPU bound is the same regardless of whether Multithreaded Rendering is taking place. The GPU always performs its tasks in a multithreaded fashion.

Low-level rendering APIs

Unity exposes a rendering API to us through their `commandBuffer` class. This allows us to control the Rendering Pipeline directly through our C# code by issuing high-level rendering commands, such as `render this object, with this Material, using this Shader, OR draw N instances of this piece of procedural geometry`. This customization is not as powerful as having direct Graphics API access, but it is a step in the right direction for Unity developers to customize unique graphical effects.

Check out the Unity documentation on `CommandBuffer` to make use of this feature at <http://docs.unity3d.com/ScriptReference/Rendering.CommandBuffer.html>.

If an even more direct level of rendering control is needed, such that we wish to make direct Graphics API calls to OpenGL, DirectX, and Metal, then be aware that it is possible to create a Native Plugin (a small library written in C++ code that is compiled specifically for the architecture of the target platform) that hooks into the Unity's Rendering Pipeline, setting up callbacks for when particular rendering events happen, similar to how hook into various callbacks of the main Unity Engine. This is certainly an advanced topic for most Unity users, but useful to know for the future as our knowledge of rendering techniques and Graphics APIs matures.

Unity provides some good documentation on generating a rendering interface in a Native Plugin

at <https://docs.unity3d.com/Manual/NativePluginInterface.html>.

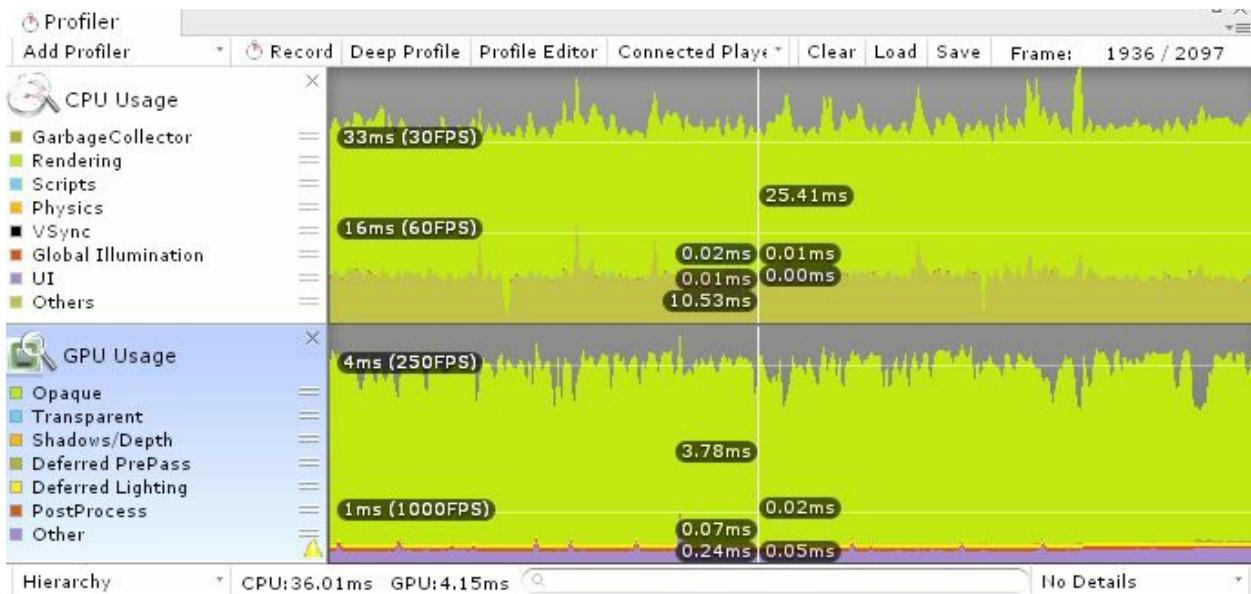
Detecting performance issues

It should be obvious that due to the number of complex processes involved, there are a lot of different ways in which the GPU can become bottlenecked. Now that we have a thorough understanding of the Rendering Pipeline and how bottlenecks may occur, let's explore how to detect these problems.

Profiling rendering issues

The Profiler can be used to quickly narrow down which of the two devices used in the Rendering Pipeline we are bottlenecked within--whether it is the CPU or GPU. We must examine the problem using both the CPU Usage and GPU Usage Areas of the Profiler window, as this can tell us which device is working the hardest.

The following screenshot shows Profiler data for a CPU-bound application. The test involved creating thousands of simple cube objects, with no batching or Shadowing techniques taking place. This resulted in an extremely large Draw Call count (around 32,000) for the CPU to generate commands for, but giving the GPU relatively little work to do due to the simplicity of the objects being rendered:



This example shows that the CPU's Rendering task is consuming a large amount of cycles (around 25 ms per-frame), whereas the GPU is processing for less than 4 milliseconds, indicating that the bottleneck resides in the CPU. Note that this Profiling test was performed against a standalone app, not within the Editor. We now know that our rendering is CPU bound and can begin to apply some CPU-saving performance improvements (being careful

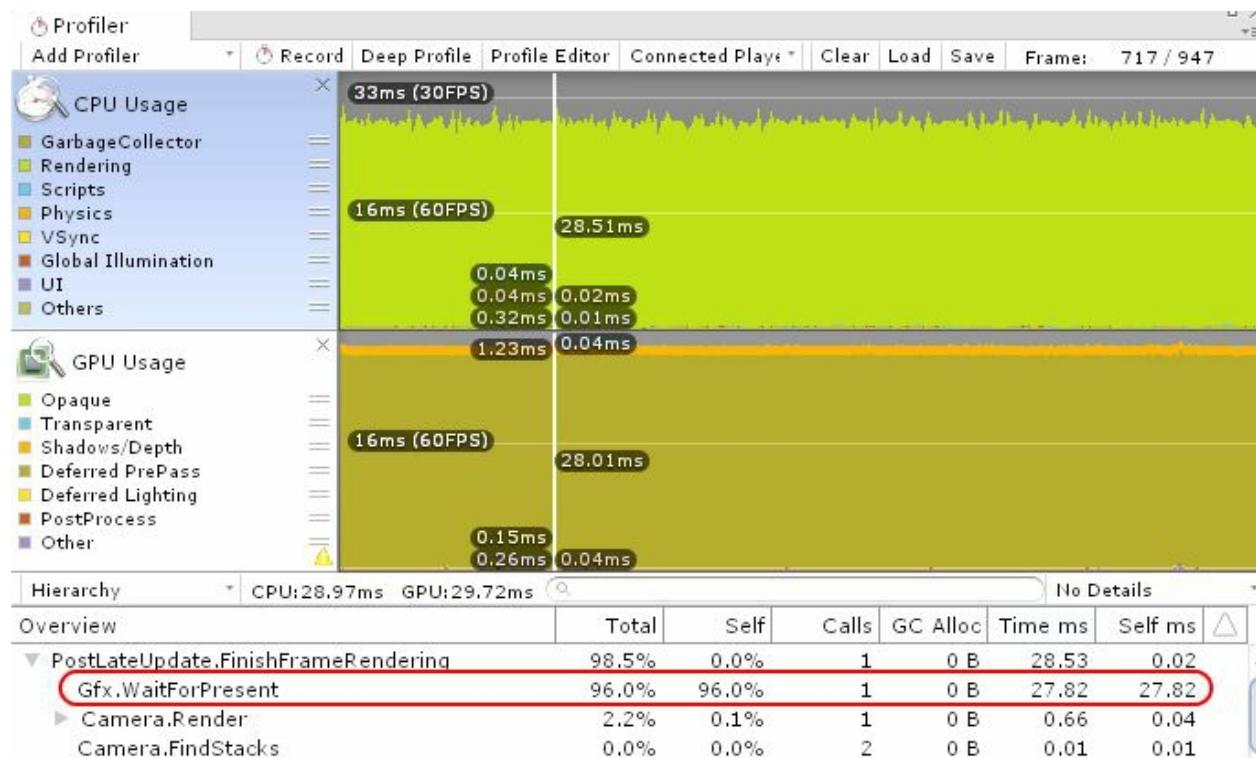
not to introduce rendering bottlenecks elsewhere by doing so).

Meanwhile, Profiling a GPU-bound application via the Profiler is a little trickier. This time, the test involves creating a simple object requiring minimal Draw Calls, but using a very expensive Shader that samples a texture thousands of times to create an absurd amount of activity in the Back End.



To perform fair GPU-bound Profiling tests, you should ensure that you disable Vertical Sync through Edit | Project Settings | Quality | Other | V Sync Count, otherwise it is likely to pollute our data.

The following screenshot shows Profiler data for this test when it is run in a standalone application:



As we can see in the preceding screenshot, the rendering task of the CPU Usage Area matches closely with the total rendering costs of the GPU Usage Area. We can also see that the CPU and GPU time costs at the bottom of the image are relatively similar (about 29 milliseconds each). This is somewhat

confusing as we seem to be bottlenecked equally in both devices, where we would expect the GPU to be working much harder than the CPU.

In actuality, if we drill down into the Breakdown View of the CPU Usage Area using the Hierarchy Mode, we will note that most of the CPU time is spent on the task labeled Gfx.WaitForPresent. This is the amount of time that the CPU is wasting while it waits for the GPU to finish the current frame. Hence, we are in fact bottlenecked by the GPU despite appearing as though we are bound by both. Even if Multithreaded Rendering is enabled, the CPU must still wait for the Rendering Pipeline to finish before it can begin the next frame.



Gfx.WaitForPresent is also used to signal that the CPU is waiting on Vertical Sync to complete, hence the need to disable it for this test.

Brute-force testing

If we're pouring over our Profiling data and still not sure where the source of the problem resides, or we're GPU-bound and need to determine where we're bottlenecked in the Rendering Pipeline, we should try the brute-force method, that is, cull a specific activity from the Scene and check whether it results in greatly improved performance. If a small change results in a big speed improvement, then we have a strong clue about where the bottleneck lies. There's no harm in this approach if we eliminate enough unknown variables to ensure that the data is leading us in the right direction.

The obvious brute-force test for CPU bounding will be to reduce Draw Calls to check whether performance suddenly improves. However, this is often not possible since, presumably, we've already been reducing our Draw Calls to a minimum through techniques such as Static Batching, Dynamic Batching, and Atlasing. This would mean that we have very limited scope for reducing them further.

What we can do, however, is intentionally increase our Draw Call count by a small number, either by introducing more objects or disabling Draw Call-saving features such as Static and Dynamic Batching and observe whether the situation gets significantly worse than before. If so, then we have evidence that we're either very close to being CPU bound or have already become so.

There are two good brute-force tests we can apply to a GPU-bound application to determine whether we're bound by Fill Rate or by Memory Bandwidth: reducing screen resolution or reducing texture resolution, respectively.

By reducing screen resolution, we will ask the Rasterizer to generate significantly fewer fragments and transpose them over a smaller canvas of pixels for the Back End to process. This will reduce the Fill Rate consumption of the application, giving this key part of the Rendering Pipeline some additional breathing room. Ergo, if performance suddenly improves with a screen resolution reduction, then Fill Rate should be our primary

concern.



A reduction from a resolution of 2560 x 1440 to 800 x 600 is an improvement factor of about eight, which is often more than enough to reduce Fill Rate costs enough to make the application perform well again.

Similarly, if we're bottlenecked on Memory Bandwidth, then reducing texture quality is likely to result in a significant performance improvement. By doing so, we have shrunk the size of our textures, greatly reducing the Memory Bandwidth costs of our Fragment Shaders, allowing the GPU to fetch the necessary textures much quicker. Globally reducing texture quality can be achieved by going to Edit | Project Settings | Quality | Texture Quality and setting the value to Half Res, Quarter Res, or Eighth Res.

An application bound by the CPU has ample opportunities for performance enhancements through practically every performance-enhancing tip in this book. If we free up CPU cycles from other activities, then we can afford to render more objects through more Draw Calls, keeping in mind, of course, that each will cost us more activity in the GPU. There are, however, additional opportunities to make some indirect improvements in Draw Call count while we try to improve other parts of the Rendering Pipeline. This includes Occlusion Culling, tweaking our Lighting and Shadowing behavior, and modifying our Shaders. These will be explained in the following sections as we investigate various performance enhancements.

Meanwhile, we will probably need to apply a little brute-force testing and guesswork to determine how a GPU-bound application is bottlenecked. Most applications are bottlenecked by Fill Rate or Memory Bandwidth, so we should start there. It is rare to find performance bottlenecks in the Front End, at least on desktop applications, so it is worth checking only after we've verified that the other sources are not the problem. Vertex Shaders are often trivial compared to Fragment Shaders, and so the only real opportunity to cause problems with Front End processing is either to push too much geometry or to have overly complex Geometry Shaders.

Ultimately, this investigation should help us determine whether we are CPU-bound or GPU-bound, and in the latter case, whether we are bound by the

Front End or Back End, and again in the latter case, whether we are bound by Fill Rate or Memory Bandwidth. With this knowledge, there are a number of techniques we can apply to improve performance.

Rendering performance enhancements

We should now have all of the information we need to make sense of performance bottlenecks so that we can start to apply fixes. For the remainder of this chapter, we will cover a series of techniques to improve Rendering Pipeline performance for CPU-bound and GPU-bound applications.

Enable/Disable GPU Skinning

The first tip involves a setting that eases the burden on the CPU or GPU Front End at the expense of the other, that is, GPU Skinning. Skinning is the process where mesh vertices are transformed based on the current location of their animated bones. The animation system, working on the CPU, transforms the object's bones that are used to determine its current pose, but the next important step in the animation process is wrapping the mesh vertices around those bones to place the mesh in the final pose. This is achieved by iterating over each vertex and performing a weighted average against the bones connected to those vertices.

This vertex processing task can either take place on the CPU or within the Front End of the GPU, depending on whether the GPU Skinning option is enabled. This feature can be toggled under Edit | Project Settings | Player Settings | Other Settings | GPU Skinning. Enabling this option pushes skinning activity to the GPU, although bear in mind that the CPU must still transfer the data to the GPU and will generate instructions on the Command Buffer for the task, so it doesn't remove the CPU's workload entirely. Disabling this option eases the burden on the GPU by making the CPU resolve the mesh's pose before transferring mesh data across and simply asking the GPU to draw it as is. Obviously, this feature is useful if we have lots of animated meshes in our Scenes and can be used to help either bounding case by pushing the work onto the device that is least busy.

Reduce geometric complexity

This tip concerns the GPU Front End. We have already covered some techniques on mesh optimization in [Chapter 4, Kickstart Your Art](#), which can help reduce our mesh's vertex attributes. As a quick reminder, it is not uncommon to use a mesh that contains a lot of unnecessary UV and Normal vector data, so our meshes should be double-checked for this kind of superfluous fluff. We should also let Unity optimize the structure for us, which minimizes cache misses as vertex data is read within the Front End.

The goal is to simply reduce actual vertex counts. There are three solutions to this. First, we can simplify the mesh by either having the art team manually tweak and generate meshes with lower polycounts or using a *mesh decimation* tool to do it for us. Second, we could simply remove meshes from the Scene, but this should be a last resort. The third option is to implement automatic culling through features such as **Level of Detail (LOD)**, which will be explained later in this chapter.

Reduce Tessellation

Tessellation through Geometry Shaders can be a lot of fun, as it is a relatively underused technique that can really make our graphical effects stand out from among the crowd of games that use only the most common effects. However, it can contribute enormously to the amount of processing work taking place in the Front End.

There aren't really any simple tricks we can exploit to improve Tessellation, besides improving our Tessellation algorithms or easing the burden caused by other Front End tasks to give our Tessellation tasks more room to breathe. Either way, if we have a bottleneck in the Front End and are making use of Tessellation techniques, we should double-check that they are not consuming the lion's share of the Front End's budget.

Employ GPU Instancing

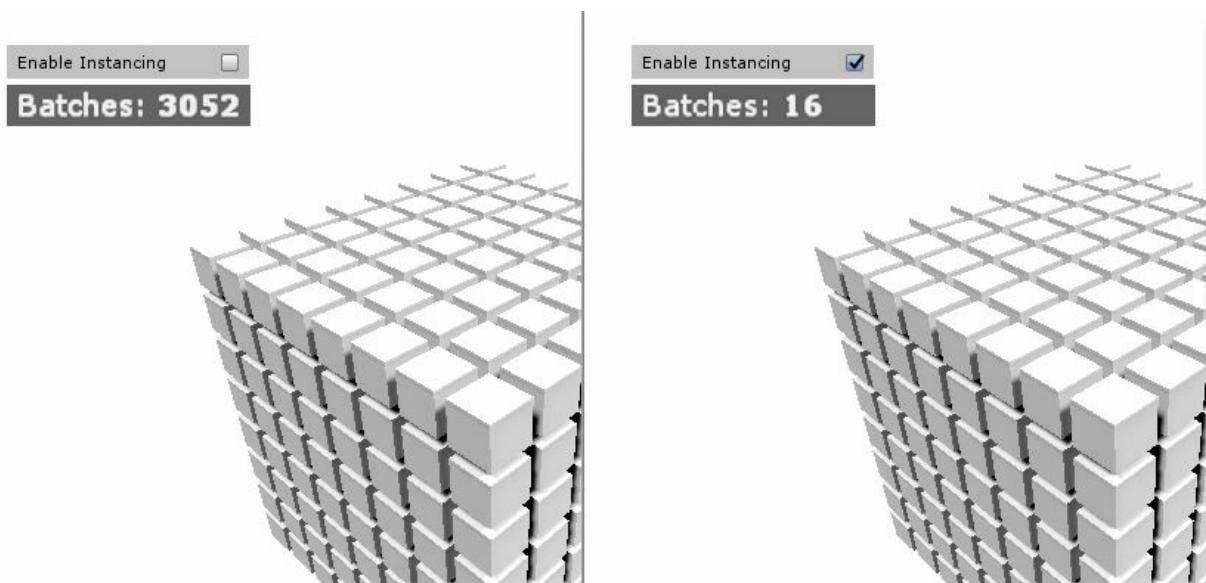
GPU Instancing is a means to render multiple copies of the same mesh quickly by exploiting the fact that they will have identical Render States, hence require minimal Draw Calls. This is practically identical to Dynamic Batching, except that it is not an automatic process. In fact, we can think of Dynamic Batching as *poor-man's GPU Instancing* since GPU Instancing can enable even better savings and allows for more customization by allowing parameterized variations.

GPU Instancing is applied at the Material level with the Enable Instancing checkbox, and variations can be introduced by modifying Shader code. This way, we can give different instances different rotations, scales, colors, and so on. This is useful for rendering Scenes such as forests and rocky areas where we want to render hundreds or thousands of different copies of a mesh with some slight variation.



Note that Skinned Mesh Renderers cannot be instanced for similar reasons that they cannot be Dynamically Batched, and not all platforms and APIs support GPU Instancing.

The following screenshot shows the benefits of GPU Instancing on a group of 512 cube objects (with some extra Lighting and Shadowing applied to increase the total Draw Call count):



This system is much more versatile than Dynamic Batching since we have more control of how objects are batched together. Of course, there are more opportunities for mistakes if we batch things in inefficient ways, so we should be careful to use it wisely.

Check out the Unity documentation for more information on GPU Instancing at <https://docs.unity3d.com/Manual/GPUInstancing.html>.

Use mesh-based Level Of Detail (LOD)

LOD is a broad term referring to the dynamic replacement of features based on their distance from the Camera and/or how much space they take up in the Camera's view. Since it can be difficult to tell the difference between a low- and high-quality object at great distances, there is very little reason to render the high-quality version, and so we may as well dynamically replace distant objects with something more simplified. The most common implementation of LOD is mesh-based LOD, where meshes are dynamically replaced with lower detailed versions as the Camera gets farther and farther away.

Making use of mesh-based LOD can be achieved by placing multiple objects in the Scene and making them children of a `GameObject` with an attached `LODGroup` Component. The LOD Group's purpose is to generate a bounding-box from these objects and decide which object should be rendered based on the size of the bounding-box within the Camera's field of view. If the object's bounding-box consumes a large area of the current view, then it will enable the mesh(es) assigned to lower LOD Groups, and if the bounding-box is very small, it will replace the mesh(es) with those from higher LOD Groups. If the mesh is too far away, it can be configured to hide all child objects. So, with the proper setup, we can have Unity replace meshes with simpler alternatives, or cull them entirely, which eases the burden on the rendering process.

Check out the Unity documentation for more detailed information on the mesh-based LOD feature at <http://docs.unity3d.com/Manual/LevelOfDetail.html>.

This feature can cost us a large amount of development time to fully implement; artists must generate lower polygon count versions of the same object, and level designers must generate LOD Groups, configure them, and test them to ensure that they don't cause jarring transitions as the Camera moves closer or farther away.

Note that some game development middleware companies offer



third-party tools for automated LOD mesh generation. These might be worth investigating to compare their ease of use versus quality loss versus cost effectiveness.

Mesh-based LOD will also cost us in disk footprint as well as RAM and CPU; the alternative meshes need to be bundled, loaded into RAM, and the `LODGroup` Component must routinely test whether the Camera has moved to a new position that warrants a change in LOD level. The benefits on the Rendering Pipeline are rather impressive, however. Dynamically rendering simpler meshes reduces the amount of vertex data we need to pass and potentially reduces the number of Draw Calls, Fill Rate, and Memory Bandwidth needed to render the object.

Due to the number of sacrifices needed for mesh-based LOD to function, developers should avoid preoptimizing by automatically assuming that mesh-based LOD will help them. Excessive use of the feature will lead to burdening other parts of our application's performance and chew up precious development time, all for the sake of paranoia. It should only be used if we start to observe problems in the Rendering Pipeline, and we've got CPU, RAM, and development time to spare.

Having said that, Scenes that feature large, expansive views of the world and have lots of Camera movement, might want to consider implementing this technique very early, as the added distance and massive number of visible objects will likely exacerbate the vertex count enormously. As a counter example, Scenes that are always indoors or feature a Camera with a viewpoint looking down at the world will find little benefit in this technique since objects will tend to be at a similar distance from the Camera at all times. Examples include **Real-Time Strategy (RTS)** and **Multiplayer Online Battle Arena (MOBA)** games.

Culling Groups

Culling Groups are a part of the Unity API that effectively allows us to create our own custom LOD system as a means of coming up with our own ways of dynamically replacing certain gameplay or rendering behaviors. Some examples of things we might want to apply LOD to include replacing animated characters with a version with fewer bones, applying simpler Shaders, skipping Particle System generation at great distances, simplifying AI behavior, and so on.

Since the Culling Group system at its most basic level simply tells us whether objects are visible to the Camera, and how big they are, it also has other uses in the realm of Gameplay, such as determining whether certain enemy spawn points are currently visible to the player or whether a player is approaching certain areas. There are a wide range of possibilities available with the Culling Group system that makes it worth considering. Of course, the time spent to implement, test, and redesign Scenes to exploit can be significant.

Check out the Unity documentation for more information on Culling Groups at <https://docs.unity3d.com/Manual/CullingGroupAPI.html>.

Make use of Occlusion Culling

One of the best ways to reduce both Fill Rate consumption and Overdraw is to make use of Unity's Occlusion Culling system. The system works by partitioning the world into a series of small cells and flying a virtual Camera through the Scene, making note of which cells are invisible from other cells (are *occluded*) based on the size and position of the objects present.



Note that this is different from the technique of Frustum Culling, which culls objects outside the current Camera view. Frustum Culling is always active and automatic. Objects culled by this process are, therefore, automatically ignored by the Occlusion Culling system.

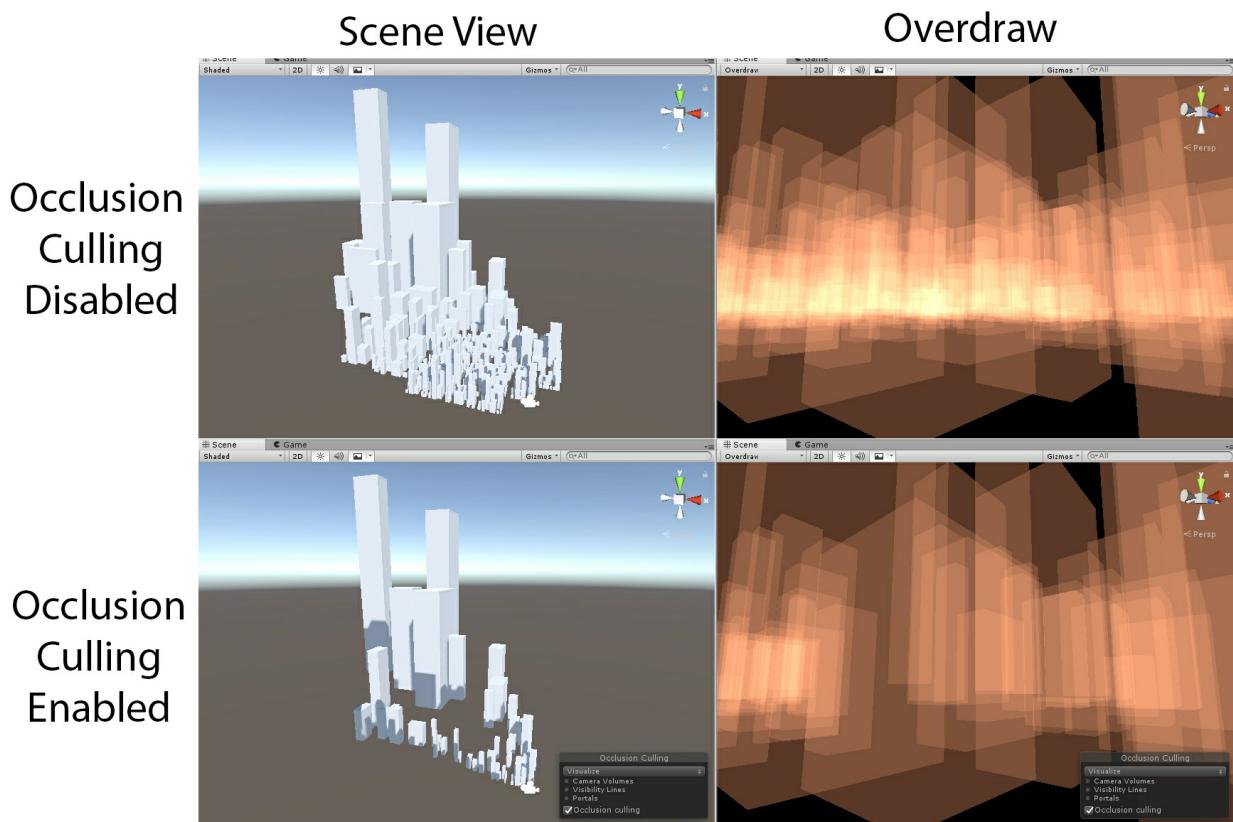
Occlusion Culling data can only be generated for objects properly labeled Occluder Static and/or Occludee Static under the StaticFlags dropdown. Occluder Static is the general setting for static objects we expect to be so large that they will both occlude and be occluded by other objects, such as sky scrapers or mountains, which can hide other objects behind them, as well as be hidden behind each other, and so on. Occludee Static is a special case for things such as transparent objects that always require other objects behind them to be rendered, but they themselves need to be hidden if something large blocks their visibility.



Naturally, because Static flags must be enabled for Occlusion Culling, this feature will not work for dynamic objects.

The following screenshot shows how effective Occlusion Culling can be at reducing the number of rendered objects from our Scene from an external point of view for the sake of demonstration. From the point of view of the Main Camera, the two situations appear identical.

The Rendering Pipeline is not wasting time rendering objects that are obscured by closer ones:



Enabling the Occlusion Culling feature will cost additional disk space, RAM, and CPU time. Extra disk space is required to store the occlusion data, extra RAM is needed to keep the data structure in memory, and there will be a CPU processing cost to determine which objects are being occluded in each frame. The Occlusion Culling data structure must be properly configured to create cells of the appropriate size for our Scene and the smaller the cells, the longer it takes to generate the data structure. However, if it is configured correctly for the Scene, Occlusion Culling can provide both Fill Rate savings through reduced Overdraw and Draw Call savings by culling nonvisible objects.



Note that even though an object may be culled by occlusion, its Shadows must still be calculated, so we won't save any Draw Calls or Fill Rate from those tasks.

Optimizing Particle Systems

Particle Systems are useful for a huge number of different visual effects, and usually the more particles they generate, the better the effect looks. However, we will need to be responsible about the number of particles generated and the complexity of Shaders used since they can touch on all parts of the Rendering Pipeline; they generate a lot of vertices for the Front End (each particle is a quad) and could use multiple textures, which consume Fill Rate and Memory Bandwidth in the Back End, so they can potentially cause an application to be bound anywhere if used irresponsibly.

Reducing Particle System density and complexity are fairly straightforward-- use fewer Particle Systems, generate fewer particles, and/or use fewer special effects. Atlasing is also another common technique to reduce Particle System performance costs. However, there is an important performance consideration behind Particle Systems that is not too well known and happens behind the Scenes, and that is the process of automatic Particle System culling.

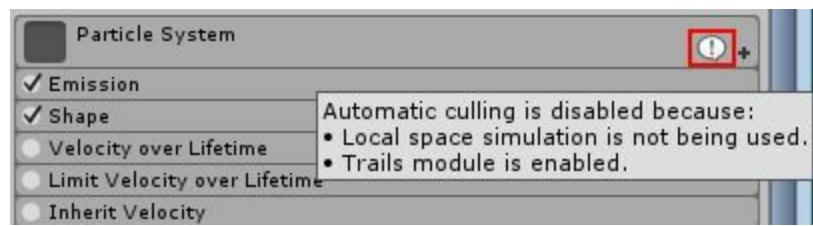
Make use of Particle System Culling

Unity Technologies have released an excellent blog post covering this topic, which can be found at <https://blogs.unity3d.com/2016/12/20/unitytips-particlesystem-performance-culling/>.

The basic idea is that all Particle Systems are either predictable or not (deterministic versus nondeterministic), depending on various settings. When a Particle System is predictable and not visible to the main view, then the entire Particle System can be automatically culled away to save performance. As soon as a predictable Particle System comes back into view, Unity can figure out exactly how the Particle System is meant to look at that moment as if it had been generating particles the entire time it wasn't visible. So long as the Particle System generates particles in a very procedural way, then the state is immediately solvable mathematically.

However, if any setting forces the Particle System to become unpredictable or *nonprocedural*, then it would have no idea what the current state of the Particle System needs to be, had it been hidden previously, and will hence need to render it fully every frame regardless of whether or not it is visible. Settings that break a Particle System's predictability include, but are not limited to making the Particle System render in world-space, applying external forces, collisions, and Trails, or using complex Animation Curves. Check out the blog post mentioned previously for a rigorous list of nonprocedural conditions.

Note that Unity provides a useful warning on Particle Systems when something would cause it to break automatic culling, as shown in the following screenshot:



Avoid recursive Particle System calls

Many methods available to a `ParticleSystem` Component are recursive calls. Calling them will iterate through each child of the Particle System, which then calls `GetComponent<ParticleSystem>()` on each child, and if the Component exists, it will call the appropriate method. This then repeats for each child `ParticleSystem` beneath the original parent, its grandchildren, and so on. This can be a huge problem with deep hierarchies of Particle Systems, which is sometimes the case with complex effects.

There are several `ParticleSystem` API calls affected by this behavior, such as `Start()`, `Stop()`, `Pause()`, `Clear()`, `Simulate()`, and `isAlive()`. We obviously cannot avoid these methods entirely since they represent the most common methods we would want to call on a Particle System. However, each of these methods has a `withChildren` parameter that defaults to `true`. By passing `false` in place of this parameter (for example, by calling `clear(false)`), it disables the recursive behavior and will not call into its children. Hence, the method call will only affect the given Particle System, thus reducing the overhead cost of the call.

This is not always ideal since we do often want all children of the Particle System to be affected by the method call. Another approach is to, therefore, cache the `ParticleSystem` Components in the same way we learned in [Chapter 2, Scripting Strategies](#), and iterate through them manually ourselves (making sure that we pass `false` for the `withChildren` parameter each time).



Note that there is a bug in Unity versions 5.4 to Unity 2017.1, where additional memory is allocated each time `stop()` and `simulate()` are called (even if the Particle System has already been stopped). This bug is fixed in Unity 2017.2.

Optimizing Unity UI

Unity's first few attempts at built-in UI Systems were not particularly successful; it is often quickly supplanted by products on the Asset Store. However, the latest generation of their solution (simply called Unity UI) has become a much more popular solution, so many developers are starting to rely on it for their UI needs so much so, in fact, that Unity Technologies bought the company behind the Text Mesh Pro asset in early 2017 and have merged it into the Unity UI as a built-in feature.

Let's explore a few techniques we can use to improve the performance of Unity's built-in UI.

Use more Canvases

A canvas Component's primary task is to manage the meshes that are used to draw the UI elements beneath them in the Hierarchy window and issues the Draw Calls necessary to render those elements. An important task of the Canvas is to batch these meshes together (which can only happen if they share the same Material) to reduce Draw Calls. However, when changes are made to a Canvas, or any of its children, this is known as *dirtying* the Canvas. When a Canvas is *dirty*, it needs to regenerate meshes for all of the UI elements beneath it before it can issue a Draw Call. This regeneration process is not a simple task and is a common source of performance problems in Unity projects, because unfortunately there are many things that can cause the Canvas to be made dirty. Even changing a single UI element within a Canvas can cause this to occur. There are so many things that cause dirtying, and so few that don't (and usually only in certain circumstances) that it's best to simply err on the side of caution and assume that any change will cause this effect.



Perhaps the only notable action that doesn't cause dirtying is changing a color property of a UI element.

If we find our UI causes a large spike in CPU usage any time something changes (or sometimes literally every frame if they're being changed every frame), one solution we can apply is to simply use more Canvases. A common mistake is to build the entire game's UI in a single Canvas and keep it this way as the game code and its UI continues to become more complex.

This means that it will need to check every element every time anything changes in the UI, which can become more and more disastrous on performance as more elements are crammed into a single Canvas. However, each Canvas is independent and does not need to interact with other Canvases in the UI and so by splitting up the UI into multiple Canvases we can separate the workload and simplify the tasks required by any single Canvas.



Ensure that you add a `GraphicsRaycaster` Component to the same `GameObject` as the child `Canvas` so that its own child elements can still be interacted with. Conversely, if none of the `Canvas`' child elements are interactable, then we can safely remove any `GraphicsRaycaster` Components from it to reduce performance costs.

In this case, even though an element still changes, fewer other elements will need to be regenerated in response, reducing the performance cost. The downside of this approach is that elements across different Canvases will not be batched together, so we should try to keep similar elements with the same Material grouped together within the same Canvas, if possible.



It's also possible to make a `Canvas` a child of another `Canvas`, for the sake of organization, and the same rules apply. If an element changes in one `Canvas`, the other will be unaffected.

Separate objects between static and dynamic canvases

We should strive to try and generate our Canvases in a way that groups elements based on when they get updated. We should think of our elements as fitting within one of three groups: *Static*, *Incidental Dynamic*, and *Continuous Dynamic*. Static UI elements are those that never change, and good examples of these are background images, labels, and so on. Dynamic elements are those that can change, where Incidental Dynamic objects are those UI elements that only change in response to something, such as a UI button press or a hover action, whereas Continuous Dynamic objects are those UI elements that update regularly, such as animated elements.

We should try to split UI elements from these three groups into three different Canvases for any given section of our UI, as this will minimize the amount of wasted effort during regeneration.

Disable Raycast Target for noninteractive elements

UI elements have a Raycast Target option, which enables them to be interacted with by clicks, taps, and other user behavior. Each time one of these events takes place, the `GraphicsRaycaster` Component will perform pixel-to-bounding-box checks to figure out which element has been interacted with and is a simple iterative `for` loop. By disabling this option for noninteractive elements, we're reducing the number of elements that the `GraphicsRaycaster` needs to iterate through, saving performance.

Hide UI elements by disabling the parent Canvas Component

The UI uses a separate Layout System to handle regeneration of certain element types, which operates in a similar way as dirtying a Canvas. `UIImage`, `UIText`, and `LayoutGroup` are all examples of Components that fall under this system. There are many things that can cause a Layout System to become dirty, most obvious of which are enabling and disabling such elements. However, if we want to disable a portion of the UI, we can avoid these expensive regeneration calls from the Layout System by simply disabling the `Canvas` Component they are children of. This can be done by setting the `Canvas` Component's `enabled` property to `false`. The drawback of this approach is that if any child objects that have some `Update()`, `FixedUpdate()`, `LateUpdate()`, or Coroutine code, then we would need to also disable them manually, otherwise they will continue to run. By disabling the `Canvas` Component, we're only stopping the UI from being rendered and interacted with, and we should expect various update calls to continue to happen as normal.

Avoid Animator Components

Unity's `Animator` Components were never intended to be used with the latest version of its UI System, and their interaction with it is a naive implementation. Each frame, the `Animator` will change properties on UI elements that causes their Layouts to be dirtied and cause regeneration of a lot of internal UI information. We should avoid using `Animators` entirely, and instead perform tweening ourselves or use a utility asset intended for such operations.

Explicitly define the Event Camera for World Space Canvases

Canvases can be used for UI interactions in both 2D and 3D. This is determined by whether the Canvas has its Render Mode setting configured to Screen Space (2D) or World Space (3D). Any time a UI interaction takes place, the `Canvas` Component will check its `eventCamera` property (exposed as Event Camera in the Inspector window) to figure out which Camera to use. By default, a 2D Canvas will set this property to the Main Camera, but a 3D Canvas leaves it set to `null`. This is unfortunate because each time the Event Camera is needed, it will still use the Main Camera, but do so by calling `FindObjectWithTag()`. Finding objects by Tag isn't as bad of a performance cost as using the other variations of `Find()`, but its performance cost scales linearly with the more Tags we use in a given project. To make matters worse, the Event Camera is accessed fairly often during a given frame for a World Space Canvas, which means leaving this property `null` will cause a huge performance hit for no real benefit. We should manually set this property to the Main Camera for all of our World Space Canvases.

Don't use alpha to hide UI elements

Rendering a UI element with an alpha value of `0` in its `color` property will still cause a Draw Call to be issued. We should favor changing the `IsActive` property of a UI element in order to hide it when necessary. Another alternative is to use Canvas Groups via `CanvasGroup` Components, which can be used to control the alpha transparency of all child elements beneath them. Setting the `alpha` value of a Canvas Group to `0` will cull away its child objects, and, therefore, no Draw Calls will be issued.

Optimizing ScrollRects

`ScrollRect` Components are UI elements that are used to scroll through a list of other UI elements and are fairly common in mobile applications.

Unfortunately, the performance of these elements scales very poorly with size since the Canvas needs to regenerate them regularly. There are a number of things we can do to improve the performance of our `ScrollRect` Components.

Make sure to use a RectMask2D

It's possible to create scrolling UI behavior by simply placing other UI elements with a lower `depth` value than the `scrollRect` elements. However, this is bad practice since there will be no culling taking place in the `scrollRect`, and every element will need to be regenerated for each frame that the `scrollRect` is moving. If we haven't already, we should use a `RectMask2D` Component to clip and cull child objects that are not visible. This Component creates a region of space whereby any child UI elements within it will be culled away if they are outside the bounds of the `RectMask2D` Component. The cost of determining whether to cull an object compared to the savings of rendering too many invisible ones is typically worth it.

Disable Pixel Perfect for ScrollRects

Pixel Perfect is a setting on a `canvas` Component that forces its child UI elements to be drawn with direct alignment to the pixels on the screen. This is often a requirement for art and design, as the UI elements will appear much sharper than if it was disabled. While this alignment behavior is a relatively expensive operation, it is effectively mandatory that it will be enabled for the majority of our UI to keep things crisp and clear. However, for animating and fast-moving objects, it can be somewhat pointless due to the motion involved. Disabling Pixel Perfect for `scrollRect` elements is a good way to make some impressive savings. However, since the Pixel Perfect setting affects the entire Canvas, we should make sure to enable the `scrollRect` element as a child object beneath a separate Canvas so that other elements will maintain their pixel-aligned behavior.



Different kinds of animated UI elements actually look better with Pixel Perfect disabled. Be sure to do some testing, as it can save quite a bit of performance.

Manually stop ScrollRect motion

The Canvas will always need to regenerate the entire `scrollRect` element even if the velocity is moving by a fraction of a pixel each frame. We can manually freeze its motion once we detect that its velocity is below a certain threshold using `scrollRect.velocity` and `scrollRect.StopMovement()`. This can help reduce the frequency of regeneration a great deal.

Use empty `UIText` elements for full-screen interaction

A common implementation in most UIs is to activate a large, transparent interactable element that covers the entire screen, forcing the player to handle a popup before proceeding, while still allowing the player to see what's going on behind it (as a means of not ripping the player out of the game experience entirely). This is often done with a `UIImage` element, but unfortunately this can break batching operations and transparency can be a problem on mobile devices.

A hacky way around this problem is to use a `UIText` element with no Font or Text defined. This creates an element that doesn't need to generate any renderable information and only handles bounding-box checks for interaction.

Check the Unity UI source code

Unity provides the code for its UI system in a bitbucket repository found at <https://bitbucket.org/Unity-Technologies/ui>.

If we're having significant problems with the performance of our UI, it's possible to look into the source code to figure out exactly what might be going on and hopefully discover ways to get around the problem.

A more drastic measure, but a potential option, could be to actually modify the UI code, compile it, and add it to our project manually.

Check the documentation

The tips mentioned previously are some of the more obscure, undocumented, or critical performance optimization tips for the UI System. There are a number of great resources on the Unity website that explain how the UI System works and how best to optimize it, which is far too large to fit in this book verbatim.

Start with the following page and work your way through them for many more helpful UI optimization tips: <https://unity3d.com/learn/tutorials/temas/best-practices/guide-optimizing-unity-ui>.

Shader optimization

Fragment Shaders are the primary consumers of Fill Rate and Memory Bandwidth. The costs depend on their complexity--how much texture sampling takes place, how many mathematical functions are used, and many more factors. The GPU's parallel nature (sharing small pieces of the overall job between hundreds of threads) means that any bottleneck in a thread will limit how many fragments can be pushed through that thread during a frame.

The classic analogy is a vehicle assembly line. A complete vehicle requires multiple stages of manufacture to complete. The critical path to completion might involve stamping, welding, painting, assembly, and inspection, where each step is completed by a single team. For any given vehicle, no stage can begin before the previous one is finished, but whatever team handled the stamping for the last vehicle can begin stamping for the next vehicle as soon as it has finished. This organization allows each team to become masters of their particular domain rather than trying to spread their knowledge too thin, which would likely result in less consistent quality in the batch of vehicles.

We can double the overall output by doubling the number of teams, but if any team gets blocked, then precious time is lost for any given vehicle, as well as all future vehicles that would pass through the same team. If these delays are rare, then they can be negligible in the grand scheme, but if not, and even if one stage takes several minutes longer than normal each and every time it must complete the task, then it can become a bottleneck that threatens the release of the entire batch.

The GPU parallel processors work in a similar way: each processor thread is an assembly line, each processing stage is a team, and each fragment is the thing that needs to be built. If the thread spends a long time processing a single stage, then time is lost on each fragment. This delay will multiply such that all future fragments coming through the same thread will be delayed. This is a bit of an oversimplification, but it often helps to paint a picture of how quickly some poorly optimized Shader code can chew up our Fill Rate and how small improvements in Shader optimization provide big benefits in

Back End performance.

Shader programming and optimization is a very niche area of game development. Their abstract and highly specialized nature requires a very different kind of thinking to generate high-quality Shader code compared to a typical gameplay or Engine code. They often feature mathematical tricks and back-door mechanisms for pulling data into the Shader, such as precomputing values and putting them in texture files. Because of this, and the importance of optimization, Shaders tend to be very difficult to read and reverse-engineer.

Consequently, many developers rely on prewritten Shaders, visual Shader creation tools from the Asset Store, such as Shader Forge, or Amplify Shader Editor. This simplifies the act of initial Shader code generation, but might not result in the most efficient form of Shaders. Whether we're writing our own Shaders, or we're relying on prewritten/pregenerated Shaders, we might find it worthwhile to perform some optimization passes over them using some tried-and-true techniques.

Consider using Shaders intended for mobile platforms

The built-in mobile Shaders in Unity do not have any specific restrictions that force them to only be used on mobile devices. They are simply optimized for minimum resource usage (and tend to feature some of the other optimizations listed in this section).

Desktop applications are perfectly capable of using these Shaders, but they tend to feature a loss of graphical quality. It only becomes a question of whether the loss of graphical quality is acceptable. So, consider doing some testing with the mobile equivalents of common Shaders to check whether they are a good fit for your game.

Use small data types

GPUs can calculate with smaller data types more quickly than larger types (particularly on mobile platforms), so the first tweak we can attempt is replacing our `float` data types (32-bit, floating-point) with smaller versions such as `half` (16-bit, floating-point) or even `fixed` (12-bit fixed point). The size of the data types listed previously will vary depending on what floating-point formats the target platform prefers. The sizes listed are the most common. The optimization comes from the relative size between formats since there are fewer bits to process.

Color values are good candidates for precision reduction, as we can often get away with less precise color values without much noticeable loss in coloration. However, the effects of reducing precision can be very unpredictable for graphical calculations. So, changes such as these can require some testing to verify that the reduced precision is costing too much graphical fidelity.



Note that the effects of these tweaks can vary enormously between one GPU architecture and another (for example, AMD versus Nvidia versus Intel), and even GPU brands from the same manufacturer. In some cases, we can make some decent performance gains for a trivial amount of effort. In other cases, we might see no benefit at all.

Avoid changing precision while swizzling

Swizzling is the Shader programming technique of creating a new vector (an array of values) from an existing `vector` by listing the components in the order in which we wish to copy them into the new structure.

Here are some examples of swizzling:

```
float4 input = float4(1.0, 2.0, 3.0, 4.0); // initial test value (x, y, z, w)

// swizzle two components
float2 val1 = input.yz; // val1 = (2.0, 3.0)

// swizzle three components in a different order
float3 val2 = input.zyx; // val2 = (3.0, 2.0, 1.0)

// swizzle the same component multiple times
float4 val3 = input.yyy; // val3 = (2.0, 2.0, 2.0, 2.0)

// swizzle a scalar multiple times
float sclr = input.w; // sclr = (4.0)
float3 val4 = sclr.xxx; // val4 = (4.0, 4.0, 4.0)
```

We can use both the `xyzw` and `rgba` representations to refer to the same components, sequentially. It does not matter whether it is a `color` or `vector`; they just make the Shader code easier to read. We can also list components in any order we like to fill in the desired data, repeating them if necessary.

Converting from one precision type to another in a Shader can be a costly operation, but converting the precision type while simultaneously swizzling can be particularly painful. If we have mathematical operations that use swizzling, ensure that they don't also convert the precision type. In these cases, it would be wiser to simply use the high-precision data type from the very beginning or reduce precision across the board to avoid the need for changes in precision.

Use GPU-optimized helper functions

The Shader compiler often performs a good job of reducing mathematical calculations down to an optimized version for the GPU, but compiled custom code is unlikely to be as effective as both the Cg library's built-in helper functions and the additional helpers provided by the Unity Cg included files. If we are using Shaders that include custom function code, perhaps we can find an equivalent helper function within the Cg or Unity libraries that can do a better job than our custom code can.

These extra `include` files can be added to our Shader within the `CGPROGRAM` block, as follows:

```
CGPROGRAM
// other includes
#include "UnityCG.cginc"
// Shader code here
ENDCG
```

Example Cg library functions to use are `abs()` for absolute values, `lerp()` for linear interpolation, `mul()` for multiplying matrices, and `step()` for step functionality. Useful `UnityCG.cginc` functions include `WorldSpaceViewDir()` for calculating the direction toward the Camera and `Luminance()` for converting a color to grayscale.

Check out http://http.developer.nvidia.com/CgTutorial/cg_tutorial_appendix_e.html for a full list of Cg standard library functions.

Check out the Unity documentation for a complete and up-to-date list of possible `include` files and their accompanying helper functions at <http://docs.unity3d.com/Manual/SL-BuiltinIncludes.html>.

Disable unnecessary features

Perhaps we can make savings by simply disabling Shader features that aren't vital. Does the Shader really need transparency, Z-writing, alpha-testing, and/or alpha blending? Will tweaking these settings or removing these features give us a good approximation of our desired effect without losing too much graphical fidelity? Making such changes is a good way of making Fill Rate cost savings.

Remove unnecessary input data

Sometimes, the process of writing a Shader involves a lot of back and forth experimentation in editing code and viewing it in the Scene. The typical outcome of this process is that input data that was needed when the Shader was going through early development is now surplus fluff once the desired effect has been obtained, and it's easy to forget what changes were made when/if the process drags on for a long time. However, these redundant data values can cost the GPU valuable time, as they must be fetched from memory even if they are not explicitly used by the Shader. So, we should double-check our Shaders to ensure that all of their input geometry, vertex, and fragment data are actually being used.

Expose only necessary variables

Exposing unnecessary variables from our Shader to the accompanying Material can be costly, as the GPU can't assume these values are constant, which means the compiler cannot compile away these values. This data must be pushed from the CPU with every pass since they can be modified at any time through a Material object's methods such as `setColor()` and `setFloat()`. If we find that, toward the end of the project, we always use the same value for these variables then they should be replaced with a constant in the Shader to remove such excess runtime workload. The only cost is obfuscating what could be critical graphical effect parameters, so this should be done very late in the process.

Reduce mathematical complexity

Complicated mathematics can severely bottleneck the rendering process, so we should do whatever we can to limit the damage. It is entirely possible to store a map of complex mathematical function outputs by precalculating them and placing them as floating-point data in a texture file. A texture file is, after all, just a huge blob of floating-point values that can be indexed quickly with three dimensions: x , y , and color (`rgba`). We can feed this texture into the Shader and sample the pregenerated table in the Shader at runtime instead of completing a complex calculation at runtime.

We may not see any improvement with functions such as `sin()` and `cos()` since they've been heavily optimized to make use of GPU architecture, but complex methods such as `pow()`, `exp()`, `log()`, and our own custom mathematical calculations can only be optimized so much and would be good candidates for simplification. This is assuming that we can easily index the result from the texture with x and y coordinates. If complex calculations are required to generate those coordinates, then it may not be worth the effort.

This technique will cost us additional graphics memory to store the texture at runtime and some Memory Bandwidth, but if the Shader has already been receiving a texture (which they are, in most cases), but the alpha channel is not being used, then we could sneak the data in through the texture's alpha channel, costing us literally no performance since that data has already been passed through anyway. This will involve hand-editing our art assets to include such data in any unused color channel(s), possibly requiring coordination between programmers and artists, but is a very good way of saving Shader processing costs with no runtime sacrifices.

Reduce texture sampling

Texture sampling is at the core of all Memory Bandwidth costs. The fewer textures we use, and the smaller we make them, the better. The more we use, the more cache misses we are likely to invoke, and the larger they are, the more Memory Bandwidth is consumed transferring them into the Texture Cache. Such situations should be simplified as much as possible to avoid severe GPU bottlenecks.

Even worse, sampling textures in a nonsequential order would likely result in some very costly cache misses for the GPU to suffer through. So, if this is being done, then the texture should be reordered so that it can be sampled in a more sequential order. For example, if we're sampling by inverting the x and y coordinates (for example, `tex2D(y, x)` instead of `tex2D(x, y)`), the texture lookup would iterate through the texture vertically, then horizontally, inflicting a cache-miss almost every iteration. A lot of performance could be saved by simply rotating the texture file data and performing a sample in the correct order (`tex2D(x, y)`).

Avoid conditional statements

When conditional statements are run through a modern day CPU, they undergo a lot of clever predictive techniques to make use of *instruction-level parallelism*. This is a feature where the CPU attempts to predict which direction a conditional statement will go in before it has actually been resolved and speculatively begins processing the most likely result of the conditional using any free cores that aren't being used to resolve the conditional (fetching some data from memory, copying some floating-point values into unused registers, and so on). If it turns out that the decision is wrong, then the current result is discarded and the proper path is taken instead. So long as the cost of speculative processing and discarding false results is less than the time spent waiting to decide the correct path, and it is right more often than it is wrong, then this is a net gain for the CPU's speed.

However, this feature is less beneficial for GPU architecture because of its parallel nature. The GPU's cores are typically managed by some higher-level construct that instructs all cores under its command to perform the same machine-code-level instruction simultaneously, such as a huge stamping machine that stamps sheets of metal in groups simultaneously. So, if the Fragment Shader requires a `f1oat` to be multiplied by 2, then the process will begin by having all cores copy data into the appropriate registers in one coordinated step. Only when all cores are finished copying to the registers will the cores be instructed to begin the second step: multiplying all registers by 2 all in a second simultaneous action.

Thus, when this system stumbles into a conditional statement, it cannot resolve the two statements independently. It must determine how many of its child cores will go down each path of the conditional, grab the list of required machine code instructions for one path, resolve them for all cores taking that path, and repeat these steps for each path until all possible paths have been processed. So, for an `if-else` statement (two possibilities), it will tell one group of cores to process the `true` path, then ask the remaining cores to process the `false` path. Unless every core takes the same path, it must process both paths every time.

So, we should avoid branching and conditional statements in our Shader code. Of course, this depends on how essential the conditional is to achieving the graphical effect we desire. However, if the conditional is not dependent on per-pixel behavior, then we would often be better off absorbing the cost of unnecessary mathematics than inflicting a branching cost on the GPU.

Reduce data dependencies

The compiler will try its best to optimize our Shader code into the more GPU-friendly low-level language so that it is not waiting on data to be fetched when it could be processing some other task. For example, the following poorly optimized code could be written in our Shader:

```
float sum = input.color1.r;
sum = sum + input.color2.g;
sum = sum + input.color3.b;
sum = sum + input.color4.a;
float result = calculateSomething(sum);
```

This code has a data dependency such that each calculation cannot begin until the last finishes due to the dependency on the `sum` variable. However, such situations are often detected by the Shader compiler and optimized into a version that uses instruction-level parallelism. The following code is the high-level code equivalent of the resulting Machine Code after the previous code is compiled:

```
float sum1, sum2, sum3, sum4;
sum1 = input.color1.r;
sum2 = input.color2.g;
sum3 = input.color3.b;
sum4 = input.color4.a;
float sum = sum1 + sum2 + sum3 + sum4;
float result = CalculateSomething(sum);
```

In this case, the compiler would recognize that it can fetch the four values from memory in parallel and complete the summation once all four have been fetched independently via thread-level parallelism. This can save a lot of time relative to performing the four fetches one after another.

However, long chains of data dependency that cannot be compiled away can absolutely murder Shader performance. If we create a strong data dependency in our Shader's source code, then it has no freedom to make any optimizations. For example, the following data dependency would be painful on performance, as one step literally cannot be completed without waiting on another to fetch data, since sampling each texture requires sampling another

texture beforehand, and the compiler cannot assume that the data hasn't changed in the meantime.

The following code represents a very strong data dependency between instructions, since each relies on texture data being sampled from the previous instruction:

```
| float4 val1 = tex2D(_tex1, input.texcoord.xy);  
| float4 val2 = tex2D(_tex2, val1.yz); // requires data from _tex1  
| float4 val3 = tex2D(_tex3, val2.zw); // requires data from _tex2
```

Strong data dependencies such as these should be avoided whenever possible.

Surface Shaders

Unity's Surface Shaders are a simplified form of Fragment Shaders, allowing Unity developers to get to grips with Shader programming in a more simplified fashion. The Unity Engine takes care of converting our Surface Shader code for us, abstracting away some of the optimization opportunities we have just covered. However, it does provide some miscellaneous values that can be used as replacements, which reduce accuracy but simplify the mathematics in the resulting code. Surface Shaders are designed to handle the general case fairly efficiently, but optimization is best achieved with a personal touch by writing our own Shaders.

The `approxview` attribute will approximate the view direction, saving costly operations. The `halfasview` attribute will reduce the precision of the view vector, but beware of its effect on mathematical operations involving multiple precision types. The `noforwardadd` attribute will limit the Shader to only considering a single Directional Light, reducing Draw Calls, since the Shader will render in only a single pass, and Lighting complexity. Finally, the `noambient` attribute will disable ambient Lighting in the Shader, removing some extra mathematical operations that we may not need.

Use Shader-based LOD

We can force Unity to render distant objects using simpler Shaders, which can be an effective way of saving Fill Rate, particularly if we're deploying our game onto multiple platforms or supporting a wide range of hardware capability. The `LOD` keyword can be used in the Shader to set the onscreen size factor that the Shader supports. If the current LOD level does not match this value, it will drop to the next fallback Shader and so on until it finds the Shader that supports the given size factor. We can also change a given Shader object's LOD value at runtime using the `maximumLOD` property.

This feature is similar to the mesh-based LOD covered earlier and uses the same LOD values for determining object form factor, so it should be configured as such.

Check out <https://docs.unity3d.com/Manual/SL-ShaderLOD.html> in the Unity documentation for more information on Shader-based LOD.

Use less texture data

This approach is simple, straightforward, and always a good idea to consider. Reducing texture quality, either through resolution or bit rate, is not ideal for graphical quality, but we can sometimes get away with using 16-bit textures without any noticeable degradation.

Mip Maps (explored in [chapter 4](#), *Kickstart Your Art*) are another excellent way of reducing the amount of texture data being pushed back and forth between VRAM and the Texture Cache. Note that the Scene window has a Mipmaps Shading Mode, which will highlight textures in our Scene blue or red, depending on whether the current texture scale is appropriate for the current Scene window's Camera position and orientation. This will help identify what textures are good candidates for further optimization.

Test different GPU Texture Compression formats

As you learned in [Chapter 4](#), *Kickstart Your Art*, there are different texture compression formats, which can reduce our application's disk footprint (executable file size), runtime CPU, and RAM usage. These compression formats are designed to support GPU architecture for the given platform. There are many different formats, such as DXT, PVRTC, ETC, and ASTC, but only a handful of these are available on a given platform.

By default, Unity will pick the best compression format determined by the Compression setting for a texture file. If we drill down into platform-specific options for a given texture file, then different compression type options will be available, listing the different texture formats the given platform supports. We may be able to find some space or performance savings by overriding the default choices for compression.

Although, beware that if we're at the point where individually tweaking Texture Compression techniques is necessary, then hopefully we have already exhausted all other options for reducing Memory Bandwidth. By going down this road, we would be committing ourselves to supporting many different devices each in their own specific way. Many developers would prefer to keep things simple with a general solution instead of personal customization and time-consuming handiwork for small performance gains.

Check out the Unity documentation for an overview of all of the different texture formats available and which formats Unity prefers by default at <https://docs.unity3d.com/Manual/class-TextureImporterOverride.html>.

In older versions of Unity, all formats were exposed for Advanced texture types, but if the platform did not support the given type, it would be handled at the software level. In other words, the CPU would need to stop and recompress the texture to the desired format the GPU wants, as opposed to the GPU



taking care of it with a specialized hardware chip. Unity Technologies decided to remove this capability in more recent versions so that we can't accidentally cause these problems.

Minimize texture swapping

This one is fairly straightforward. If Memory Bandwidth is a problem, then we need to reduce the amount of texture sampling we're doing. There aren't really any special tricks to exploit here since Memory Bandwidth is all about throughput, so the primary metric under consideration is the volume of data we're pushing.

One way to reduce volume is to simply lower texture resolution and hence quality. This is obviously not ideal, so another approach is to find clever ways to reuse textures on different meshes, but using different Material and Shader properties. For instance, a properly darkened brick texture may appear to look like a stone wall instead. Of course, this will require different Render States and hence we won't save on Draw Calls, but it could reduce Memory Bandwidth consumption.



Did you ever notice how clouds and bushes looked exactly the same in Super Mario Bros but with different colours? This is the same concept.

There could also be ways to combine textures into Atlases to reduce the number of swaps needed. If there are a group of textures that are always used together at similar times, then they could potentially be merged together. This could save the GPU from having to pull in separate texture files over and over again during the same frame.

Finally, removing textures from the application entirely is always a last resort option we could employ.

VRAM limits

One last consideration related to textures is how much VRAM we have available. Most texture transfer from CPU to GPU occurs during initialization, but can also occur when a nonexistent texture is first required by the current view. This process is normally asynchronous and will result in a blank texture being used until the full texture is ready for rendering (refer to [Chapter 4](#), *Kickstart Your Art*, to note that this assumes read/write access is disabled for the texture). As such, we should avoid introducing new textures at runtime too frequently.

Preload textures with hidden GameObjects

The blank texture that is used during asynchronous texture loading can be jarring when it comes to game quality. We would like a way to control and force the texture to be loaded from disk to RAM and then to VRAM before it is actually needed.

A common workaround is to create a hidden `GameObject` that uses the texture and place it somewhere in the Scene on the route that the player will take toward the area where it is actually needed. As soon as the player looks at that object, the texture is needed by the Rendering Pipeline (even if it's technically hidden), it will begin the process of copying the data from RAM to VRAM. This is a little clunky, but easy to implement and works sufficiently well in most cases.

We can also control such behavior via Script code by changing a Material's `texture` property:

```
GetComponent<Renderer>().material.texture = textureToPreload;
```

Avoid texture thrashing

In the rare event that too much texture data is loaded into VRAM, and the required texture is not present, the GPU will need to request it from RAM and overwrite one or more existing textures to make room for it. This is likely to worsen over time as the memory becomes fragmented, and it introduces a risk that the texture just flushed from VRAM needs to be pulled again within the same frame. This will result in a serious case of memory *thrashing* and should be avoided at all costs.

This is less of a concern on modern consoles such as the PS4, Xbox One, and WiiU since they share a common memory space for both CPU and GPU. This design is a hardware-level optimization, given the fact that the device is always running a single application, and almost always rendering 3D graphics. However, most other platforms must share time and space with multiple applications, where a GPU is merely an optional device and is not always present. They, therefore, feature separate memory spaces for the CPU and GPU, and we must ensure that the total texture usage at any given moment remains below the available VRAM of the target hardware.

Note that this *thrashing* is not precisely the same as hard disk thrashing, where memory is copied back and forth between main memory and virtual memory (the swap file), but it is analogous. In either case, data is being unnecessarily copied back and forth between two regions of memory because too much data is being requested in too short a time period for the smaller of the two memory regions to hold it all.



Thrashing such as this can be a common cause of dreadful rendering performance when games are ported from modern consoles to desktop platforms and should be treated with care.

Avoiding this behavior may require customizing texture quality and file sizes on a per-platform and per-device basis. Be warned that some players are likely to notice these inconsistencies if we're dealing with hardware from the same console or desktop GPU generation. As many of us will know, even

small differences in hardware can lead to a lot of apples-versus-oranges comparisons, but hardcore gamers tend to expect a similar level of quality across the board.

Lighting optimization

We covered the theory of Lighting behavior earlier in this chapter, so let's run through some techniques we can use to improve Lighting costs.

Use real-time Shadows responsibly

As mentioned previously, Shadowing can easily become one of the largest consumers of Draw Calls and Fill Rate, so we should spend the time to tweak these settings until we get the performance and/or graphical quality we need. There are multiple important settings for Shadowing that can be found under Edit | Project Settings | Quality | Shadows. As far as the Shadows option is concerned, Soft Shadows are expensive, Hard Shadows are cheap, and No Shadows are free. Shadow Resolution, Shadow Projection, Shadow Distance, and Shadow Cascades are also important settings, which affect the performance of our Shadows.

Shadow Distance is a global multiplier for runtime Shadow rendering. There is little point in rendering Shadows at a great distance from the Camera, so this setting should be configured specific to our game and how much Shadowing we expect to witness during gameplay. It is also a common setting that is exposed to the user in an options screen, so they can choose how far to render Shadows to get the game's performance to match their hardware (at least on desktop machines).

Higher values of Shadow Resolution and Shadow Cascades will increase our Memory Bandwidth and Fill Rate consumption. Both of these settings can help curb the effects of artifacts generated by Shadow rendering, but at the cost of much larger Shadowmap texture sizes, costing increased Memory Bandwidth and VRAM.



The Unity documentation contains an excellent summary on the topic of the aliasing effect of Shadowmaps and how the Shadow Cascades feature helps to solve the problem at <http://docs.unity3d.com/Manual/DirLightShadows.html>.

It's worth noting that Soft Shadows do not consume any more memory or CPU overhead relative to Hard Shadows, as the only difference is a more

complex Shader. This means that applications with enough Fill Rate to spare can enjoy the improved graphical fidelity of Soft Shadows.

Use Culling Masks

A `Light` Component's Culling Mask property is a Layer-based mask that can be used to limit the objects that will be affected by the given Light. This is an effective way of reducing Lighting overhead, assuming that the Layer interactions also make sense with how we are using Layers for physics optimization. Objects can only be a part of a single Layer, and reducing physics overhead probably trumps Lighting overhead in most cases; thus, if there is a conflict, then this may not be the ideal approach.



Note that there is limited support for Culling Masks when using Deferred Shading. Due to the way it treats Lighting in a very global fashion, only four Layers can be disabled from the mask, limiting our ability to optimize its behavior.

Use baked Lightmaps

Baking Lighting and Shadowing into a Scene is significantly less processor-intensive than generating them at runtime. The downside is the added application disk footprint, memory consumption, and potential for Memory Bandwidth abuse. Ultimately, unless a game's Lighting effects are being handled exclusively through the legacy Vertex Lit Shading format or through a single `DirectionalLight`, then it should probably include Lightmapping somewhere to make some huge budget savings on Lighting calculations. Relying entirely on real-time Lighting and Shadows is a recipe for disaster due to the performance costs they are likely to inflict.

There are several metrics that can affect the cost of Lightmapping, however, such as their resolution, compression, whether we are using Pre-computed Realtime GI, and of course, the number of objects in our Scene. The Lightmapper generates textures that span all of the objects marked Lightmap Static in the Scene, and hence the more we have, the more texture data must be generated for them. This would be an opportunity to make use of additive or subtractive Scene loading to minimize how many objects need to be processed each frame. This, of course, pulls in even more Lightmap data while more than one Scene is loaded, so we should expect a big bump in memory consumption each time this happens, only to have it freed once the old Scene is unloaded.

Optimizing rendering performance for mobile devices

Unity's ability to deploy to mobile devices has contributed greatly to its popularity among hobbyist, small, and mid-size development teams. As such, it would be prudent to cover some approaches that are more beneficial for mobile platforms than for desktop and other devices.



Note that any, or all, of the following approaches may eventually become obsolete, at least for newer devices. The capabilities of mobile devices have advanced blazingly fast, and the following techniques as they apply to mobile devices merely reflect conventional wisdom from the last half decade or so. We should test the assumptions behind these approaches to check whether the limitations of mobile devices still fit the mobile marketplace.

Avoid Alpha Testing

Mobile GPUs haven't quite reached the same levels of chip optimization as desktop GPUs, and Alpha Testing remains a particularly costly task on mobile devices. In most cases, it should simply be avoided in favor of Alpha Blending.

Minimize Draw Calls

Mobile applications are more often bottlenecked on Draw Calls than on Fill Rate. Not that Fill Rate concerns should be ignored (nothing should, ever!), but this makes it almost necessary for any mobile application of reasonable quality to implement mesh combining, Batching, and Atlasing techniques from the very beginning. Deferred Rendering is also the preferred technique, as it fits well with other mobile-specific concerns, such as avoiding transparency and having too many animated characters, but of course not all mobile devices and Graphics APIs support it.

Check out the Unity documentation for more information on which platforms/APIs support Deferred Shading at <https://docs.unity3d.com/Manual/RenderingPaths.html>.

Minimize Material count

This concern goes hand in hand with the concepts of Batching and Atlasing. The fewer Materials we use, the fewer Draw Calls required. This strategy will also help with concerns relating to VRAM and Memory Bandwidth, which tend to be very limited on mobile devices.

Minimize texture size

Most mobile devices feature a very small Texture Cache relative to desktop GPUs. There are very few devices on the market still supporting OpenGL ES 1.1 or lower, such as the iPhone 3G, but these devices could only support a maximum texture size of 1024 x 1024. Devices supporting OpenGL ES 2.0, such as everything from the iPhone 3GS to the iPhone 6S, can support textures up to 2048 x 2048. Finally, devices supporting OpenGL ES 3.0 or greater, such as devices running iOS 7, can support textures up to 4096 x 4096.



There are way too many Android devices to list here, but the Android developer portal gives a handy breakdown of OpenGL device support. This information is updated regularly to help developers determine supported APIs in the Android market at <https://developer.android.com/about/dashboards/index.html>

Double-check the device hardware we are targeting to be sure that it supports the texture file sizes we wish to use.^[P] However, later-generation devices are never the most common devices in the mobile marketplace. If we wish our game to reach a wide audience (increasing its chances of success), then we must be willing to support weaker hardware.

Note that textures that are too large for the GPU will be downscaled by the CPU during initialization. This wastes valuable loading time and is going to leave us with unintended loss of quality due to an uncontrolled reduction in resolution. This makes texture reuse of paramount importance for mobile devices due to the limited VRAM and Texture Cache sizes available.

Make textures square and power-of-two

We have already covered this topic in [Chapter 4, Kickstart Your Art](#), but it is worth revisiting the subject of GPU-level Texture Compression. The GPU will find it difficult, or simply be unable to compress the texture if it is not in a square format, so make sure that you stick to the common development convention and keep things square and sized to a power-of-two.

Use the lowest possible precision formats in Shaders

Mobile GPUs are particularly sensitive to precision formats in its Shaders, so the smallest formats should be used such as `half`. On a related note, precision format conversion should be avoided at all costs for the same reason.

Summary

If you've made it this far without skipping ahead, then congratulations are in order. That was a lot of information to absorb for just one subsystem of the Unity Engine, but then it is clearly the most complicated of them all, requiring a matching depth of explanation. Hopefully, you've learned a lot of approaches to help you improve your rendering performance and enough about the Rendering Pipeline to know how to use them responsibly.

By now, we should be used to the idea that, with the exception of algorithm improvements, every performance enhancement we implement will come with some related cost that we must be willing to bear for the sake of removing one bottleneck. We should always be ready to implement multiple techniques until we've squashed them all, and potentially spend a lot of additional development time to implement and test some performance enhancing features.

In the next chapter, let's bring performance optimization into the modern era by exploring some performance improvements we can apply to Virtual Reality and Augmented Reality projects.

Virtual Velocity and Augmented Acceleration

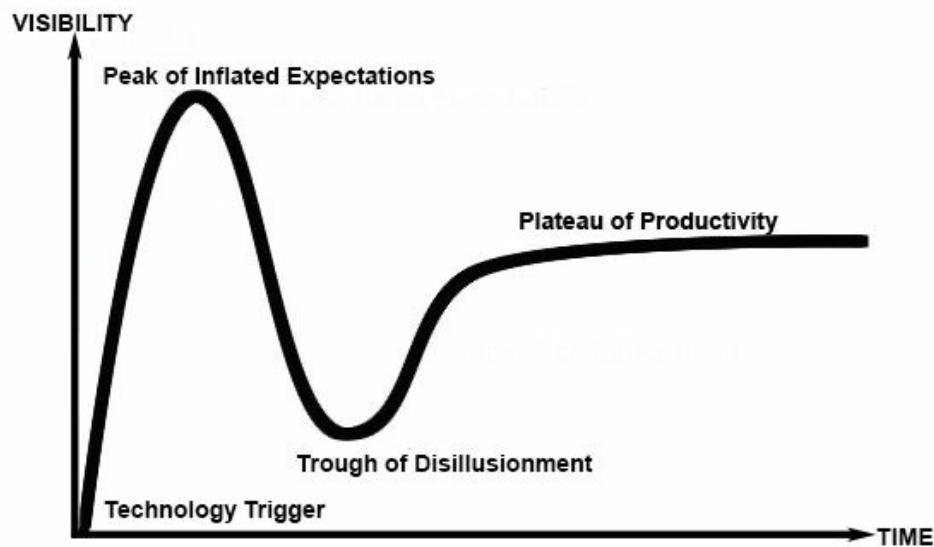
Two whole new entertainment mediums have entered the world stage in the forms of **Virtual Reality (VR)** where users are transported into a virtual space through the use of a **Head Mounted Device (HMD)**, and **Augmented Reality (AR)** where virtual elements are superimposed on top of a display showing the real world. For the sake of brevity, these two terms are often combined into the singular term **XR**. There is also the **Mixed Reality (MR)** (also known as **Hybrid Reality** or **HR**), where an application mixes the real and virtual worlds together, which encompasses all of previously mentioned formats, while also including **Augmented Virtuality** where real world objects are scanned and superimposed inside a mostly virtual world.

The markets for these media formats have sprung up very fast and are continuing to grow rapidly, with huge investments from the technology industry's biggest players. Naturally, Game Engines such as Unity jumped on the bandwagon quickly, providing ample support for most of the top contending platforms, such as Google's Cardboard, HTC's Vive, Oculus' Rift, Microsoft's HoloLens and Samsung's GearVR platforms, as well as the more recent entries such as Apple's ARKit, Google's AR Core, Microsoft's Windows Mixed Reality platform, PTC's (originally Qualcomm's) Vuforia, and Sony's PlayStation VR.

XR offers a whole new realm for developers and creatives to explore. This includes entertainment products such as games and 360-degree Videos (or *360 Video* for short), where a series of cameras are bundled together, each facing a different direction--the various captures from those cameras are stitched together and later played back as a movie in a VR headset with visibility in all directions. Creative industry tools are also common in XR, such as 3D modeling software, workflow visualizations, and quality-of-life gadgets. There are very few rules that have been set in stone, so there are plenty of opportunities to innovate, contribute to this new wave of technology and become the one to create those rules. This has led to a lot of buzz and

excitement as people explore what is possible and try to make their mark on the future of entertainment and interactive experiences.

Of course, almost every new and budding technology goes through the Hype Cycle (from the Gartner Hype Cycle available at <https://www.gartner.com/technology/research/methodologies/hype-cycle.jsp>). The Hype Cycle starts with its honeymoon period of excessive hype, and where early adopters will evangelize its benefits. Later, there is an eventual cooling of emotions as it enters the *trough of disillusionment* since it hasn't quite hit the mainstream, and its benefits are not taking hold just yet. This continues until it either the technology fails to capture hearts and minds, thus falling out of existence, or takes a firm hold and continues steady adoption. The following image shows the essentials of the Gartner Hype Cycle:



Arguably, XR has recently been passing through this final phase and is starting to enjoy much better support and much higher quality experiences than it did during the early days, although it is true that the adoption rate of XR has been slower than initially predicted. It remains to be seen whether XR will grow into a multi-billion dollar industry or fade into a niche market of gadgets. Consequently, developing within this new medium is not without its risks, and we can find industry analysts who will agree with our opinions, regardless of where we stand on the future of XR. One thing is for certain, though, that every time someone experiences firsthand what VR and AR are capable of, they're blown away by the level of immersion and the medium's

ability to convincingly transport them into another world. This level of immersion and interactivity is unparalleled, and teases many more possibilities as support for the platforms mature and technology continues to advance.

In this chapter, we will explore the following topics:

- What concerns to keep in mind when developing VR or AR projects in Unity and what must be avoided
- Performance enhancements specific to the XR medium

XR Development

Developing an XR product in Unity involves importing one of several **XR Software Development Kits (SDKs)** into our Unity project and making some specialized API calls to configure and use the platform at runtime. Each SDK is different in its own way and offers a different set of features. For instance, the Oculus Rift and HTC Vive SDKs provide APIs to control VR HMDs and their respective Controllers, whereas Apple's ARKit provides utilities to determine spatial positioning and superimpose objects on the display. Unity technologies have been working hard to create APIs that support all of these variations, so the APIs for XR development in Unity have changed a lot over the past couple of years.

The early days of Unity VR development meant dropping Native Plugins into our Unity projects, importing SDKs directly from an external developer portal, and all kinds of annoying grunt-work in setup, and applying updates manually. Since then, however, Unity has incorporated several of these SDKs directly into the Editor. In addition, since AR has become more popular recently, the main API has been renamed from `UnityEngine.VR` to `UnityEngine.XR` in Unity 2017.2.0 and greater, in order to be inclusive of several AR SDKs. Importing various XR SDKs and configuring them can be managed through the `Edit | Project Settings | Project | XR Settings` area.

The development experience of working on XR products is a bit of a mixed bag right now. It involves working on some top-of-the-line hardware and software, which means that there are constant changes, redesigns, breakages, patches, bugs, crashes, compatibility issues, performance problems, rendering artifacts, lack of feature parity between platforms, and so on. All of these problems serve to slow down our progress, which makes gaining a competitive advantage extraordinarily difficult in the XR field. On the bright side, pretty much everyone is having the same issues so they get a lot of attention from their developers, making them easier to develop with all the time. Lessons are learned; APIs are cleaned up; and new features, tools, and optimizations are made available with every passing update.

Performance problems limit an XR product's success perhaps more so than non-XR projects due to the current state of the medium. Firstly, our users will be spending significant money to purchase VR HMDs and sensor equipment, or AR-capable hardware. Both of these platforms can be very resource intensive, requiring similarly expensive graphics hardware to support them. This typically leads users to expect a much higher level of quality compared to typical games, in order to make the investment feel worthwhile, or to put it another way, this makes poor user experiences understandably less forgivable due to the monetary investment required by the user. Secondly, perhaps more so for VR projects than AR ones, poor application performance can lead to serious physical user discomfort, quickly turning even the staunchest advocate into a detractor. Thirdly, the XR platform's primary draw is its immersiveness, and nothing breaks that faster than frame drops, flickering, or any kind of application breakdown that forces the user to remove their headset or reboot the app.

Ultimately, we must be prepared to profile our XR applications early and to make sure we aren't exceeding our runtime budget, as it will be stretched thin by the complex and resource-intensive nature of the technology behind these mediums.

Emulation

Unity offers several emulation options, particularly for AR app development with HoloLens, and ARKit, called **Holographic Remoting / Simulation** or ARKitRemote. We should be careful and use these emulation features as they're intended to be used--as convenience tools for development workflow to check that the essentials are working correctly. They are no substitute for the real thing when it comes to benchmarking and should not be trusted to give accurate profiling results. This feature can be reached through the Window | Holographic Emulation window.

User comfort

Unlike typical games and apps, VR apps need to consider user comfort as a metric to optimize. Dizziness, motion sickness, eye strain, headaches, and even physical injuries from loss of balance have unfortunately been all too common for early VR adopters, and the onus is on us to limit these negative effects for users. In essence, content is just as important to user comfort as the hardware is, and we need to take the matter seriously if we are building for the medium.

Not everyone experiences these issues, and there are a lucky few who have experienced none of them. However, the overwhelming majority of users have reported these problems at one point or another. Also, just because our game doesn't trigger these problems in ourselves when we're testing them doesn't mean they won't trigger them in someone else. In fact, we will be the most biased test subject for our game due to familiarity. Without realizing it, we might start to predict our way around the most nauseating behavior our app generates, making it an unfair test compared to a new user experiencing the same situation. This unfortunately raises the costs of VR app development further, as a lot of testing with different unbiased individuals is required if we want to figure out whether our experience will cause discomfort, which may be needed each time we make significant changes that affect motion and frame rate.



There are a number of things that users can do to improve their VR comfort, such as starting with small sessions and working their way up to get practice in balancing and training their brain to expect the mismatched motion. A more drastic option is to take motion sickness medication or drink a little ginger tea beforehand to settle the stomach. However, we will hardly convince users to try our app if we promise it'll only take a few sessions of motion sickness before it starts to get enjoyable.

There are three main discomfort effects users can experience in VR:

- Motion sickness
- Eye strain
- Disorientation

The first problem, nausea, caused by motion sickness, typically happens when there is a sensory disconnect between where the user's eyes think the horizon is and what their other senses are telling their brain, such as the inner ear's sense of balance. The second problem, eye strain, comes from the fact that the user is staring at a screen mere inches from their eyes, which tends to lead to a lot of eye strain and ultimately headaches after prolonged use.

Finally, disorientation typically occurs because a user in VR is sometimes standing within a confined space, so if a game features any kind of acceleration-based motion, the user will instinctively try to offset that acceleration by adjusting their balance, which can lead to disorientation, falling over, and the user hurting themselves if we are not careful about ensuring the user experiences smooth and predictable motion.



Note that the term acceleration is used intentionally since it is a vector, which means it has both magnitude and direction. Any kind of acceleration can cause disorientation, which includes not only accelerating forward, backward, and sideways, but also acceleration in a rotational fashion (turning around), falling, jumping, and so on.

Another potential problem for VR apps is the possibility of invoking seizures. VR is in the unique position of being able to blast images into the user's eyes at a close range, which opens up some risks that we unintentionally cause vulnerable users to seize if rendering behavior breaks down and starts flickering. These are all things we need to keep in mind during development, that need to be tested for and fixed sooner rather than later.

Perhaps, the most important performance metric to reach in a VR app is having a high value for FPS, preferably, 90 FPS or more, as this will generate a smooth viewing experience since there will be very little disconnection between the user's head motion, and the motion of the world. Any period of extended frame drops or having an FPS value consistently below this value are likely to cause a lot of problems for our users, making it critical that our

application performs well at all times. In addition, we should be very careful about how we control the user's viewpoint. We should avoid changing an HMD's field of view ourselves (let the user dictate the direction it is facing), generating acceleration over long periods, or causing uncontrolled world rotation and horizon motion since these are extremely likely to trigger motion sickness and balance problems for the user.

A strict rule, that is not up for debate, is that we should never apply any kind of gain, multiplier effect, or acceleration effect to the positional tracking of an HMD in the final build of our product. Doing so for the sake of testing is fine, but if a real user moves their head two inches to the side, then it should feel like it moved the same relative distance inside the application and should stop the moment their head stops. Doing otherwise is not only going to cause a disconnect between where the player's head feels like it should be and where it is, but may cause some serious discomfort if the camera becomes offset with respect to the player's orientation and the hinge of their neck.

It is possible to use acceleration for the motion of the player character, but it should be incredibly short and rapid before the user starts to self-adjust too quickly. It would be wisest to stick to motion that relies on constant velocities and/or teleportation.



Placing banked turns in racing games seems to improve user comfort a great deal since the user naturally tilts their head and adjusts their balance to match the turn.

All of the above rules apply just as well to 360 Video content as they do to VR games. Frankly, there have been an embarrassing number of 360 Videos released to market, which are not taking the above into account--they feature too many jerking movements, lack of camera stabilization, manual viewport rotation, and so on. These hacks are often used to ensure the user is facing in the direction we intend. However, more effort must be spent to do so without hacking in such nausea-inducing behaviour. Humans are naturally very curious about things that move. If they notice something moving in the corner of their eye, then they will most likely turn to face it. This can be used to great effect to keep the user facing in the direction we intend as they watch the video.



Laziness is not the way to go when generating VR content. Don't just slap a 360 Camera on top of a dirt rally car and hack an unexpected camera rotation into the video to keep the action in the center. Motion needs to be smooth and predictable. During production, we need to constantly keep in mind where we expect the user to be looking so that we capture action shots correctly.

Fortunately, for the 360 Video format it seems as though industry standard frame rates such as 24FPS or 29.97FPS do not have a disastrous effect on user comfort, but note that this frame rate applies to video playback only. Our rendering FPS is a separate FPS value, and dictates how smooth positional head tracking will be. The rendering FPS must always be very high to avoid discomfort (ideally 90 FPS).

Other problems arise with building VR apps--different HMDs and controllers support different inputs and behavior, making feature-parity across VR platforms difficult. A problem called **Stereo-Fighting** can occur if we try to merge 2D and 3D content together, where 2D objects appear to be rendering deep inside 3D objects since the eyes can't distinguish the distance correctly; this is typically a big problem for the User Interface of VR applications and 360 Video playback, which tends to be a series of flat panels superimposed over a 3D background. Stereo-Fighting does not usually lead to nausea, but it can cause additional eye strain.

Although the effects of discomfort are not quite as pronounced in the AR platform, it's still important not to ignore it. Since AR apps tend to consume a lot of resources, low frame rate applications can cause some discomfort. This is especially true if an AR app makes use of superimposing objects onto a camera image (which is the majority of them), where there will probably be a disconnect in the frame rate between the background camera image and the objects we're superimposing over it. We should try to synchronize these frame rates to limit that disconnect.

Performance enhancements

That's enough talk about the industry and XR development. Let's cover some performance enhancements that can be applied to XR projects.

The kitchen sink

Since AR and VR apps are built using the same Engine, the same subsystems, assets, tools, and utilities as any other Unity game, literally every other performance enhancement mentioned in this book can help VR and AR apps in some fashion, and we should try them all before getting too in-depth with XR-specific enhancements. This is reassuring, as there are a lot of potential performance enhancements we could apply. The downside is that we may need to apply many of them to reach the level of performance we need for our app.

The biggest threat to a VR app's performance is GPU Fill Rate, which is already one of the more likely bottlenecks in any other game, but significantly more so for VR since we will always be trying to render a high resolution image to a much larger Frame Buffer (since we're effectively rendering the Scene twice--once for each eye). AR apps are typically going to find extreme consumption in both the CPU and the GPU since AR platforms make heavy use of the GPU's parallel pipeline to resolve spatial locality of objects and perform tasks such as image recognition, as well as needing a lot of Draw Calls to support those activities.

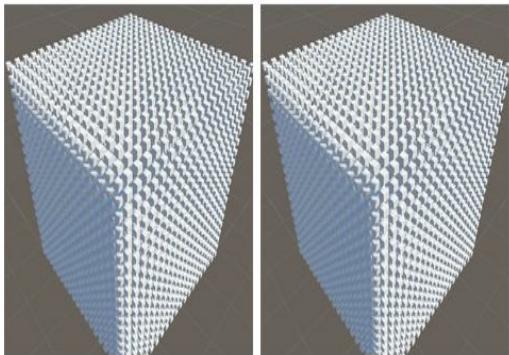
Of course, certain performance enhancing techniques are not going to be particularly effective in XR. Occlusion Culling in a VR app may be difficult to set up since the user can look under, around, and sometimes through objects in the Scene (although it can still be enormously beneficial). Meanwhile, AR apps normally render objects at reachable distances, LOD enhancements may be fairly pointless to setup. We must use our better judgement to determine whether a performance optimization technique is worth implementing *before* we start implementing it since many of them take a lot of time to implement and support.

Single-Pass versus Multi-Pass Stereo Rendering

For VR apps, Unity provides two rendering modes; Multi-Pass and Single-Pass. This can be configured under Edit | Project Settings | Player | XR Settings | Stereo Rendering Method. Multi-Pass rendering will render the Scene to two different images, which are displayed separately for each eye. Single-Pass Stereo Rendering combines both images into a single double-width Render Texture, where only the relevant half is displayed to each eye.

Multi-Pass Stereo Rendering is the default case. The advantage of Single-Pass rendering is that it provides significant savings in the CPU work in the main thread by reducing Draw Call setup and in the GPU since less texture swapping needs to occur. Of course, the GPU will need to work just as hard to render the objects since each object is still rendered twice from two different perspectives (there are no freebies here). The disadvantage is that this effect is currently only usable when using OpenGL ES3.0 or higher and hence is not available on all platforms. In addition, its effects on the Rendering Pipeline require extra care and effort, particularly surrounding any Shaders making use of screen-space effects (effects which only use data already drawn to the Frame Buffer). With Single-Pass Stereo Rendering enabled, Shader code can no longer make the same assumptions about the incoming screen space information. The following image shows how screen space coordinates vary between Multi-Pass and Single-Pass Stereo Rendering

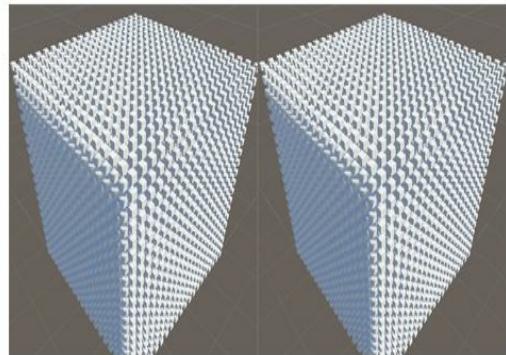
Multi-Pass Stereo Rendering



x 0.0 0.5 1.0 0.0 0.5 1.0

Two Frame Buffers. Unique screenspace coordinates

Single-Pass Stereo Rendering



x 0.0 0.5 1.0

Shared Frame Buffer. Shared screenspace coordinates

The Shader is always informed of the screen space coordinates relative to the entire output Render Texture rather than just the portion it is interested in. For example, we would normally expect an x value of 0.5 to correspond to the horizontal half-way point of the screen, which would be the case when we use Multi-Pass Stereo Rendering, however, if we use Single-Pass Stereo Rendering, then an x value of 0.5 would correspond to the half-way point between the rendering of both eyes (the right-edge of the left eye, or the left-edge of the right eye).



Unity provides some useful helper methods for screen space conversion for Shaders, which can be found at <https://docs.unity3d.com/Manual/SinglePassStereoRendering.html>.

Another problem to worry about is Post-Processing effects. We essentially always pay double the cost for any Post-Processing effect applied to the Scene in VR since it needs to be evaluated once for each eye. Single-Pass Stereo Rendering can reduce the Draw Calls needed to set up our effect, but we can't blindly apply a Post-Processing effect to both images simultaneously. Consequently, Post-Processing effect Shaders must also be

adjusted to ensure that they render to the correct half of the output Render Texture. Without doing this, a Post-Processing effect will be stretched over both eyes, twice, which might look incredibly bizarre for effects such as lens flares.

The Single-Pass Stereo Rendering feature has been around since Unity 5.4, but is still currently labeled Preview in the Editor due to the current lack of support for all platforms. We can expect it to be rolled out to more platforms eventually, but for platforms supporting it we will need to perform some profiling and sensible sanity checks on our screen space Shaders to ensure that we are making positive gains from enabling this option.

Apply anti-aliasing

This tip is less of a performance enhancement and more of a requirement. Anti-aliasing significantly improves the fidelity of XR projects since objects will blend together better and appear less pixelated, improving immersion, which can cost a lot of Fill Rate. We should enable this feature early and try to reach our performance goals with the assumption that it is simply always there, only disabling it as an absolute last resort.

Prefer Forward Rendering

The advantage of Deferred Rendering is the ability to resolve many Light sources with minimal Draw Calls. Unfortunately, if we follow the preceding advice and apply anti-aliasing effects, this must be done as a Post-Processing screen space Shader when Deferred Rendering is used. This can cost a considerable amount of performance compared to how the same technique is applied as a multi-sampling effect in Forward Rendering, potentially making Forward Rendering the more performant between the two options.

Image effects in VR

The effects applied by Normal Maps tend to break down easily in VR, where the texture appears *painted on* to the surface instead of giving the illusion of depth. Normal Maps normally break down very quickly with viewing angles that are very oblique (shallow) with the surface, which is not particularly common in a typical game. However, in VR, since most HMDs allow users to move their heads around in a 3D space via positional tracking (which, granted, not all of them do), they will quickly find positions that break the effect for any objects close to the Camera. Normal Maps have been known to improve the quality of high polygon count objects in VR, but it rarely provides a benefit for those with a low polygon count, so we should perform a little testing to make sure that any visual improvement is worth the costs in Memory Bandwidth.

Ultimately, we cannot rely on Normal Mapping to provide a quick and cheap increase in graphical fidelity for low polygon count objects that we might expect from a non-VR Scene, so testing is required to find out whether the illusion is working as intended. Displacement maps, Tessellation, and/or Parallax Mapping should be used instead to create a more believable appearance of depth. Unfortunately, all of these techniques are more expensive than a typical Normal Map, but it is a burden we must suffer to achieve good graphical quality in VR.

Other Post-Processing effects such as Depth of Field, blurring, and lens flares are effects that look good in a typical 3D game, but are generally not effects we witness in the real world and will seem out of place in VR (at least until eye-tracking support is available) and should generally be avoided.

Backface culling

Backface culling (removing faces from objects that will never be visible) can be tricky for VR and AR projects since the player's viewing angle could potentially come from any direction for objects near the camera. Assets near the camera should be a fully closed shape if we want to avoid immersion-breaking viewpoints. Also, we should think carefully about applying backface culling for distant objects, particularly if the user travels by teleportation since it can be tricky to completely restrict a user's location. Ensure that you test your game world's bounding volumes to ensure that the user cannot escape.

Spatialized audio

The audio industry is abuzz with new techniques to present audio experiences for VR (or more technically, old techniques that have finally found good use) in the form of spatial audio. Audio data for these formats no longer represents audio data from specific channels, but instead contain data for certain audio harmonics that are merged together at runtime to create a more believable audio experience depending on the current camera viewport, particularly vertical orientations. The keyword from the previous sentence being *runtime*, meaning that this effect has a continuous nontrivial cost associated with it. These techniques will require CPU activity, but may also use GPU acceleration to generate their effects, so we should double-check the behavior of both devices if we're experiencing performance problems when we're making use of spatial audio.

Avoid camera physics collisions

In VR and AR, it's possible for the user to move the camera through objects, which can lead to breaking immersion. Although it may be tempting to add physics Colliders to such surfaces in order to prevent the camera moving through it, this will cause disorientation in VR since the Camera will not move in unison with the user's movements. This could also break the positional tracking calibration of an AR app. A better approach is to either allow the user to see into objects or to maintain a safe buffer zone between the Camera and such surfaces. If we don't allow the player to teleport too close to them in the first place, then there's no risk of sticking their head through walls.

This will save on performance due to a reduced number of Colliders, but should be followed as more of a quality-of-life issue. We shouldn't be too concerned about risking immersion-breaking behavior by doing this, as research has shown users tend to avoid looking into objects once they realize they can do it. They may experience a moment of confusion or hilarity when it happens initially, but fortunately people tend to want to remain in the immersive experience we've created and will tend to avoid putting their head through walls. Although, this changes if doing so awards some kind of gameplay advantage, such as seeing through a wall in order to observe which way enemies are about to come from in a strategy game, so we should develop our Scenes with that in mind.

Avoid Euler angles

Avoid using Euler angles for any kind of orientation behavior. Quaternions are designed to be much better for representing angles (the only downside is that it is more abstract and harder to visualize when debugging) and maintaining accuracy whenever there are changes, while also avoiding the dreaded *Gimbal Lock*. Using Euler angles for calculations could eventually lead to inaccuracies after there are a lot of rotation changes, which is incredibly likely since the user's viewpoint will change by tiny amounts many times per second in both VR and AR.



Gimbal Lock is a problem that can occur with Euler angles. Since Euler angles represent orientation via three axes, and there are overlaps when one of these axes is rotated 90 degrees, we could accidentally lock them together, becoming mathematically inseparable and causing future orientation changes to affect both axes simultaneously. Of course, a human being can figure out how to rotate the object in order to solve this problem, but Gimbal Lock is purely a mathematical problem. The classic example is the orientation bubble in a fighter jet. The pilot never has problems with Gimbal Lock, but the orientation instruments in their heads-up display could become inaccurate because of it. Quaternions solve this problem by including a fourth value that effectively allows overlapping axes to still be distinguished from one another.

Exercise restraint

Performance targets for VR apps are very difficult to reach. It is, therefore, important to recognize when we are simply trying to cram too much quality into our app than is tolerable for the current generation of XR devices and typical user hardware. The last resort is always to cull objects from our Scenes until we reach our performance goals. We should be more willing to do so for an XR app than a non-XR one, since the costs of poor performance often far outweigh the gains of higher quality. We must exercise restraint from adding more detail to our Scenes if it has become apparent that the rendering budget has been exhausted. This can be difficult to admit with immersive VR content, where we want to create as much compelling immersion as we can, but until the technology catches up with the capabilities, we need to remain frugal.

Keep up to date with the latest developments

Unity provides a list of useful articles containing VR design and optimization tips, which will likely get updated as the medium and market matures and new techniques are discovered. It can be kept more up-to-date than this tome ever could be, so check them out from time to time to catch the latest tips. The articles in question can be found at <https://unity3d.com/learn/tutorials/topics/virtual-reality>.

We should also keep an eye on Unity blogs to make sure that we don't miss anything important with regard to XR API changes, performance enhancements, and performance suggestions.

Summary

Hopefully, this brief guide helps us improve the performance of our XR applications. The reassuring news is that we have many performance optimization options to choose from since Unity XR apps are built upon the same underlying platform we've been exploring throughout this book. Less reassuring is the fact that we might have to test and implement *all of them* in order to stand a chance of reaching our quality goals. We can expect hardware to get more powerful over time, prices to come down, and adoption to increase as a result. However, until then, we need to pull out all of the stops if we're going to compete in the tech world's latest craze.

In the next chapter, we'll dig into Unity's underlying Engine along with the various frameworks, layers, and languages that it is built from. In essence, we will take a look at our Script code in a more advanced light and investigate some methods to improve our CPU and memory management across the board.

Masterful Memory Management

Memory efficiency is an important element of performance optimization. It's possible for games of limited scope, such as hobby projects and prototypes, to get away with ignoring memory management; they will tend to waste a lot of resources, and potentially leak memory, but this won't be a problem if we limit its exposure to friends and coworkers. However, anything we want to release professionally needs to take this subject seriously. Unnecessary memory allocations lead to poor user experience due to excessive garbage collection (costing precious CPU time), and memory leaks which will lead to crashes. None of these situations are acceptable in modern game releases.

Using memory efficiently with Unity requires a solid understanding of the underlying Unity Engine, Mono platform, the C# language. Also, if we're making use of the new IL2CPP scripting backend then it would be wise to become familiar with its inner workings. This can be a bit of an intimidating place for some developers since many pick Unity3D for their game development solution primarily to avoid the kind of low-level work that comes from Engine development and memory management. We'd prefer to focus on higher-level concerns related to gameplay implementation, level design, and art asset management, but, unfortunately, modern computer systems are complex tools and ignoring low-level concerns for too long could potentially lead to disaster.

Understanding what is happening with memory allocations and C# language features, how they interact with the Mono platform and how Mono interacts with the underlying Unity Engine are absolutely paramount to making high-quality, efficient script code. So, in this chapter, you will learn about all of the nuts and bolts of the underlying Unity Engine: the Mono platform, the C# language, IL2CPP, and the .NET Framework.

Fortunately, it is not necessary to become absolute masters of the C# language to use it effectively. This chapter will boil these complex subjects down to a more digestible form and is split into the following subjects:

- The Mono platform:
 - Native and Managed Memory Domains
 - Garbage collection
 - Memory Fragmentation
- IL2CPP
- How to profile memory issues
- Various memory-related performance enhancements:
 - Minimizing garbage collection
 - Using Value types and Reference types properly
 - Using strings responsibly
 - A multitude of potential enhancements related to the Unity Engine
 - Object and Prefab Pooling

The Mono platform

Mono is a magical sauce, mixed into the Unity recipe, which gives it a lot of its cross-platform capability. Mono is an open source project that built its own platform of libraries based on the API, specifications, and tools from Microsoft's .NET Framework. Essentially, it is an open source recreation of the .NET Library, was accomplished with little-to-no access to the original source code, and is fully compatible with the original library from Microsoft.

The goal of the Mono project is to provide cross-platform development through a framework that allows code, written in a common programming language, to run against many different hardware platforms, including Linux, MacOS, Windows, ARM, PowerPC, and more. Mono even supports many different programming languages. Any language that can be compiled into .NET's **Common Intermediate Language (CIL)** is sufficient to integrate with the Mono platform. This includes C# itself, but also several other languages such as F#, Java, Visual Basic .NET, PythonNet, and IronPython. Of course, the only three exposed to us through Unity are C#, Boo, and UnityScript.



Note that the Boo language has already been deprecated in previous versions of Unity, and the UnityScript language will start becoming phased out in future versions. The blog post at <https://blogs.unity3d.com/2017/08/11/unityscripts-long-ride-off-into-the-sunset/> explains the reasoning behind these changes.

A common misconception about the Unity Engine is that it is built on top of the Mono platform. This is untrue, as its Mono-based layer does not handle many important game tasks such as audio, rendering, physics and keeping track of time. Unity Technologies built a Native C++ backend for the sake of speed and allows its users control of this Game Engine through Mono as a scripting interface. As such, Mono is merely an ingredient of the underlying Unity Engine. This is equivalent to many other Game Engines, which run C++ under the hood, handling important tasks such as rendering, animation, and resource management, while providing a higher-level scripting language

for Gameplay logic to be implemented. As such, the Mono platform was chosen by Unity Technologies to provide this feature.



Native Code is common vernacular for code that is written specifically for the given platform. For instance, writing code to create a window object or interface with networking subsystems in Windows would be completely different to code performing the tasks for a Mac, Unix, Playstation 4, XBox One, and so on.

Scripting languages typically abstract away complex memory management through automatic garbage collection and provide various safety features, which simplify the act of programming at the expense of runtime overhead. Some scripting languages can also be interpreted at runtime, meaning that they don't need to be compiled before execution. The raw instructions are converted dynamically into Machine Code and executed the moment they are read during runtime; of course, this often makes the code relatively slow. The last feature, and probably the most important one, is that they allow simpler syntax of programming commands. This usually improves development workflow immensely, as team members without much experience using languages such as C++ can still contribute to the code base. This enables them to implement things such as Gameplay logic in a simpler format at the expense of a certain amount of control and runtime execution speed.

Note that such languages are often called *Managed Languages*, which feature *Managed Code*. Technically, this was a term coined by Microsoft to refer to any source code that must run inside their **Common Language Runtime (CLR)** environment, as opposed to code that is compiled and run *Natively* through the target OS.

However, because of the prevalence and common features that exist between the CLR and other languages that feature their own similarly designed runtime environments (such as Java), the term *Managed* has since been hijacked. It tends to be used to refer to any language or code that depends on its own runtime environment and that may, or may not, include automatic garbage collection. For the rest of this chapter, we will adopt this definition and use the term Managed to refer to code that both depends on a separate runtime environment in order to execute and is being monitored by automatic garbage collection.

The runtime performance cost of Managed Languages is always greater than the equivalent Native Code, but it is becoming less significant every year. This is partly due to gradual optimizations in tools and runtime environments, and partly due to the computing power of the average device gradually becoming greater. Although, the main point of controversy with using Managed Languages still remains their automatic memory management. Managing memory manually can be a complex task that can take many years of difficult debugging to be proficient at, but many developers feel that Managed Languages solve this problem in ways that are too unpredictable, risking too much product quality. Such developers might cite that Managed Code will never reach the same level of performance as Native Code, and hence it is foolhardy to build high-performance applications with them.

This is true to an extent, as Managed Languages invariably inflict runtime overheads, and we lose partial control over runtime memory allocations. This would be a deal-breaker for high-performance server architecture; however, for game development, it becomes a balancing act since not all resource usage will necessarily result in a bottleneck, and the best games aren't necessarily the ones that use every single byte to their fullest potential. For example, imagine a user interface that refreshes in 30 microseconds via Native Code versus 60 microseconds in Managed Code due to an extra 100 percent overhead (an extreme example). The Managed Code version is still fast enough such that the user will never be able to notice the difference, so is there really any harm in using Managed Code for such a task?

In reality, at least for game development, working with Managed Languages often just means that developers have a unique set of concerns to worry about compared to Native Code developers. As such, choosing to use a Managed Language for game development is partly a matter of preference and partly a compromise of control over development speed.

Let's revisit a topic we touched upon in earlier chapters, but didn't quite flesh out: the concept of Memory Domains in the Unity Engine.

Memory Domains

Memory space within the Unity Engine can be essentially split into three different Memory Domains. Each Domain stores different types of data and takes care of a very different set of tasks.

The first Memory Domain--the Managed Domain--should be very familiar. This Domain is where the Mono platform does its work, where any `MonoBehaviour` scripts and custom C# classes we write will be instantiated at runtime, and so we will interact with this Domain very explicitly through any C# code we write. It is called the *Managed Domain* because this memory space is automatically managed by a Garbage Collector.

The second Domain--the Native Domain--is more subtle since we only interact with it indirectly. Unity has an underlying Native Code foundation, which is written in C++ and compiled into our application differently, depending on which platform is being targeted. This Domain takes care of allocating internal memory space for things such as asset data (for example, textures, audio files, and meshes) and memory space for various subsystems such as the Rendering Pipeline, Physics System, and User Input System. Finally, it includes partial Native *representations* of important Gameplay objects such as `GameObjects` and Components so that they can interact with these internal systems. This is where a lot of built-in Unity classes keep their data, such as the `Transform` and `Rigidbody` Components.

The Managed Domain also includes wrappers for the very same object representations that are stored within the Native Domain. As a result, when we interact with Components such as `Transform`, most instructions will ask Unity to dive into its Native Code, generate the result there, and then copy it back to the Managed Domain for us. This is where the *Native-Managed Bridge* between the Managed Domain and Native Domains derives from, which was briefly mentioned in previous chapters. When both Domains have their own representations for the same entity, crossing the bridge between them requires a memory context-switch that can potentially inflict some fairly significant performance hits on our game. Obviously, crossing

back and forth across this bridge should be minimized as much as possible due to the overhead involved. We covered several techniques for this in [Chapter 2, Scripting Strategies](#).

The third and final Memory Domains are those of external libraries, such as DirectX and OpenGL libraries, as well as any custom libraries and plugins we include in our project. Referencing these libraries from our C# code will cause a similar memory context switch and subsequent cost.

Memory in most modern Operating Systems (OS) splits runtime memory space into two categories: the stack and the heap. The stack is a special reserved space in memory, dedicated to small, short-lived data values, which are automatically deallocated the moment they go out of scope, hence why it is called the stack. It literally operates like a stack data structure, pushing and popping data from the top. The stack contains any local variables we declare and handles the loading and unloading of functions as they're called. These function calls expand and contract through what is known as the call stack. When the call stack is done with the current function, it jumps back to the previous point on the call stack and continues from where it left off. The start of the previous memory allocation is always known, and there's no reason to perform any clean-up operations since any new allocations can simply overwrite the old data. Hence, the stack is relatively quick and efficient.

The total stack size is usually very small, usually on the order of Megabytes. It's possible to cause a stack overflow by allocating more space than the stack can support. This can occur during exceptionally large call stacks (for example, an infinite loop) or having a large number of local variables, but in most cases, causing a stack overflow is rarely a concern despite its relatively small size.

The heap represents all remaining memory space, and it is used for the overwhelming majority of memory allocation. Since we want most of the memory allocated to persist longer than the current function call, we couldn't allocate it on the stack since it would just get overwritten when the current function ends. So, instead, whenever a data type is too big to fit in the stack or must persist outside the function it was declared in, it is allocated on the heap. There's nothing physically different between the stack and the heap; they're both just memory spaces containing bytes of data that exist in RAM,

which have been requested and set aside for us by the OS. The only difference is in when, where, and how they are used.

In Native Code, such as code written in languages such as C++, these memory allocations are handled very manually in that we are responsible for ensuring that all pieces of memory we allocate are properly and explicitly deallocated when they are no longer needed. If this is not done properly, then we could easily and accidentally introduce memory leaks since we are likely to keep allocating more and more memory space from RAM that is never cleaned up until there is no more space to allocate, and the application crashes.

Meanwhile, in Managed Languages, this process is automated through the Garbage Collector. During initialization of our Unity app, the Mono platform will request a given chunk of memory from the OS and use it to generate a heap memory space that our C# code can use (often known as the *Managed Heap*). This heap space starts off fairly small, less than 1 Megabyte, but will grow as new blocks of memory are needed by our script code. This space can also shrink by releasing it back to the OS if Unity determines that it's no longer needed.

Garbage collection

The Garbage Collector (hereafter referred to as the *GC*) has an important job, which is to ensure that we don't use more Managed Heap memory than we need, and that memory that is no longer needed will be automatically deallocated. For instance, if we create a `GameObject`, and then later destroy it, the GC will flag the memory space used by the `GameObject` for eventual deallocation later. This is not an immediate process, as the GC only deallocates memory when necessary.

When a new memory request is made, and there is enough empty space in the Managed Heap to satisfy the request, the GC simply allocates the new space and hands it over to the caller. However, if the Managed Heap does not have room for it, then the GC will need to scan all the existing memory allocations for anything that is no longer being used and cleans them up first. It will only expand the current heap space as a last resort.

The GC in the version of Mono that Unity uses is a type of Tracing Garbage Collector, which uses a Mark-and-Sweep strategy. This algorithm works in two phases: each allocated object is tracked with an additional bit. This flags whether the object has been marked or not. These flags start off set to `false` to indicate that it has not yet been marked.

When the collection process begins, it marks all objects that are still reachable to the program by setting their flags to `true`. Either the reachable object is a direct reference, such as static or local variables on the stack, or it is an indirect reference through the fields (member variables) of other directly or indirectly accessible objects. In essence, it is gathering a set of objects that are still referenceable to our application. Everything that is not still referenceable would be effectively invisible to our application and can be deallocated by the GC.

The second phase involves iterating through this catalog of references (which the GC will have kept track of throughout the lifetime of the application) and determining whether or not it should be deallocated based on its

marked status. If the object is marked, then it is still being referenced by something else, and so the GC leaves it alone. However, if it is not marked, then it is a candidate for deallocation. During this phase, all marked objects are skipped over, but not before setting their flag back to `false` for the first phase of the next garbage collection scan.



In essence, the GC maintains a list of all objects in memory, while our application maintains a separate list containing only a portion of them. Whenever our application is done with an object, it simply forgets it exists, removing it from its list. Hence, the list of objects that can be safely deallocated would be the difference between the GC's list, and our application's list.

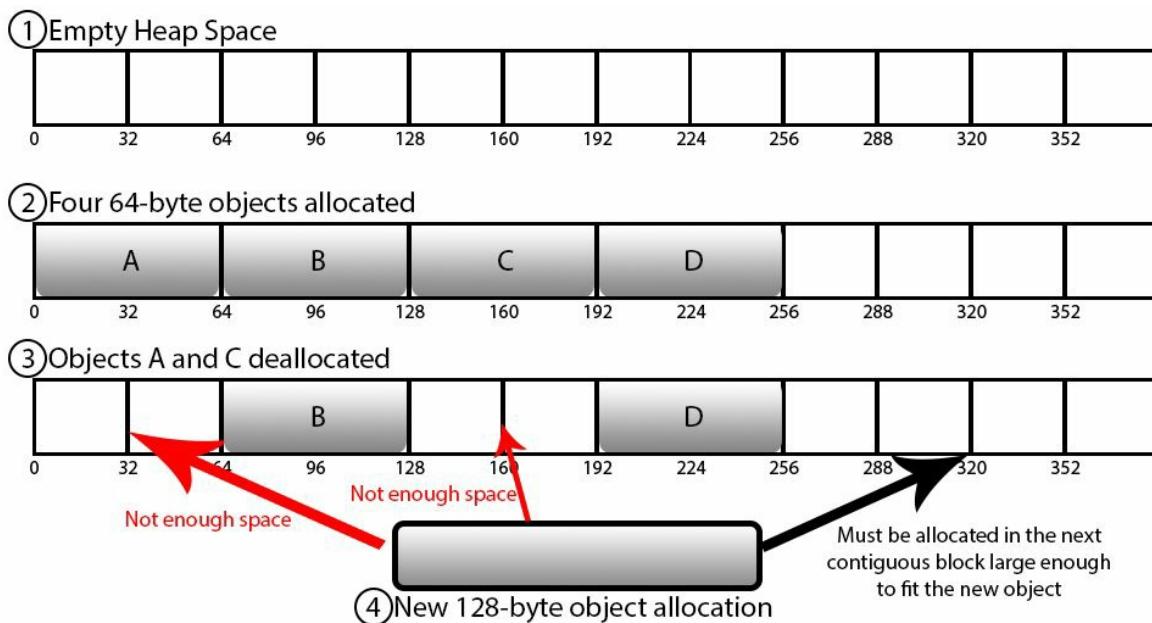
Once the second phase ends, all unmarked objects are deallocated to free space, and then the initial request to create the object is revisited. If the GC has freed up enough space for the object, then it is allocated within that newly-freed space and returned to the caller. However, if it is not, then we hit the last-resort situation and must expand the Managed Heap by requesting it from the OS, at which point the object space can finally be allocated and returned to the caller.

In an ideal world, where we only keep allocating and deallocating objects but only a finite number of them exist at once, the heap would maintain a roughly constant size because there's always enough space to fit the new objects we need. However, all objects in an application are rarely deallocated in the same order they were allocated and even more rarely do they all have the same size in memory. This leads to *Memory Fragmentation*.

Memory Fragmentation

Fragmentation occurs when objects of different sizes are allocated and deallocated in alternating orders and if lots of small objects are deallocated, following by lots of large objects being allocated.

This is best explained through an example. The following shows four steps we take in allocating and deallocating memory in a typical heap memory space:



The memory allocation takes place, as follows:

1. We start with an empty heap space.
2. We then allocate four objects on the heap, **A**, **B**, **C**, and **D**, each 64-bytes in size.
3. At a later time, we deallocate two of the objects, **A** and **C**, freeing up 128 bytes.
4. We then try to allocate a new object that is 128-bytes large.

Deallocating objects **A** and **C** technically frees 128-bytes worth of space, but since the objects were not contiguous (adjoining neighbors) in memory, we

cannot allocate an object larger than both individual spaces there. New memory allocations must always be contiguous in memory; therefore, the new object must be allocated in the next available contiguous 128-byte space available in the Managed Heap. We now have two empty 64-byte holes in our memory space, which will never be reused unless we allocate objects sized 64-bytes or smaller.

Over long periods of time, our heap memory can become riddled with more, smaller empty spaces such as these, as objects of different sizes are deallocated, and then the system later tries to allocate new objects within the smallest available space that it can fit within, leaving some small remainder that becomes harder to fill. In the absence of background techniques that automatically clean up this Fragmentation, this effect would occur in literally any memory space--RAM, heap space, and even hard drives--which are just larger, slower, and more permanent memory storage areas (this is why it's a good idea to defragment our hard drives from time to time).

Memory Fragmentation causes two problems. Firstly, it effectively reduces the total usable memory space for new objects over long periods of time, depending on the frequency of allocations and deallocations. This is likely to result in the GC having to expand the heap to make room for new allocations. Secondly, it makes new allocations take longer to resolve due to the extra time it takes to find a new memory space large enough to fit the object.

This becomes important when new memory allocations are made in a heap since the location of available space becomes just as important as how much free space is available. There is no way to split an object across partial memory locations, so the GC must either continue searching until it finds a large enough space or the entire heap size must be increased to fit the new object, costing even more time after it just spent a bunch of time doing an exhaustive search.

Garbage collection at runtime

So, in a worst-case scenario, when a new memory allocation is being requested by our game, the CPU would have to spend cycles completing the following tasks before the allocation is finally completed:

1. Verify that there is enough contiguous space for the new object.
2. If there is not enough space, iterate through all known direct and indirect references, marking everything they connect to as reachable.
3. Iterate through all of these references again, flagging unmarked objects for deallocation.
4. Iterate through all flagged objects to check whether deallocating some of them would create enough contiguous space for the new object.
5. If not, request a new memory block from the OS in order to expand the heap.
6. Allocate the new object at the front of the newly allocated block and return it to the caller.

This can be a lot of work for the CPU to handle, particularly if this new memory allocation is an important object such as a Particle Effect, a new character entering the Scene, or a cutscene transition,. Users are extremely likely to note moments where the GC is freezing gameplay to handle this extreme case. To make matters worse, the garbage collection workload scales poorly as the allocated heap space grows since sweeping through a few Megabytes of space will be significantly faster than scanning several Gigabytes of space.

All of this makes it absolutely critical to control our heap space intelligently. The lazier our memory usage tactics are, the worse the GC will behave in an almost exponential fashion, as we are more and more likely to hit this worst-case scenario. So, it's a little ironic that despite the efforts of Managed Languages to make the memory management problem easier, Managed Language developers still find themselves being just as, if not more, concerned with memory consumption than developers of Native applications. The main difference is in the types of problems they're trying to solve.

Threaded garbage collection

The GC runs on two separate threads: the main thread and what is called the *Finalizer Thread*. When the GC is invoked, it will run on the main thread and flag heap memory blocks for future deallocation. This does not happen immediately. The Finalizer Thread, controlled by Mono, can have a delay of several seconds before the memory is finally freed and available for reallocation.

We can observe this behavior in the Total Allocated block (the green line, with apologies to that 5 percent of the population with deuteranopia/deuteranomaly) of the Memory Area within the Profiler window. It can take several seconds for the total allocated value to drop after a garbage collection has occurred. Owing to this delay, we should not rely on memory being available the moment it has been deallocated, and as such, we should never waste time trying to eke out every last byte of memory that we believe should be available. We must ensure that there is always some kind of buffer zone available for future allocations.

Blocks that have been freed by the GC may sometimes be given back to the OS after some time, which would reduce the reserved space consumed by the heap and allow the memory to be allocated for something else, such as another application. However, this is very unpredictable and depends on the platform being targeted, so we shouldn't rely on it. The only safe assumption to make is that as soon as the memory has been allocated to Mono, it's then reserved and is no longer available to either the Native Domain or any other application running on the same system.

Code compilation

When we make changes to our C# code, it is automatically compiled when we switch back from our favorite IDE (which is typically either MonoDevelop or the much more feature-rich Visual Studio) to the Unity Editor. However, the C# code is not converted directly into Machine Code, as we would expect static compilers to do if we are using languages such as C++.

Instead, the code is converted into an intermediate stage called **Common Intermediate Language (CIL)**, which is an abstraction above Native Code. This is how .NET can support multiple languages--each uses a different compiler, but they're all converted into CIL, so the output is effectively the same regardless of the language that we pick. CIL is similar to Java bytecode, upon which it is based, and the CIL code is entirely useless on its own, as CPUs have no idea how to run the instructions defined in this language.

At runtime, this intermediate code is run through the Mono **Virtual Machine (VM)**, which is an infrastructure element that allows the same code to run against multiple platforms without the need to change the code itself. This is an implementation of the .NET **Common Language Runtime (CLR)**. If we're running on iOS, we run on the iOS-based Virtual Machine infrastructure, and if we're running on Linux, then we simply use a different one that is better suited for Linux. This is how Unity allows us to write code once, and it works magically on multiple platforms.

Within the CLR, the intermediate CIL code will actually be compiled into Native Code on demand. This immediate Native compilation can be accomplished either by an **Ahead-Of-Time (AOT)** or **Just-In-Time (JIT)** compiler. Which one is used will depend on the platform that is being targeted. These compilers allow code segments to be compiled into Native Code, allowing the platform's architecture to complete the written instructions without having to write them ourselves. The main difference between the two compiler types is when the code is compiled.

AOT compilation is the typical behavior for code compilation and happens early (*ahead of time*) either during the build process or in some cases during app initialization. In either case, the code has been precompiled and no further runtime costs are inflicted due to dynamic compilation since there are always Machine Code instructions available whenever the CPU needs them.

JIT compilation happens dynamically at runtime in a separate thread and begins just prior to execution (*just in time* for execution). Often, this dynamic compilation causes the first invocation of a piece of code to run a little (or a lot) more slowly, because the code must finish compiling before it can be executed. However, from that point forward, whenever the same code block is executed, there is no need for recompilation, and the instructions run through the previously compiled Native Code.

A common adage in software development is that 90 percent of the work is being done by only 10 percent of the code. This generally means that JIT compilation turns out to be a net positive on performance than if we simply tried to interpret the CIL code directly. However, because the JIT compiler must compile code quickly, it is not able to make use of many optimization techniques that static AOT compilers are able to exploit.



Not all platforms support JIT compilation, but some scripting functionalities are not available when using AOT. Unity provides a complete list of these restrictions at <https://docs.unity3d.com/Manual/ScriptingRestrictions.html>.

IL2CPP

A few years ago, Unity Technologies was faced with a choice to either continue to support the Mono platform, which Unity was finding more and more difficult to keep up with, or implement their own scripting backend. They chose the latter option and multiple platforms now support **IL2CPP**, which is short for **Intermediate Language to C++**.



The Unity Technologies' initial post about IL2CPP, the reasoning behind the decision, and its long-term benefits can be found at <https://blogs.unity3d.com/2014/05/20/the-future-of-scripting-in-unity/>.

IL2CPP is a scripting backend designed to convert Mono's CIL output directly into Native C++ code. This leads to improved performance since the application will now be running Native Code. This ultimately gives Unity Technologies more control of runtime behavior since IL2CPP provides its own AOT compiler and VM, allowing custom improvements to subsystems such as the GC and compilation process. IL2CPP does not intend to completely replace the Mono platform, but it is an additional tool we can enable, which improves part of the functionality that Mono provides.

Note that IL2CPP is automatically enabled for iOS and WebGL projects. For other platforms that support it, IL2CPP can be enabled under Edit | Project Settings | Player | Configure | Scripting Backend.



A list of platforms currently supporting IL2CPP can be found at <https://docs.unity3d.com/Manual/IL2CPP.html>.

Profiling memory

There are two issues we are concerned about when it comes to memory management: how much we're consuming, and how often we're allocating new blocks. Let's cover each of these topics separately.

Profiling memory consumption

We do not have direct control over what is going on in the Native Domain since we don't have the Unity Engine source code and hence can't add any code that will interact with it directly. We can, however, control it indirectly by means of various script-level functions that serve as interaction points between Managed and Native Code. There are technically a variety of memory allocators available, which are used internally for things such as `GameObjects`, `Graphics` objects, and the Profiler, but these are hidden behind the Native-Managed Bridge.

However, we can observe how much memory has been allocated and reserved in this Memory Domain via the Memory Area of the Profiler window. Native memory allocations show up under the values labeled Unity, and we can even get more information using Detailed Mode and sampling the current frame:

```
Used Total: 101.2 MB [Unity: 68.1 MB] Mono: 7.8 MB GfxDriver: 15.8 MB FMOD: 1.3 MB Video: 224 B Profiler: 9.5 MB  
Reserved Total: 241.4 MB [Unity: 199.0 MB] Mono: 10.7 MB GfxDriver: 15.8 MB FMOD: 1.3 MB Video: 224 B Profiler: 16.0 MB  
Total System Memory Usage: 0.78 GB
```

Under the Scene Memory section of Breakdown View, we can observe that `MonoBehaviour` objects always consume a constant amount of memory, regardless of their member data. This is the memory consumed by the Native representation of the object.



Note that memory consumption in Edit Mode is always wildly different to that of a stand-alone version due to various debugging and editor hook data being applied. This adds a further incentive to avoid using Edit Mode for benchmarking and instrumentation purposes.

We can also use the `Profiler.GetRuntimeMemorySize()` method to get the Native memory allocation size of a particular object.

Managed object representations are intrinsically linked to their Native representations. The best way to minimize our Native memory allocations is

to simply optimize our Managed memory usage.

We can verify how much memory has been allocated and reserved for the Managed Heap using the Memory Area of the Profiler window, under the values labeled Mono, as follows:

```
Used Total: 101.2 MB Unity: 68.1 MB Mono: 7.8 MB GfxDriver: 15.8 MB FMOD: 1.3 MB Video: 224 B Profiler: 9.5 MB  
Reserved Total: 241.4 MB Unity: 199.0 MB Mono: 10.7 MB GfxDriver: 15.8 MB FMOD: 1.3 MB Video: 224 B Profiler: 16.0 MB  
Total System Memory Usage: 0.78 GB
```

We can also determine the current used and reserved heap space at runtime using the `Profiler.GetMonoUsedSize()` and `Profiler.GetMonoHeapSize()` methods, respectively.

Profiling memory efficiency

The best metric we can use to measure the health of our memory management is simply watching the behavior of the GC. The more work it's doing, the more waste we're generating and the worse our application's performance is likely to become.

We can use both the CPU Usage Area (the GarbageCollector checkbox) and Memory Area (the GC Allocated checkbox) of the Profiler window to observe the amount of work the GC is doing and the time it is doing it. This can be relatively straightforward for some situations, where we only allocated a temporary small block of memory or we just destroyed a `GameObject`.

However, root-cause analysis for memory efficiency problems can be challenging and time-consuming. When we observe a spike in the GC's behavior, it could be a symptom of allocating too much memory in a previous frame and merely allocating a little more in the current frame, requiring the GC to scan a lot of fragmented memory, determine whether there is enough space, and decide whether to allocate a new block. The memory it cleaned up could have been allocated a long time ago, and we may only be able to observe these effects when our application runs over long periods of time and could even happen when our Scene is sitting relatively idle, giving no obvious cause for the GC to suddenly trigger. Even worse, the Profiler can only tell us what happened in the last few seconds or so, and it won't be immediately obvious what data was being cleaned up.

We must be vigilant and test our application rigorously, observing its memory behavior while simulating a typical play session if we want to be certain we are not generating memory leaks or creating a situation where the GC has too much work to complete in a single frame.

Memory management performance enhancements

In most Game Engines, we would have the luxury of being able to port inefficient Managed Code into faster Native Code if we were hitting performance issues. This is not an option unless we invest serious cash in obtaining the Unity source code, which is offered as a license separate from the Free/Personal/Pro licensing system, and on a per case, per title basis. We could also purchase a license of Unity Pro with the hope of using Native Plugins, but doing so rarely leads to a performance benefit since we must still cross the Native-Managed Bridge to invoke function calls inside of it. Native Plugins are normally used to interface with systems and libraries that are not built specifically for C#. This forces the overwhelming majority of us into a position of needing to make our C# script-level code as performant as possible ourselves.

With this in mind, we should now have enough understanding of Unity Engine internals and memory spaces to detect and analyze memory performance issues and understand and implement enhancements for them. So, let's cover some performance enhancements we can apply.

Garbage collection tactics

One strategy to minimize garbage collection problems is concealment by manually invoking the GC at opportune moments, when we're certain the player would not notice. Garbage collection can be manually invoked by calling `System.GC.Collect()`.

Good opportunities to invoke a collection may be while loading between levels, when gameplay is paused, shortly after a menu interface has been opened, during cutscene transitions, or any break in Gameplay that the player would not witness, or care about, a sudden performance drop. We could even use the `Profiler.GetMonoUsedSize()` and `Profiler.GetMonoHeapSize()` methods at runtime to determine whether a garbage collection needs to be invoked in the near future.

We can also cause the deallocation of a handful of specific objects. If the object in question is one of the Unity object wrappers, such as a `GameObject` or `MonoBehaviour` Component, then the finalizer will first invoke the `Dispose()` method within the Native Domain. At this point, the memory consumed by both the Native and Managed Domains will then be freed. In some rare instances, if the Mono wrapper implements the `IDisposable` Interface Class (that is, it has a `Dispose()` method available from script code), then we can actually control this behavior and force the memory to be freed instantly.

There are a number of different object types in the Unity Engine (most of which are introduced in Unity 5 or later), which implement the `IDisposable` Interface Class, as follows: `NetworkConnection`, `WWW`, `UnityWebRequest`, `UploadHandler`, `DownloadHandler`, `VertexHelper`, `CullingGroup`, `PhotoCapture`, `VideoCapture`, `PhraseRecog` and more.

These are all utility classes for pulling in potentially large datasets where we might want to ensure immediate destruction of the data it has acquired, since they normally involve allocating several buffers and memory blocks in the

Native Domain in order to accomplish their tasks. If we kept all of this memory for a long time, it would be a colossal waste of precious space. So, by calling their `Dispose()` method from script code, we can ensure that the memory buffers are freed promptly and precisely when they need to be.

All other Asset objects offer some kind of unloading method to clean up any unused asset data, such as `Resources.UnloadUnusedAssets()`. Actual asset data is stored within the Native Domain, so the GC technically isn't involved here, but the idea is basically the same. It will iterate through all Assets of a particular type, check whether they're no longer being referenced, and, if so, deallocate them. However, again, this is an asynchronous process, and we cannot guarantee exactly when the deallocation will occur. This method is automatically called internally after a Scene is loaded, but this still doesn't guarantee instant deallocation. The preferred approach is to use `Resources.UnloadAsset()` instead, which will unload one specific asset at a time. This method is generally faster since time will not be spent iterating through an entire collection of asset data, in order to figure out what is unused.

However, the best strategy for garbage collection will always be avoidance; if we allocate as little heap memory and control its usage as much as possible, then we won't have to worry about the GC inflicting frequent, expensive performance costs. We will cover many tactics for this throughout the remainder of this chapter.

Manual JIT compilation

In the event that JIT compilation is causing a runtime performance loss, be aware that it is actually possible to force JIT compilation of a method at any time via **Reflection**. Reflection is a useful feature of the C# language that allows our code base to explore itself introspectively for type information, methods, values, and metadata. Using Reflection is often a very costly process. It should be avoided at runtime or, at the very least, only used during initialization or other loading times. Not doing so can easily cause significant CPU spikes and Gameplay freezing.

We can manually force JIT compilation of a method using Reflection to obtain a function pointer to it:

```
var method = typeof(MyComponent).GetMethod("MethodName");
if (method != null) {
    method.MethodHandle.GetFunctionPointer();
    Debug.Log("JIT compilation complete!");
}
```

The preceding code only works on `public` methods. Obtaining `private` or `protected` methods can be accomplished through the use of `BindingFlags`:

```
using System.Reflection;
// ...
var method = typeof(MyComponent).GetMethod("MethodName",
    BindingFlags.NonPublic | BindingFlags.Instance);
```

This kind of code should only be run for very targeted methods where we are certain that JIT compilation is causing CPU spikes. This can be verified by restarting the application and profiling a method's first invocation versus all subsequent invocations. The difference will tell us the JIT compilation overhead.



Note that the official method for forcing JIT compilation in the .NET library is `RuntimeHelpers.PrepareMethod()`, but this is not properly implemented in the current version of Mono that comes with Unity (Mono Version 2.6.5). The aforementioned workaround should be used until Unity has pulled in a more

recent version of the Mono project.

Value types and Reference types

Not all memory allocations we make within Mono will go through the heap. The .NET Framework (and, by extension, the C# language, which merely implements the .NET specification) has the concept of Value types and Reference types, and only the latter needs to be marked by the GC while it is performing its Mark-and-Sweep algorithm. Reference types are expected to (or need to) last a long time in memory due to their complexity, their size, or how they're used. Large datasets, and any kind of object instantiated from a class, is a Reference type. This also includes arrays (regardless of whether it is an array of Value types or Reference types), delegates, all classes, such as MonoBehaviour, GameObject, and any custom classes we define.

Reference types are always allocated on the heap, whereas Value types can be allocated either on the stack or the heap. Primitive data types such as bool, int, and float are examples of Value types. These values are typically allocated on the stack, but as soon as a Value type is contained within a Reference type, such as a class or an array, then it is implied that it is either too large for the stack or will need to survive longer than the current scope and must be allocated on the heap, bundled with the Reference type it is contained within.

All of this can be best explained through examples. The following code will create an integer as a Value type that exists on the stack only temporarily:

```
public class TestComponent {
    void TestFunction() {
        int data = 5; // allocated on the stack
        DoSomething(data);
    } // integer is deallocated from the stack here
}
```

As soon as the start() method ends, the integer is deallocated from the stack. This is essentially a free operation since, as mentioned previously, it doesn't bother doing any cleanup; it just moves the stack pointer back to the previous

memory location in the call stack (back to whichever function called `TestFunction()` on the `TestComponent` object). Any future stack allocations simply overwrite the old data. More importantly, no heap allocation took place to create the data, so the GC has no need to track its existence.

However, if we created an integer as a member variable of the `MonoBehaviour` class definition, then it is now contained within a Reference type (a `class`) and must be allocated on the heap along with its container:

```
public class TestComponent : MonoBehaviour {
    private int _data = 5;
    void TestFunction() {
        DoSomething(_data);
    }
}
```

The `_data` integer is now an additional piece of data that consumes space in the heap alongside the `TestComponent` object it is contained within. If the `TestComponent` is destroyed, then the integer is deallocated along with it, but not before then.

Similarly, if we put the integer into a normal C# `class`, then the rules for Reference types still apply and the object is allocated on the heap:

```
public class TestData {
    public int data = 5;
}

public class TestComponent {
    void TestFunction() {
        TestData dataObj = new TestData(); // allocated on the heap
        DoSomething(dataObj.data);
    } // dataObj is not immediately deallocated here, but it will
      // become a candidate during the next GC sweep
}
```

So, there is a big difference between creating a temporary Value type within a `class` method versus storing long-term Value type as a member field of a `class`. In the former case, we're storing it in the stack, but in the latter case, we're storing it within a Reference type, which means it can be referenced elsewhere. For example, imagine that `DoSomething()` has stored the reference to `dataObj` within a member variable:

```
public class TestComponent {
    private TestData _testDataObj;
```

```

void TestFunction() {
    TestData dataObj = new TestData(); // allocated on the heap
    DoSomething(dataObj.data);
}

void DoSomething (TestData dataObj) {
    _testDataObj = dataObj; // a new reference created! The referenced
    // object will now be marked during Mark-and-Sweep
}
}

```

In this case, we would not be able to deallocate the object pointed to `dataObj` as soon as the `TestFunction()` method ends because the total number of things referencing the object would go from 2 to 1. This is not 0, and hence the GC would still mark it during Mark-and-Sweep. We would need to set the value of `_testDataObj` to `null` or make it reference something else, before the object is no longer reachable.

Note that a Value type must have a value and can never be `null`. If a stack-allocated Value type is assigned to a Reference type, then the data is simply copied. This is true even for arrays of Value types:

```

public class TestClass {
    private int[] _intArray = new int[1000]; // Reference type
                                              // full of Value types
    void StoreANumber(int num) {
        _intArray[0] = num; // store a Value within the array
    }
}

```

When the initial array is created (during object initialization), 1,000 integers will be allocated on the heap set to a value of 0. When the `StoreANumber()` method is called, the value of `num` is merely copied into the *zeroth* element of the array rather than storing a reference to it.

The subtle change in the referencing capability is what ultimately decides whether something is a Reference type or a Value type, and we should try to use Value types whenever we have the opportunity so that they generate stack allocations instead of heap allocations. Any situation where we're just sending around a piece of data that doesn't need to live longer than the current scope is a good opportunity to use a Value type instead of a Reference type. Ostensibly, it does not matter if we pass the data into another method of the same class or a method of another class--it still remains a

Value type that will exist on the stack until the method that created it goes out of the scope.

Pass by value and by reference

Technically, something is duplicated every time a data value is passed as an argument from one method to another, and this is true whether it is a Value type or a Reference type. When we're passing the object's data, this is known as *passing by value*. When we're simply copying a reference to something else, it is called *passing by reference*.

An important difference between Value types and Reference types is that a Reference type is merely a pointer to another location in memory that consumes only 4 or 8 bytes in memory (32 bit or 64 bit, depending on the architecture), regardless of what it is actually pointing to. When a Reference type is passed as an argument, it is only the value of this pointer that gets copied into the function. Even if the Reference type points to a humongous array of data, this operation will be very quick since the data being copied is very small.

Meanwhile, a Value type contains the full and complete bits of data stored within a concrete object. Hence, all of the data of a Value type will be copied whenever they are passed between methods or stored in other Value types. In some cases, it can mean that passing a large Value type as arguments around too much can be more costly than just using a Reference type and letting the GC take care of it. For most Value types, this is not a problem since they are comparable in size to a pointer, but this becomes important when we begin to talk about the `struct` type in the next section.

Data can also be passed around by reference using the `ref` keyword, but this is very different from the concept of Value and Reference types, and it is very important to keep them distinct in our mind when we try to understand what is going on under the hood. We can pass a Value type by value, or by reference, and we can pass a Reference type by value, or by reference. This means that there are four distinct data passing situations that can occur, depending on which type is being passed and whether the `ref` keyword is being used or not.

When data is passed by reference (even if it is a Value type), then making any changes to the data will change the original. For example, the following code would print the value as 10:

```
void Start() {
    int myInt = 5;
    DoSomething(ref myInt);
    Debug.Log(String.Format("Value = {0}", myInt));
}

void DoSomething(ref int val) {
    val = 10;
}
```

Removing the `ref` keyword from both places would make it print the value 5 instead (and removing it from only one of them would lead to a compiler error since the `ref` keyword needs to be present in both locations or neither). This understanding will come in handy when we start to think about some of the more interesting data types we have access to, namely, structs, arrays, and strings.

Structs are Value types

The `struct` type is an interesting special case in C#. A `struct` object can contain `private`, `protected`, and `public` fields, have methods, and can be instantiated at runtime, just like a `class` type. However, there is a fundamental difference between the two: a `struct` type is a Value type, and a `class` type is a Reference type. Consequently, this leads to some important differences between the two, namely, that a `struct` type cannot support inheritance, their properties cannot be given custom default values (member data always defaults to values such as `0` or `null` since it is a Value type), and their default constructors cannot be overridden. This greatly restricts their usage compared to classes, so simply replacing all classes with structs (under the assumption that it will just allocate everything on the stack) is not as easy as it sounds.

However, if we're using a class in a situation whose only purpose is to send a blob of data to somewhere else in our application, and it does not need to last beyond the current scope, then we might be able to use a `struct` type instead, since a `class` type would result in a heap allocation for no particularly good reason:

```
public class DamageResult {
    public Character attacker;
    public Character defender;
    public int totalDamageDealt;
    public DamageType damageType;
    public int damageBlocked;
    // etc.
}

public void DealDamage(Character _target) {
    DamageResult result = CombatSystem.Instance.CalculateDamage(this, _target);
    CreateFloatingDamageText(result);
}
```

In this example, we're using a `class` type to pass a bunch of data from one subsystem (the combat system) to another (the UI system). The only purpose of this data is to be calculated and read by various subsystems, so this is a good candidate to convert into a `struct` type.

Merely changing the `DamageResult` definition from a `class` type to a `struct` type

could save us quite a few unnecessary garbage collections since it would be allocated on the stack as a Value type instead of the heap as a Reference type:

```
public struct DamageResult {  
    // ...  
}
```

This is not a catch-all solution. Since structs are Value types, the entire blob of data will be duplicated and provided to the next method in the call stack, regardless of how large or small it is. So, if a struct object is passed by value between five different methods in a long chain, then five different stack copies will occur at the same time. Recall that stack deallocations are effectively free, but stack allocations (which involves copying of data) is not. This data copying is pretty much negligible for small values, such as a handful of integers or floating-point values, but passing around ridiculously large datasets through structs over and over again is obviously not a trivial task and should be avoided.

We can work around this problem by passing the `struct` object by reference using the `ref` keyword to minimize the amount of data being copied each time (just a single pointer). However, this can be dangerous since passing by reference allows any subsequent methods to make changes to the `struct` object in which case it would be prudent to make its data values `readonly`. This means that the values can only be initialized in the constructor, and never again, even by its own member functions, which prevents accidental changes as it's passed through the chain.

All of the above is also true when structs are contained within Reference types, as follows:

```
public struct DataStruct {  
    public int val;  
}  
  
public class StructHolder {  
    public DataStruct _memberStruct;  
    public void StoreStruct(DataStruct ds) {  
        _memberStruct = ds;  
    }  
}
```

To the untrained eye, the preceding code appears to be attempting to store a

stack-allocated struct (`ds`) within a Reference type (`structHolder`). Does this mean that a `structHolder` object on the heap can now reference an object on the stack? If so, what will happen when the `storeStruct()` method goes out of scope and the `struct` object is (effectively) erased? It turns out that these are the wrong questions.

What's actually happening is that while a `DataStruct` object (`_memberStruct`) has been allocated on the heap within the `structHolder` object, it is still a Value type and does not magically transform into a Reference type when it is a member variable of a Reference type. So, all of the usual rules for Value types apply. The `_memberStruct` variable cannot have a value of `null`, and all of its fields will be initialized to `0` or `null` values. When `storeStruct()` is called, the data from `ds` will be copied into `_memberStruct` in its entirety. There are no references to stack objects taking place, and there is no concern about lost data.

Arrays are Reference types

The purpose of arrays is to contain large datasets, which make them difficult to be treated as a Value type since there's probably not enough room on the stack to support them. Therefore, they are treated as a Reference type so that the entire dataset can be passed around via a single reference (if it were a Value type, we would need to duplicate the entire array every time it is passed around). This is true irrespective of whether the array contains Value types or Reference types.

This means that the following code will result in a heap allocation:

```
TestStruct[] dataObj = new TestStruct[1000];
for(int i = 0; i < 1000; ++i) {
    dataObj[i].data = i;
    DoSomething(dataObj[i]);
}
```

However, the following, functionally equivalent, code would not result in any heap allocations since the `struct` objects being used are Value types, and hence, it would be created on the stack:

```
for(int i = 0; i < 1000; ++i) {
    TestStruct dataObj = new TestStruct();
    dataObj.data = i;
    DoSomething(dataObj);
}
```

The subtle difference in the second example is that only one `TestStruct` exists on the stack at a time, whereas the first example needs to allocate 1,000 of them via an array. Obviously, these methods are kind of ridiculous as they're written, but they illustrate an important point to consider. The compiler isn't smart enough to automatically find these situations for us and make the appropriate changes. Opportunities to optimize our memory usage through Value type replacements will be entirely down to our ability to detect them, and understand why conversions from Reference types to Value types will result in stack allocations, rather than heap allocations.

Note that when we allocate an array of Reference types, we're creating an array of references, which can each reference other locations on the heap. However, when we allocate an array of Value types, we're creating a packed list of Value types on the heap. Each of these Value types will be initialized with a value of `0` (or equivalent) since they cannot be `null`, while each reference within an array of Reference types will always initialize to `null` since no references have been assigned, yet.

Strings are immutable Reference types

We briefly touched upon the subject of strings in [Chapter 2, Scripting Strategies](#), but now it's time to go into more detail about why proper string usage is extremely important.

Strings are essentially arrays of characters, and so they are considered Reference types, and follow all of the same rules as other Reference types; they will be allocated on the heap, and a pointer is all that is copied from one method to the next. Since a string is effectively an array, this implies that the characters it contains must be contiguous in memory. However, we often find ourselves expanding, contracting, or combining strings to create other strings. This can lead us to making some faulty assumptions about how strings work. We might assume that because strings are such common, ubiquitous objects, performing operations on them is fast and cheap. Unfortunately, this is incorrect. Strings are not made to be fast. They are only made to be convenient.

The `string` object class is immutable, which means they cannot be changed after they've been allocated. Therefore, when we change a string, we are actually allocating a whole new string on the heap to replace it, where the contents of the original will be copied and modified as needed into a whole new character array, and the original `string` object reference now points to a completely new `string` object. In which case, the old `string` object might no longer be referenced anywhere, will not be marked during Mark-and-Sweep, and will eventually be purged by the GC. As a result, lazy string programming can result in a lot of unnecessary heap allocations and garbage collection.

A good example to illustrate how strings are different than normal Reference types is the following code:

```
| void TestFunction() {  
|     string testString = "Hello";
```

```
    DoSomething(testString);
    Debug.Log(testString);
}

void DoSomething(string localString) {
    localString = "World!";
}
```

If we were under the mistaken assumption that strings worked just like other Reference types, then we might be forgiven for assuming that the log output of the following to be `World!`. It appears as though `testString`, a Reference type, is being passed into `DoSomething()`, which would change what `testString` is referencing to, in which case, the log statement will print out the new value of the string.

However, this is not the case, and it will simply print out `Hello`. What is actually happening is that the `localString` variable, within the scope of `DoSomething()`, starts off referencing the same place in memory as `testString` due to the reference being passed by value. This gives us two references pointing to the same location in memory as we would expect if we were dealing with any other Reference type. So far, so good.

However, as soon as we change the value of `localString`, we run into a little bit of a conflict. Strings are immutable, and we cannot change them, so, therefore, we must allocate a new string containing the value `World!` and assign its reference to the value of `localString`; now, the number of references to the `Hello` string returns back to one. The value of `testString`, therefore, has not been changed, and that is still the value that will be printed by `Debug.Log()`. All we've succeeded in doing by calling `DoSomething()` is creating a new string on the heap that gets garbage-collected and doesn't change anything. This is the textbook definition of wasteful.

If we change the method definition of `DoSomething()` to pass the string by reference via the `ref` keyword, the output would indeed change to `World!`. Of course, this is also what we would expect to happen with a Value type, which leads a lot of developers to incorrectly assume that strings are Value types. However, this is an example of the fourth and final data-passing case, where a Reference type is being passed by reference, which allows us to change what the original reference is referencing.

So, let's recap:

- If we pass a Value type by value, we can only change the value of a copy of its data
- If we pass a Value type by reference, we can change the value of the original data passed in
- If we pass a Reference type by value, we can make changes to the object that the original reference is referencing
- If we pass a Reference type by reference, we can change the object or dataset the original reference is referencing

If we find functions that seem to generate a lot of GC allocations the moment they are called, then we might be causing undue heap allocations due to a misunderstanding of the preceding rules.

String concatenation

Concatenation is the act of appending strings to one another to form a larger string. As you've learned, any such cases are likely to result in excess heap allocations. The biggest offender in a string-based memory waste is concatenating strings using the `+` operator and `+=` operators, because of the allocation chaining effect they cause.

For example, the following code tries to combine a group of `string` objects together to print some information about a combat result:

```
void CreateFloatingDamageText(DamageResult result) {  
    string outputText = result.attacker.GetCharacterName() + " dealt " + result.totalDamageDealt.ToString() + " " + result.damageType.ToString() + " damage to " + result.defender.GetCharacterName() + " (" + result.damageBlocked.ToString() + " blocked);  
}  
// ...
```

An example output of this function might be a string that reads as follows:

```
| Dwarf dealt 15 Slashing damage to orc (3 blocked)
```

This function features a handful of string literals (hardcoded strings that are allocated during application initialization) such as `" dealt "`, `" damage to "`, and `" blocked"`, which are simple constructs for the compiler to pre-allocate for us. However, because we are using other local variables within this combined string, it cannot be compiled away at build time, and, therefore, the complete string be regenerated dynamically at runtime each time the function is called.

A new heap allocation will be generated each time a `+`, or `+=`, operator is executed; only a single pair of strings will be merged at a time, and it allocates a new `string` object each time. Then, the result of one merger will be fed into the next and merged with the next `string` and so on until the final `string` object has been built.

So, the previous example will result in nine different strings being allocated

all in one statement. All of the following strings would be allocated to satisfy this instruction, and all would eventually need to be garbage collected (note that the operators are resolved from right to left):

```
"3 blocked"
" (3 blocked)"
"Orc (3 blocked)"
" damage to Orc (3 blocked)"
"Slashing damage to Orc (3 blocked)"
" Slashing damage to Orc (3 blocked)"
"15 Slashing damage to Orc (3 blocked)"
" dealt 15 Slashing damage to Orc (3 blocked)"
"Dwarf dealt 15 Slashing damage to Orc (3 blocked)"
```

That's 262 characters being used, instead of 49. In addition, because a character is a 2-byte data type (for Unicode strings), that's 524-bytes of data being allocated when we only need 98-bytes. Chances are that if this code exists in the code base once, it exists all over the place; so, for an application that's doing a lot of lazy string concatenation like this, that is a ton of memory being wasted on generating unnecessary strings.



Note that big, constant string literals can be safely combined using the + and += operators. The compiler knows that you will eventually need the full string and pregenerates the string automatically. This helps us to make huge block of text more readable within the code base, but only if they will result in a constant string.

Better approaches for generating strings are to use either the `StringBuilder` class or one of several string class methods for string formatting.

StringBuilder

Conventional wisdom says that if we roughly know the final size of the resultant string, then we can allocate an appropriate buffer ahead of time and save ourselves undue allocations. This is the purpose of the `StringBuilder` class. It is effectively a mutable (changeable) string-based object that works like a dynamic array. It allocates a block of space, which we can copy future `string` objects into, and allocates additional space whenever the current size is exceeded. Of course, expanding the buffer should be avoided as much as possible by predicting the maximum size we will need and allocating a sufficiently sized buffer ahead of time.

When we use a `StringBuilder`, we can retrieve the resultant `string` object by calling the `ToString()` method. This still results in a memory allocation for the completed `string`, but, at the very least, we only allocated one large string as opposed to dozens of smaller strings, had we used the `+` or `+=` operators.

For the previous example, we might allocate a `StringBuilder` buffer of 100 characters to make room for long character names and damage values:

```
using System.Text;
// ...
StringBuilder sb = new StringBuilder(100);
sb.Append(result.attacker.GetCharacterName());
sb.Append(" dealt " );
sb.Append(result.totalDamageDealt.ToString());
// etc.
string result = sb.ToString();
```

String formatting

If we don't know the final size of the resultant string, then using a `StringBuilder` class is unlikely to generate a buffer that fits the result size exactly. We will either end up with a buffer that's too large (wasted space), or, worse, a buffer that's too small, which must keep expanding as we generate the complete string. In this scenario, it might be best to use one of the various string class formatting methods.

There are three string class methods available for generating strings: `string.Format()`, `string.Join()`, and `string.Concat()`. Each operates slightly differently, but the overall output is the same. A new `string` object is allocated, containing the contents of the `string` objects we pass into them, and it is all done in a single action, which reduces excess string allocations.



Unfortunately, regardless of the approach we use, if we're converting other objects into additional string objects (such as the calls to generate the strings for "orc", "Dwarf", or "slashing" in the preceding example), then this will allocate an additional string object on the heap. There is nothing we can do about this allocation, except perhaps cache the result so that we don't need to recalculate it each time its needed.

It can be surprisingly hard to say which one of these string generation approaches would be more beneficial in a given situation, as there are a lot of silly little nuances involved that tends to explode into religious debate (just do a Google search for *c sharp string concatenation performance*, and you'll see what I mean), so the simplest approach is to implement one or the other using the conventional wisdom described previously. Whenever we run into bad performance with one of the string-manipulation methods, we should also try the other to check whether it results in a performance improvement. The best way to be certain is to profile them both for comparison, and then pick the best option of the two.

Boxing

Everything in C# is an object (caveats apply), meaning that they derive from the `System.Object` class. Even primitive data types such as `int`, `float`, and `bool` are implicitly derived from `System.Object`, which is itself a Reference type. This is a special case, which allows them access to helper methods such as `Tostring()` so that they can customize their string representation, but without actually turning them into Reference types. Whenever one of these Value types is implicitly treated in such a way that it must act like an object, the CLR automatically creates a temporary object to store, or *box*, the value inside so that it can be treated as a typical Reference type object. As we should expect, this results in a heap allocation to create the containing vessel.



Note that boxing is not the same thing as using Value types as member variables of Reference types. Boxing only takes place when Value types are treated as Reference types via conversion or casting.

For example, the following code will cause the integer variable `i` to be boxed inside object `obj`:

```
| int i = 128;
| object obj = i;
```

The following code will use the object representation `obj` to replace the value stored within the integer, and *unbox* it back into an integer, storing it in `i`. The final value of `i` would be 256:

```
| int i = 128;
| object obj = i;
| obj = 256;
| i = (int)obj; // i = 256
```

These types can be changed dynamically. The following is perfectly legal C# code, where we override the type of `obj`, converting it into a `float`:

```
| int i = 128;
| object obj = i;
| obj = 512f;
```

```
| float f = (float)obj; // f = 512f
```

The following is also legal--conversion into a `bool`:

```
| int i = 128;
| object obj = i;
| obj = false;
| bool b = (bool)obj; // b = false
```

Note that attempting to unbox `obj` into a type that isn't the most recently assigned type would result in an `InvalidOperationException`:

```
| int i = 128;
| object obj = i;
| obj = 512f;
| i = (int)obj; // InvalidOperationException thrown here since most recent conversic
```

All of this can be a little tricky to wrap our head around until we remember that, at the end of the day, everything is just bits in memory and that we are free to interpret them any way we like. After all, data types like `int`, `float`, and so on are just an abstraction over binary lists of `0` and `1`. What's important is knowing that we can treat our primitive types as objects by boxing them, converting their types, and then unboxing them into a different type at a later time, but each time we do this results in a heap memory allocation.



Note that it's possible to convert a boxed object's type using one of the many `System.Convert.To...()` methods.

Boxing can be either implicit, as shown in the preceding examples, or explicit, by typecasting to `System.Object`. Unboxing must always be explicit by typecasting back to its original type. Whenever we pass a Value type into a method that uses `System.Object` as arguments, boxing will be applied implicitly.

Methods such as `String.Format()`, which take `System.Object` as arguments, are one such example. We typically use them by passing in Value types, such as `int`, `float`, and `bool`, to generate a string with. Boxing is automatically taking place in these situations, causing additional heap allocations that we should be aware of. `Collections.Generic.ArrayList` is another such example since `ArrayList` always contains converts its inputs into `System.Object` references,

regardless of what types are stored within.

Any time we use a function definition that takes `System.Object` as arguments, and we're passing in Value types, we should be aware that we're implicitly causing heap allocations due to boxing.

The importance of data layout

The importance of how our data is organized in memory can be surprisingly easy to forget about, but can result in a fairly big performance boost if it is handled properly. Cache misses should be avoided whenever possible, which means that in most cases, arrays of data that are contiguous in memory should be iterated over sequentially as opposed to any other iteration style.

This means that data layout is also important for garbage collection since it is done in an iterative fashion, and if we can find ways to have the GC skip over problematic areas, then we can potentially save a lot of iteration time.

In essence, we want to keep large groups of Reference types separated from large groups of Value types. If there is even one Reference type within a Value type, such as a `struct`, then the GC considers the entire object, and all of its data members, indirectly referenceable objects. When it comes time to `Mark-and-Sweep`, it must verify all fields of the object before moving on. However, if we separate the various types into different arrays, then we can make the GC skip the majority of the data.

For instance, if we have an array of `struct` objects that looks like the following code, then the GC will need to iterate over every member of every `struct`, which could be fairly time-consuming:

```
public struct MyStruct {
    int myInt;
    float myFloat;
    bool myBool;
    string myString;
}
MyStruct[] arrayOfStructs = new MyStruct[1000];
```

However, if we reorganize all pieces of this data into multiple arrays of each type, then the GC will ignore all of the primitive data types and only check the `string` objects. The following code will result in much a faster garbage collection sweep:

```
int[] myInts = new int[1000];
```

```
| float[] myFloats = new float[1000];  
| bool[] myBools = new bool[1000];  
| string[] myStrings = new string[1000];
```

The reason this works is because we're giving the GC fewer indirect references to check. When the data is split into separate arrays (Reference types), it finds three arrays of Value types, marks the arrays, and then immediately moves on because there's no reason to mark the contents of an array of Value types. It must still iterate through all of the `string` objects within `myStrings` since each is a Reference type and it needs to verify that there are no indirect references within it. Technically, the `string` objects cannot contain indirect references, but the GC works at a level where it only knows whether the object is a Reference type or Value type and, therefore, can't tell the difference between a `string` and `class`. However, we have still spared the GC from needing to iterate over an extra 3,000 pieces of data (the 3,000 values in `myInts`, `myFloats`, and `myBools`).

Arrays from the Unity API

There are several instructions within the Unity API that result in heap memory allocations, which we should be aware of. This essentially includes everything that returns an array of data. For example, the following methods allocate memory on the heap:

```
| GetComponents<T>(); // (T[])
| Mesh.vertices; // (Vector3[])
| Camera.allCameras; // (Camera[])
```

Each and every time we call a Unity API method that returns an array, will cause a whole new version of that data to be allocated. Such methods should be avoided whenever possible or at the very least called once and cached so that we don't cause memory allocations more often than necessary.

There are other Unity API calls where we provide an array of elements to a method, and it writes the necessary data into the array for us. One such example is providing a `Particle[]` array to a `ParticleSystem` to get its `Particle` data. The benefit of these types of API calls is that we can avoid reallocating large arrays, whereas the downside is that the array needs to be large enough to fit all of the objects. If the number of objects we need to acquire keeps increasing, then we may find ourselves reallocating larger arrays. In the case of a `ParticleSystem`, we need to be certain we create an array large enough to contain the maximum number of `Particle` objects it generates at any given time.



Unity Technologies have hinted in the past that they may eventually change some of the API calls that return arrays into the form that requires an array to be provided. The API of the latter form can be confusing for new programmers on first glance; however, unlike the first form, it allows responsible programmers to use memory much more efficiently.

Using InstanceIDs for dictionary keys

As mentioned in [Chapter 2](#), *Scripting Strategies*, dictionaries are used to map associations between two different objects, which are very quick at telling us whether a mapping exists, and if so, what that mapping is. It's common practice to map a `MonoBehaviour` or `Scriptable Object` reference as the key of a dictionary, but this causes some problems. When the dictionary element is accessed, it will need to call into several derived methods of `UnityEngine.Object`, which both of these object types derive from. This makes element comparison and mapping acquisition relatively slow.

This can be improved by making use of `Object.GetInstanceID()`, which returns an integer representing a unique identification value for that object that never changes and is never reused between two objects during the entire lifecycle of the application. If we cache this value in the object somehow and use it as the key in our dictionary, then the element comparison will be around two to three times faster than if we used the object reference directly.

However, there are caveats to this approach. If the instance ID value is not cached (we keep calling `Object.GetInstanceID()` each time we need to index into our dictionary) and we are compiling with Mono (and not IL2CPP), then element acquisition could end up being slow. This is because it will call some thread-unsafe code in order to acquire the instance ID, in which case, the Mono compiler cannot optimize the loop, and, therefore, causes some additional overhead by comparison to caching the instance ID value. If we are compiling with IL2CPP, which doesn't have this problem, then the benefits are still not as great (only around 50 percent faster) than if we had simply cached the value beforehand. Therefore, we should aim to cache the integer value in some way so that we avoid having to call `Object.GetInstanceID()` too often.

foreach loops

The `foreach` loop keyword is a bit of a controversial issue in Unity development circles. It turns out that a lot of `foreach` loops implemented in Unity C# code will incur unnecessary heap memory allocations during these calls, as they allocate an `Enumerator` object as a `class` on the heap, instead of a `struct` on the stack. It all depends on the given collection's implementation of the `GetEnumerator()` method.

It turns out that every single collection that has been implemented in the version of Mono that comes with Unity (Mono version 2.6.5) will create classes instead of structs, which results in heap allocations. This includes, but is not limited to, `List<T>`, `LinkedList<T>`, `Dictionary<K,V>`, and `ArrayList`.



Note that it is safe to use foreach loops on typical arrays. The Mono compiler secretly converts foreach over arrays into simple for loops.

The cost is fairly negligible, as the heap allocation cost does not scale with the number of iterations. Only one `Enumerator` object is allocated, and reused over and over again, which only costs a handful of bytes of memory overall. So, unless our `foreach` loops are being invoked for every update (which is typically dangerous in, and of, itself), the costs will be mostly negligible on small projects. The time taken to convert everything to a `for` loop may not be worth the time. However, it's definitely something to keep in mind for the next project we begin to write.

If we're particularly savvy with C#, Visual Studio, and manual compilation of the Mono assembly, then we can have Visual Studio perform code compilation for us and copy the resulting assembly DLL into the `Assets` folder, which will fix this mistake for the generic collections.

Note that performing `foreach` over a `Transform` Component is a typical shortcut to iterating over a `Transform` Component's children. Let's consider the following example:

```
| foreach (Transform child in transform) {  
|     // do stuff with 'child'  
| }
```

However, this results in the same heap allocations mentioned above. As a result, that coding style should be avoided in favor of the following code:

```
| for (int i = 0; i < transform.childCount; ++i) {  
|     Transform child = transform.GetChild(i);  
|     // do stuff with 'child'  
| }
```

Coroutines

As mentioned before, starting a Coroutine costs a small amount of memory to begin with, but note that no further costs are incurred when the method calls `yield`. If memory consumption and garbage collection are significant concerns, we should try to avoid having too many short-lived Coroutines and avoid calling `StartCoroutine()` too much during runtime.

Closures

Closures are useful, but dangerous tools. Anonymous methods and lambda expressions are not always Closures, but they can be. It all depends on whether the method uses data outside of its own scope and parameter list or not.

For example, the following anonymous function would not be a Closure, since it is self-contained and functionally equivalent to any other locally defined function:

```
| System.Func<int,int> anon = (x) => { return x; };  
| int result = anon(5); // result = 5
```

However, if the anonymous function pulled in data from outside itself, it becomes a Closure, as it *closes the environment* around the required data. The following would result in a Closure:

```
| int i = 1024;  
| System.Func<int,int> anon = (x) => { return x + i; };  
| int result = anon(5);
```

In order to complete this transaction, the compiler must define a new custom class that can reference the *environment* where the data value *i* would be accessible. At runtime, it creates the corresponding object on the heap and provides it to the anonymous function. Note that this includes Value types (as per the above example), which were originally on the stack, possibly defeating the purpose of them being allocated on the stack in the first place. So, we should expect each invocation of the second method to result in heap allocations and inevitable garbage collection.

The .NET library functions

The .NET library offers a huge amount of common functionalities that help solve numerous problems that programmers may come across during day-to-day implementation. Most of these classes and functions are optimized for general use cases, which may not be optimal for a specific situation. It may be possible to replace a particular .NET library class with a custom implementation that is more suited to our specific use case.

There are also two big features in the .NET library that often become big performance hogs whenever they're used. This tends to be because they are only included as a quick-and-dirty solution to a given problem without much effort put into optimization. These features are **LINQ** and **Regular Expressions**.

LINQ provides a way to treat arrays of data as miniature databases and perform queries against them using a SQL-like syntax. The simplicity of its coding style and complexity of the underlying system (through its usage of Closures) implies that it has a fairly large overhead cost. LINQ is a handy tool, but is not really intended for high-performance, real-time applications, such as games, and does not even function on platforms that do not support JIT compilation, such as iOS.

Meanwhile, Regular Expressions through the `Regex` class allow us to perform complex string parsing to find substrings that match a particular format, replace pieces of a string, or construct strings from various inputs. Regular Expression is another very useful tool, but tends to be overused in places where it is largely unnecessary or in so-called *clever* ways to implement a feature such as text localization, when straightforward string replacement would be far more efficient.

Specific optimizations for both of these features go far beyond the scope of this book, as they could fill an entire book by themselves. We should either try to minimize their usage as much as possible, replace their usage with something less costly, bring in a LINQ or Regex expert to solve the problem

for us, or do some Googling on the subject to optimize how we're using them.



One of the best ways to find the correct answer online is to simply post the wrong answer. People will either help us out of kindness or will take such a great offense to our implementation that they will consider it their civic duty to correct us. Just be sure to do some kind of research on the subject first. Even the busiest of people are generally happy to help if they can see that we've put in our fair share of effort beforehand.

Temporary work buffers

If we get into the habit of using large, temporary work buffers for one task or another, then it just makes sense that we should look for opportunities to reuse them, instead of reallocating them over and over again, as this lowers the overhead involved in allocation and garbage collection (often called *memory pressure*). It might be worthwhile to extract such functionality from case-specific classes into a generic God Class that contains a big work area for multiple classes to reuse.

Object Pooling

Speaking of temporary work buffers, Object Pooling is an excellent way of both minimizing and establishing control over our memory usage by avoiding deallocation and reallocation. The idea is to formulate our own system for object creation, which hides away whether the object we're getting has been freshly allocated or has been recycled from an earlier allocation. The typical terms to describe this process are to *spawn* and *despawn* the object rather than creating and deleting them in memory. When an object is despawned, we're simply hiding it, making it lay dormant until we need it again, at which point it is *respawned* from one of the previously despawned objects and used in place of an object we might have otherwise newly allocated.

Let's cover a quick implementation of an Object Pooling System.

An important feature of this system is to allow the pooled object to decide how to recycle itself when the time comes. The following Interface Class called `IPoolableObject` will satisfy this requirement nicely:

```
public interface IPoolableObject{
    void New();
    void Respawn();
}
```

This Interface Class defines two methods: `New()` and `Respawn()`. These should be called when the object is first created and when it has been respawned, respectively.

The following `ObjectPool` class definition is a fairly simple implementation of the Object Pooling concept:

```
using System.Collections.Generic;

public class ObjectPool<T> where T : IPoolableObject, new() {
    private Stack<T> _pool;
    private int _currentIndex = 0;

    public ObjectPool(int initialCapacity) {
        _pool = new Stack<T>(initialCapacity);
        for(int i = 0; i < initialCapacity; ++i) {
```

```

        Spawn (); // instantiate a pool of N objects
    }
    Reset ();
}

public int Count {
    get { return _pool.Count; }
}

public void Reset() {
    _currentIndex = 0;
}

public T Spawn() {
    if (_currentIndex < Count) {
        T obj = _pool.Peek ();
        _currentIndex++;
        IPoolableObject po = obj as IPoolableObject;
        po.Respawn();
        return obj;
    } else {
        T obj = new T();
        _pool.Push(obj);
        _currentIndex++;
        IPoolableObject po = obj as IPoolableObject;
        po.New();
        return obj;
    }
}
}

```

This class allows the ObjectPool to be used with any object type so long as it fits the following two criteria: it must implement the `IPoolableObject` Interface Class, and the derived class must allow for a parameter-less constructor (specified by the `new()` keyword in the class declaration).

An example poolable object would look like so: it must implement two `public` methods, `New()` and `Respawn()`, which are invoked by the `ObjectPool` class at the appropriate times:

```

public class TestObject : IPoolableObject {
    public void New() {
        // very first initialization here
    }
    public void Respawn() {
        // reset data which allows the object to be recycled here
    }
}

```

Finally, consider an example usage to create a pool of 100 `TestObject` objects:

```
| private ObjectPool<TestObject> _objectPool = new ObjectPool<TestObject>(100);
```

The first 100 calls to `spawn()` on the `_objectPool` will cause the objects to be respawned, providing the caller with a unique instance of the object each time. If there are no more objects to provide (we have called `spawn()` more than 100 times), then we will allocate a new `TestObject` object and push it onto the Stack. Finally, if `Reset()` is called on the `_objectPool`, it will begin again from the start, recycling objects and providing them to the caller.

Note that we are using the `Peek()` method on the `stack` object so that we don't remove the old instance from the Stack. We want the `ObjectPool` to maintain references to all of the objects we create.

Also, note that this Pooling solution will not work for classes we haven't defined and cannot derive from `IPoolableObject`, such as `Vector3` and `Quaternion`. This is normally dictated by the `sealed` keyword in the class definition. In these cases, we would need to define a containing class:

```
public class PoolableVector3 : IPoolableObject {
    public Vector3 vector = new Vector3();
    public void New() {
        Reset();
    }
    public void Respawn() {
        Reset();
    }
    public void Reset() {
        vector.x = vector.y = vector.z = 0f;
    }
}
```

We could extend this system in a number of ways, such as defining a `Despawn()` method to handle destruction of the object, making use of the `IDisposable` Interface Class and using blocks when we wish to automatically spawn and despawn objects within a small scope, and/or allowing objects instantiated outside the pool to be added to it.

Prefab Pooling

The previous Pooling solution is useful for typical C# objects, but it won't work for specialized Unity objects, such as `GameObject` and `MonoBehaviour`. These objects tend to consume a large chunk of our runtime memory, can cost us a great deal of CPU usage when they're created and destroyed, and tend to risk a large amount of garbage collection at runtime. For instance, during the lifecycle of a small RPG game, we might spawn a thousand Orc creatures, but at any given moment, we may only need at a maximum, 10 of them. It would be nice if we could perform similar Pooling as before, but, for Unity Prefabs, to save on a lot of unnecessary overhead creating and destroying 990 Orcs we don't need.

Our goal is to push the overwhelming majority of object instantiation to Scene initialization rather than letting them get created at runtime. This can provide some big runtime CPU savings and avoids a lot of spikes caused by object creation/destruction and garbage collection at the expense of Scene loading times and runtime memory consumption. As a result, there are quite a few Pooling solutions available on the Asset Store for handling this task, with varying degrees of simplicity, quality, and feature sets.



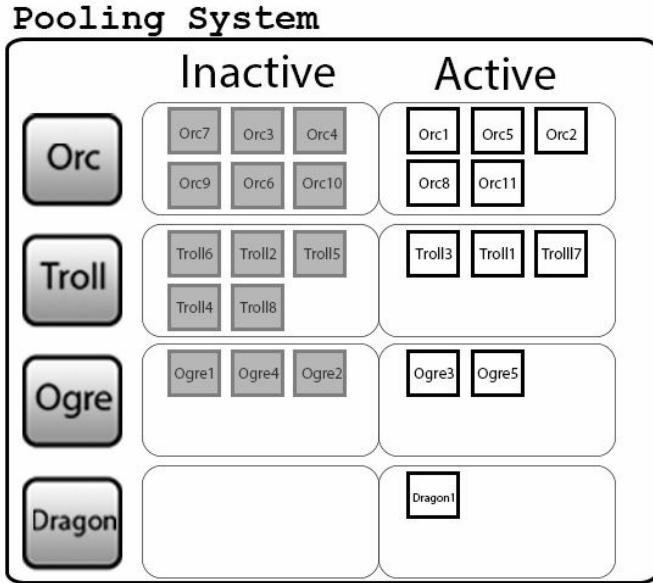
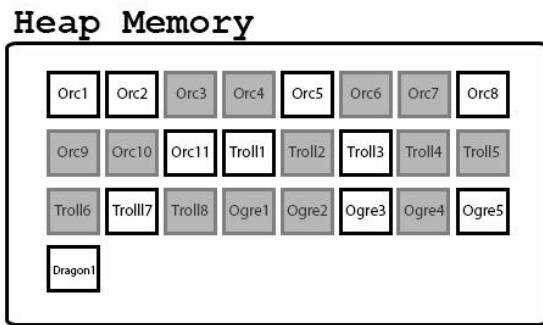
It is often recommended that Pooling should be implemented in any game that intends to be deployed on mobile devices, due to the greater overhead costs involved in the allocation and deallocation of memory compared to desktop applications.

However, creating a Pooling solution is an interesting topic, and building one from scratch is a great way of getting to grips with a lot of important internal Unity Engine behavior. Also, knowing how such a system is built makes it easier to extend if we wish it to meet the needs of our particular game, rather than relying on a prebuilt solution.

The general idea of Prefab Pooling is to create a system that contains lists of active and inactive `GameObjects` that were all instantiated from the same Prefab reference. The following diagram shows how the system might look after

several spawns, despawns, and respawns of various objects derived from four different Prefabs (**Orc**, **Troll**, **Ogre**, and **Dragon**):

11 Orcs (5 active, 6 inactive)
8 Trolls (3 active, 5 inactive)
5 Ogres (2 active, 3 inactive)
1 Dragon (1 active)



Note that the Heap Memory area in the previous screenshot represents the objects as they exist in memory, while the Pooling System area represents the Pooling System's references to those objects.

In this example, several instances of each Prefab were instantiated (11 Orcs, 8 Trolls, 5 Ogres, and 1 Dragon). Currently, only 11 of these objects are active, while the other 14 have been previously despawned and are inactive. Note that the despawned objects still exist in memory, although they are not visible and cannot interact with the game world until they have been respawned.

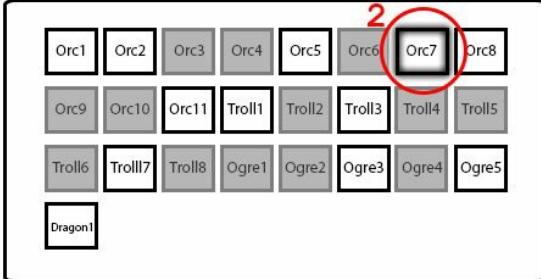
Naturally, this costs us a constant amount of heap memory at runtime in order to maintain the inactive objects, but when a new object is instantiated, we can reuse one of the existing inactive objects rather than allocating more memory in order to satisfy the request. This saves significant runtime CPU costs during object creation and destruction and avoids garbage collection.

The following diagram shows the chain of events that needs to occur when a new Orc is spawned:

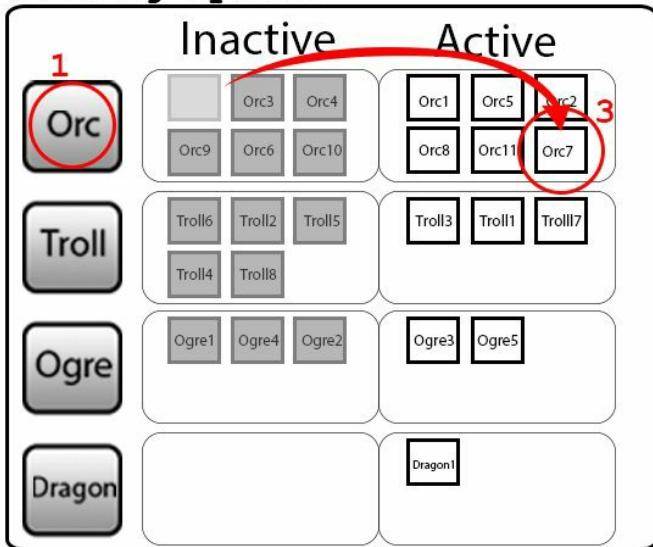
New Orc is spawned

- Determine which Pool corresponds to the given Prefab
- The first inactive Orc in the Inactive Group (Orc7) is activated - the corresponding object in the Heap is therefore activated
- Newly-spawned Orc is moved to Active group

Heap Memory



Pooling System



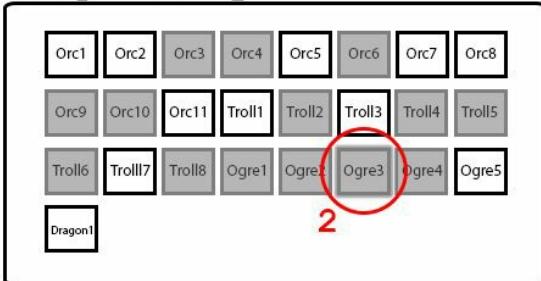
The first object in the **Inactive** Orc pool (**Orc7**) is reactivated and moved into the **Active** pool. We now have six active Orcs and five inactive Orcs.

The following figure shows the order of events when an **Ogre** object is despawned:

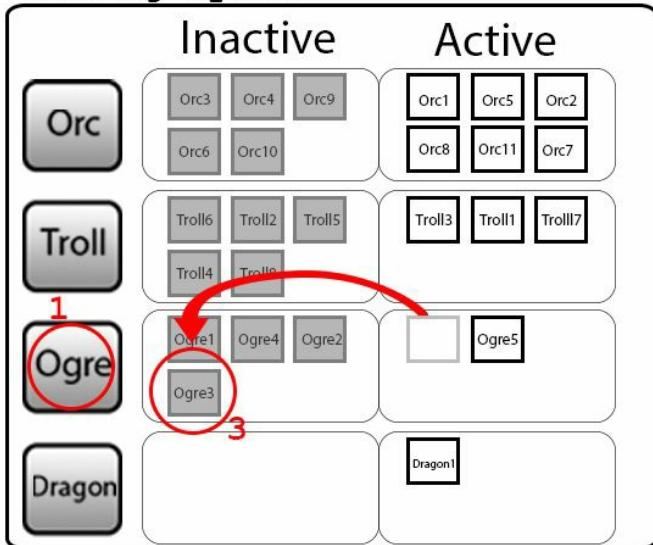
Ogre3 is despawned

- Determine which Pool corresponds to the given Object
- Deactivate Ogre3 - the corresponding object in the Heap is therefore deactivated
- Move Ogre3 to Inactive group

Heap Memory



Pooling System



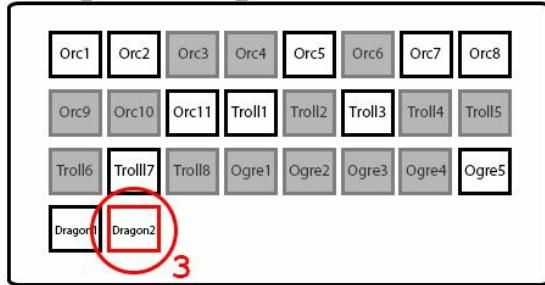
This time, the object is deactivated and moved from the Active pool into the Inactive pool, leaving us with one active Ogre and four inactive Ogres.

Finally, the following diagram shows what happens when a new object is spawned, but there are no inactive objects to satisfy the request:

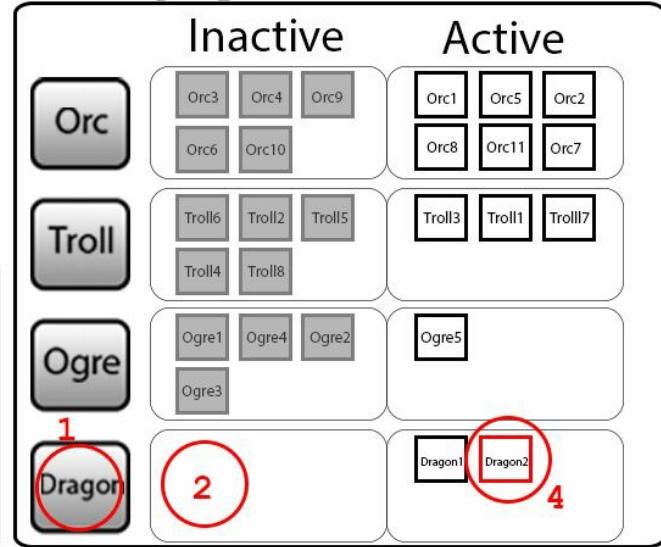
New Dragon is spawned

1. Determine which Pool corresponds to the given Prefab
2. Inactive group is empty, so a new Instance of Dragon must be created
3. Instantiate a new Dragon from the Prefab on the heap
4. Add the newly -created Dragon to the Active list

Heap Memory



Pooling System



In this scenario, more memory must be allocated to instantiate the new Dragon object since there are no Dragon objects in its Inactive pool to reuse. Therefore, in order to avoid runtime memory allocations for our `GameObject`s, it is critical that we know beforehand how many we will need and that there is sufficient memory space available to contain them all at once. This will vary depending on the type of object in question and requires occasional testing and sanity-checking to ensure that we have a sensible number of each Prefab instantiated at runtime.

With all of this in mind, let's create a Pooling System for Prefabs.

Poolable Components

Let's first define an `interface` class for a Component that can be used in the Pooling System:

```
| public interface IPoolableComponent {  
|     void Spawned();  
|     void Despawned();  
| }
```

The approach for `IPoolableComponent` will be very different from the approach taken for `IPoolableObject`. The objects being created this time are `GameObjects`, which are a lot trickier to work with than standard objects because of how much of their runtime behavior is already handled through the Unity Engine and how little low-level access we have to it.

`GameObjects` do not give us access to an equivalent `New()` method that we can invoke any time the object is created, and we cannot derive from the `GameObject` class in order to implement one. `GameObjects` are created either by placing them in a Scene or by instantiating them at runtime through `GameObject.Instantiate()`, and the only inputs we can apply are an initial position and rotation. Of course, their Components have an `Awake()` callback that we can define, which is invoked the first time the Component is brought to life, but this is merely a compositional object--it's not the actual parent object we're spawning and despawning.

So, because we have control on only a `GameObject` class's Components, it is assumed that the `IPoolableComponent` Interface Class is implemented by at least one of the Components that is attached to the `GameObject` we wish to pool.

The `Spawned()` method should be invoked on every implementing Component each time the pooled `GameObject` is respawned, while the `Despawned()` method gets invoked whenever it is despawned. This gives us entry points to control the data variables and behavior during the creation and destruction of the parent `GameObject`.

The act of despawning a `GameObject` is trivial: turn its `active` flag to `false`

through `SetActive()`. This disables the `collider` and `Rigidbody` for physics calculations, removes it from the list of renderable objects, and essentially takes care of disabling all interactions with all built-in Unity Engine subsystems in a single stroke. The only exception is any Coroutines that are currently invoking on the object, since as you learned in [Chapter 2, Scripting Strategies](#), Coroutines are invoked independently of any `update()` and `GameObject` activity. We will, therefore, need to call `StopCoroutine()` or `StopAllCoroutines()` during the despawning of such objects.

In addition, Components typically hook into our own custom gameplay subsystems as well, so the `Despawn()` method gives our Components the opportunity to take care of any custom cleanup before shutting down. For example, we would probably want to use `Despawn()` to deregister the Component from the Messaging System we defined in [Chapter 2, Scripting Strategies](#).

Unfortunately, successfully respawning the `GameObject` is a lot more complicated. When we respawn an object, there will be many settings that were left behind when the object was previously active, and these must be reset in order to avoid conflicting behaviors. A common problem with this is the `Rigidbody`'s `linearVelocity` and `angularVelocity` properties. If these values are not explicitly reset before the object is reactivated, then the newly respawned object will continue moving with the same velocity the old version had when it was despawned.

This problem becomes further complicated by the fact that built-in Components are `sealed`, which means that they cannot be derived from. So, to avoid these issues, we can create a custom Component that resets the attached `Rigidbody` whenever the object is despawned:

```
public class ResetPooledRigidbodyComponent : MonoBehaviour, IPoolableComponent
[SerializeField] Rigidbody _body;
public void Spawed() { }
public void Despawed() {
    if (_body == null) {
        _body = GetComponent<Rigidbody>();
        if (_body == null) {
            // no Rigidbody!
            return;
        }
    }
    _body.velocity = Vector3.zero;
```

```
| } } _body.angularVelocity = Vector3.zero;  
| }
```

Note that the best place to perform the cleanup task is during despawning, because we cannot be certain in what order the `GameObject` class's `IPoolableComponent` Interface Classes will have their `Spawned()` methods invoked. It is unlikely that another `IPoolableComponent` will change the object's velocity during despawning, but it is possible that a different `IPoolableComponent` attached to the same object might want to set the Rigidbody's initial velocity to some important value during its own `Spawned()` method. Ergo, performing the velocity reset during the `ResetPooledRigidbodyComponent` class's `spawned()` method could potentially conflict with other Components and cause some very confusing bugs.



In fact, creating Poolable Components that are not self-contained and tend to tinker with other Components like this, is one of the biggest dangers of implementing a Pooling System. We should minimize such design and routinely verify them when we're trying to debug strange issues in our game.

For the sake of illustration, here is the definition of a simple poolable Component making use of the `MessagingSystem` class we defined in [Chapter 2, Scripting Strategies](#). This Component automatically handles some basic tasks every time the object is spawned and despawned:

```
public class PoolableTestMessageListener : MonoBehaviour, IPoolableComponent  
public void Spawned() {  
    MessagingSystem.Instance.AttachListener(typeof(MyCustomMessage),  
                                            this.HandleMyCustomMessage);  
}  
  
bool HandleMyCustomMessage(BaseMessage msg) {  
    MyCustomMessage castMsg = msg as MyCustomMessage;  
    Debug.Log (string.Format("Got the message! {0}, {1}",  
                            castMsg._intValue,  
                            castMsg._floatValue));  
    return true;  
}  
  
public void Despawned() {  
    if (MessagingSystem.IsAlive) {  
        MessagingSystem.Instance.DetachListener(typeof(MyCustomMessage),  
                                                this.HandleMyCustomMessage);  
    }  
}
```

| }

The Prefab Pooling System

Hopefully, we now have an understanding of what we need from our Pooling System, so all that's left is to implement it. The requirements are as follows:

- It must accept requests to spawn a `GameObject` from a Prefab, an initial position, and an initial rotation:
 - If a despawned version already exists, it should respawn the first available one
 - If it does not exist, then it should instantiate a new `GameObject` from the Prefab
 - In either case, the `spawned()` method should be invoked on all `IPoolableComponent` Interface Classes attached to the `GameObject`
- It must accept requests to despawn a specific `GameObject`:
 - If the object is managed by the Pooling System, it should deactivate it and call the `Despawned()` method on all `IPoolableComponent` Interface Classes attached to the `GameObject`
 - If the object is not managed by the Pooling System, it should send an error

The requirements are fairly straightforward, but the implementation requires some investigation if we wish to make the solution performance-friendly. Firstly, a typical Singleton would be a good choice for the main entry point since we want this system to be globally accessible from anywhere:

```
| public static class PrefabPoolingSystem {}
```

The main task for object spawning involves accepting a Prefab reference and figuring whether we have any despawned `GameObjects` that were originally instantiated from the same reference. To do this, we will essentially want our Pooling System to keep track of two different lists for any given Prefab reference: a list of active (spawned) `GameObjects`, and a list of inactive (despawned) objects that were instantiated from it. This information would be best abstracted into a separate class, which we will name `PrefabPool`.

In order to maximize the performance of this system (and hence make the

largest gains possible, relative to just allocating and deallocating objects from memory all of the time), we will want to use some fast data structures in order to acquire the corresponding `PrefabPool` objects whenever a spawn or despawn request comes in.

Since spawning involves being given a Prefab, we will want a data structure that can quickly map Prefabs to the `PrefabPool` that manages them. Also, since despawning involves being given a `GameObject`, we will want another data structure that can quickly map spawned `GameObjects` to the `PrefabPool` that originally spawned them. A pair of dictionaries would be a good choice for both of these needs.

Let's define these dictionaries in our `PrefabPoolingSystem` class:

```
public static class PrefabPoolingSystem {
    static Dictionary<GameObject, PrefabPool> _prefabToPoolMap = new Dictionary<
    static Dictionary<GameObject, PrefabPool> _goToPoolMap = new Dictionary<Game
```

Next, we'll define what happens when we spawn an object:

```
public static GameObject Spawn(GameObject prefab, Vector3 position, Quaternion
    if (!_prefabToPoolMap.ContainsKey (prefab)) {
        _prefabToPoolMap.Add (prefab, new PrefabPool());
    }
    PrefabPool pool = _prefabToPoolMap[prefab];
    GameObject go = pool.Spawn(prefab, position, rotation);
    _goToPoolMap.Add (go, pool);
    return go;
}
```

The `Spawn()` method will be given a Prefab reference, an initial position, and an initial rotation. We need to figure out which `PrefabPool` the Prefab belongs to (if any), ask it to spawn a new `GameObject` using the data provided, and then return the spawned object to the requestor. We will first check our *Prefab-to-Pool* map to check whether a pool already exists for this Prefab. If not, we immediately create one. In either case, we then ask the `PrefabPool` to spawn us a new object. The `PrefabPool` will either end up respawning an object that was despawned earlier or instantiate a new one (if there aren't any inactive instances left).

This class doesn't particularly care how the `PrefabPool` creates the object. It just wants the instance generated by the `PrefabPool` class so that it can be

entered into the *GameObject-to-Pool* map and returned to the requestor.

For convenience, we can also define an overload that places the object at the world's center. This is useful for `GameObjects` that aren't visible and just need to exist in the Scene:

```
public static GameObject Spawn(GameObject prefab) {  
    return Spawn (prefab, Vector3.zero, Quaternion.identity);  
}
```



Note that no actual spawning and despawning are taking place, yet. This task will eventually be implemented within the `PrefabPool` class.

Despawning involves being given a `GameObject` and then figuring out which `PrefabPool` is managing it. This could be achieved by iterating through our `PrefabPool` objects and checking whether they contain the given `GameObject`. However, if we end up generating a lot of Prefab Pools, then this iterative process can take a while. We will always end up with as many `PrefabPool` objects as we have Prefabs (at least, so long as we manage all of them through the Pooling System). Most projects tend to have dozens, hundreds, if not thousands, of different Prefabs.

So, the *GameObject-to-Pool* map is maintained to ensure that we always have rapid access to the `PrefabPool` that originally spawned the object. It can also be used to quickly check whether the given `GameObject` is even managed by the Pooling System to begin with. Here is the method definition for the despawning method, which takes care of these tasks:

```
public static bool Despawn(GameObject obj) {  
    if (!_goToPoolMap.ContainsKey(obj)) {  
        Debug.LogError (string.Format ("Object {0} not managed by pool system!",  
            obj));  
        return false;  
    }  
  
    PrefabPool pool = _goToPoolMap[obj];  
    if (pool.Despawn (obj)) {  
        _goToPoolMap.Remove (obj);  
        return true;  
    }  
    return false;  
}
```



Note that the `Despawn()` method of both `PrefabPoolingSystem` and



PrefabPool returns a Boolean that can be used to check whether the object was successfully despawned.

As a result, thanks to the two maps we're maintaining, we can quickly access the `PrefabPool` that manages the given reference, and this solution will scale for any number of Prefab that the system manages.

Prefab pools

Now that we have a system that can handle multiple Prefab pools automatically, the only thing left is to define the behavior of the pools. As mentioned previously, we will want the `PrefabPool` class to maintain two data structures: one for active (spawned) objects that have been instantiated from the given Prefab and another for inactive (despawned) objects.

Technically, the `PrefabPoolingSystem` class already maintains a map of which Prefab is governed by which `PrefabPool`, so we can actually save a little memory by making the `PrefabPool` dependent upon the `PrefabPoolingSystem` class to give it the reference to the Prefab it is managing. Consequently, the two data structures would be the only member variables the `PrefabPool` needs to keep track of.

However, for each spawned `GameObject`, it must also maintain a list of all of its `IPoolableComponent` references in order to invoke the `Spawned()` and `Despawned()` methods on them. Acquiring these references can be a costly operation to perform at runtime, so it would be best to cache the data in a simple struct:

```
| public struct PoolablePrefabData {  
|     public GameObject go;  
|     public IPoolableComponent[] poolableComponents;  
| }
```

This struct will contain a reference to the `GameObject` and the precached list of all of its `IPoolableComponent` Components.

Now, we can define the member data of our `PrefabPool` class:

```
| public class PrefabPool {  
|     Dictionary<GameObject, PoolablePrefabData> _activeList = new Dictionary<Game  
|         Queue<PoolablePrefabData> _inactiveList = new Queue<PoolablePrefabData>();  
| }
```

The data structure for the active list should be a dictionary in order to do a quick lookup for the corresponding `PoolablePrefabData` from any given `GameObject` reference. This will be useful during object despawning.

Meanwhile, the inactive data structure is defined as a `Queue`, but it will work equally well as a `List`, `Stack`, or really any data structure that needs to regularly expand or contract, where we only need to pop items from one end of the group, since it does not matter which object it is. It only matters that we retrieve one of them. A `Queue` is useful in this case because we can both retrieve and remove the object from the data structure in a single call to `Dequeue()`.

Object spawning

Let's define what it means to spawn a `GameObject` in the context of our Pooling System: at some point, `PrefabPool` will get a request to spawn a `GameObject` from a given Prefab, at a particular position and rotation. The first thing we should check is whether or not we have any inactive instances of the Prefab. If so, then we can dequeue the next available one from the `queue` and respawn it. If not, then we need to instantiate a new `GameObject` from the Prefab using `GameObject.Instantiate()`. At this moment, we should also create a `PoolablePrefabData` object to store the `GameObject` reference and acquire the list of all `MonoBehaviours` that implement `IPoolableComponent` that are attached to it.

Either way, we can now activate the `GameObject`, set its position and rotation, and call the `Spawned()` method on all of its `IPoolableComponent` references. Once the object has been respawned, we can add it to the list of active objects and return it to the requestor.

The following is the definition of the `Spawn()` method that defines this behavior:

```
public GameObject Spawn(GameObject prefab, Vector3 position, Quaternion rotation)
    PoolablePrefabData data;

    if (_inactiveList.Count > 0) {
        data = _inactiveList.Dequeue();
    } else {
        // instantiate a new object
        GameObject newGO = GameObject.Instantiate(prefab, position, rotation) as
        data = new PoolablePrefabData();
        data.go = newGO;
        data.poolableComponents = newGO.GetComponents<IPoolableComponent>();
    }

    data.go.SetActive (true);
    data.go.transform.position = position;
    data.go.transform.rotation = rotation;

    for(int i = 0; i < data.poolableComponents.Length; ++i) {
        data.poolableComponents[i].Spawned ();
    }
    _activeList.Add (data.go, data);

    return data.go;
```

| }

Instance prespawning

Since we are using `GameObject.Instantiate()` whenever the `PrefabPool` has run out of despawned instances, this system does not completely rid us of runtime object instantiation, hence heap memory allocation. It's important to prespawn the expected number of instances that we will need during the lifetime of the current Scene so that we minimize or remove the need to instantiate more during runtime.

Note that we shouldn't prespawn too many objects. It would be wasteful to prespawn 100 explosion Particle Effects, if the most we will ever expect to see in the Scene at any given time is three or four. Conversely, spawning too few instances will cause excessive runtime memory allocations, and the goal of this system is to push the majority of allocation to the start of a Scene's lifetime. We need to be careful about how many instances we maintain in memory so that we don't waste more memory space than necessary.

Let's define a method in our `PrefabPoolingSystem` class that we can use to quickly prespawn a given number of objects from a Prefab. This essentially involves spawning `N` objects and then immediately despawning them all:

```
public static void Prespawn(GameObject prefab, int numToSpawn) {
    List<GameObject> spawnedObjects = new List<GameObject>();

    for(int i = 0; i < numToSpawn; i++) {
        spawnedObjects.Add (Spawn (prefab));
    }

    for(int i = 0; i < numToSpawn; i++) {
        Despawn(spawnedObjects[i]);
    }

    spawnedObjects.Clear ();
}
```

We would use this method during Scene initialization to prespawn a collection of objects to use in the level. Take for example the following code:

```
public class OrcPreSpawner : MonoBehaviour
    [SerializeField] GameObject _orcPrefab;
    [SerializeField] int _numToSpawn = 20;
```

```
| void Start() {  
|     PrefabPoolingSystem.Prespawn(_orcPrefab, _numToSpawn);  
| }  
| }
```

Object despawning

Finally, there is the act of despawning the objects. As mentioned previously, this primarily involves deactivating the object, but we also need to take care of various bookkeeping tasks and invoking `Despawned()` on all of its `IPoolableComponent` references.

Here is the method definition for `PrefabPool.Despawn()`:

```
public bool Despawn(GameObject objToDespawn) {
    if (!_activeList.ContainsKey(objToDespawn)) {
        Debug.LogError ("This Object is not managed by this object pool!");
        return false;
    }

    PoolablePrefabData data = _activeList[objToDespawn];

    for(int i = 0; i < data.poolableComponents.Length; ++i) {
        data.poolableComponents[i].Despawned ();
    }

    data.go.SetActive (false);
    _activeList.Remove (objToDespawn);
    _inactiveList.Enqueue(data);
    return true;
}
```

First, we verify that the object is being managed by the pool and then we grab the corresponding `PoolablePrefabData` in order to access the list of `IPoolableComponent` references. Once `Despawned()` is invoked on all of them, we deactivate the object, remove it from the active list, and push it into the inactive queue so that it can be respawned later.

Prefab pool testing

The following class definition allows us to perform a simple hands-on test with the `PrefabPoolingSystem` class; it will support three Prefabs, and prespawn five instances of each during application initialization. We can press the `1`, `2`, `3`, or `4` keys to spawn an instance of each type, and then press `Q`, `W`, `E` and `R` to despawn a random instance of each type, respectively:

```
public class PrefabPoolingTestInput : MonoBehaviour {
    [SerializeField] GameObject _orcPrefab;
    [SerializeField] GameObject _trollPrefab;
    [SerializeField] GameObject _ogrePrefab;
    [SerializeField] GameObject _dragonPrefab;

    List<GameObject> _orcs = new List<GameObject>();
    List<GameObject> _trolls = new List<GameObject>();
    List<GameObject> _ogres = new List<GameObject>();
    List<GameObject> _dragons = new List<GameObject>();

    void Start() {
        PrefabPoolingSystem.Prespawn(_orcPrefab, 11);
        PrefabPoolingSystem.Prespawn(_trollPrefab, 8);
        PrefabPoolingSystem.Prespawn(_ogrePrefab, 5);
        PrefabPoolingSystem.Prespawn(_dragonPrefab, 1);
    }

    void Update () {
        if (Input.GetKeyDown(KeyCode.Alpha1)) {SpawnObject(_orcPrefab, _orcs);}
        if (Input.GetKeyDown(KeyCode.Alpha2)) {SpawnObject(_trollPrefab, _trolls)}
        if (Input.GetKeyDown(KeyCode.Alpha3)) {SpawnObject(_ogrePrefab, _ogres);}
        if (Input.GetKeyDown(KeyCode.Alpha4)) {SpawnObject(_dragonPrefab, _dragons)}
        if (Input.GetKeyDown(KeyCode.Q)) { DespawnRandomObject(_orcs); }
        if (Input.GetKeyDown(KeyCode.W)) { DespawnRandomObject(_trolls); }
        if (Input.GetKeyDown(KeyCode.E)) { DespawnRandomObject(_ogres); }
        if (Input.GetKeyDown(KeyCode.R)) { DespawnRandomObject(_dragons); }
    }

    void SpawnObject(GameObject prefab, List<GameObject> list) {
        GameObject obj = PrefabPoolingSystem.Spawn (prefab,
                                                    5.0f * Random.insideUnitSphere
                                                    Quaternion.identity);
        list.Add (obj);
    }

    void DespawnRandomObject(List<GameObject> list) {
        if (list.Count == 0) {
            // Nothing to despawn
            return;
        }
    }
}
```

```
|     int i = Random.Range (0, list.Count);  
|     PrefabPoolingSystem.Despawn(list[i]);  
|     list.RemoveAt(i);  
| } }
```

Once we spawn more than five instances of any of the Prefabs, it will need to instantiate a new one in memory, costing us some memory allocation. However, if we observe the Memory Area in the Profiler window, while we only spawn and despawn instances that already exist, then we will notice that absolutely no new allocations take place.

Prefab Pooling and Scene loading

There is one important caveat to this system that has not yet been mentioned: the `PrefabPoolingSystem` class will outlast Scene lifetime since it is a Static Class. This means that when a new Scene is loaded, the Pooling System's dictionaries will attempt to maintain references to any pooled instances from the previous Scene, but Unity forcibly destroys these objects regardless of the fact that we are still keeping references to them (unless they were set to `DontDestroyOnLoad()`), and so the dictionaries will be full of `null` references. This would cause some serious problems for the next Scene.

We should, therefore, create a method in `PrefabPoolingSystem` that resets the Pooling System in preparation for this likely event. The following method should be called before a new Scene is loaded so that it is ready for any early calls to `Prespawn()` in the next Scene:

```
public static void Reset() {
    _prefabToPoolMap.Clear ();
    _goToPoolMap.Clear ();
}
```

Note that if we also invoke a garbage collection during Scene transitions, there's no need to explicitly destroy the `PrefabPool` objects these dictionaries were referencing. Since these were the only references to the `PrefabPool` objects, they will be deallocated during the next garbage collection. If we aren't invoking garbage collection between Scenes, then the `PrefabPool` and `PooledPrefabData` objects will remain in memory until that time.

Prefab Pooling summary

This Pooling System provides a decent solution to the problem of runtime memory allocations for `GameObjects` and Prefabs, but, as a quick reminder, we need to be aware of the following caveats:

- We need to be careful about properly resetting important data in respawned objects (such as `Rigidbody` velocity)
- We must ensure that we don't prespawn too few, or too many, instances of a Prefab
- We should be careful of the order of execution of `Spawned()` and `Despawned()` methods on `IPoolableComponent` and not assume that they will be called in a particular order
- We must call `Reset()` on `PrefabPoolingSystem` when loading a new Scene in order to clear any `null` references to objects, which may no longer exist

There are several other features that we could implement. These will be left as academic exercises if we wish to extend this system in the future:

- Any `IPoolableComponent` added to the `GameObject` after initialization will not have their `Spawned()` or `Despawned()` methods invoked since we only collect this list when the `GameObject` is first instantiated. We could fix this by changing `PrefabPool` to keep acquiring `IPoolableComponent` references every time `Spawned()` and `Despawned()` are invoked at the cost of additional overhead during spawning/despawning.
- Any `IPoolableComponent` attached to children of the Prefab's root will also not be counted. This could be fixed by changing `PrefabPool` to use `GetComponentsInChildren<T>` at the cost of additional overhead if we're using Prefabs with deep hierarchies.
- Prefab instances that already exist in the Scene will not be managed by the Pooling System. We could create a Component that needs to be attached to such objects and that notifies the `PrefabPoolingSystem` class of its existence in its `Awake()` callback, which passes the reference along to the corresponding `PrefabPool`.
- We could implement a way for `IPoolableComponent` to set a priority during

acquisition and directly control the order of execution for their `Spawned()` and `Despawned()` methods.

- We could add counters that keep track of how long objects have been sitting in the Inactive list relative to total Scene lifetime and print out the data during shutdown. This could tell us whether or not we're prespawning too many instances of a given Prefab.
- This system will not interact kindly with Prefab instances that set themselves to `DontDestroyOnLoad()`. It might be wise to add a Boolean to every `Spawn()` call to say whether the object should persist or not and keep them in a separate data structure that is not cleared out during `Reset()`.
- We could change `Spawn()` to accept an argument that allows the requestor to pass custom data to the `Spawned()` function of `IPoolableObject` for initialization purposes. This could use a system similar to how custom message objects were derived from the `Message` class for our Messaging System in [Chapter 2, Scripting Strategies](#).

IL2CPP optimizations

Unity Technologies have released a few blog posts on interesting ways to improve the performance of IL2CPP in some circumstances, but they can be difficult to manage. If you're using IL2CPP and need to eke out the last little bit of performance from our application that we can, then check out the blog series at the following links:

- <https://blogs.unity3d.com/2016/07/26/il2cpp-optimizations-devirtualization/>
- <https://blogs.unity3d.com/2016/08/04/il2cpp-optimizations-faster-virtual-method-calls/>
- <https://blogs.unity3d.com/2016/08/11/il2cpp-optimizations-avoid-boxing/>

WebGL optimizations

Unity Technologies have also released a number of blog posts covering WebGL applications, which includes some crucial information about memory management that all WebGL developers should know. These can be found at the following links:

- <https://blogs.unity3d.com/2016/09/20/understanding-memory-in-unity-webgl/>
- <https://blogs.unity3d.com/2016/12/05/unity-webgl-memory-the-unity-heap/>

The future of Unity, Mono, and IL2CPP

When the first edition of this book was written, IL2CPP was being teased in a variety of blog posts, and Unity was still running a very old version of Mono. Well, the good news is that IL2CPP has finally landed, but the bad news is Unity is *still* running a very old version of Mono. This is an unfortunate consequence of how various elements of the Mono Framework are licensed, meaning that Unity Technologies has only been able to update Mono on an infrequent basis.

After many years, Unity is still only able to make use of the .NET 3.5 library functionality, which is about a decade old at this point. This limits what kinds of .NET library classes we can use, restricts certain C# language features, and limits performance since many enhancements would have been made to the library during this time. However, to the great relief of many Unity developers, an experimental upgrade to a more recent version of Mono has been in beta for a while and is now available in Unity 2017 by switching the Edit | Project Settings | Player | Configuration | Scripting Runtime Version option to Experimental (.NET 4.6 Equivalent). This setting upgrades the Mono runtime and allows us to make use of the .NET 4.6 functionality, which is only a couple of years old at the time of writing.

This feature is still experimental, however. This could lead to obscure bugs and crashes until it's been deemed stable. At which point, the Mono upgrade will become permanent so that all users can enjoy it. For the time-being, we should only use this option if we need it for a vital bug fix or feature that we cannot work around. Also, there is never a guarantee that the version fixing your bug isn't going to introduce an even worse bug. However, at the very least, hope is on the horizon that Unity is serious about making the upgrade happen sooner rather than later, as there are regular bug fixes around this area each update.

Meanwhile, Unity's rollout of IL2CPP has been rapid and impressive. It

currently supports a large set of modern platforms and offers some decent performance improvements. Of course, again, it is slightly unstable, and we are likely to find bugs if we use the latest and greatest beta versions of Unity.

Where things will go from here is pretty predictable; since Unity is reliant on IL2CPP for certain platforms, and others gain benefit from it, IL2CPP will continue to gain more support, more fixes, and offer more performance improvements with each passing release, and the Mono upgrade will finally land at some point in the future.

Due to the advancement of time and technology, some features of Unity have become deprecated. A few examples that have been mentioned during this book include Animators (at least when used with UI elements), the Sprite Packer tool, and the UnityScript and Boo languages (which aren't completely gone yet, but it would be wise to start learning C# scripting sooner rather than later). Other features are also planned to be dropped. The most imminent among which is dropping of all support for DirectX 9 for Windows in Unity 2017.3 (check out <https://blogs.unity3d.com/2017/07/10/deprecating-directx-9/> for more information). We can expect further deprecated features as development of Unity 2017 progresses.

Of course, there are a lot of new features as well. The bulk of the features new in Unity 2017 primarily focus on helper utilities such as Timeline, Cinemachine, the Post-Processing stack, Collaborate, and built-in Analytics. These tools are helpful to game developers, but also make Unity a little more approachable to creatives among nongaming industries, so it's a safe bet that Unity Technologies is trying to entice a broader audience, so we can expect more on that front. Of course, game developers are still reaping the benefits of Unity Technologies' hard work. Since the first edition of this book around 2 years ago, Unity Technologies have released a whole new User Interface system (including the acquisition of Text Mesh Pro); the Progressive Lightmapper along with Global Illumination, Vulkan, and Metal support; huge improvements to the WebGL platform, 2D games, Particle Systems, Physics (both 3D and 2D), and Video playback; and, of course, a massive amount of support for the XR platforms and the more popular HMDs and their controllers.

Naturally, there is even more to come in the form of feature improvements

(Particles, Animation, and 2D), expanded APIs (Physics and Networking), performance enhancements, platform support (WebVR, 360 Video playback, and Apple Watch), so ensure that you check out the Unity Road Map for an idea of what Unity Technologies are working on and when we can expect them to become available at <https://unity3d.com/unity/roadmap>.

The upcoming C# Job System

A huge performance-enhancing feature that Unity Technologies have been teasing for a while is a feature dubbed as the *C# Job System*. The feature is still in active development and has not yet made it to any release versions of Unity, but it would be wise to start becoming familiar with it sooner rather than later, as it will introduce huge changes to how Unity developers will be writing high-performance code. It is plausible that the difference in quality of a game that uses this system well, versus one that doesn't, might become very noticeable, which may cause some fragmentation of the Unity development community. It is in our best interests to understand and exploit the benefits of the new Job System so that our application will have the most potential for success.

The idea of the C# Job System is to be able to create simple tasks that run in background threads to offload work from the main thread. The C# Job System will be ideal for tasks that are *embarrassingly parallel* such as having hundreds of thousands of simple AI agents operating in a Scene simultaneously and really any problem that can be boiled down to thousands of small, independent operations. Of course, it can also be used for typical multithreading behavior as well, where we perform some calculations in the background that are not needed immediately. The Job System also introduces some compiler technology improvements to get an even greater performance boost than just moving the tasks to separate threads.

A big problem with writing multithreaded code is the risk of *race conditions*, *deadlock*, and bugs that are notoriously difficult to reproduce and debug. The C# Job System aims to make this task easier than usual (but by no means trivial). One of Unity's founders, Joachim Ante, gave a presentation on the upcoming C# Job System at Unite Europe 2017 which gives a preview of the C# Job System and takes us through the ways in which how we think about programming Unity code will need to change. Of course, this won't apply to all code we write, but it should be treated as a valuable tool where it can be deployed for huge performance improvements if we understand how it works and are able to identify sensible situations in which it would help.



Race conditions are where two or more calculations are racing toward completion, but actual outcome depends on the order in which they finish. Imagine one thread trying to add 3 to a number, whereas another thread multiplies it by 4. The result will be different, depending on which operation happens first. Deadlock is a problem where two or more threads are competing for shared resources, where each needs the full resource collection to complete its task, but each has reserved a separate small portion of resources and refuses to relinquish control of them to another thread, in which case, none of the threads can get any work done because neither has the complete set it needs.

The presentation can be found at <https://www.youtube.com/watch?v=AXUvnk7Jws4>.

Summary

We've covered a humongous amount of theory and language concepts in this chapter, which have hopefully shed some light on how the internals of the Unity Engine and C# language work. These tools try their best to spare us from the burden of complex memory management, but there is still a whole host of concerns we need to keep in mind as we develop our game. Between the compilation processes, multiple Memory Domains, the complexities of Value types versus Reference types, passing by value versus passing by reference, boxing, Object Pooling, and various quirks within the Unity API, you have a lot of things to worry about. However, with enough practice, you will learn to overcome them without needing to keep referring to giant tomes such as this!

This chapter effectively concludes all of the techniques we can bestow that explicitly aim to improve application performance. However, optimizing your workflow is also enormously beneficial. As mentioned previously, the one constant cost of performance optimization work is development time. However, if you can speed up our development working, saving some time during the more tedious parts of the job, then hopefully you can save yourself enough time to actually implement as many optimization techniques we've talked about through this entire book as you can. There are a lot of neat little nuances to the Unity Engine that aren't well known or clearly documented and that only become apparent through experience with the Engine or by involving ourselves in its community. As such, the next chapter will be full of hints and tips for improving how to manage your project and Scenes more effectively and how to make the most of the Unity Editor.

Tactical Tips and Tricks

Software engineers are an optimistic bunch, and as such, we often underestimate the amount of work it takes to fully implement new features or make changes to an existing codebase. A common mistake is to only consider how long it will take to write the code required to create that feature. In which case, we are forgetting to include the time it takes for several important tasks. We often need to spend time refactoring other subsystems to support the changes we're making. This can happen either because we didn't think it would be necessary at the time, or because we thought of a better way to implement it halfway-through which can quickly turn into a rabbit-hole dive of redesign and refactoring if we don't plan sufficiently far ahead. We should also consider the time needed for testing and documentation. Even if a QA team will do a testing pass against the change after it has been implemented, we still need to run through some scenarios on our own system during implementation in order to ensure that the change actually does what's intended.

The one constant cost included in all performance optimization work is time. So, with limited time at our disposal to both implement the features we want to implement and keep everything working, an important skill to learn for any developer is workflow optimization. Better understanding of the tools we use will save us more time in the long run, and hopefully provide the extra time we need to implement everything we want to, which applies not only to the Unity Engine, but to every tool we use--IDEs, build systems, analytics systems, social media platforms, app stores, and so on.

There are a lot of little nuances to using the Unity Engine that can help improve our project workflow. However, quite a lot of the Editor's functionality is not well documented, well known, or just not something we think about until after quite some time--we realize the fact that it could have been applied perfectly to solve a particular problem we were having 6 months ago.

The Internet is crammed full of blogs, tweets, and forum posts that try to help

other Unity developers learn about these useful features, but they only tend to focus on a handful of tips at a time. There don't seem to be any online resources that group together many of them in one place. As a result, the Internet browsers of intermediate and advanced Unity developers are probably bursting at the seams with links to these tips that we bookmark for later and then completely forget about.

So, because this book is primarily for such users, I felt like it was worth including a short chapter to pool together many of these tips and tricks in one location. This chapter serves as a reference list in the hope of saving us time during future development effort.

In this chapter, we'll cover tips and tricks for the following areas:

- Convenient Editor hotkeys
- Better and faster ways to use the Unity Editor UI
- Ways to speed up or simplify our scripting practices
- Other tips that involve using external tools outside of Unity

Editor hotkey tips

The Editor is rife with hotkeys that can aid rapid development, and it's worth checking out the documentation. However, let's be honest--nobody reads the manual until they need something specific from it, and so here are some of the most useful, yet lesser-known hotkeys available when playing with the Unity Editor.



In all cases, the Windows hotkey is listed. If the MacOS hotkey requires a different set of keystrokes, then it will be shown in parentheses.

GameObjects

GameObjects can be duplicated by selecting them in the Hierarchy window and pressing *Ctrl + D* (*Command + D*). New, empty GameObjects can be created using *Ctrl + Shift + N* (*Command + Shift + N*).

Press *Ctrl + Shift + A* (*Command + Shift + A*) to quickly open the Add Component menu. From there, you can type in the name of the Component you wish to add.

Scene window

Pressing *Shift + F* or double-tapping the *F* key will follow an object in the Scene window (assuming that the Scene window is open and visible), which can be helpful for tracking high-velocity objects or figuring out why objects might be falling out of our Scene.

Holding *Alt* and left-click dragging with the mouse in the Scene window will make the Scene window Camera orbit the currently selected object (as opposed to looking around it). Holding *Alt* and right-click dragging with the mouse in the Scene window will zoom the Camera in/out (*Alt + Ctrl + left-drag*).

Holding *Ctrl* and left-click dragging will cause the selected object to snap to the grid as it moves. The same can be done for rotation by holding *Ctrl* as we adjust the rotation widgets around the object. Selecting *Edit | Snap Settings...* opens a window where we can edit the grid that objects snap to on a per-axis basis.

We can force objects to snap to each other through their vertices by holding the *V* key as we move an object around in the Scene window. By doing so, the selected object will automatically snap its vertices to the nearest vertex, of the nearest object to the mouse cursor. This is very useful for aligning Scene pieces into place, such as floors, walls, platforms, and other tile-based systems, without needing to make small manual position adjustments.

Arrays

We can duplicate array elements that have been exposed in the Inspector window by selecting them and pressing *Ctrl + D* (*Command + D*). This will copy the element and insert it into the array immediately after the current selection.

We can remove entries from an array of references (for example, an array of `GameObject`s) by pressing *Shift + Delete* (*Command + Delete*). This will strip away the element and condense the array. Note that the first press will clear the reference, setting it to `null`, but the second press will remove the element. Removing elements in arrays of primitive types (`int`, `float`, and so forth) can be accomplished by simply pressing *Delete* without the *Shift* key (*Command*) modifier held down.

While holding down the right-mouse button in the Scene window, we can use the *W*, *A*, *S*, and *D* keys to fly around with the Camera, in typical first-person Camera control style. The *Q* and *E* keys can also be used to fly up and down, respectively.

Interface

We can press *Alt* and click on any Hierarchy window arrow (the small gray arrow to the left of any parent object name) to expand the object's entire hierarchy rather than just the next level in the Hierarchy window. This works on `GameObjects` in the Hierarchy window, folders and Prefabs within the Project window, lists in the Inspector window, and so on.

We can save and restore object selections in the Hierarchy or Project windows much like a typical RTS game. Make the selection and press *Ctrl + Alt + <0-9>* (*Command + Alt + <0-9>*) to save the selection. Press *Ctrl + Shift + <0-9>* (*Command + Shift + <0-9>*) to restore it. This is exceptionally useful if we find ourselves selecting the same handful of objects over and over again while we make adjustments.

Pressing *Shift + spacebar* will expand the current window to fill the entire Editor screen. Pressing it again will shrink the window and restore it to its previous location.

Pressing *Ctrl + Shift + P* (*Command + Shift + P*) will toggle the Pause button while in Play Mode. This is usually an awkward key combination to press if we're trying to pause in a hurry, so it often helps to create a custom hotkey for pausing:

```
void Update() {
    if (Input.GetKeyDown(KeyCode.P)) {
        Debug.Break();
    }
}
```

In-editor documentation

We can quickly access the documentation of any Unity keyword or class by highlighting it in MonoDevelop and pressing *Ctrl + '*(*Command + '*). This will open the default browser and perform a search on the Unity documentation for the given keyword or class.



Note that users with European keyboards may need to also hold down the Shift key for this feature to work.

The same can be done in Visual Studio by pressing *Ctrl + Alt + M*, followed by *Ctrl + H*.

Editor UI tips

The following collection of tips relates to the Editor and its interface controls.

Script Execution Order

We can prioritize which Scripts will have their `Update()` and `FixedUpdate()` callbacks called before others by navigating to `Edit | Project Settings | Script Execution Order`. If we find ourselves trying to solve complex problems using this feature (with the exception of time-sensitive systems, such as audio processing), it implies that we've got some very fragile and tight coupling going on between our Components. From a software design perspective, this can be a warning sign that we might need to approach the problem from another angle. However, this can be helpful to use as a quick fix.

Editor files

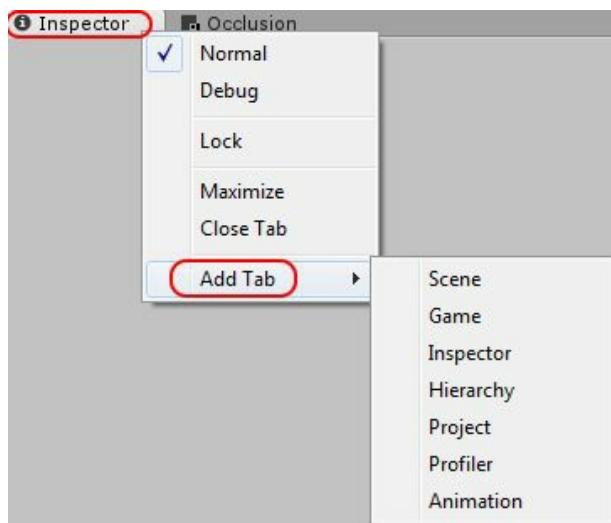
Integrating Unity projects with a Source Control solution can be a little tricky. The first step is to include *.meta* files that Unity generates for various assets; if we don't do this, then anyone pulling data into their local Unity project must regenerate their own metadata files. This could potentially cause conflicts, so it is essential that everyone uses the same versions. This can be enabled by navigating to Edit | Project Settings | Editor | Version Control | Mode | Visible Meta Files.

It can also be helpful to convert certain asset data into a text-only format, rather than binary data, to allow manual editing of data files. This turns many data files into the much more human-readable YAML format. For instance, if we're using Scriptable Objects to store custom data, we can use a text editor to search and edit these files without having to do it all through the Unity Editor and Serialization System. This can save a lot of time, especially when searching for a particular data value or performing multi-editing across different derived types. This option can be enabled by navigating to Edit | Project Settings | Editor | Asset Serialization | Mode | Force Text.

The Editor has a log file, which can be accessed by opening the Console window (where log messages are printed), left-clicking on the *hamburger* icon in the top-right corner (which looks like three thin horizontal lines), and selecting Open Editor Log. This can help us get more information about build failures. Alternatively, if we successfully built our project, it will contain a breakdown of compressed file sizes of all assets that were packed into the executable and ordered by size. This is an extremely helpful way of figuring out which assets consume the majority of our application footprint (hint: it's almost always texture files), and which files take up more space than we would expect.



Additional windows can be added to the Editor by right-clicking on the title of an existing window and selecting Add Tab. This also allows us to add duplicate windows, such as having more than one Project window or Inspector window open at a time. This can be particularly useful for moving files between different locations via multiple Project windows.



Having duplicate Inspector windows can be kind of redundant since they'll show the exact same information when we click on a new object. However, by making use of the *lock icon*, we can lock the given Inspector window to its current selection. When we select an object, all Inspector windows will update to show the object's data, except for any locked Inspector windows, which continue to show the data of the object they were locked to.



Common tricks that make use of window locking include the following:

- Using two of the same window (Inspector, Animation, and so forth) to

compare two objects side by side or more easily copy data from one object to another.

- Watching what happens to any dependent object if an object is tweaked during Play Mode.
- Selecting multiple objects in the Project window, then dragging and dropping them into a serialized array in the Inspector window without losing the original selection.

The Inspector window

We can enter calculations into numeric Inspector window fields. For example, typing `4*128` into an `int` field will resolve the value to 512, sparing us from having to pull out a calculator or do the math in our head.

Array elements can be duplicated and deleted from a list (in the same fashion as the hotkeys mentioned previously) by right-clicking on the root element and selecting Duplicate Array Element or Delete Array Element.

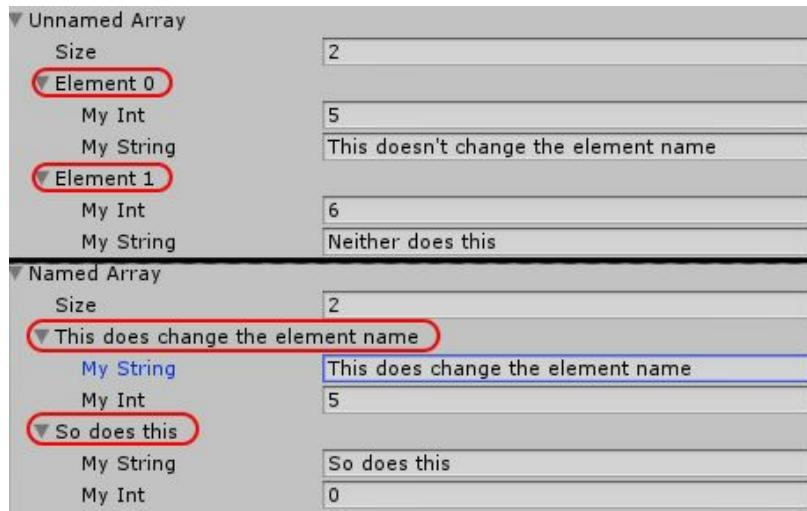
A Component's context menu can be accessed through both the small *cog* icon in the upper-right corner or by right-clicking on the name of the Component. Every Component's context menu contains a Reset option, which resets all values back to their default, sparing us from having to reset values manually. This is useful when working with `Transform` Components, as this option will set the object's position and rotation to `(0,0,0)` and its scale to `(1,1,1)`.

It's commonly known that if a `GameObject` was spawned from a Prefab, then the entire object can be reverted back to its initial Prefab state using the Revert button at the top of the Inspector window. However, it's less well known that individual values can be reverted by right-clicking on the name of the value and selecting Revert Value to Prefab. This restores the selected value while leaving the rest untouched.

The Inspector window has a Debug Mode that can be accessed by left-clicking on the *hamburger* icon next to the *lock* icon and selecting Debug. This will disable all custom Inspector window drawing from Editor scripts and the like, and instead reveal all pieces of raw data within the given `GameObject` and its Components. Even `private` data fields become visible. Although they are grayed-out and cannot be modified through the Inspector window, this still gives us a useful way of examining the `private` data and other hidden values during Play Mode. The Debug Mode of the Inspector window also reveals internal ObjectIDs, which can be useful if we're doing interesting things with Unity's Serialization System and want to

resolve conflicts. Since Editor scripts are also disabled during this mode, it can be useful for debugging such scripts by comparing its internal data to what we are trying to reveal in our Editor script.

If we have an array of data elements serialized in the Inspector window, then they are typically labeled Element N, where N represents the array index of that element, starting from 0. This can make it tricky to find a specific element if our array elements are a series of serialized classes or structs, which tend to have multiple children themselves. However, if the very first field in the object is a string, then the elements will be named after the value of the string field:



When a mesh object is selected, the Preview subsection at the bottom of the Inspector window is often fairly small, making it hard to take a look at details in the mesh and how it will look when it appears in our Scene. However, if we right-click on the top bar of the Preview subsection, it will be detached and enlarged into a separate Preview window, making it much easier to see our mesh. We don't have to worry about setting the detached window back to its original home, as if the detached window is closed, then the Preview subsection will return to the bottom of the Inspector window.

The Project window

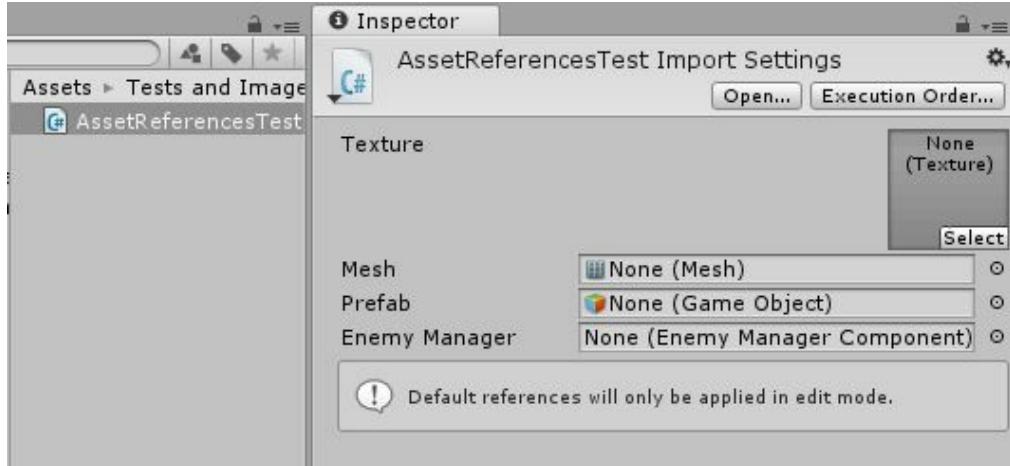
The Project window's search bar allows us to filter for objects of a particular type by clicking on the small icon to the right of the search bar. This provides a list of different types we can filter by revealing all objects of that type within the entire project. However, selecting these options simply fills the search bar with a string of the `t:<type>` format, which applies the appropriate filter.

Thus, we can simply type the equivalent strings into the search bar for the sake of speed. For instance, typing `t:prefab` will filter for all Prefabs, no matter where they are found in the Hierarchy window.

Similarly, `t:texture` will reveal textures, `t:scene` will reveal Scene files, and so on. Adding multiple search filters in the search bar will include objects of all types (it does not reveal objects that only satisfy both filters). These filters are modifiers in addition to name-based filtering, so adding a plain text string will cause a name-based search through the filtered objects. For example, `t:texture normalmap` will find all texture files that include the word *normalmap* in their name.

If we're making use of Asset Bundles and the built-in labeling system, the Project window's search bar also allows us to hunt down bundled objects by their label using `l:<label type>`.

If a `MonoBehaviour` script contains serialized references (using `[SerializeField]` or `public`) to Unity assets, such as meshes and textures, then we can assign default values directly into the script itself. Select the script file in the Project window and the Inspector window should contain a field for the asset for us to drag and drop the default assignment into.



By default, the Project window splits files and folders into two columns and treats them separately. If we prefer the Project window to have a typical hierarchical folder and file structure, then we can set it to One Column Layout in its context menu (the hamburger icon at the top right). This can be a great space saver in some Editor layouts.

Right-clicking on any object in the Project window and selecting Select Dependencies will reveal all objects on which this asset relies in order to exist, such as textures, meshes, and `MonoBehaviour` script files. For Scene files, it will list all entities referenced within that Scene. This is helpful if we're trying to perform asset cleanup.

The Hierarchy window

A lesser-known feature of the Hierarchy window is its ability to perform Component-based filtering within the currently active Scene. This can be accomplished by typing `t:<component name>`. For example, typing `t:light` inside the Hierarchy window search bar will reveal all objects in the Scene that contain a `Light` Component.

This feature is not case-sensitive, but the string we input must match the full Component name for the search to complete. Components that derive from the given type will also be revealed, so typing `t:renderer` will reveal all objects with derived Components such as `MeshRenderer` and `SkinnedMeshRenderer`.

The Scene and Game windows

The Scene window Camera is not visible from the Game window, but it is generally a lot easier to be moved around and placed through the use of the hotkeys mentioned previously. The Editor allows us to align the selected object to the same position and to rotate the Scene window Camera by navigating to GameObject | Align with View or pressing *Ctrl + Shift + F* (*Command + Shift + F*). This means that we can use the Camera controls to place the Scene window Camera where we would like our object to be and place the object there by aligning it with the Camera.

Similarly, we can align the Scene window Camera to the selected object by selecting GameObject | Align View to Selected (note that there is no hotkey for this on either Windows or MacOS). This is useful for checking whether the given object is pointing toward the right direction.

We can perform similar Component-based filtering on the Scene window, as we can with the Hierarchy window, using the `t:<component>` syntax within its search bar. This will cause the Scene window to only render objects containing the given Component (or those that derive from it). Note that this textbox is linked to the same textbox in the Hierarchy window, so anything we type in one will automatically affect the other, which is very helpful when searching for elusive objects.

At the very top right of the Unity Editor is a drop-down menu labeled Layers. This contains a Layer-based filtering and locking system for the Scene window. Enabling the eye icon for a given Layer will show/hide all objects of that Layer within the Scene window. Toggling the lock icon will allow or prevent objects of the given Layer from being selected or modified (at least through the Editor UI).

This is helpful for things such as preventing someone from accidentally selecting and moving background objects that have been already situated in the perfect location:



A commonly known and useful feature of the Editor is that `GameObjects` can be given special icons or labels to make them easier to find in the Scene window. This is particularly helpful for objects with no renderer but that we wish to find easily. For instance, objects such as Lights and Cameras have built-in icons that identify them in our Scene window more easily. However, the same gizmos can be revealed within the Game window by clicking on the Gizmos button at the top right of the Game window. The dropdown for this option determines what gizmos will be visible when this option is enabled.

Play Mode

Since Play Mode changes are not automatically saved, it is wise to modify the tint color applied during Play Mode to make it blatantly obvious which mode we're currently working with. This value can be set by navigating to `Edit | Preferences | Colors | Playmode tint`.

Changes can be saved from Play Mode simply using the clipboard. If we're tweaking an object in Play Mode and we're happy with its settings, then we can copy the object into the clipboard using `Ctrl + C (Command + C)` and paste it back into the Scene once Play Mode ends via `Ctrl + V (Command + V)`.

All settings on the object at the time of the copy will be kept. The same can be done with individual values or entire Components using the Copy Component and Paste Component options in the Component's context menu. However, the clipboard can only contain data for one `GameObject`, Component, or value at a time.

Another approach, which allows us to save the data of multiple objects during Play Mode, is to create Prefabs from them by dragging and dropping them into the Project window at runtime once we're happy with the settings. If the original object was derived from a Prefab, and we wish to update it across all instances, then we only need to overwrite the old Prefab with the new one we created by dragging and dropping the copy on top of the original. Note that this also works while Play Mode is active, but it can be dangerous since there is no dialog popup to confirm the overwrite. Be very careful not to overwrite the wrong Prefab.

We can use the Frame Skip button (the button to the right of the Pause button in the Editor) to iterate one frame at a time. This can be useful for watching frame-by-frame physics or Gameplay behavior. Keep in mind that this causes both one Fixed Update and one Update to be called per iteration, in equal counts, which may not exactly reflect actual runtime behavior where we tend to have an unequal number of calls to these callbacks.

If the Pause button is enabled when Play Mode begins, then the game will be paused just after the very first frame, giving us a chance to observe any anomalies that occurred during initialization of our Scene.

Scripting tips

The following tips are useful features to know when scripting.

General

We can modify various templates of new script, Shader, and Compute Shader files. This can be helpful to remove the empty Update stubs, which as we covered in [Chapter 2, Scripting Strategies](#), can cause unnecessary runtime overhead. These files can be found in the following locations:

- Windows: <Unity install>\Editor\Data\Resources\ScriptTemplates\
- MacOS: /Applications/Unity/Editor/Data/Resources/ScriptTemplates/

The `Assert` class allows for assert-based debugging, which some developers are more comfortable with as opposed to exception-based debugging. Check out the Unity documentation for more information on Asserts at <http://docs.unity3d.com/ScriptReference/Assertions.Assert.html>.

Attributes

Attributes are very useful meta-level tags that can be given to almost any target in C#. They are most commonly used on fields and classes, allowing us to flag them with special properties so that they can be processed differently. Intermediate and advanced Unity developers will find it worthwhile to read the C# documentation on attributes and use their imagination to come up with their own attributes that help accelerate their workflow. There are quite a few attributes built in to the Unity engine that can be exceptionally useful when used in the right place.



Advanced users will note that attributes can also be given to enums, delegates, methods, parameters, events, modules, and even assemblies.

Variable attributes

The `[Range]` attribute can be added to an integer or floating-point field to convert it to a slider in the Inspector window. We can give minimum and maximum values, thus limiting the range that the value can contain.

Normally, if a variable is renamed, even if we do a refactor through our IDE (whether its MonoDevelop or Visual Studio), then the values are lost as soon as Unity recompiles the `MonoBehaviour` and makes the appropriate changes to any instances of the Component. However, the `[FormerlySerializedAs]` attribute is incredibly helpful if we want to rename a variable that has been previously serialized, as it will copy the data from the variable named within the attribute into the given variable during compilation time. No more lost data due to renaming stuff!

Note that it is not safe to remove the `[FormerlySerializedAs]` attribute after the conversion is completed unless the variable has been manually changed and resaved into every relevant Prefab since the attribute was included. The `.prefab` data file will still contain the old variable name, and so it still needs the `[FormerlySerializedField]` attribute to figure out where to place the data the next time the file is loaded (for example, when the Editor is closed and reopened). Thus, this is a helpful attribute, but extended use does tend to clutter up our codebase a lot.

Class attributes

The `[SelectionBase]` attribute will mark any `GameObject` the Component is attached to as the root of selection for the Scene window. This is especially useful if we have meshes that are children of other objects, as we might want the parent object to be selected with the first click, instead of the object with the `MeshRenderer` Component.

If we have Components with a strong dependency, we can use the `[RequireComponent]` attribute to force level designers to attach vital Components to the same `GameObject`. This ensures that any dependencies that our codebase relies on will be satisfied by designers, without having to write out a whole bunch of documentation for them.

The `[ExecuteInEditMode]` attribute will force the object's `Update()`, `OnGUI()`, and `OnRenderObject()` callbacks to be called even during Edit Mode. However, there are caveats, as follows:

- The `Update()` method is only called if something changes in the Scene, such as moving the Camera around or changing an object property
- `OnGUI()` is only called during Game window events, not for other window such as the Scene window
- `OnRenderObject()` is called during any repaint event for the Scene and Game windows

However, this attribute gives such objects a different set of event hooks and entry points compared to typical Editor scripts, so it still has its uses.

Logging

We can add rich text tags to debug strings. Tags such as `<size>`, `` (bold), `<i>` (italics), and `<color>` work on debug strings. This can be helpful for differentiating the different kinds of log messages and highlighting specific elements, as follows:

```
| Debug.Log ("<color=red>[ERROR]</color>This is a <i>very</i> <size=14><b>speci  
! [ERROR]This is a very specific kind of log message
```

The `MonoBehaviour` class has a `print()` method for convenience, which does the same thing as `Debug.Log()`.

It can help to create a custom logger class, which automatically appends `\n\n` to the end of every log message. This will push away the unnecessary `UnityEngine.Debug.Log(Object)` clutter that tends to fill the Console window.

Useful links

Unity Technologies provide many useful tutorials on the usage of various scripting features, which primarily target beginner- and intermediate-level developers. The tutorials can be found at <https://unity3d.com/learn/tutorials/topics/scripting>.

There's a helpful post on Unity Answers, which provides a reference list that covers many of the different scripting and compilation errors we might run across during development, found at <http://answers.unity3d.com/questions/723845/what-are-the-c-error-messages.html>.

Nested Coroutines is an interesting and useful area of scripting that is not well documented. However, the following third-party blog post, which covers a lot of the interesting details, should be considered when working with Nested Coroutines:

<http://www.zingweb.com/blog/2013/02/05/unity-coroutine-wrapper>

We can figure out when a particular feature was added to the Unity API by checking the API history page at http://docs.unity3d.com/ScriptReference/40_history.html. This page currently shows API history only up to version 5.0.0. Hopefully, Unity Technologies will update this page someday, because it can sometimes be useful to know what features were added when we're trying to support multiple versions of Unity simultaneously.

Custom Editor scripts and menu tips

While it's common knowledge that we can create an Editor menu item in an Editor script with the `[MenuItem]` attribute, a lesser-known ability is being able to set custom hotkeys for menu items. For example, we can make the "K" key trigger our menu item method by defining the `[MenuItem]` attribute ending with `_k` as follows:

```
| [MenuItem("My Menu/Menu Item _k")]
```

We can also include modifier keys such as *Ctrl* (*Command*), *Shift*, and *Alt* using the %, #, and & characters, respectively.

`[MenuItem]` also has two overloads, which allows us to set two additional parameters: a Boolean that determines whether the menu item requires a validation method, and an integer that determines the menu item's priority in the Hierarchy window.

Check out the documentation for `[MenuItems]` for a complete list of available hotkey modifiers, special keys, and how to create validation methods at <http://docs.unity3d.com/ScriptReference/MenuItem.html>.

It is possible to *ping* an object in the Hierarchy window, similar to what happens when we click on a `GameObject` reference in the Inspector window by calling `EditorGUIUtility.PingObject()`.

The original implementation of the `Editor` class, and the way that most people learned how to write Editor scripts, originally involved writing all logic and content drawing in the same class. However, the `PropertyDrawer` class is an effective way of delegating Inspector window drawing to a different class from the main `Editor` class. This effectively separates input and validation behavior from display behavior, allowing more fine-tuned control of rendering on a per-field basis and more effective reuse of code. We can even use `PropertyDrawer` to override default Unity drawing for built-in objects, such

as Vector and Quaternion.

PropertyDrawer makes use of the SerializedProperty class to accomplish serialization of individual fields, and they should be preferred when writing Editor scripts, since they make use of built-in undo, redo, and multi-edit functionality. Data validation can be a little problematic, and the best solution is to use `onValidate()` calls in *setter* properties. A session at Unite 2013 by Unity Technologies developer Tim Cooper explains the benefits and pitfalls of various serialization and validation approaches in great detail at https://www.youtube.com/watch?v=0zc_hXzp_KU.

We can add entries to Component context menus and even the context menus of individual fields with the `[ContextMenu]` and `[MenuItem]` attributes. This allows an easy way to customize Inspector window behavior for our Components without needing to write broad Editor classes or custom Inspector windows.

Advanced users may find it useful to store custom data within Unity metadata files through the `AssetImporter.userData` variable. There are also a multitude of opportunities to make use of Reflection of the Unity codebase. Ryan Hippel's session at Unite 2014 outlines a huge number of neat little hacks and tricks one can achieve with Reflection in the Unity Editor at <https://www.youtube.com/watch?v=SyR40YZpVqQ>.

External tips

The following tips and tricks relate to topics outside the Unity Editor itself that can help Unity development workflow enormously.

The Twitter hashtag `#unitytips` is a great resource for useful tips and tricks for Unity development and is, in fact, where many of the tips in this chapter originate from. However, hashtags are difficult to filter for tips you haven't seen before, and it tends to be abused for marketing. A great resource that pulls together a bundle of weekly tips from `#unitytips` can be found at <http://devdog.io/blog>.

Googling Unity-related problems or concerns can go a lot faster if we start the search with `site:unity3d.com`, which will filter all results so that only those under the `unity3d.com` domain will appear.

If the Unity Editor crashes, for whatever reason, then we can potentially restore our Scene by renaming the following file to include the `.unity` extension (for Scene files), and copying it into our `Assets` folder:

```
| \<project folder>\Temp\_EditModeScene
```

If we're developing on Windows, then there's very little reason not to use Visual Studio at this point. MonoDevelop has been dragged along kicking and screaming for many years, and many developers switched over to the more feature-rich Visual Studio Community Edition for most of their development workflow needs, particularly with incredibly helpful plugins such as Resharper.

There is a great resource for game programming patterns (or, rather, typical programming patterns explained in a way that is pertinent to game development), and it's completely free and available online. The following guide includes more information on several of the design patterns and game features we explored in this book, such as the Singleton Pattern, Observer Pattern, the Game Loop, and doubling-up on Frame Buffers:

<http://gameprogrammingpatterns.com/contents.html>

Keep an eye on any session videos that come from Unite conferences, whenever they happen (or better yet, try to attend them). There's usually a couple of panels at each conference held by Unity employees and experienced developers who will share lots of cool and interesting things they've been able to accomplish with the Engine and Editor. In addition to this, make sure that you involve yourself in the Unity community, through the forums on <https://unity3d.com>, Twitter, Reddit, Stack Overflow, Unity Answers, or at whatever social gathering places pop out of the woodwork in the coming years.

Every single tip that was included in this book wasn't conjured out of thin air. It started out as an idea or tidbit of knowledge that someone shared somewhere at some point and somehow eventually found their way to this book's author. So, the best way to keep up to date on the best tips, tricks, and techniques is to keep our fingers on the pulse of where Unity is heading by staying involved in its community.

Other tips

Finally, the following section contains tips that didn't quite fit into other categories.

It's always a good idea to organize our Scenes using empty `GameObjects` as a parent for a group of objects while naming it to something sensible for that group. The only drawback to this method is that the empty object's `Transform` is included during position or rotation changes and gets included during recalculations, and as we know, reparenting a `GameObject` to another `Transform` has its own costs. Proper object referencing, `Transform` change caching, and/or use of `localPosition/localRotation` solve some of these problems adequately. In almost all cases, the benefits to workflow from Scene organization are significantly more valuable than such trivial performance losses.

Animator Override Controllers were introduced way back in Unity v4.3, but tend to be forgotten or rarely mentioned. They are an alternative to standard Animation Controllers that allow us to reference an existing Animation Controller, and then override specific animation states to use different animation files. This allows for much faster workflows since we don't need to duplicate and tweak Animation Controllers multiple times; we only need to change a handful of animation states.

When Unity is launched, it automatically opens the Project Wizard, allowing us to open a recent project. However, if we prefer the default behavior from Unity 4, which is to automatically open the previous project, we can edit this behavior by navigating to `Edit | Preferences | General | Load Previous Project` on startup. Note that if the Project Wizard is enabled, we can open multiple instances of Unity Editor simultaneously (although not to the same project).

The amazing customizability of the Unity Editor and its ever-growing feature set means that there are tons of little opportunities to improve workflows and more are being discovered or invented every single day. The Asset Store marketplace is absolutely rife with different products that try to solve some

kind of problem that modern developers are having trouble with, which makes it a great place to browse if we're looking for ideas or, if we're willing, drop some money to save us a ton of hassle.

Since these assets tend to sell to a broad audience, this tends to keep prices low, and we can pick up some amazingly useful tools and scripts for surprisingly little cost. In almost all cases, it would take us a significant number of hours to develop the same solution ourselves. If we consider our time as valuable, then scanning the Asset Store on occasion can be a very cost-effective approach to development.

Summary

This brings us to the book's conclusion. Hopefully, you have enjoyed the ride. To reiterate, perhaps, the most important tip in this book, always make sure that you verify the source of the performance bottleneck via benchmarking before making a single change. The last thing we want to waste time on is chasing ghosts in the codebase, when 5 minutes of profiler testing can save us an entire day of work. In a lot of cases, the solution requires a cost-benefit analysis to determine whether we're not sacrificing too much in any other area at the risk of adding further bottlenecks. Make sure that you have a reasonable understanding of the root cause of the bottleneck to avoid putting other performance metrics at risk. To also reiterate the second most important tip in this book, always profile and test after making changes to ensure that it had the intended effect.

Performance enhancement is all about problem solving, which can be a lot of fun since due to the complexity of modern computer hardware, small tweaks can yield big rewards. There are many techniques that can be implemented to improve application performance or speed up our workflows. Some of these are hard to fully realize without the experience and skills necessary to spend a reasonable amount of time implementing them. In most cases, the fixes are relatively simple if we simply take the time to find and understand the source of the problem. So, go forth and use your repository of knowledge to make your games the best they can be.