



**Faculty of Electrical Engineering
Department of Computer Graphics and Interaction**

Implementation of rendering system in Rust

Eduard Lavuš

**Supervisor: doc. Ing. Jiří Bittner, Ph.D.
Field of study: Open Informatics
Subfield: Computer Games and Graphics
Date: 2020-05-22**

I. OSOBNÍ A STUDIJNÍ ÚDAJE

Příjmení: **Lavuš** Jméno: **Eduard** Osobní číslo: **474497**
Fakulta/ústav: **Fakulta elektrotechnická**
Zadávající katedra/ústav: **Katedra počítačové grafiky a interakce**
Studijní program: **Otevřená informatika**
Studijní obor: **Počítačové hry a grafika**

II. ÚDAJE K BAKALÁŘSKÉ PRÁCI

Název bakalářské práce:

Implementace zobrazovacího systému v jazyce Rust

Název bakalářské práce anglicky:

Implementation of rendering system in Rust

Pokyny pro vypracování:

Provedte rešerši metod používaných pro zobrazování virtuálních scén v současných herních enginech. Vyberte důležitou podmnožinu zmapovaných metod a implementujte ji ve vlastním zobrazovacím systému založeném na jazyce Rust. Pro implementaci využijte 3D rozhraní Vulkan. V práci rozeberte výhody a nevýhody jazyka Rust ve srovnání s jazykem C++ pro implementaci dané úlohy. Vytvořte demonstrační aplikaci, která ukáže možnosti vytvořeného systému na nejméně třech scénách různé složitosti. Součástí aplikace bude vyhodnocení rychlosti zobrazování (benchmark) a identifikace úzkých hrdel výpočtu.

Seznam doporučené literatury:

- [1] Jason Gregory. Game Engine Architecture (3rd edition). CRC Press, 2018.
- [2] Tomas Akenine-Moller et al. Real-Time Rendering (4th edition). CRC Press, 2018.
- [3] Lagarde, S., and C. D. Rousiers. 'Moving frostbite to physically based rendering.' SIGGRAPH 2014 Conference, Vancouver. 2014.
- [4] Daniel Šimek. Rozšiřitelný zobrazovací řetězec založený na odloženém stínování. Diplomá práce ČVUT FEL, 2013.
- [5] Tomáš Dřínovský. Nepřímé osvětlení pomocí trasování kuželů. Diplomá práce ČVUT FEL, 2013.

Jméno a pracoviště vedoucí(ho) bakalářské práce:

doc. Ing. Jiří Bittner, Ph.D., Katedra počítačové grafiky a interakce

Jméno a pracoviště druhé(ho) vedoucí(ho) nebo konzultanta(ky) bakalářské práce:

Datum zadání bakalářské práce: **11.02.2020**

Termín odevzdání bakalářské práce: _____

Platnost zadání bakalářské práce: **30.09.2021**

doc. Ing. Jiří Bittner, Ph.D.
podpis vedoucí(ho) práce

podpis vedoucí(ho) ústavu/katedry

prof. Mgr. Petr Páta, Ph.D.
podpis děkana(ky)

III. PŘEVZETÍ ZADÁNÍ

Student bere na vědomí, že je povinen vypracovat bakalářskou práci samostatně, bez cizí pomoci, s výjimkou poskytnutých konzultací. Seznam použité literatury, jiných pramenů a jmen konzultantů je třeba uvést v bakalářské práci.

Datum převzetí zadání

Podpis studenta

Acknowledgements

Thanky.

Declaration

I declare that this work is all my own work and I have cited all sources I have used in the bibliography.

.....

Abstract

English abstract

Keywords: key word here

Abstrakt

Slovenský abstrakt

Kľúčové slová: zámok kľúč hrad

Preklad názvu: Implementácia zobrazovacieho systému v jazyku Rust

Contents

1	Introduction	1
1.1	Vulkan API	1
2	Related work	3
2.1	V-EZ	3
2.2	gfx-hal	3
2.3	Vulkano	3
2.4	Tephra	4
2.4.1	Summary	4
3	Design	5
3.1	Rust	5
3.1.1	Safety and speed	6
3.1.2	Cargo	6
3.2	Object lifetime management	7
3.3	Memory management	7
3.4	Synchronization	8
3.5	Validations	8
3.6	Windows	8
4	Implementation	9
4.1	Bindings	9
4.2	Cargo features	9
4.3	Details	10
4.3.1	Vrc, Vutex etc.	10
4.4	Generics	10
4.4.1	Storing generic parameters	11
4.4.2	Dynamic generics	11
4.4.3	Generics in Vulkayes	12
4.4.4	Deref	14
4.4.5	User code a.k.a. dyn FnMut	14
4.5	Swapchain recreate	14
5	Evaluation	17
5.1	User code	17
5.2	Benchmark	18
5.3	Safety	21
6	Conclusion	23
	Bibliography	25

1 Introduction

Since its release in 2016, Vulkan API[1] has been gaining traction as a go-to API for high-performance realtime 3D applications across all platforms. The main reason for this, apart from being cross-platform, is that Vulkan is designed as to be low-level, close to metal and with minimal overhead. This, in contrast to Khronos' older API OpenGL, leaves most of the overhead, but also complexity, to the user of the API. The user can then make decisions on where to sacrifice performance for added usability or vice versa.

This project aims to design a flexible, usable and performant wrapper on top of Vulkan API in the Rust language. It aims to provide statically upholdable invariants that are easy to break in C language. It aims to add minimal required overhead to ensure basic memory safety that is the core concept of the Rust language. The name is a play on the Rust library the project is inspired by, the `Vulkano`[2] library.

1.1 Vulkan API

Vulkan API, originally released in 2016[3], is a specification of a an open API for high-efficiency, cross-platform access to graphics and compute on modern GPUs.

It is designed to minimize the overhead between the user application and the hardware device. Vulkan achieves this by staying low level a explicitly requiring all relevant state to be reference by the user application, minimizing required lookups and orchestration on the driver side. This allows the user application to optimize for their specific usecase instead of relying on the driver to guess the correct strategy.

One of the reasons for Vulkans popularity is that it was designed in an intense collaboration between leading hardware, game engine and platform vendors[3]. This resulted in a lot of vendors having zero-day support for the specification in their drivers and software and it being immediatelly adopted as a native rendering platform on many platforms.

The openness of Vulkan also goes hand-in-hand its cross-platform capabilities. Vulkan is available on all three major desktop platforms (Linux, macOS, Windows) and both major smart-phone platforms (Android, iOS), but also on many smaller and embedded platforms. This allows applications to easily target multiple platforms with minimal variance in the rendering code. It also prevents vendor locks as seen with DirectX or Metal APIs. Lastly, it allows the community of both professionals and hobbyists to participate in the standard itself and improve it.

Khronos Group, the industry consortium responsible for Vulkan API, has been continuously improving the API and releasing updates. The API is currently on version 1.2[4], which brough important updates that have been requested by the community. This proves that Vulkan aims to improve alongside the industry and provide support and improvements into the foreseeable future.

2 Related work

There are already many libraries aiming to provide similar abstraction over Vulkan. Some of the most prominent and closest to this work are mentioned below.

2.1 V-EZ

“V-EZ is an open source, cross-platform (Windows and Linux) wrapper intended to alleviate the inherent complexity and application responsibility of using the Vulkan API. V-EZ attempts to bridge the gap between traditional graphics APIs and Vulkan by providing similar semantics to Vulkan while lowering the barrier to entry and providing an easier to use API.”[5]

This ease of use does come at a price, however. The design of V-EZ leaves no room for the user to properly express their intent at critical points of execution. This leads to unnecessary slowdowns and hashmap lookups which outweigh most of the benefits gained by simplified API.

Last commit to V-EZ was on 2018-10-05[5].

2.2 gfx-hal

gfx-hal or graphics hardware abstraction layer[6] is a project aimed at abstracting graphics computations not only from hardware, but also from low-level APIs like Vulkan or OpenGL. It is, in a sense, lower level than Vulkayes aims to be. The abstraction over multiple APIs, while very useful for most common usages, can hurt usability in niche cases where a specific extension or feature is only available in one API.

In contrast, Vulkayes aims to provide a *transparent* abstraction over Vulkan API. This allows users to use any features available to them by the API even if the abstraction doesn’t implement it directly.

2.3 Vulkano

Vulkano[2] aims to provide complete validation and synchronization guarantees for the user. This proved to be too limiting and the original developer eventually left the project. Since then, not much work has been done.

Vulkayes originally started as a fork of Vulkano, however, over time, it grew into a rewrite because of many questionable design choices taken in Vulkano. Vulkano makes heavy use of dynamic dispatch, which impacts performance. Its API also promises thorough validation checks, however at the expense of API flexibility, which makes it less likely to be widely adopted. For example, it is still impossible to upload mipmaps to Vulkano’s `ImmutableImage` (which is intended as one-time write image abstraction, e.g. for textures in games).

2.4 Tephra

Tephra[7] is a very recent work with very similar aims to Vulkayes. It can be thought of as a C++ version of Vulkayes. It takes a fresh look at the existing solutions and comes up with a transparent and flexible API for handling Vulkan.

However, many of the design considerations taken in Tephra revolve around safety and sanity of C++ language itself. This is of questionable importance and puts unnecessary strain on the library designer. Overall, most of the well designed concepts in Tephra have to be weighed against the unfriendliness of the language.

2.4.1 Summary

Table 2.1: *Related work summary*

Library	Status	Language	Goals
V-EZ	Abandoned	C++	Usability over performance
gfx-hal	Active	Rust	Hardware abstractions
Vulkano	Abandoned	Rust	Safety over performance
Tephra	Unknown	C++	Performance and usability
Vulkayes	Active	Rust	Performance, usability and increased safety

In summary, many projects aiming at simplifying the Vulkan API have been either abandoned or are too broad in scope to consider them production-ready. In the end, Tephra comes out as the closest and most practically usable work. However, the C++ language is itself a complex and hard to master system that places many requirements on the user of the library.

In contrast, Rust, and by extension Vulkayes, aims to offload as much off the user as possible without unnecessary and performance-reducing restrictions.

3 Design

The API was designed to fulfill three goals:

1. Be transparent - The API must allow falling back to pure Vulkan if a certain feature is not supported or implemented in the API.
2. Be fast - The API must carefully manage abstraction costs and minimize overhead.
3. Be flexible - The API must be easy to use in different contexts. It must not force the user unreasonably to change their code to fit the API.

3.1 Rust

Where performance is critical, programmers often fall back to the “classical” languages such as C and C++. These languages, however, are often burdened by legacy, backwards compatibility and outdated design concepts.

C is a very simple and fast language. However, programming industry has changed quite a lot since its first appearance 48 years ago[8]. Concepts common at the time in programming, such as easy low-level memory access and easy mapping to machine instruction, are hardly transferrable to today's high-level requirements of programming.

C++ attempted to extend C with a useful standard library of data types, algorithms and other features. This made C++ a much better candidate at creating complex performance-critical applications. However, stemming from C, it still carries the burden of past decisions. Writing sound code often requires the programmer to be *more* expressive and pay more attention to intricacies of the language. This comes at an expense in code quality, readability and sometimes programmer sanity.

The Rust programming language became a natural choice for this project because goals 2. and 3. are already core concepts of the language itself. Unlike C, it has extensive standard library and was designed for high-level programming. Unlike C++, higher code safety requires *less* work from the programmer. That is, safety is enforced by the language features in form of static analysis:

Ownership Rust implements a very simple but powerful ownership model. Values are always moveable. You cannot prevent the compiler from moving your value. However, the language is smart about this. Moving a value does not just create a bitwise copy, it also moves the ownership. Ownership has serious consequences: the owner has to clean up. Values that have non-trivial destructors should run those destructors at some point.

In C++ the only difference between a copy and a move is that the new value has a chance to take apart the old value. For example, for heap allocated types, this means the new value will take the heap memory (pointer) from the old value. The destructor, however, is still run for both the values, as if it was simply copied.

In contrast, Rust statically prevents use of moved-out variables. Once you move a value out of a variable, that variable now acts as if it was uninitialized, it cannot be used anymore and its destructor is not called. The destructor is only called for the “new” value once it goes out of scope. Moreover, this move is often optimizable by the compiler and thus is almost or entirely free.

Borrow checker The Rust borrow checker tracks borrowed values. A value is borrowed when a reference to it is created. A reference can either be immutable or mutable. There can only ever be one mutable reference and it also cannot coexist with any immutable references. This completely prevents all read-write race conditions *statically*.

Borrow checking also prevents problems such as use-after-free or iterator invalidation. These problems can be considered single-thread race conditions. A reference is created, then the original referred value is destroyed and then the reference is used (to read or write). Such a reference is called dangling. Rust statically prevents the existence of dangling references. When a value is borrowed, it must outlive any references taken from it. This is done using lifetimes.

Lifetimes Lifetimes are how Rust tracks borrows. Each borrow (a reference) has a lifetime associated with it. The borrow cannot be used for longer than that. For example, if a value is created in a certain scope then a reference to it cannot escape that scope since it could lead to use after free. Additionally, programmers can use these lifetimes too, as generic arguments, to express concepts like borrowing subfields or narrowing array views.

3.1.1 Safety and speed

Of course, some of the lowest-level code cannot be created in this somewhat restricted environment. The abstraction has to be built somehow. This is where **unsafe** Rust comes in. Instead of specifying additional safety features, Rust programmers have to explicitly ask to disable existing features. Code blocks marked `unsafe` are free to work with dangling pointers, have data races or cause other unsoundness, just like C++ normally does.

The implementation of the Rust standard library has empirically proven that the system truly only needs unsafe blocks few and far between. Indeed, only the most basic building blocks have to rely on unsafe operations, while all the other parts can just rely on the soundness of these simple code snippets that can easily be checked and verified over and over by quanta of programmers to ensure they truly are sound. This safety system reduces possible failures to a few narrow blocks of code, instead of leaving the programmer with having to find the bug in all of their code.

All of this is done at compile time and thus has no runtime cost. All code is as fast as the same C++ code would be, but safe.

3.1.2 Cargo

Cargo^[9] is Rust's package manager. It takes care of indexing and retrieving dependencies, compiling them and publishing your own libraries and binaries to the registry. Cargo also takes care of project configuration. In C/C++ codebases it is common to either invoke the compiler directly, or to use build tools such as `make` or `CMake`. Cargo is similar to those build tools, but it's a component of Rust ecosystem and is targeted at Rust only.

Being a part of Rust itself, cargo is able to provide lots of useful abstraction over the rust compiler. The configuration file `Cargo.toml` is filled with useful project information such as the project name, author and short description. The file also contains technical information, such as the targeted language edition, compiler and optimization flags, all of the dependencies (and how/where to look for them) and project features. Platform-specific configuration is also possible.

Features defined in `Cargo.toml` are project-unique strings. These strings can then be used from within the codebase to conditionally compile part of the code, similar to C preprocessor `#ifdef` statements. Contrary to the C preprocessor, however, these strings are defined in one central place and can even define dependency chains, so that certain features might require other features or additional dependencies. This is often used when developing on top of platform-dependent code to provide uniform interface to the user.

3.2 Object lifetime management

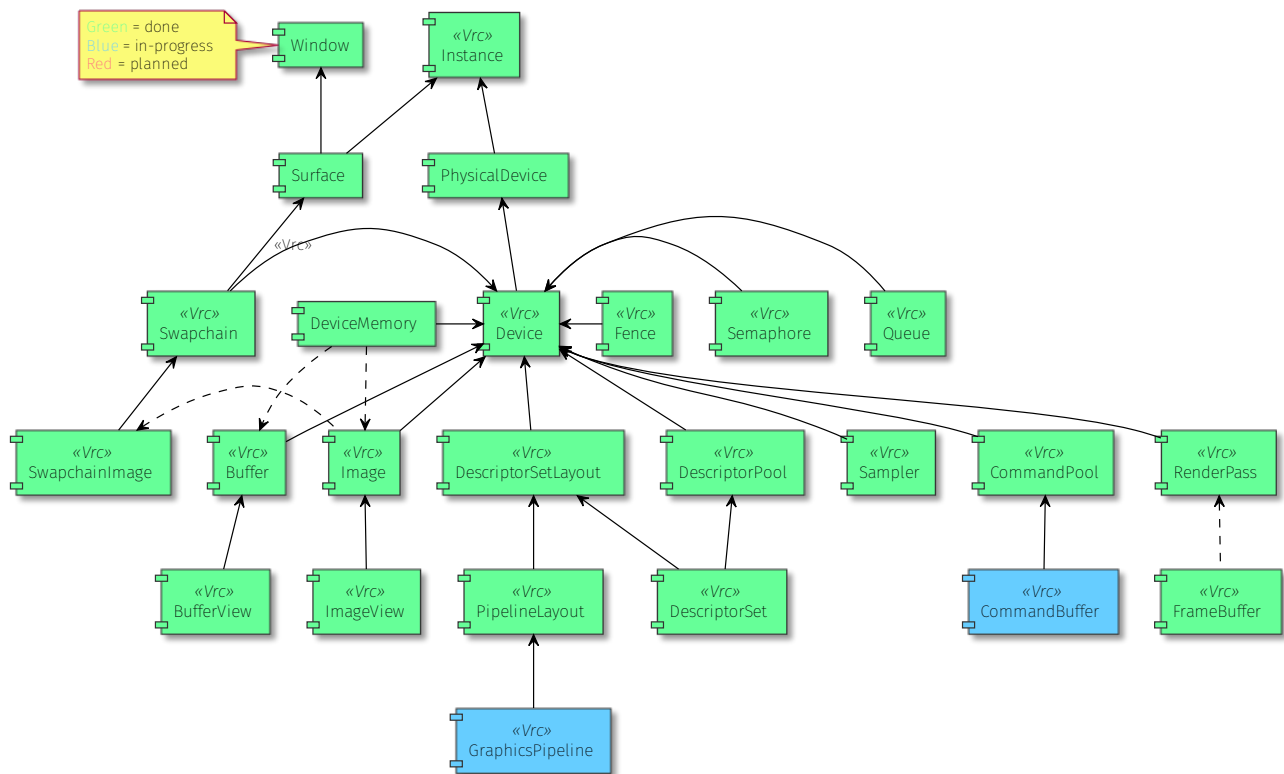


Figure 3.1: Object Dependency Graph

Objects in Vulkan have certain lifetime dependencies - some objects must outlive others - displayed in the diagram[TODO Figure number thing]. Some dependencies are simpler and always apply, others are more complex and conditional. Most of these dependencies in Vulkayes are handled using reference counting. Reference counting is a programming concept where data is shared among multiple actors using some kind of reference (pointer). The pointed-to memory, apart from storing the data object itself, also stores a count of existing references to that memory. This provides an easy way to clean up resources when they are no longer used, all automatically at runtime, with overhead only during the creation and destruction of the resource itself, not during usage. The Rust safety system also prevents the pointed-to memory to be freed or otherwise deinitialized, ensuring safety.

3.3 Memory management

Talk about how device memory management is done through user-supplied memory allocator

3.4 Synchronization

Talk about how some objects are internally synchronized

Talk about how GPU synchronization is left for future work

3.5 Validations

Talk about how only implicit validation are guaranteed, but some explicit validations are implicitly handled by api design and type system

3.6 Windows

Talk about how windows are handle, what is a surface and a swapchain and how they are supported

4 Implementation

4.1 Bindings

Vulkan API is an interface specified in the C programming language. C language is the de-facto standard in cross-language APIs. This means the system of bindings is available in almost any practical language, including Rust. Vulkayes relies on the ash[10] crate to provide these binding and some syntactic sugar on top. This library uses the Vulkan API Registry, canonical machine-readable definition of the API, to generate bindings from Rust to C automatically.

4.2 Cargo features

An important part of any flexible project is to give the user as much control as possible, so the library will fit their usecase. One way to achieve this in Rust are cargo features already mentioned in sec. 3.1.2.

The most important features defined and exposed in Vulkayes are described below.

naive_device_allocator This feature conditionally compiles a very naive device memory allocator into the project. Device memory allocation is a complex topic and applications are required to provide their own allocators to fit their own needs. One popular allocator is the Vulkan Memory Allocator[11], but it is a big dependency that might not be easily accessible for certain usecases. Vulkayes supports integration with VMA (and other allocators) seamlessly, but also provides the naive allocator as a simple no-dependency alternative for quick prototyping and debugging.

multi_thread One of the biggest selling points of Vulkan are its multi-threading capabilities. Since the user is in charge of synchronizing the resources, they can design their application to fit their needs. Single-threaded applications require no synchronization, while multi-threaded applications should allow for the full power of multi-threading to be leveraged.

Safe Rust statically prevents data races using the built-in Send and Sync traits. These traits are automatically implemented (or not implemented) by the compiler to mark types as “capable of being sent between threads” and “capable of being borrowed across threads” respectively. The user is free to unsafely implement these traits back if the compiler decides to not implement them, provided that the user takes the burden of preventing data races upon themselves.

By default, object wrappers in Vulkayes are no Send nor Sync, simply because they use the Rc type, which is a shared pointer wrapper type that uses non-atomic loads and writes to count. By turning this feature on, all usages of Rc across the crate are switched to Arc, which is atomically counted and thus implemented Send and Sync safely.

Additionally, single-threaded Vulkayes replaces use of mutexes with simple wrapper types that emulate the mutex API, but do not implement Send/Sync and do not do any synchronization. This makes the API of both single-threaded and multi-threaded Vulkayes uniform. The main reason for this feature is performance, since atomic operations and synchronization is costly compared to non-atomic counterparts.

runtime_implicit_validations Vulkayes aims to increase safety of Vulkan calls as much as possible without any performance impact. The idea is to always guarantee that the implicit

validations defined in Vulkan spec are fulfilled and the explicit validations are only fulfilled when they can be easily derived from the existing API design.

This proved to not be always possible, so a small portion (tbl. 5.5) of implicit validations requires some runtime checking to ensure their fulfillment. These validation, producing runtime overhead, are conditionally compiled using this feature to ensure that the user can always opt-out to achieve greater performance.

4.3 Details

4.3.1 Vrc, Vutex etc.

Talk about Vrc and Vutex aliases

4.4 Generics

Generics are a very powerful tool in programming. They help avoid a common problem in libraries: “What if my object doesn’t cover all usecases”. Generics provide a way for the library user to specify their own object with their own implementation and it only has to conform to some predefined bounds. In Rust, this is done by specifying trait bounds:

```
trait BoundTrait {
    fn required_method(&self) -> u32;
}

fn generic_function<P: BoundTrait>(generic_parameter: P) -> u32 {
    generic_parameter.required_method()
}
```

In this code snippet, the P parameter of the generic_function is generic. The user can then do this:

```
struct Foo;
impl BoundTrait for Foo {
    fn required_method(&self) -> u32 {
        0
    }
}

struct Bar(u32);
impl BoundTrait for Bar {
    fn required_method(&self) -> u32 {
        self.0
    }
}
```

Now both the Foo struct and the Bar implement the trait BoundTrait and can be used to call generic_function:

```
let foo = Foo;
generic_function(foo);
```



```
let bar = Bar(1);
generic_function(bar);
```

This usage is zero-cost because the functions are monomorphised at the compile time for each calling type.

4.4.1 Storing generic parameters

Using generic parameters is one thing, but storing them is harder. Generic parameters can have different sizes that are not known at the definition time:

```
struct Holder<B: BoundTrait> {
    item: B
}

let a = Holder { item: Foo };
let b = Holder { item: Bar(1) };
```

In this snippet, it is unknown at the definition time how big the `Holder` struct will be in memory. Instead, it is decided at the use time. That is, the variable `a` possibly takes less space on the stack than the variable `b`. The size of a type is a function of its fields, if the field is generic, it can't be known up front.

Generic parameters are a part of the type. Two `Holders` with different generic parameters cannot be stored together in an uniform collection (like `Vec`). The only way to achieve that is by using dynamic dispatch.

4.4.2 Dynamic generics

Dynamically dispatched generics can be used to mix and match different implementations of traits in the same place. It works by taking a pointer to the generic parameter and then “forgetting” the type of that parameter, only remembering the bounds. In rust, this is handled by trait objects in the form of `dyn BoundTrait`. This is an unsized (size isn't known at compile time) type and it cannot be stored directly on the stack or in uniform collections either. It needs to be behind some kind of pointer, whether it be a reference, `Box`, `Rc/Arc` or a raw pointer. This pointer will be a so-called “fat” pointer.

For example, to store any kind of `BoundTrait` implementor in a `Vec`, it can be written like this:

```
let a = Foo;
let b = Bar(1);

let vec: Vec<Box<dyn BoundTrait>> = vec![
    Box::new(a) as Box<dyn BoundTrait>,
    Box::new(b) as Box<dyn BoundTrait>
];
```

The downside of this is the access speed. Accessing methods on the object has to go through one more level of indirection than normally and also prevents certain powerful compiler optimizations. Thus is it undesirable to use dynamic dispatch when it is not necessary.

4.4.3 Generics in Vulkayes

Generics are used in key places across Vulkayes. One example are device memory allocators, another would be image views. They are described in detail below.

Device memory allocator generics Device memory allocators have one of the biggest impact on performance of Vulkan. There is no default memory allocator in Vulkan. Instead, memory has to be allocated manually from the device. That operation, however, can be slow. That is why it is recommended by the Vulkan specification to allocate memory in bigger chunks (about 128 to 256 MB) at once and then distribute and reuse the memory as best as possible in the user code.

For Vulkayes, this means it is required to support user-defined allocators. This is the perfect usecase for generics. An image, which needs some kind of memory backing to operate, has a simplified constructor like this:

```
trait DeviceMemoryAllocation {
    // Allocation trait methods
}

trait DeviceMemoryAllocator {
    type Allocation: DeviceMemoryAllocation;

    fn allocate(&self) -> Self::Allocation;
}

struct Image {
    // Image fields
    memory: ??
}

impl Image {
    pub fn new<A: DeviceMemoryAllocator>(
        // Other fields
        memory_allocator: &A
    ) -> Self {
        // Initialization code
    }
}
```

The `memory_allocator` parameter can be any user-defined type that implement the `DeviceMemoryAllocator` trait (and thus is capable of distributing memory given some requirements). However, given the requirements of Vulkan specification, we need to ensure that the memory outlives all usages of the image. This implies we need to store some kind of handle to the allocated memory, which can be any type implementing `DeviceMemoryAllocation` (as can be seen in the `DeviceMemoryAllocator` traits associated type `Allocation`).

Storing this memory thus has the same implications as mentioned above. We could make the `Image` struct generic over the memory it stores. This would however mean that the memory generic parameter would have to be present on anything that can possibly store the image, including swapchain images, image views, command buffers and so on. This could prevent us in the future from creating a command buffer and recording into it operations on images with possibly different memory allocations (for example, because one is a sparse image and

the other is fully-backed).

Since this is very limiting, the memory inside an image can be stored using dynamic generics. So the ?? in the above code snippet would be replaced with `Box<dyn DeviceMemoryAllocation>`.

This would be ideal for images, where the memory does not need to be accessed until it is to be deallocated (barring linearly tiled images). For buffers, however, this is a common use case. Buffers are often used as staging. Data is uploaded into a buffer from the host and then copied using device operations into an image backed by fast device-local memory. The upload of data is done by mapping the memory into host memory using Vulkan provided mechanism and then writing to it as if it was normal host memory.

Mappable memory generics Some use cases for mapped memory are performance-critical. For example, vertex animating data is done by continuously changing vertex buffer data according to the animation properties. This means the mapped memory has to be accessed every frame. This is where dynamic dispatch cost would be substantial, it is best to avoid it.

One of the ways to avoid this cost is to simply push it back. There are only 3 places where the generics are truly needed:

- The memory map function
- The memory unmap function
- The cleanup function

No other place of the memory handling needs custom user coding. This means it is enough to store 3 generic user-provided functions. In Rust, this can be done using the `Fn` family of traits. For example, instead of `Box<dyn DeviceMemoryAllocation>` for the cleanup function we will use `Box<dyn FnOnce(&Vrc<Device>, vk::DeviceMemory, vk::DeviceSize, NonZeroU64)>` inside a concrete struct `DeviceMemoryAllocation`. The cleanup function can be simply `FnOnce`, which can only ever be called once, while the map and unmap functions might need to be called multiple times and have to be `FnMut`.

Image view generics Image views are another object in Vulkano that has to deal with generics. Image view can wrap any type that can “act like” an image and create a view into some kind of subrange. This can be expressed using the `ImageTrait` like so:

```
struct ImageView {
    // Image view fields
    image: ??
}

impl ImageView {
    pub fn new<I: ImageTrait>(
        image: I
    ) -> Self {
        // Initialization code
    }
}
```

As mentioned above, this is very limiting because of the generic parameter. Unlike the above case, however, the image field needs to be accessed considerably more often.

The following table shows a benchmark of so-called mixed dispatch, where an enum is used to provide common possible values for a given generic type and the last variant, which is the only

one truly generic, is provided as a `Box<dyn Trait>` to allow using dynamic dispatch where the set of provided types is not extensive enough.

benchmark	avg. black box	avg. no black box
Enum::Foo	499.01 ps	251.31 ps
Enum::Bar	499.47 ps	252.67 ps
Enum::Dyn	1.3018 ns	1.2512 ns
Foo	499.36 ps	260.76 ps
Bar	499.03 ps	252.18 ps
Qux	313.34 ps	250.41 ps
dyn Qux	1.5104 ns	1.5028 ns

As can be seen from the table, accessing a value through a dynamic dispatch is at least twice as slow as accessing it through static dispatch, and this is with optimizations prevented by using the concept of a black box from the Rust stdlib.

Non-black boxed benchmarks show that the optimizations provided by the compiler for statically dispatched values can further reduce the overhead of static dispatch, while the dynamic dispatch stays mostly the same.

// TODO: Reference to the benchmark code

4.4.4 Deref

Talk about Deref trait usage

4.4.5 User code a.k.a. `dyn FnMut`

Talk about usage of `dyn`

4.5 Swapchain recreate

Swapchain is an object in Vulkan that facilitates image presentation onto surfaces. Surfaces are an abstraction over regions of the physical display, intended mainly for windowing systems and compositors. A swapchain is created for a combination of a surface and a device.

Requirement for our Swapchain object are:

1. Only one swapchain can exist for one surface.
2. Allow user to retrieve the surface when the swapchain is no longer in use.
3. Allow user to recreate the swapchain, transferring the ownership of the surface to the new instance, retiring the old swapchain.
4. Keep retired swapchain alive until all its acquired images are no longer in use.

Satisfying all three conditions as they are is not trivial, mainly because the the first two conditions lead to the requirement of dropping the swapchain once the surface is moved out of it, however, the fourth condition requires us to keep it alive. This can also create problems where for some reason the retired swapchain outlives the active one. In such cases, the surface can happen to be dropped before the retired swapchain, which is incorrect.

To satisfy all 4 conditions, we first have to rewrite them into terms that can be expressed in the language.

1. The creation of a swapchain requires full ownership of the surface, thus our constructor has to take surface by value.
2. The swapchain has to have a method that consumes the swapchain and returns the surface by value.
3. The new, recreated swapchain has to take the old swapchain by value and extract the surface from it using method from 2.
4. The swapchain has to be reference counted to outlive all its images.

Now it is much clearer why these requirements are hard to satisfy - 4. requires that the swapchain reference counted and its lifetime is guarded dynamically, however, 2. and 3. require for the lifetime of the swapchain to end immediately rather than sometime in the future. We need to rewrite the requirements to work with reference counting.

Adapting 2. is implementationally trivial. We must rely on the user to first drop all outstanding shared pointers except for one and then use that one to retrieve the swapchain back as an owned value.

Adapting 3. however, is much harder to implement as we can't expect the user to wait until all outstanding operations on the current swapchain are done until creating a new one since that would limit the functionality too much. Instead, we need to make sure that the surface is alive for the longer of the two lifetimes. This is exactly what reference counting does. By reference counting the surface inside a swapchain but still requiring an owned value for swapchain creation, we can make sure that no two active swapchains are ever created for one surface while still leaving the possibility of retrieving the surface after all but one of the shared pointers are dropped.

The resulting API thus looks like this:

```
pub struct Swapchain {
    surface: Vrc<Surface>,
    // Other fields
}

impl Swapchain {
    pub fn new(
        surface: Surface,
        // Other parameters
    ) -> Vrc<Self> {
        Vrc::new(
            Swapchain {
                surface: Vrc::new(surface),
                // Other fields
            }
        )
    }

    pub fn recreate(
        self: &Vrc<Self>,
        // Other parameters
    ) -> Vrc<Self> {
        Vrc::new(
```

```

        Swapchain {
            surface: self.surface.clone(),
            // Other fields
        }
    )
}

pub fn surface(&self) -> &Vrc<Self> {
    &self.surface
}
}

```

This satisfies all the rules:

1. We cannot retrieve the surface back from the swapchain without destroying the shared pointer, which dynamically ensures there are no other instances.
2. The swapchain returns a reference to the reference counted surface, which can be destroyed to gain the surface after dropping all swapchains and swapchain images in the same way as above.
3. Both the new and the old swapchain contain a reference to the surface and thus will keep it alive for as long as is needed.
4. Swapchain is reference counted and can be kept alive by the images.

5 Evaluation

5.1 User code

One of the main concerns when designing a library is the user code. How the user code will look like, if it will be readable and comfortable to write.

Below is an example of the code with same functionality from the original examples and from the current ones. The code after is three times shorter than the original code while exposing the same functionality and providing static validation guarantees.

```
let (vertex_buffer, vertex_buffer_memory) = {
    let create_info = vk::BufferCreateInfo {
        size: std::mem::size_of_val(
            &data::VERTICES
        ) as vk::DeviceSize,
        usage: vk::BufferUsageFlags::VERTEX_BUFFER,
        sharing_mode: vk::SharingMode::EXCLUSIVE,
        ..Default::default()
    };
    let buffer = unsafe {
        device
            .create_buffer(&create_info, None)
            .expect("Could not create vertex buffer")
    };

    let memory_req = unsafe {
        device.get_buffer_memory_requirements(buffer)
    };
    let memory_index = memory::find_memory_type_index(
        &memory_req,
        &device_memory_properties,
        vk::MemoryPropertyFlags::HOST_VISIBLE
        | vk::MemoryPropertyFlags::HOST_COHERENT
    )
    .expect("Unable to find suitable memory type");

    let allocate_info = vk::MemoryAllocateInfo {
        allocation_size: memory_req.size,
        memory_type_index: memory_index,
        ..Default::default()
    };
    let memory = unsafe {
        device
            .allocate_memory(&allocate_info, None)
            .expect("Could not allocate memory")
    };
    unsafe {
        device
            .bind_buffer_memory(buffer, memory, 0)
            .expect("Could not bind memory");
    }

    (buffer, memory)
};

let vertex_buffer = {
    Buffer::new(
        device.clone(),
        std::num::NonZeroU64::new(
            std::mem::size_of_val(&data::VERTICES) as u64
        ).unwrap(),
        vk::BufferUsageFlags::VERTEX_BUFFER,
        SharingMode::from(queue.as_ref()),
        BufferAllocatorParams::Some {
            allocator: &device_memory_allocator,
            requirements: vk::MemoryPropertyFlags::HOST_VISIBLE
                | vk::MemoryPropertyFlags::HOST_COHERENT
        },
        Default::default()
    )
    .expect("Could not create vertex buffer")
};
```

Overall, the code for benchmarking the spinning teapots written in pure ash has 1400 lines of Rust code. The code with Vulkayes with same semantics and even improved static validation guarantess has 942 lines. This is a difference of 458 lines of code. These numbers are clear indicators of the improvement in developer experience by using correctly designed wrappers.

Another interesting code example is the uniform buffer usage:

```

unsafe {
    *uniform_buffer_memory_ptr = frame_state;
}
let flush_ranges = [
    vk::MappedMemoryRange::builder()
        .memory(uniform_buffer_memory)
        .size(
            std::mem::size_of::<data::UniformData>()
            as vk::DeviceSize
        )
        .build()
];
unsafe {
    device
        .flush_mapped_memory_ranges(&flush_ranges)
        .expect("Could not flush uniform data memory");
}

```

```

uniform_buffer
    .memory().unwrap()
    .map_memory_with(|mut mem| {
        mem.write_slice(&[frame_state], 0, Default::default());
        mem.flush().expect("Could not flush uniform data memory");
        MappingAccessResult::Continue
    })
    .unwrap();

```

The ash code is twice as long and in some cases possibly even unsafe. Vulkayes API guarantees proper locking and borrowing, provides simplified way to flush the memory and prevents unaligned writes which on some platforms might cause hard errors and abort the process. The checking for correctness, however, does have some runtime cost. One of the guarantees of safe Rust is memory safety and Vulkayes is targeting safe Rust. That is why the `mem.write_slice` method call above does more than just write to a pointer. There is logic to check the align of the pointer and make sure all writes are either properly aligned, or an unaligned instruction is used.

5.2 Benchmark

A benchmark of ash vs Vulkayes was designed to compare the speed of Vulkayes against a “control sample” of ash. This benchmark measures several stages of a common rendering loop between ash and Vulkayes. Since Vulkayes is mostly safety and usability wrapper, not much runtime overhead is added, at least not in the critical hot-paths used in rendering loops. Some specific areas, however, such as memory mapping and writing need special handling to ensure safety, as mentioned in sec. 5.1.

The rendering loop was split into 8 stages:

preinit In this stage frame specific variables are calculated, such as the data dependent on the elapsed time. This stage represents the user logic that is not being benchmarked.

acquire In this stage the present image is acquired from Vulkan implementation. This stage is heavily dependent on the Vulkan implementation and is not being benchmarked.

uniform In this stage uniform data specific for the frame is written into device visible mapped memory and flushed. Some absolute overhead is expected because Vulkayes does checks to ensure memory safety.

command In this stage command buffer is recorded by binding necessary state and submitting draw commands for each teapot, along with push constants. Minimal overhead is expected as only one mutex needs to be locked.

submit In this stage the previously recorded command buffer is submitted for execution to Vulkan. This operation is costly on its own, but only minimal overhead is expected.

present In this stage the acquired image is submitted for presentation and a happens-before relationship is established using semaphores and fences so that submitted execution is guar-

anteed to finish before presentation begins. Again, overhead of a mutex is expected.

wait In this stage the application waits on the presentation fence. This ensures that all measurements done inside one loop are correctly assigned to that loop. In real life applications, this wait should not happen and each frame should asynchronously finish in the background while the user logic computes data for the next frame (akin to the preinit stage). This stage represents the cost of the submitted operations on the GPU, but might also invoke some implementation-dependent synchronization the application is not aware of. Timings of this stage are thus not considered.

update In this stage the update function is called on the window and all outstanding windowing events are handled. This stage represents the update of the windowing system events and a window redraw request and is not being benchmarked.

The benchmarks were run on three hardware and software configurations, note that only the relevant stages are present:

Table 5.1: *Average median time (n = 99000): macOS 10.15.3 (19D76), Quad-Core Intel Core i5, Intel Iris Plus Graphics 655, Vulkan 1.2.135*

Stage	ash	vy_ST	vy_MT
uniform	1.5 us	2.37 us (157%)	2.39 us (159%)
command	23.66 us	24.43 us (103%)	26.51 us (112%)
submit	169.43 us	171.11 us (101%)	170.99 us (101%)
present	32.76 us	33.36 us (102%)	34.14 us (104%)

Table 5.2: *Average median time (n = 99000): Linux 5.4.35_1, Intel i5-7300HQ, TODO*

Stage	ash	vy_ST	vy_MT
uniform	717.0 ns	1.53 us (214%)	1.38 us (193%)
command	39.37 us	40.03 us (102%)	39.78 us (101%)
submit	39.57 us	37.36 us (94%)	38.64 us (98%)
present	25.34 us	26.07 us (103%)	26.45 us (104%)

Table 5.3: *Average median time (n = 99000): Linux 5.4.35_1, Intel i5-7300HQ, NVIDIA GeForce GTX 1050 Mobile, Vulkan v1.2.137*

Stage	ash	vy_ST	vy_MT
uniform	1.42 us	2.15 us (152%)	2.24 us (158%)
command	28.51 us	27.99 us (98%)	28.8 us (101%)
submit	13.15 us	13.76 us (105%)	14.38 us (109%)
present	27.08 us	26.63 us (98%)	27.29 us (101%)

As can be seen, all three tested systems exhibit similar trends. The command stage is on par with pure ash benchmark, the only possible overhead is one mutex lock, which will only have an effect on multi-thread feature in the worst case.

The submit stage also closely follows the ash baseline. This stage potentially locks great num-

ber of mutexes, so could be a potential performance bottleneck on the multi-thread feature. However, the intention of an explicit submit operation in Vulkan API is to reduce overhead of submitting smaller batches of work in favor of bigger ones, where the overhead is less noticeable. Thus, for real life applications where the command buffer size will be much bigger, it is expected to be manageable.

The present stage, similaliry, does not exhibit any noticeable slowdown. The reasoning is the same as for the submit stage. Additionally, the present stage may also include the vertical synchronization delay if enabled, and will thus shadow smaller overhead factors such as locking mutexes.

Finally, the uniform stage exhibits the most interesting results. The accesses performed in Vulkayes are 1.5 to 2 times as slow as when performed by ash. This seems like a lot, but it is important to mention that the absolute difference between the median points is in range of 1 micro second and the overhead is of constant nature.

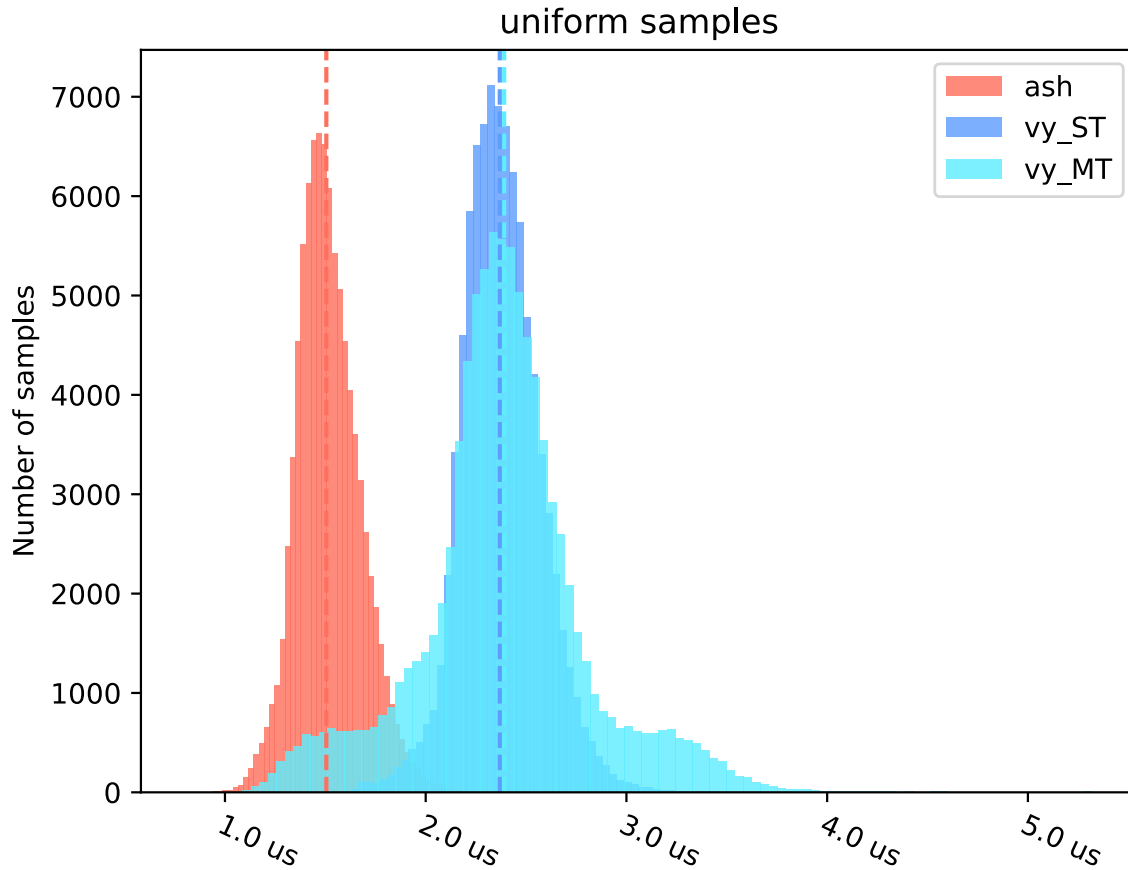


Figure 5.1: Histogram of uniform stage of the benchmarks ($n = 99000$). It's clear that ash is faster than both single- and multi-threaded Vulkayes. However, the overhead is constant.

Furthermore, tbl. 5.4 demonstrates doing 1000 writes into the mapped memory instead of 1 each frame. In fact, Vulkayes is even slightly faster in this case because it decides on the most efficient strategy for the write, which becomes efficient with larger number of writes.

Table 5.4: Average median time ($n = 99000$): macOS 10.15.3 (19D76), Quad-Core Intel Core i5, Intel Iris Plus Graphics 655, Vulkan 1.2.135

Stage	ash_u1000	vy_ST_u1000	vy_MT_u1000
uniform	45.16 us	40.61 us (90%)	42.14 us (93%)

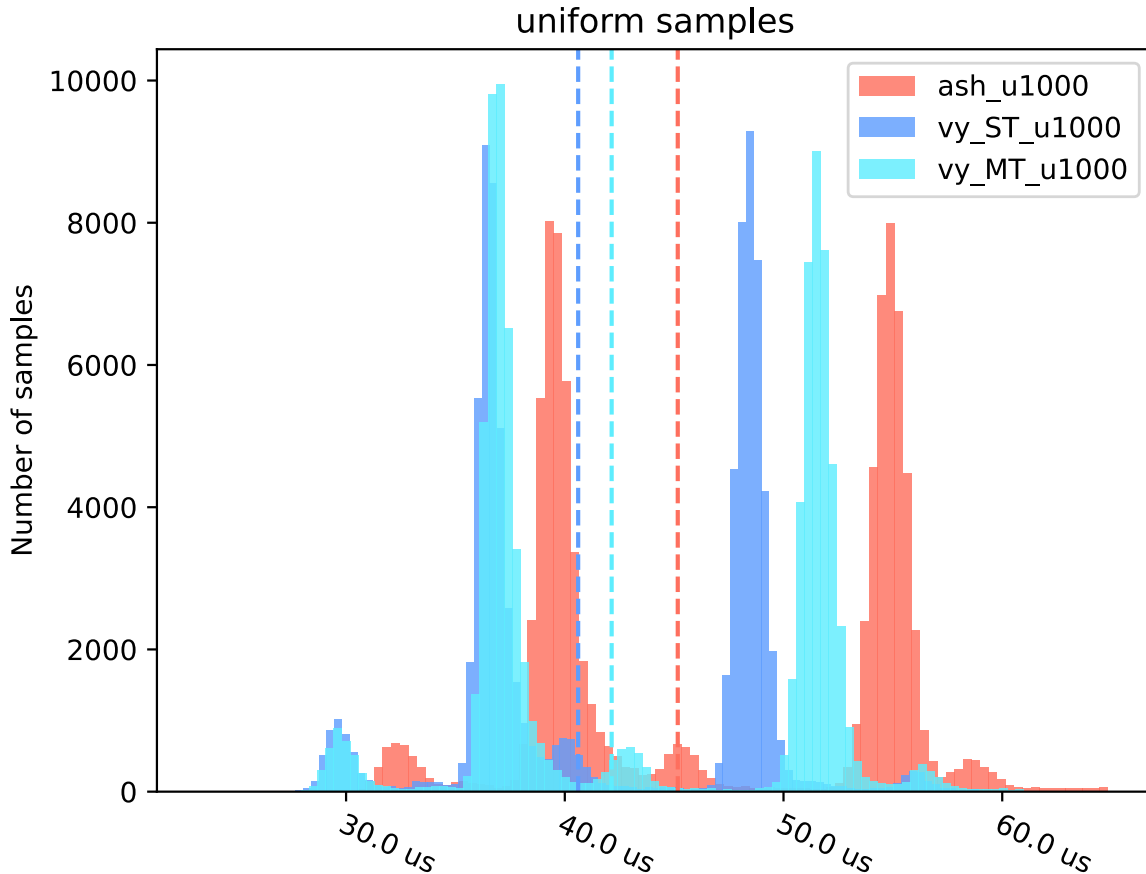


Figure 5.2: Histogram of uniform stage of the benchmarks ($n = 99000$) with 1000 writes instead of 1. The overhead displayed in previous bench is overshadowed by the gains of proper writing strategy.

5.3 Safety

One of the main goals of Vulkayes is increasing safety. This mainly includes memory safety. Vulkayes, being a safe wrapper, provides safe abstraction in the types it wraps in both the Rust way and the Vulkan API way.

Table 5.5: Vulkan API validations status in the project.

Category	Statically solved	Dynamically solved	Left to user	Total
Implicit	317	28	2	347
Creation	91	0	314	405
Usage	29	3	122	154

Category	Statically solved	Dynamically solved	Left to user	Total
Total	437	30	439	906

In tbl. 5.5 it can be seen that the goal was achieved almost perfectly. Only two implicit validations are left to the user. This decision wasn't made lightly, but it was chosen as the most sensible one given the current limitations of the stable version of language. A small number of implicit validations couldn't be solved statically. These validations are instead checked at runtime, but only conditionally under the `runtime_implicit_validations` Cargo feature. All other implicit validations were successfully solved statically. More details about the specific validations can be found in the appendix.

Additionally, a significant amount of explicit validations, categorized under creation and usage, have been solved statically as a consequence of the natural API design and/or the implicit validations. Overall this means increased safety for the user of the API at no runtime cost.

6 Conclusion

The core Vulkayes library is successful at reducing the complexity of creating and using Vulkan types, as well as correctly destroying them at appropriate times and checking basic safety requirements. Benchmarks show that this added complexity is mostly compile-time and scales well into the runtime where applicable. Additionally, safety is guaranteed at a certain level that should provide the user of the API with certain amount of confidence that their application will not sefault. Overall, the Vulkayes project is a good step towards a flexible and transparent Vulkan API in the Rust ecosystem.

However, there still remains a lot of work to be done to create an API with a application design advantage as well. Designing synchronization in Vulkan by hand is error prone due to high complexity and Vulkayes should be extended with user-friendly API that is capable of lifting the burden off the user onto the implementation. This and other improvements to Vulkayes are left for future work.

Bibliography

- [1] “Vulkan® 1.2.136 - A Specification.” [Online]. Available: <https://www.khronos.org/registry/vulkan/specs/1.2-extensions/html/vkspec.html>. [Accessed: 05-Apr-2020]
- [2] “Vulkano.” [Online]. Available: <https://github.com/vulkano-rs/vulkano>. [Accessed: 10-Apr-2020]
- [3] “Khronos Releases Vulkan 1.0 Specification.” [Online]. Available: <https://www.khronos.org/news/press/khronos-releases-vulkan-1-0-specification>. [Accessed: 29-Apr-2020]
- [4] “Khronos Releases Vulkan 1.2 Specification.” [Online]. Available: <https://www.khronos.org/news/press/khronos-group-releases-vulkan-1.2>. [Accessed: 29-Apr-2020]
- [5] “V-EZ.” [Online]. Available: <https://github.com/GPUOpen-LibrariesAndSDKs/V-EZ>. [Accessed: 10-Apr-2020]
- [6] “gfx-hal.” [Online]. Available: <https://github.com/gfx-rs/gfx>. [Accessed: 10-Apr-2020]
- [7] R. Galajda, “Designing a modern high-level graphics API,” Czech Technical University in Prague. Computing and Information Centre., 2020.
- [8] “The Development of the C Language.” [Online]. Available: <https://www.bell-labs.com/usr/dmr/www/chist.html>. [Accessed: 10-Apr-2020]
- [9] “The Cargo Book.” [Online]. Available: <https://doc.rust-lang.org/cargo/>. [Accessed: 30-Apr-2020]
- [10] “ash.” [Online]. Available: <https://github.com/MaikKlein/ash>. [Accessed: 29-Apr-2020]
- [11] “Vulkan Memory Allocator.” [Online]. Available: <https://github.com/GPUOpen-LibrariesAndSDKs/VulkanMemoryAllocator>. [Accessed: 30-Apr-2020]

