

FRANKFURT UNIVERSITY OF APPLIED SCIENCES

MASTER'S THESIS

Containerized multi-level deployment for a distributed adaptive microservice application

Author:

Tim WIGMANN

Supervisors:

Prof. Dr. Peter THOMA

Prof. Dr. Eicke GODEHARDT

*A thesis submitted in fulfilment of the requirements
for the degree of Master of Science*

in the course

Allgemeine Informatik Master

March 17th, 2023

Declaration of Authorship

I, Tim WISMANN, declare that this thesis titled, 'Containerized multi-level deployment for a distributed adaptive microservice application' and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

Date:

“Thanks to my solid academic training, today I can write hundreds of words on virtually any topic without possessing a shred of information, which is how I got a good job in journalism.”

Dave Barry

FRANKFURT UNIVERSITY OF APPLIED SCIENCES

Abstract

Faculty 2 - Computer Science and Engineering

Allgemeine Informatik Master

Master of Science

Containerized multi-level deployment for a distributed adaptive microservice application

by Tim WILDMANN

The Thesis Abstract is written here (and usually kept to just this page). The page is kept centered vertically so can expand into the blank space above the title too...

Acknowledgements

The acknowledgements and the people to thank go here, don't forget to include your project advisor...

Contents

Declaration of Authorship	i
Abstract	iii
Acknowledgements	iv
Contents	v
List of Figures	vii
List of Tables	viii
Abbreviations	ix
1 Introduction	1
1.1 Scope	1
1.2 Intended audience	1
1.3 Outline	1
2 Background and related work	2
2.1 Application under study	2
2.2 Baseline architecture	2
2.3 Problem statement	7
2.4 Limitations	8
2.5 Related Work	8
3 System design	9
3.1 Orchestration engine	9
3.1.1 Hyper-V Replication	10
3.1.2 Docker Swarm	10
3.1.3 Kubernetes	10
3.2 Kubernetes	10
3.2.1 Entities	11
3.2.2 Services	11
3.2.3 Pod life cycle	13
3.2.4 Cluster networking	13
3.3 Container environment	15

3.3.1	ContainerD	15
3.3.2	Docker Container Runtime Interface	16
3.4	Container Image	17
3.4.1	Base image	17
3.4.2	Custom image	18
3.5	Target architecture	18
4	Implementation	19
4.1	Containerization of Services	19
4.1.1	Code changes on core application	19
4.1.2	Container definition	22
4.2	Cluster Setup	22
4.2.1	Creating the master node	22
4.2.1.1	Installing a Container Network Interface	23
4.2.1.2	Adding the proxy daemonsets	23
4.2.2	Creating the worker node	23
4.3	Automatic setup	24
5	Results	25
5.1	Containerization	25
5.1.1	Container manifest	25
5.1.2	Windows Base Image	26
5.2	Discovered issues and pitfalls	26
6	Discussion	27
6.1	Analysis	27
7	Conclusion and future work	28
7.1	Future work	28
7.2	Conclusion	28
A	Overview about submitted code	29
	Acronyms	31
	Bibliography	33

List of Figures

2.1	The OpenTwin (OT) project overview.	3
2.2	The OT login screen.	3
2.3	A opened project inside OT with a few created geometric models and subtracted computation.	4
2.4	Communication overview and service organization for OT main services. In 1.1 the Local Session Service (LSS) registers at Global Session Service (GSS). As soon as the User Interface (UI) front end connects to the GSS (2.1), service information is exchanged (2.2) and the user is authenticated (2.3). As a consequence, the GSS creates a new session and tells the LSS to spawn new compute services. From now on the UI front end communicates directly with the Compute services via the Relay Service over a web socket connection.	6
2.5	Service initialization of OT processes. In the beginning, the main services GSS, Authorization Service (AUTH) and an optional LSS are initialized. While the GSS checks the reachability of AUTH, the LSS registers itself at the GSS. After starting the UI front end, the service information is requested from a GSS and the user is authenticated. After success, the UI front end connects to the LSS and requests a new session. As consequence, the LSS spawns the compute services and connects them to the UI front end via a Relay Service. (Ping messages are omitted.)	7
3.1	Core and compute services for Kubernetes[1]	12
3.2	Abstraction layers for ContainerD on Windows. The image shows the technology stack from the Docker and Kubernetes (K8s) command line to the Container layer.[2]	16

List of Tables

2.1	List of compute services and their corresponding tasks.	5
3.1	List of K8s states during pod life cycle[3].	13
A.1	List of files in the code repository and their usage description.	30

Abbreviations

LAH List Abbreviations Here

Chapter 1

Introduction

This thesis is about the deployment of a microservice application.

1.1 Scope

1.2 Intended audience

1.3 Outline

Chapter 2

Background and related work

2.1 Application under study

OpenTwin ([OT](#)) is a open-source simulation platform developed by the university of Applied Sciences in Frankfurt, Germany. It covers features like computer aided design and meshing and also is a physics simulation (having solvers for Finite Integration and PHREEC). The projects can be administered by a user and group management (see [Figure 2.1](#)). Furthermore, all changes on a project are version-controlled. The application is designed in a way, that only a local thin-client needs to run on the users computer. After entering the login credentials (see [Figure 2.2](#)), the client securely connects to a centralized service platform where the computation is made. The results and even the [UI](#) information is sent back to the client application. This has the benefit, that also weak computers can run the application.

[Figure 2.3](#) shows the application itself with a loaded project and a simple geometric model.

The development team consists of a small core team and several student groups during the semester.

2.2 Baseline architecture

The current system design consists of multiple levels. It is a multi-process application based on the programming languages C++ and Rust. The source code is mainly aligned to be built on Microsoft Windows® ([Windows](#)). A port to Unix based systems is currently in

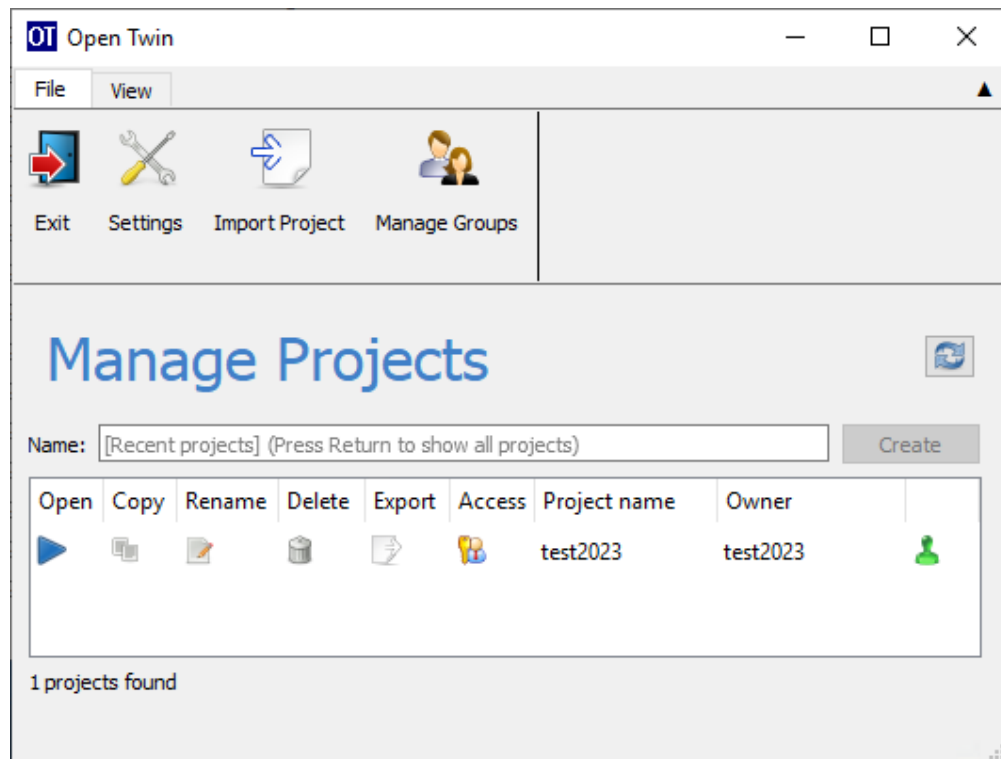


FIGURE 2.1: The OT project overview.



FIGURE 2.2: The OT login screen.

work. Therefore parts of the code base are aligned for multiple system architectures already, but the application is not yet able to be compiled for Linux.

Each microservice of the application is included dynamically and linked as a Dynamic Link Library (DLL) file. For starting the microservice environment, a central executable ("open_

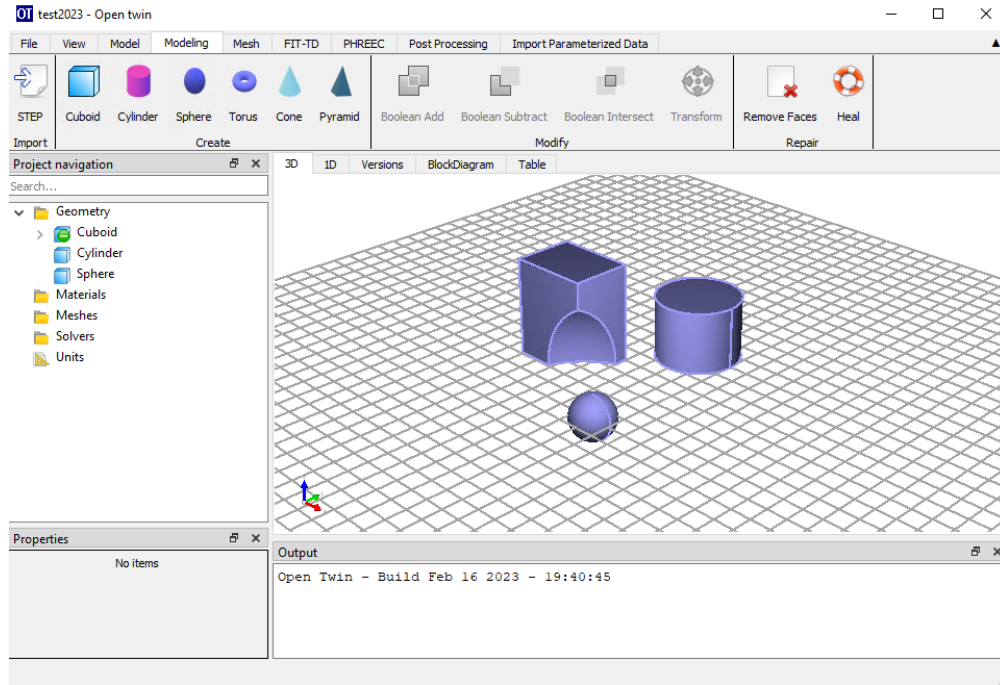


FIGURE 2.3: A opened project inside OT with a few created geometric models and subtracted computation.

twin.exe”) is started with the corresponding arguments for the services (like the server’s binding address, port numbers, and passwords) (see Listing 2.1) and the path to the DLL file itself. The UI front end, which is started by the user directly, is compiled in its own executable (“uiFrontend.exe”).

```

1 open_twin.exe GlobalSessionService.dll \
2   "0" "127.0.0.1:8091" "tls@127.0.0.1:27017" "127.0.0.1:8092"

```

LISTING 2.1: Command line of Open Twin Service start

For conveniently running the services with all their necessary arguments, batch files were provided that read environment variables and convert them into runtime arguments for the service executable. Therefore, if the services are started locally, the user runs a batch file that sets up the environment for the network binding details, path to certificates and encrypted database credentials.

The system consists of the following micro services that are permanently accessible: Global Session Service (GSS), Authorization Service (AUTH) and the database. The database is running on MongoDB¹. Another Service is the Local Session Service (LSS) that spawns the so called compute services. Those are services for running the actual computation after opening

¹MongoDB: <https://www.mongodb.com/>

a project that can dynamically spawn and exit. A partial list of compute services and their corresponding tasks can be found in Table 2.1. Each service runs in its own operating system process.

Name	Task
CartesianMeshService	If demanded, it converts a continuous geometry into a discrete Cartesian mesh.
FITTDService	If demanded, it runs a solver algorithm for transverse electromagnetic simulation based on the finite integration technique (FIT).
KrigingService	If demanded, it runs a kriging interpolation of result data.
LoggerService	A background service, that accepts logging messages from other services.
ModelingService	Performs calculations for the creation, modeling and boolean combination of geometric data.
PHREECService	If demanded, it runs simulation based on PHREEC.
TetMeshService	If demanded, it meshes a form with an tetrahedral mesh.
VisualizationService	Runs the graphical calculations for displaying the geometric and data based results on the UI.

TABLE 2.1: List of compute services and their corresponding tasks.

As shown in Figure 2.4, the services can be separated by their network space. Not all services require a public available network address. While GSS, AUTH and database are globally accessible via a fixed network address, the LSS can theoretically run on a dedicated host and is only communicated to other parties after it has registered itself to the GSS. The services, spawned by LSS do not require a public address space either. All communication between the UI front end and the compute services is achieved via a relay service and a web socket communication channel.

The whole process of the LSS registration and connection of the UI front end to the compute services is depicted in Figure 2.5. Once started, the user can login. In order to connect to the database, the following steps are performed:

1. The UI front end requests further service information from the publicly available GSS. The address for this service is provided by the user. The GSS responds with Uniform Resource Locators (URLs) to the database and the AUTH.
2. The UI front end connects to the AUTH using the authentication information provided by the user.

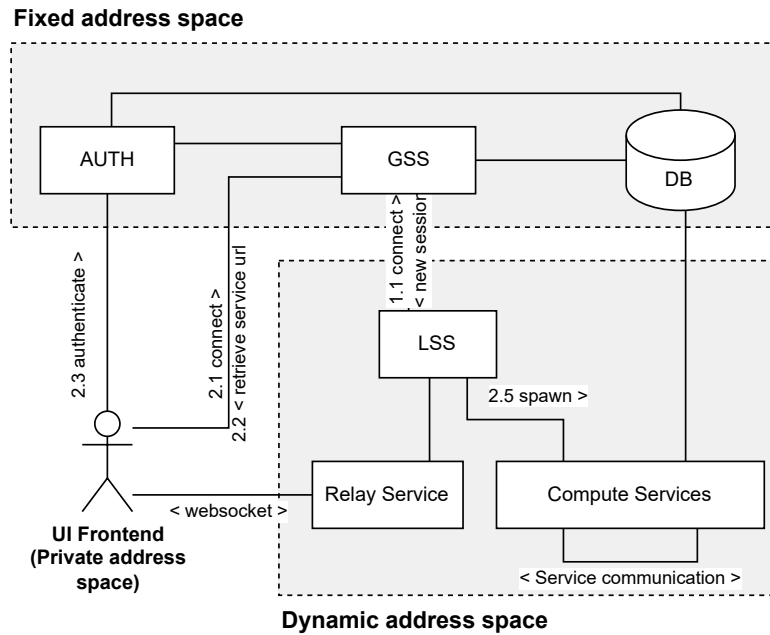


FIGURE 2.4: Communication overview and service organization for OT main services. In 1.1 the LSS registers at GSS. As soon as the UI front end connects to the GSS (2.1), service information is exchanged (2.2) and the user is authenticated (2.3). As a consequence, the GSS creates a new session and tells the LSS to spawn new compute services. From now on the UI front end communicates directly with the Compute services via the Relay Service over a websocket connection.

3. If the AUTH replies with a positive authentication, the UI front end connects to the database and lists the projects.
4. Once a project is opened or created, the UI front end requests a new session from the GSS. The GSS replies with the connection URLs of the LSS. The LSS has been registered to the GSS during its initialization.
5. The UI front end then connects to the LSS and requests a new session. As a result, the LSS spawns new application service processes and replies with the respective service URLs.
6. From now on, the UI front end communicates with the application services via the Relay service over a web socket.

The traffic between services is encrypted using mutual Transfer Layer Security (mTLS) technology. While regular Transfer Layer Security (TLS) ensures the authenticity of the server by using Certificates and the chain of trust, it does not verify the identity of the client. This is the benefit of mTLS. In mTLS, both sides, client and server has to verify their identity by providing a certificate inherited from a common root authority.

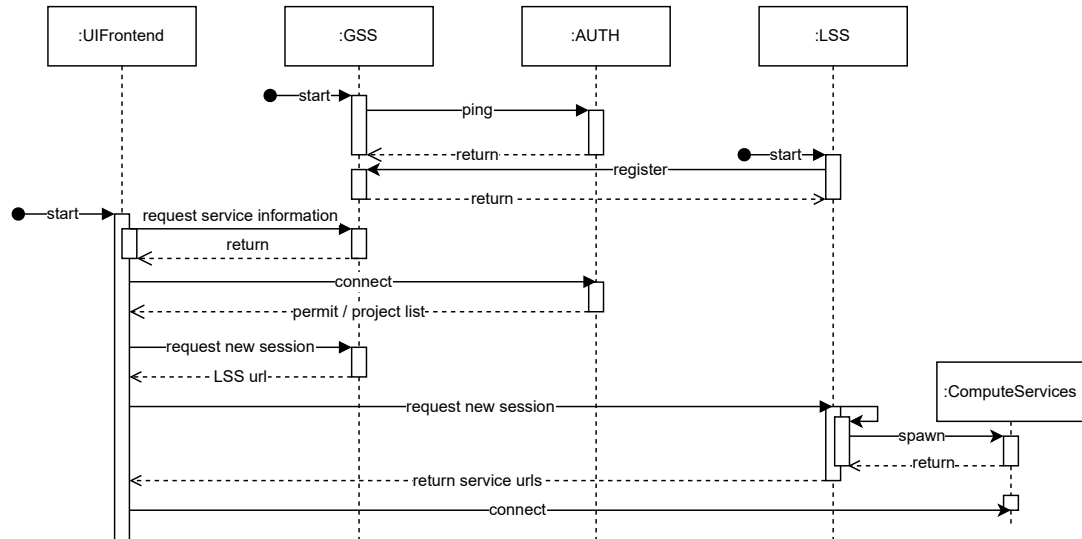


FIGURE 2.5: Service initialization of OT processes. In the beginning, the main services GSS, AUTH and an optional LSS are initialized. While the GSS checks the reachability of AUTH, the LSS registers itself at the GSS. After starting the UI front end, the service information is requested from a GSS and the user is authenticated. After success, the UI front end connects to the LSS and requests a new session. As consequence, the LSS spawns the compute services and connects them to the UI front end via a Relay Service. (Ping messages are omitted.)

2.3 Problem statement

Even though, the application is clearly based on a microservice architecture and it is able to run on a distributed system, it is not designed for an automated cluster yet. It consists of multiple processes where many of them have to run on the same system and need a full working operating system as baseline. Containerization of the system has never been tested and needs to be introduced. First, the cluster engine needs to be set up for Windows compute nodes to allow Windows containers to run inside the cluster. Additionally, it needs to have full network capabilities as well as inter connectivity between the several services. Next, the application needs to run inside containers. Therefore, container images must be created and provided for the cluster engine. Additionally, the automatic extension of services requires communication between the cluster orchestration management and the applications running on the nodes. A feature that needs to be introduced later on.

Regarding logging, while the front end application does, the microservices currently do not produce log files. Instead, only a few sub processes write the information on its standard output stream. In some cases, the error information given by exceptions is dropped. Furthermore, proper exit codes in error cases are not returned. That is, if the application exits there is currently no way to detect if the process terminated normally or crashed as part of an error.

2.4 Limitations

Due to the limited amount of time, not all code changes are applied. On the one hand, this involves the adaption for automatic extension of services. On the other hand, it implies the changes required to make the application more fault-tolerant. The changes that would be necessary, would be too extensive. Therefore, they are only made to the main processes.

As a first case study, the application is not fully containerized. Since the network connectivity is known to cause troubles in [Windows](#) container networks, there is more investigation required later on. As part of this study, only the main services are containerized and the cluster is set up to investigate the behavior in cluster environments. The actual distribution of an full functional cluster network can be part of further studies later on.

2.5 Related Work

Chapter 3

System design

Various applications for realizing the architecture have been compared. In the following sections the different options that were taken into account are presented.

3.1 Orchestration engine

Orchestration engines aggregate the processes and tools that are used to distribute services across multiple machines. Further, multiple replications are provided to maintain reliability. In addition, some solutions offer load balancing of incoming requests and network interconnection. What all of these engines have in common is that a group of virtual machines or containers, known as "nodes", are managed from a central spot. An administrator directs what application is run on the cluster. Based on the application's metadata, the orchestration engine then decides where to run the application by selecting a node inside that cluster.

3.1.1 Hyper-V Replication

Microsoft [Windows](#) supports a replication mechanism for virtual machines hosted by Hyper-V. The existing virtual machines are mirrored to secondary virtual machine host servers which increases scalability and reliability. The replications are replicated to a secondary Hyper-V host server, enabling process continuity and recovery on outages. Although there are benefits, like scalability and recovery, Hyper-V is mainly designed for virtual machines. Therefore, the cluster management solution is not applicable on this use case.

3.1.2 Docker Swarm

"Docker Swarm" is a cluster and orchestration engine for the container service "Docker". The offered extension mode has more features compared to the Hyper-V replication and is specialized for containers. For example, Load Balancing, increased fault tolerance and automatic service discovery. A highlighted feature among Docker Swarm is the decentralized design. That means, manager and application service can both run on any node within the cluster. Since it comes with Docker, no additional installation is required if Docker is already installed on the system. However, since it is bound to the Docker Application Programming Interface (API), using this orchestration technology involves the risk of inflexibility later on ("vendor lock-in").

3.1.3 Kubernetes

Kubernetes (K8s) is a orchestration engine similar to "Docker Swarm". Load balancing, auto-scaling and automatic service discovery are also offered. However, K8s additionally comes with the ability to rollback to a previous version in a product lifecycle and has built-in support for auto-scaling. However, K8s has more sophisticated configuration options which makes it harder to configure in the beginning.

The engine of choice was K8s because of its rich feature set. Also studies showed that K8s outperforms Docker Swarms when it comes to performance. For example, Marathe et. al. [4] compared a simple web server service deployed on a Docker Swarm cluster with a K8s cluster. The results showed better performance for K8s in terms of memory consumption and CPU usage. Another study of Kang et. al. [5] compared the performance of Docker Swarm and K8s in a limited computing environment on Raspberry Pi boards. They also concluded that K8s outperforms Docker Swarm if used with a high amount (=30) of service containers on 3 Pi boards [5]. Since they focused on container distribution and management methods this might get handy in the use case scenario under study.

3.2 Kubernetes

Since K8s is the chosen orchestration engine, the following sections are taking a deeper look inside its architecture.

3.2.1 Entities

There are many entities for objects inside the cluster. For description of those entities the configuration language YAML is used. Some of the most widely used entities are described in the following paragraphs.

Pod A pod represents a set of running containers on a node. Each pod has additional information stored, such as Health state, the cluster internal network Internet Protocol (IP) address or the amount of replications.

Deployment Deployments are used to define declarative states for Pods. This allows to maintain consecutive versions of the pod and upgrade them during runtime.

Daemon set These ensure that multiple (or all) nodes run a certain pod[6]. Common use cases are tasks for all nodes or running the network overlay pod.

User This entity describes a user that can access the K8s cluster and API services. Users can be part of a group and permission roles.

Node A node represents a physical machine inside the cluster. Nodes can run multiple pods.

3.2.2 Services

K8s comes with a set of core services (see Figure 3.1) that ensure the life span of scheduled containers, and the application services that offer the actual application.

In the following paragraphs, the crucial services are described in detail. Since every service is a pod, they can have multiple replicas. Only the core service have to run on a dedicated Linux node the so called "control plane node". The other services can run on nodes for executing the applications and perform computations ("compute node").

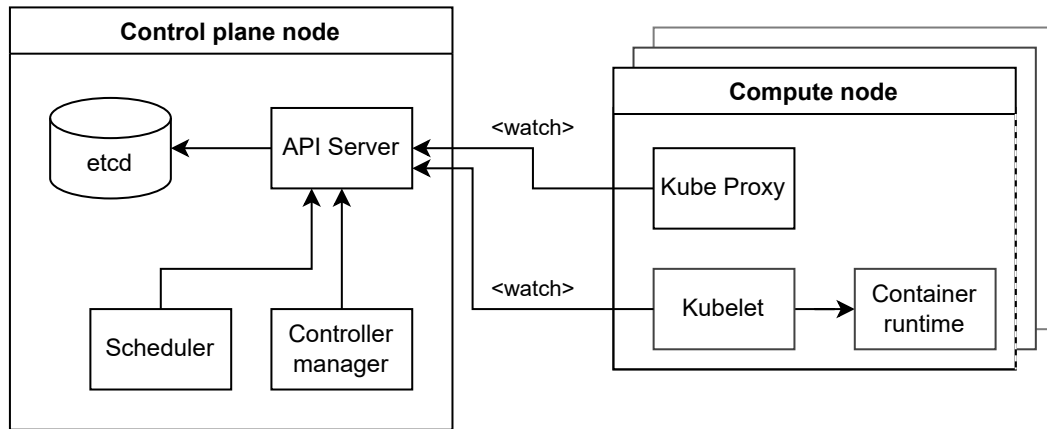


FIGURE 3.1: Core and compute services for Kubernetes[1]

etcd The etcd¹ database server is a key-value store designed for distributed systems[1]. That means it could run with multiple replications and would still be able to keep a persistent storage synchronized across multiple instances. It contains the applied configuration of several cluster entities (e.g. User configurations, deployments, pod configurations).

API server This is a RESTful web server that serves the Kubernetes [API](#) via Hypertext transfer protocol ([HTTP](#))[7]. It is the central joint between the services and establishes communication between users, external components and other core services. It makes the objects stored in etcd accessible over an Open [API](#) specification[1, 8] and allows observing changes on the entities. The Command line interface ([CLI](#)) tools "kubectl" and "kubeadm" both interact with the [API](#) server.

Kubelet Kubelet is the service on the operating system level that maintains the pod life cycle and ensures the runtime of a container inside a pod. Furthermore, it manages the registration of the node to the control plane and reports its health and pod status to the [API](#) server.

Kube Proxy The Kube-Proxy runs as a separate pod on every compute node. It maintains the connectivity between the services and pods[1]. For a given [IP](#) address and port combination it assures the connection to the corresponding pod. If multiple pods can offer a service, the proxy also acts as a load balancer[1].

¹etcd: <https://etcd.io/>

Scheduler The scheduler is responsible for distributing services on the cluster and determining which node to choose during runtime. It reads conditions for scheduling (e.g. hardware resources, operating system, labels) from the [API](#) server and decides which node matches the configuration[1].

Controller manager While the [API-Server](#) is responsible for storing data in etcd and announcing changes to the clients, the Controller manager and its parts try to achieve a described target state[1]. The controller manager consists of several controllers for replications, daemon sets, deployments, volumes, and so on.

3.2.3 Pod life cycle

Similar to the underlying application container, Pods in [K8s](#) have a ephemeral lifetime[3]. After creation on the cluster, a unique identifier is assigned before a pod gets scheduled to an available node[3]. The pod keeps alive until its termination or deletion[3]. For distinguishing different kind of states of a pod life cycle, [K8s](#) defines the pod states as described in [Table 3.1](#).

State	Description
Pending	The pod has been set up, the container and pod is currently initialized.
Running	The pod is bound to a node, the container is running.
Succeeded	The container terminated with a zero exit code.
Failed	The container terminated with a non-zero exit code or was terminated by the system.
Unknown	The pod state could not be obtained.

TABLE 3.1: List of [K8s](#) states during pod life cycle[3].

A terminated pod automatically gets restarted based on a configured restart policy. As the [K8s](#) documentation states, "the kubelet restarts them [the containers of a pod] with an exponential back-off delay (10s, 20s, 40s, ...), that is capped at five minutes"[3]. Furthermore, it is explained that the back-off time gets reset, once a container keeps running for 10 minutes[3].

3.2.4 Cluster networking

Cluster networking is achieved using two components: The network plugin and the Container Network Interface ([CNI](#)). Pods receive their own [IP](#) address and can communicate with other

Pods. However, this is not a functionality which is achieved by Kubernetes directly. By using a [CNI](#) the automated generation of network addresses and their inclusion is achieved when new containers are created or destroyed. It is crucial that pods share the same subnet across all the nodes in a cluster and Network Address Translation (NAT) is avoided[1].

Network plugins do implement the [CNI](#). They usually come with a manifest for a daemon set that introduces a network agent on all nodes inside the cluster to support the network communication. For setting up the network interface, namespace and its [IP](#) address, a dedicated container image is used. This is called the "pause container" image.

Even though, the team behind [K8s](#) do not recommend any specific network plugin, there are only a few common network plugins widely used. For this case, Flannel is used as network plugin, since it is the described plugin used in the documentations for setting up [K8s](#) with [Windows](#) containers[9, 10]. Hence, support for this [CNI](#) plugin in relation with [Windows](#) containers is larger.

Flannel Flannel² was originally developed as part from Fedora CoreOS³[11]. It works with various backends for transferring packets in the internal network. Two possible backends are virtual extensible Local Area Network ("vxlan") and host gateway ("host-gw"). While "host-gw" needs an existing infrastructure and performs routing on the layer 3 network level, VXLAN is more flexible and could also be used in cloud environments[12]. VXLAN is an overlay protocol and encapsulates layer 2 Ethernet frames within datagrams[11]. It is similar to regular VLAN, but offers more than 4,096 network identifiers[11]. Thus, VXLAN is a good choice for highly scalable systems.

Calico Compared to Flannel, Calico⁴ is stated to be more performative, flexible and powerful[11, 13]. Calico comes with a sophisticated access control system[13] and more configuration options. However, its advanced configuration makes it hard to maintain long-term.

²Flannel: <https://github.com/flannel-io/flannel>

³CoreOS: <https://getfedora.org/en/coreos>

⁴Tigera's Calico: <https://www.tigera.io/project-calico/>

3.3 Container environment

The ecosystem around containerization defines terminology that needs to be looked at before going into details for [K8s](#). First of all, the Container Runtime Interface (CRI) defines the interface between [K8s](#) and container runtime. Most of the container runtimes follow the design principles defined by the Open Container Initiative (OCI)⁵ for describing images and containers. The actual container runtime runs the isolation layer between the physical host machine and the [K8s](#) cluster by using containerization of processes. This is what can be selected when working with [K8s](#).

While [K8s](#) used to support Docker as their standard container runtime, they announced it to be deprecated in 2020, and finally removed the support in February 2022[14, 15]. The teams behind [K8s](#) decided to drop the hard coded support for Docker and offer ContainerD instead. However, the specification for ContainerD's "Containerfile" has only minor differences compared to Docker's "Dockerfile". Thus, ContainerD files are fully compatible to docker files.

Some of the container runtimes offered by [K8s](#) are not available for [Windows](#) hosts. For example, Linux containers (LXC)⁶ use process groups, control groups (cgroups) and name spaces on the operating system level. The CRI from the Open Container Initiative (CRI-O)⁷ is another alternative offered for [K8s](#) on Linux systems. Since those are not available in [Windows](#), they are not further considered.

At the current time being, container networking with ContainerD is not well-established on [Windows](#) [16–19] even though the docker runtime is already removed in current versions of [K8s](#) [14]. However, these are the only two working container backends for [Windows](#) containers. Therefore ContainerD as container backend was chosen.

3.3.1 ContainerD

Containerd is a native version of a container runtime. Newer versions of Docker on Linux, are running Containerd under the hood for process isolation. On [Windows](#), ContainerD uses slim host process isolation. The process isolation with ContainerD consists of multiple abstraction layers (shown in [Figure 3.2](#)). Its back end contacts the containerd-shim which is maintaining

⁵OCI: <https://opencontainers.org/>

⁶LXC: <https://linuxcontainers.org/>

⁷CRI-O: <https://cri-o.io/>

an abstraction layer for communication for the underlying layers (depending on Linux and Windows). Below that, Windows offers a custom fork of the CLI *runc*, so called *runhcs*[2]. Using *runc*, new containers can be created by running a simple command[2]. The layer for *runhcs* connects to the Host Compute Service (HCS) which is another abstraction layer of Windows for providing a stable API to the low level functionality of the operating system[20]. Containerd does not come with any mechanisms for networking. Instead, this is in responsibility of the HCS.

The developers of K8s marked the Docker CRI as deprecated in version v1.20[14]. Since version v1.23 of K8s, Docker was fully removed which lead to ContainerD being the only available CRI for Windows containers.

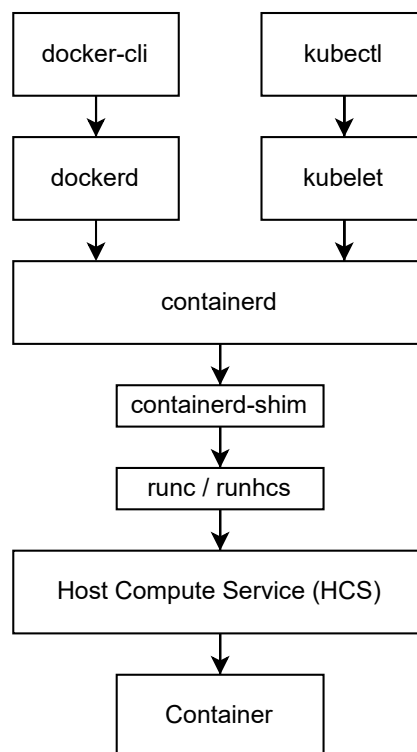


FIGURE 3.2: Abstraction layers for ContainerD on Windows. The image shows the technology stack from the Docker and K8s command line to the Container layer.[2]

3.3.2 Docker Container Runtime Interface

The Docker CRI (so called "Docker shim") is using the internal mechanisms from Docker to run containers. Older versions in Linux were using control group isolation.

For [Windows](#), there are two different modes available. The first option is using the process isolation mode offered by ContainerD. It can be enabled by switching to "Windows Containers". It is the default mode for Windows Server systems. However, older versions of Docker and Docker on Windows 10 and above has the opposite behavior[21]. On client versions of Windows, a dedicated hypervisor-isolated virtualization is the default option[21]. This means, during the installation of Docker on [Windows](#), the underlying Windows container host creates a separate Hyper-V Virtual Machine (VM) to run container images. However, this is not a regular Hyper-V VM. Instead, it is a purpose-built VM, often referred to as utility VM or UVM that can't be managed directly and is fully controlled by the [Windows](#) container runtime[21]. For networking, the internal mechanisms from Docker are used. All running containers are deployed inside the dedicated Hyper-V VM. Therefore, this is a mixture of process isolation and full isolation using virtualization.

However, the hypervisor approach still has the disadvantage of using large resources for containers, even though they are running in one virtual machine. In addition, containers running in hypervisor isolation take longer to start up than those running in process isolation[21].

3.4 Container Image

The container image consists of a base part and a part for custom configuration. Both are further explained in the following sections.

Container images are described, using the Containerfile⁸ format. There are container images for each process of the system architecture, each of them having their own Containerfile definition with different command line arguments and environment variables.

3.4.1 Base image

The container images to use for running the OpenTwin processes need to run a [Windows](#) base image. Beside the full [Windows](#) images, Microsoft® offers the more common images "Windows Server Core" and "Windows Nanoserver"[22]. They significantly differ in the download size, their on-disk footprint and the features supported[22]. As Microsoft states, "Nanoserver was built to provide just enough [API](#) surface to run apps that have a dependency on .NET core or

⁸Containerfile: <https://www.mankier.com/5/Containerfile>

other modern open source frameworks. PowerShell, Windows Management Instrumentation, and the [Windows](#) servicing stack are absent from the Nanoserver image”[22].

The design of containerization of [OT](#) envisages the usage of ”Server Core” as base image. Even though the ”Server Core” image is not the smallest base image, it provides full functionality for the required technologies for the current use case.

3.4.2 Custom image

On top of the base image, customizations and the actual application are applied. The binary files are included in the container image and added during build. The common [CRI](#)s only forward the output of processes with process id 1 to the host machine. Furthermore, this is also the only process the [CRI](#) is waiting for, to keep the container alive. Thus, instead of using the provided batch files to start the application services, the OpenTwin process is called directly with the appropriate command line arguments as command for the container. Therefore, the environment variables needs to be set up as part of the container file. The root certificate (certificate authority) is passed as file mount into the container later on.

3.5 Target architecture

The application needs to be distributed on multiple systems. Kubernetes supports application rollout only as container images. To be able to distribute the services on a cluster management tool a containerization of the application is necessary.

Chapter 4

Implementation

4.1 Containerization of Services

Container images are provided for running the application inside the [CRI](#). Definition of the container manifest is done in the "Containerfile"¹ format.

4.1.1 Code changes on core application

For containerizing the [OT](#) application, several changes were applied to its code base. This improves error handling and error tracing of the application and therefore simplifies development of the cluster. In the following sections, the several code changes are described in detail.

Introduction of exit codes

The microservice [DLL](#) files have return codes for different error cases. This is, for instance a code 200 in the [AUTH](#) for connection issues to the database. As with process exit codes, a zero return code indicates a successful termination. Even though the main executable "open_twin.exe" retrieves the return code, it did not convert these codes into proper process exit codes. Exit-Codes are a crucial part of the [K8s](#)'s life cycle management as described in [subsection 3.2.3](#).

The surrounding lines of code in the main executable are shown in [Listing 4.1](#).

¹Containerfile: <https://www.mankier.com/5/Containerfile>

```

85 let _result = initialize(siteid_c_str.as_ptr(), service_c_str.as_ptr(), db_c_str.as_ptr(),
    dir_c_str.as_ptr());

```

LISTING 4.1: Former code snippet from main executable for calling the microservice library code in Rust (/Microservices/OpenTwin/src/main.rs)

A new condition for non-zero exit codes was added as shown in [Listing 4.2](#). The variable `_result` contains the returned exit code of the library [DLL](#), whereas `lib_path` contains the path to the library [DLL](#). As can be seen, it is used to describe the affected service name in a error message (line 87).

```

86 if _result != 0 {
87     eprintln!("Library/Service initialization ({}:init()) failed with exit code {}",
        lib_path, _result);
88     std::process::exit(_result);
89 }

```

LISTING 4.2: Code changes in Rust main executable for additional treatment of exit codes (/Microservices/OpenTwin/src/main.rs)

Debug verbosity in launcher

By default Rust programs show a window for console output no matter if it was built in release mode or with debug parameters. However, the main executable uses conditional compilation to set configuration attributes about the windows subsystem. Rust provides a compilation option "debug_assertions" that is set to "true" for compilations without code optimization[23]. Therefore it is set to "true" if the application runs in debug mode. As shown in [Listing 4.3](#) with conditional compilation, this build option is checked and console output is only shown for debug builds.

```

2 #![cfg_attr(not(debug_assertions), windows_subsystem = "windows")]

```

LISTING 4.3: Conditional compilation for disabling console output in non-debug builds (/Microservices/OpenTwin/src/main.rs)

Even if output is shown on the console window for debug builds only, the error messages are not able to be read conveniently in cases where the application crashes. This aggravates debugging of regular application errors outside of a containerized environment. To avoid this

behavior, the launcher batch files were adapted. For this, a new parameter was invented to the batch file for running the batch file in verbose mode.

```

18 IF "%~1"==" /V" (
19     REM OT is opening console windows in debug build, so we want to pause them at the end
20     SET pause_prefix=cmd.exe /S /C "
21     SET pause_suffix=" ^& pause
22     ECHO ON
23 )

```

LISTING 4.4: Additional command argument for preventing close of window after termination
(/Microservices/Launcher/OpenTwin_session.bat)

As first step, a new command line argument has been introduced. If "/V" is appended to the start of the launcher batch file it will run with higher verbosity. As can be seen in [Listing 4.4](#), with "/V" appended, two new variables "pause_prefix" and "pause_suffix" are set (line 20-21). Furthermore, command output is enabled to debug the launcher file itself (line 22).

```

34 START "AUTHORIZATION SERVICE" %pause_prefix%open_twin.exe AuthorisationService.dll
    "%OPEN_TWIN_SERVICES_ADDRESS%:%OPEN_TWIN_AUTH_PORT%" "%OPEN_TWIN_MONGODB_ADDRESS%"
    "%OPEN_TWIN_MONGODB_PWD%"%pause_suffix%

```

LISTING 4.5: Additional command extension for preventing close of window after termination
(/Microservices/Launcher/OpenTwin_session.bat)

As can be seen in [Listing 4.5](#), the newly defined variables "pause_prefix" and "pause_suffix" are appended to the command of starting a service. This ensures that a service process is started in a new window and the command is followed by the [Windows](#) "pause" command to stop and wait for user interaction.

Enhanced logging and error tracing

For improving the error tracing, the overall log amount has been increased. This involves enabling the logging of the central logging functions inside the library "OpenTwinCommunication". The logging mechanism did not use logging to standard output. Instead, calling the respective functions for logs were just ignored. This was changed and logging to standard output has been introduced. Additionally, more calls to the logging mechanism, including caller information, have been added. For instance, the [AUTH](#) now shows error messages for caught exceptions and errors during database initialization. Also, the general error tracing has been

improved. In the main services [GSS](#), [LSS](#) and [AUTH](#), exception handling has been added, where it was not present before. Furthermore, the formatting of exception messages was improved and more information has been added.

For the current time being, there is an ongoing work to replace the logging to standard output by forwarding the log lines to a central logging library.

Certificate changes

- Certificate template was extended by localhost/127.0.0.1
- OpenTwin executable (rust server) listens on all interface (because of Docker); Port parsing - Bugfix for OpenGL error in uiFrontend

4.1.2 Container definition

4.2 Cluster Setup

The following section describes the setup of different machines in the cluster, so called nodes. While the master node refers to the [K8s](#) Control-Plane node which is responsible for distribution of the workers, the worker nodes are the actual machines that are executing the applications. During development the cluster was set up on virtual machines completely, due to the lack of physical hardware.

4.2.1 Creating the master node

For setting up the master node on Linux a system based on Debian Bullseye 11.5 has been used. After installing and setting up the operating system, the swap mechanism needs to be permanently turned off. This is done by editing the file system table (fstab) in file `/etc/fstab` respectively by commenting out the swap partitions and masking the systemd swap units.

After installing the pre-requisite packages, a containerd config file needs to be created. For this, the command from [Listing 4.6](#) is applied.

```
sudo sysctl net.bridge.bridge-nf-call-iptables=1
echo 1 > /proc/sys/net/ipv4/ip_forward
```

```
sudo containerd config default | sudo tee /etc/containerd/config.toml &>/dev/null
```

LISTING 4.6: Bash command for setting up containerd config

Afterwards the systemd `cgroup` is added to the runtime options of containerd and the its service is restarted. After setting up the prerequisites, the cluster can be initialized by running the command line tool as shown in Listing 4.7 with the appropriate configuration as parameter.

```
sudo kubeadm init --config config.yaml
```

LISTING 4.7: Bash command for setting up the cluster

4.2.1.1 Installing a Container Network Interface

After successfully running the initialization, the cluster overlay network `flannel` needs to be setup. This is required for working with Windows worker nodes. To setup `flannel` the respective pod description can be directly downloaded from the vendor². In the configuration the VNI (4096) and port (4789) for Flannel on Windows were set. Afterwards the configuration has been applied on the cluster. After linking `kubectl` to the local control plane node, the successful setup of the cluster can be checked with the `kubectl` command.

4.2.1.2 Adding the proxy daemonsets

For using Windows worker nodes in a Kubernetes cluster, additional configmaps and daemonsets need to be applied on the cluster. Those are used for setting up a proxy for `flannel`.

```
curl -L https://github.com/kubernetes-sigs/sig-windows-tools/releases/latest \
  /download/kube-proxy.yml | sed 's/VERSION/v1.25.3/g' | kubectl apply -f -
kubectl apply -f https://github.com/kubernetes-sigs/sig-windows-tools/releases \
  /latest/download/flannel-overlay.yml
```

4.2.2 Creating the worker node

On the Windows worker node, the prerequisites were installed first. While installing the prerequisite `crictl` it was added to the `PATH` environment variable. After the setup, the node preparation scripts of the Kubernetes Special Interest Group (SIG) were retrieved. Before running

²<https://raw.githubusercontent.com/flannel-io/flannel/master/Documentation/kube-flannel.yml>

the scripts, the [CNI](#) version needs to be aligned to the same version written on the master node. Therefore the appearances of v0.2.0 have been replaced with v0.3.0 respectively. Afterwards the installation script was executed. It sets up the [NAT](#) configuration on the worker node and registers containerd as a service. For having a valid image at a later point in time, the value `sandbox_image` in `cri` in containerd's configuration file (`config.toml`) needs to be replaced with a newer version.

After successfully setting up the [NAT](#) and installing containerd as a service, the node will be finally prepared to host tasks. For this, another powershell script "PrepareNode"³ from the Kubernetes [SIG](#) was run. After running the script the resulting "StartKubelet" file needs to be changed to drop invalid arguments.

Furthermore, the following lines were added to the Kubelet configuration:

```
enforceNodeAllocatable: []  
cgroupsPerQOS: false  
enableDebuggingHandlers: true
```

Since the configuration changes needs to be served only to Windows machines (config value are valid for Windows only) we (!!WE!!) need to manually change the configuration on the nodes locally. This

After successful run of the preparation script the node was ready to join the cluster.

4.3 Automatic setup

³<https://github.com/kubernetes-sigs/sig-windows-tools/releases/download/v0.1.5/PrepareNode.ps1>

Chapter 5

Results

The development of the [K8s](#) cluster and the containerization of [OT](#) hides several pitfalls. These are discussed in this chapter.

5.1 Containerization

The issues that require consideration during the containerization of the application are discussed in this section.

5.1.1 Container manifest

Even though the format of the container manifests "Containerfile" is compatible to the proprietary "Dockerfile" format from Docker, the [CRIs](#) do not follow the specification everywhere[24]. This is an issue while writing a Containerfile for the ContainerD [CRI](#). Especially in cases where line breaks in the Containerfile might be necessary to shorten long lines and increase readability. ContainerD is treating line breaks paths in string notation different compared to paths in JSON array notation and is not following the specification[25].

```
ENTRYPOINT open_twin.exe \  
Service.dll
```

LISTING 5.1: Containerfile entrypoint specification across multiple lines in text format.

[Listing 5.1](#) shows an example of the problem. While the entry point in Docker is interpreted as *open_twin.exe Service.dll*, the interpreter in ContainerD only reads the first line as entry point

and ignores the line break character "ä" in *open_twin.exe*. This causes troubles during build and execution of a container image. The container image is built in Docker, whereas it does not run in ContainerD. To overcome this flaw, the *ENTRYPOINT* definition has to be written in JSON notation. If defined as in [Listing 5.2](#) the interpreter

```
ENTRYPOINT [ "open_twin.exe",  
"Service.dll"]
```

LISTING 5.2: Containerfile entrypoint specification across multiple lines in JSON format.

5.1.2 Windows Base Image

While it is normal to have a requirement for the same operating system kernel when running container images, it is a uncommon requirement to have the exact same build version.

- Windows base image has to have same build number than host computer (not the same in linux containers)
- OpenTwin needs to be compiled on the system it is running. Developers need to fix this issue.

5.2 Discovered issues and pitfalls

- K8s Documentation redirects to a different page when searching for tutorial for adding windows nodes
- Error messages of Windows hcs shim are not explanatory
- Windows needs to have certain Update installed to run flannel container
- Scripts and docs are maintained by a relatively small community (SIG windows tools)
- Cluster networking throws weird error messages when performed inside hyper-v vm ("Directory not found")
- Docker support was removed, containerd is not fully supported yet (without bugs) in Windows

Chapter 6

Discussion

6.1 Analysis

Chapter 7

Conclusion and future work

7.1 Future work

- Linux port and cluster based on linux - Automated image building using Packer.io (for multiple platforms) - Ranger for streamlining kubernetes deployment - Images verkleinern - nur die Dateien ins Image bundlen, die für den entsprechenden Service notwendig sind.

7.2 Conclusion

Appendix A

Overview about submitted code

The created code and configuration file was uploaded to the project GitHub repository. The repository has private access only. Access to the repository can be granted by request.

The GitHub repository is located at: <https://github.com/pth68/SimulationPlatform>

The location of the additional files that are relevant for creation of the cluster are located in the **"Distribution"** sub directory of the repository. [Table A.1](#) describes the submitted files (relative from the "Distribution" directory) and their intention.

Path	Filename	Intention
Container/	build-image.ps1	Automation script for building, saving, and re-importing container images for ContainerD.
	setup-node.ps1	Script for automated installation of Kubernetes node prerequisites and setup by treating the flaws of a manual installation.
	authorisation.Containerfile	Containerfile for OpenTwin's Authorisation service.
	globalsession.Containerfile	Containerfile for OpenTwin's global session service.
	mongodb.Containerfile	Containerfile for the containerized MongoDB server.
	session.Containerfile	Containerfile for OpenTwin's local session service.
	compose.yaml	A container compose file for testing the setup of containerized services.
Container/mongodb/	0-init.js	A initialization script for setting up the admin user and roles of the OpenTwin database.
	mongodb.conf	The MongoDB Server configuration file.
Controlplane/	kubeadm_config.yaml	The configuration which should be passed during initialization of the cluster using kubeadm.
	open_twin.yaml	The K8s configuration file for deploying the cluster.
Kubernetes/	open_twin.yaml	The Kubernetes configuration file for deploying the cluster.

TABLE A.1: List of files in the code repository and their usage description.

Acronyms

OT OpenTwin

Windows Microsoft® Windows®

cgroup control group

SIG Special Interest Group

DLL Dynamic Link Library

URL Uniform Resource Locator

UI User Interface

TLS Transfer Layer Security

mTLS mutual Transfer Layer Security

GSS Global Session Service

AUTH Authorization Service

LSS Local Session Service

HTTP Hypertext transfer protocol

API Application Programming Interface

CLI Command line interface

K8s Kubernetes

IP Internet Protocol

NAT Network Address Translation

CNI Container Network Interface

CRI Container Runtime Interface

VM Virtual Machine

HCS Host Compute Service

Bibliography

- [1] Marko Lukša. *Kubernetes in action: Anwendungen in Kubernetes-Clustern bereitstellen und verwalten*. Hanser eLibrary. Hanser, München, 2018. ISBN 9783446456020. doi: 10.3139/9783446456020. URL <https://www.hanser-elibrary.com/doi/book/10.3139/9783446456020>.
- [2] Scooley. Windows container platform: Microsoft learn, 2022. URL <https://learn.microsoft.com/en-us/virtualization/windowscontainers/deploy-containers/containerd>.
- [3] Kubernetes. Pod lifecycle: Documentation, 2023-02-17. URL <https://kubernetes.io/docs/concepts/workloads/pods/pod-lifecycle/>.
- [4] Nikhil Marathe, Ankita Gandhi, and Jaimeel M. Shah. Docker swarm and kubernetes in cloud computing environment. In *Proceedings of the International Conference on Trends in Electronics and Informatics (ICOEI 2019)*, pages 179–184, Piscataway, NJ, 2019. IEEE. ISBN 978-1-5386-9439-8. doi: 10.1109/ICOEI.2019.8862654.
- [5] Byungseok Kang, Jaeyeop Jeong, and Hyunseung Choo. Docker swarm and kubernetes containers for smart home gateway. *IT Professional*, 23(4):75–80, 2021. ISSN 1520-9202. doi: 10.1109/MITP.2020.3034116.
- [6] Kubernetes. Daemonset: Kubernetes documentation, 2022-08-31. URL <https://kubernetes.io/docs/concepts/workloads/controllers/daemonset/>.
- [7] Kubernetes. The kubernetes api: Documentation, 2022-10-24. URL <https://kubernetes.io/docs/concepts/overview/kubernetes-api/>.
- [8] OpenAPI Initiative. Version 3.1.0 of the openapi-specification: Github - oai/openapi-specification, 2023-02-10. URL <https://github.com/OAI/OpenAPI-Specification/blob/main/versions/3.1.0.md>.

- [9] GitHub - Kubernetes SIG Windows Tools. Guide for adding windows node: Documentation, 2023-02-13. URL <https://github.com/kubernetes-sigs/sig-windows-tools/blob/727707fa7d83d401956b467a5ce41700cce7d9a3/guides/guide-for-adding-windows-node.md>.
- [10] Kubernetes. Adding windows nodes: Kubernetes - documentation (archived / v1.23), 2022-04-19. URL <https://v1-23.docs.kubernetes.io/docs/tasks/administer-cluster/kubeadm/adding-windows-nodes/>.
- [11] Suse Rancher Community. Comparing kubernetes cni providers: Flannel, calico, canal, and weave | suse communities, 2023-02-12. URL https://www.suse.com/c/rancher_blog/comparing-kubernetes-cni-providers-flannel-calico-canal-and-weave/.
- [12] GitHub - Flannel.io. Flannel backends documentation, 2023-02-12. URL <https://github.com/flannel-io/flannel/blob/master/Documentation/backends.md>.
- [13] Tigera. Project calico | tigera, 2023-02-10. URL <https://www.tigera.io/project-calico/>.
- [14] Kubernetes. Don't panic: Kubernetes and docker, 2020. URL <https://kubernetes.io/blog/2020/12/02/dont-panic-kubernetes-and-docker/>.
- [15] Kubernetes. Kubernetes is moving on from dockershim: Commitments and next steps, 2022. URL <https://kubernetes.io/blog/2022/01/07/kubernetes-is-moving-on-from-dockershim/>.
- [16] GitHub. Guide for adding windows node: Rbac config not found: Issue #261 - kubernetes-sigs/sig-windows-tools, 2023-02-02. URL <https://github.com/kubernetes-sigs/sig-windows-tools/issues/261>.
- [17] GitHub. Windows node with containerd can't run flannel and kubeproxy daemonsets: Issue #128 - kubernetes-sigs/sig-windows-tools, 2023-02-02. URL <https://github.com/kubernetes-sigs/sig-windows-tools/issues/128>.
- [18] amalinsky. Windows containers issue: Windows server 2022 container on windows 10 (10.0.19044) host., 2022. URL <https://github.com/microsoft/Windows-Containers/issues/258>.

- [19] GitHub. kube-flannel for linux fails to (re)start after rbac for kube-flannel for windows is applied: Issue #277 - kubernetes-sigs/sig-windows-tools, 2023-02-02. URL <https://github.com/kubernetes-sigs/sig-windows-tools/issues/277>.
- [20] Microsoft. Introducing the host compute service (hcs): Blog post, archived technet article, 2017. URL <https://techcommunity.microsoft.com/t5/containers/introducing-the-host-compute-service-hcs/ba-p/382332>.
- [21] Vinicius Ramos Apolinario. *Windows Containers for IT Pros: Transitioning Existing Applications to Containers for on-Premises, Cloud, or Hybrid*. Apress L. P, Berkeley, CA, 2021. ISBN 9781484266861. URL <https://ebookcentral.proquest.com/lib/kxp/detail.action?docID=6508410>.
- [22] Mattbriggs - Microsoft. Windows container base images: Microsoft documentation, 2023-02-14. URL <https://learn.microsoft.com/en-us/virtualization/windowscontainers/manage-containers/container-base-images>.
- [23] Conditional compilation: The rust reference, 2023-02-09. URL <https://doc.rust-lang.org/reference/conditional-compilation.html>.
- [24] failed to create containerd task: hcsshim::createcomputesystem kube-proxy: The directory name is invalid. URL <https://stackoverflow.com/questions/74799620/kubernetes-windows-worker-node-addition-failed-to-create-containerd-task-hcss>.
- [25] ManKier. Containerfile format: automate the steps of creating a container image, 2023-02-22. URL <https://www.mankier.com/5/Containerfile>.