

Master Thesis

Name: Fabian Burth

Topic: Design and Implementation of a Security and Compliance Data Lake

Place of work: SAP SE, Walldorf

Supervisor: Prof. Dr.-Ing. Vogelsang

Co-examiner: Prof. Dr. Körner

Deadline: 16/02/2023

Karlsruhe, 17/08/2022

The Chairman of the examination committee



Prof. Dr. Heiko Körner

Statutory declaration

I declare that I have composed the master thesis myself and without use of any other than the cited sources and aids. Sentences or parts of sentences quoted literally are marked as such; other references with regard to the statement and scope are indicated by full details of the publications concerned. The thesis in the same or similar form has not been submitted to any examination body and has not been published. This thesis was not yet, even in part, used in another examination or as a course performance. Furthermore I declare that the submitted written (bound) copies of the master thesis and the version submitted in digital format are consistent with each other in contents.

(Place, Date)

(Fabian Burth)

Abstract

The importance of software is constantly growing and so is its complexity. Thus, large-scale enterprise software systems are usually not developed from scratch. Commonly required functionality like logging or serialization is often imported through already existing *open source software*. Hence, modern software systems are composed of several components. These components usually have individual version updates, vulnerabilities, and licenses. Version updates of a component may introduce new vulnerabilities and licenses. Vulnerabilities in one specific component may compromise the security of the whole software system. Violating license agreements may lead to litigations due to copyright infringement. Furthermore, it must be considered that each component itself may be composed of several other components.

So, developing and maintaining secure large-scale enterprise software systems poses a considerable challenge to software companies. The state of the art approach to mitigate these risks is by incorporating compliance scanning tools into the CI/CD pipeline. These tools analyze the components within a software system and retrieve information like composition, vulnerabilities, and licenses. But this approach has its flaws. The components within a software system are commonly built through multiple different CI/CD pipelines maintained by individual development teams. These individual development teams frequently come up with technology specific point-to-point solutions of integrating the scanning tools into their CI/CD pipeline. Therefore, the information about the software systems of a company is frequently distributed over multiple development teams and therein again over multiple technology specific solutions.

The main objective of this thesis is the design and prototypical implementation of a central application to store this distributed information and to provide a consolidated view on all the software systems of a company. Through an API, the application should enable consumers to query this information to answer questions that might not even have been known at the time the compliance scan was performed. A recent and popular example for such a question is “Which systems contain Log4j?”.

Abstract (German)

Die Relevanz und Komplexität von Software nimmt kontinuierlich zu. Daher werden Unternehmens-Software-Systeme in der Regel nicht von Grund auf entwickelt. Standardfunktionalitäten wie Logging oder Serialisierung werden häufig durch bereits bestehende Open-Source-Software importiert. Moderne Software-Systeme setzen sich dementsprechend aus diversen kleineren Software-Komponenten zusammen. Diese Komponenten haben dabei individuelle Versionswechsel, Verwundbarkeiten und Lizenzen. Verwundbarkeiten in einer einzelnen Komponente können das gesamte System compromittieren. Verstöße gegen Lizenzvereinbarungen können zu teuren Gerichtsverfahren führen. Versionswechsel von Komponenten können weitere neue Verwundbarkeiten und Lizenzen einführen. Außerdem muss berücksichtigt werden, dass jede Komponente wiederum selbst aus weiteren Komponenten bestehen kann.

Die Entwicklung und Wartung sicherer Unternehmens-Software-Systeme stellt daher eine wesentliche Herausforderung für Softwareunternehmen dar. Der State-of-the-Art Ansatz, diese Risiken zu mitigieren, besteht im Integrieren von Compliance-Scanning Tools in die CI/CD-Pipeline. Diese Tools analysieren die Software-Komponenten eines Software-Systems und liefern Informationen zur Zusammensetzung, zu Verwundbarkeiten und zu Lizenzen. Dabei gibt es aber auch Probleme. Die Komponenten, aus denen sich ein Software-System zusammensetzt, werden häufig durch mehrere verschiedene CI/CD Pipelines erstellt, die von unterschiedlichen Entwicklungsteams gepflegt werden. Diese Entwicklungsteams entwickeln dabei oft technologiespezifische Punkt-zu-Punkt Lösungen, um die Scanning Tools in ihre CI/CD Pipeline einzubauen. Daher sind die Information über die Software Systeme eines Unternehmens oft über verschiedene Entwicklungsteams und in den Entwicklungsteams wiederum über verschiedene Technologien verteilt.

Das Hauptziel dieser Arbeit ist das Design und die Implementierung einer Applikation zur zentralen Speicherung dieser verteilten Informationen, die eine konsolidierte Sicht auf die Software Systeme eines Unternehmens ermöglicht. Durch eine API soll die Applikation die Abfrage von Informationen ermöglichen. Dabei sollen Fragen beantwortet werden können, die zum Zeitpunkt des Scans potentiell noch nicht bekannt waren. Ein aktuelles und populäres Beispiel für eine solche Frage ist: "Welche Systeme enthalten Log4j?".

Contents

Statutory declaration	ii
Abstract	iii
Abstract (German)	iv
Contents	v
List of Figures	viii
List of Tables	ix
List of Listings	x
Abbreviations	xii
1 Introduction	1
1.1 Motivation	1
1.2 Goals	2
1.3 Scope	3
1.4 Environment	3
1.5 Structure of the Thesis	4
2 Foundations	6
2.1 Software Identification	6
2.2 Open Source Software and Licensing	7
2.3 Vulnerability Management	8
2.4 Software Bill of Materials	12
2.5 Regulations	16
3 State of the Art	18
3.1 Software Development Life Cycle	18
3.2 Integrating Security and Compliance Measures	19

3.3	Producing Software Bill of Materials	21
3.4	Limitations	21
4	Reference Architecture	24
4.1	Open Component Model	24
4.1.1	Model Elements	25
4.1.2	Capabilities of the Open Component Model	30
4.2	Abstract Solution Idea	31
4.2.1	Abstract Data Lake Design	32
4.2.2	Reference Integration Architecture	33
4.2.3	Holistic Solution Architecture	35
5	System Design	38
5.1	Requirements	38
5.1.1	Functional Requirements	38
5.1.2	Non-functional Requirements	42
5.2	Data Model	42
5.2.1	Meta Data Model	44
5.2.2	Insights into the Development Process	47
5.2.3	Application of the Meta Data Model	48
5.2.4	Artifact Identity Problem	50
5.3	Database	53
5.3.1	Requirements for the Database	54
5.3.2	Relational Databases	55
5.3.3	Document Databases	65
5.3.4	Graph Databases	70
5.3.5	Database Selection	76
5.4	Programming Language	77
5.5	Application Architecture	77
5.5.1	Types API & Services	80
5.5.2	Ingestion API & Services	91
5.5.3	Query API & Services	92
5.5.4	Limitations and Problems	101
6	Result	104
6.1	Presentation	104
6.1.1	Application Infrastructure	104
6.1.2	Application Configuration & Data	105

6.1.3	Demonstration	107
6.2	Evaluation	111
6.3	Research Question and Retrospective	114
6.4	Outlook	117
Bibliography		119
A Additional Information		127
A.1	Black Duck Binary Analysis Result	127
A.2	Relational Model of the Example Data Model Instance	128
A.3	Relation of Number of Keys and Number of Leaves in a B-tree	129
B Data Definition Statements for Sample Data		131
B.1	Data Definition Statement for SQL Sample Data	131
B.2	Data Definition Statement for MongoDB Sample Data	132
B.3	Data Definition Statement for Neo4J Sample Data	133

List of Figures

Figure 2.1	CVE Record	9
Figure 2.2	CWE Hierarchy	10
Figure 2.3	C Project Directory Structure	14
Figure 2.4	ECCN Structure	16
Figure 3.1	CI/CD Pipeline as Stage-Gate System	19
Figure 4.1	Abstract Data Lake Design	33
Figure 4.2	Data Lake Integration	34
Figure 4.3	Holistic Solution Architecture	36
Figure 5.1	Meta Data Model	43
Figure 5.2	Data Model Instance	49
Figure 5.3	Local Artifact Identities	53
Figure 5.4	B ⁺ -tree as 4-way Search Tree	57
Figure 5.5	Package Versions and Dependencies as Table	61
Figure 5.6	Sales Data Model	66
Figure 5.7	Graph Representations	71
Figure 5.8	Graph Database Representation	72
Figure 5.9	Package Versions and Dependencies as Graph	74
Figure 5.10	Application Architecture	78
Figure 6.1	Neo4j Graph	105
Figure 6.2	Neo4j Component Tree	106
Figure 6.3	GraphQL Query Vulnerabilities in Component	109
Figure 6.4	GraphQL Query Components with Vulnerability	110
Figure 6.5	GraphQL Query Packages in Component	111
Figure A.1	Snippet of BDBA Analysis Result	127

List of Tables

Table 2.1 Minimum Elements of a SBOM	12
Table 5.1 Requirements for the Security and Compliance Data Lake . . .	41

List of Listings

2.1	Package URL Components	6
2.2	Package URL Examples	6
2.3	CPE Formatted String Binding	11
2.4	SPDX Document	14
4.1	Component Descriptor	26
5.1	SQL Query – Package Version Dependencies	61
5.2	SQL Query – Package Version Reciprocal Dependencies	62
5.3	SQL Query – Package Version Dependencies (transitive)	63
5.4	SQL Query – Package Version Reciprocal Dependencies (transitive) . .	63
5.5	Document Database – Data as Nested JSON Document	66
5.6	Document Database – MongoDB Join Example	68
5.7	Cypher Query – Package Version Dependencies (transitive)	74
5.8	Cypher Query – Package Version Reciprocal Dependencies (transitive) .	75
5.9	Example OpenAPI Document	80
5.10	Entity Type Schema Schema	84
5.11	Request Body for Component Version Schema Creation	85
5.12	Component Version Schema	87
5.13	Types API Endpoints	89
5.14	Ingestion API Endpoint	91
5.15	REST API Endpoints	93
5.16	REST API Query Endpoints	94
5.17	Ambiguous REST API URL	95
5.18	REST API Source Endpoint	95
5.19	Querying for Entity ID	95
5.20	Exposing Cypher through API	97
5.21	GraphQL Component Version Schema	98
5.22	GraphQL Query	99
5.23	Nested GraphQL Query	100
B.1	SQL Example Data Creation Statements	131
B.2	Document Database Example Data Creation	132

Abbreviations

AI Artificial Intelligence

B2B Business to Business

B2C Business to Consumer

BDBA Black Duck Binary Analysis

blob Binary Large Object

BOM Bill of Materials

CCL Commerce Control List

CD Continuous Delivery

CI/CD Continuous Integration/Continuous Delivery

CI Continuous Integration

CNA CVE Numbering Authority

CPE Common Platform Enumeration

CPU Central Processing Unit

CTE Common Table Expression

CVE Common Vulnerabilities and Exposures

CVSS Common Vulnerability Scoring System

CWE Common Weakness Enumeration

DDL Data Definition Language

DRAM Dynamic Random Access Memory

ECCN Export Control Classification Number

ERM Entity-Relationship Model

ERP Enterprise Resource Planning

FIRST Forum of Incident Response and Security Teams

Gi Gibibyte

GSD Global Security Database

IoT Internet of Things

MUX Multiplexer

NIST National Institute of Standards and Technology

NTIA National Telecommunications and Information Administration

NVD National Vulnerability Database

OAI OpenAPI Initiative

OAS Open API Specification

OCI Open Container Initiative

OCM Open Component Model

OSS Open Source Software

OSV Open Source Vulnerability

PoC Proof of Concept

purl Package URL

RAM Random Access Memory

RCE Remote Code Execution

RPC Remote Procedure Call

SaaS Software as a Product

SaaS Software as a Service

SBOD Software Bill of Delivery

SBOM Software Bill of Materials

SCA Software Composition Analysis

SCDL Security and Compliance Data Lake

SDLC Software Development Life Cycle

SPDX Software Package Data eXchange

SWID Software Identification

URI Uniform Resource Identifier

URL Uniform Resource Locator

WSL2 Windows Subsystem for Linux 2

Chapter 1

Introduction

"[T]he trust we place in our digital infrastructure should be proportional to how trustworthy and transparent that infrastructure is"

Executive Order on Improving the Nation's Cybersecurity [1]

1.1 Motivation

The introductory quote above initially sounds pretty intuitive and self evident. Probably everyone in the IT industry would be able to agree on this. Yet, the panic and public outcry unleashed in the software industry after the vulnerability in Log4j, a popular and widely used logging library, was discovered goes to show how far the statement strives from reality. The vulnerability is rated with the highest severity possible since it enables *Remote Code Execution (RCE)*, which, in other words, allows an attacker to run any code on the machine using Log4j [2]. After the discovery, IT specialists all over the world had to identify which of their applications and systems were using vulnerable versions of the library before they could even start dealing with the vulnerability. This increases the response time and subsequently the risk exposure significantly. But what is it, that makes this such a difficult task?

Modern software systems are composed of numerous software components, such as Log4j, and these components may themselves be composed of other software components and so on. This adds several layers of complexity obfuscating the dependencies of a software system. The manufacturing industry has been dealing with such issues in their supply chains for years. But these companies usually maintain close relationships and contracts with their suppliers, perhaps even exchanging *bills of materials (BOM)*. Thus, the companies can easily keep track of each individual part in their products by accumulating the BOMs of their suppliers and they will also be informed about problems affecting this parts by their suppliers. This is different from the software industry. Especially when it comes to *open source software (OSS)*, companies do not have a contract with the supplier. On the contrary, they frequently may not even know the maintainer of the software component. Subsequently, they

also will not be informed about issues such as vulnerabilities with the components.

The software industry has come up with approaches to deal with this issue and implemented measures to proactively monitor applications and detect known vulnerabilities in their components. With growing complexity of software and changing development and deployment landscapes, these approaches are not sufficient and fail to answer questions such as which systems use vulnerable versions of Log4j. Why and how they fail in these situations will be further examined in the course of this thesis. But as "[s]oftware is eating the world"[3] and thus, as companies, devices of our everyday life, cars and even medical and military devices rely on software, as the impact of software on our privacy and also on our physical security increases, so does the responsibility of companies providing software and the gap between trust and transparency becomes less acceptable.

In an effort to address this issue, the US government has published an Executive Order on Improving the Nation's Cybersecurity, which will require every company supplying software to the government to provide an *Software Bill of Materials (SBOM)* [1, 4]. As previously discussed, this is only possible if every company in the supply chain provides an SBOM for their software components. Therefore, the downstream impact of this Executive Order will most likely affect the entire industry.

1.2 Goals

The goals of this thesis aim to provide a research-based practical contribution to closing the gap between the trust and transparency placed into our digital infrastructure [1].

To make this more concrete, the main goal of this thesis is to develop a *proof of concept (POC)* implementation of a *Security and Compliance Data Lake (SCDL)*. As the name suggests, this is an application for storing information about software components. Such information may be about the occurrence of known vulnerabilities, as already indicated in the previous paragraphs, but also about licenses or dependencies. This kind of information about software components will be called metadata from here on.

The goal is to design and develop a central application for storing and querying software metadata, which is able to cope with the complexity and requirements of modern development and deployment landscapes. Therefore, improving the transparency of the entire software supply chain and enabling companies to answer questions such as: Which of our applications, systems or even landscapes contain a certain vulnerability?

Thereby, this thesis also contributes to the knowledge concerning metadata stores and the corresponding API design by answering the following research question:

What are the challenges arising during the development of a database application for the central metadata management of enterprise software systems and how may these challenges be solved?

In addition, the application may support companies in fulfilling the requirements resulting from the recently announced Executive Order.

1.3 Scope

As already stated in the goal section, this thesis is about designing and implementing a central application for software metadata management and its integration into modern development and deployment landscapes.

It is not about developing new ways of vulnerability detection. It is rather about monitoring systems and keeping track of already known vulnerabilities. This is just as important since a large amount of security incidents are caused by attackers exploiting such known vulnerabilities [5]. Therefore, neither does this thesis compete with current security testing techniques which are used to find vulnerabilities in applications [6], nor does it discuss these in detail. From the perspective of the *Security and Compliance Data Lake*, each security testing technique applied may be viewed as a potential data source. But storing known vulnerabilities is still just one, although probably the most popular, use case of this central metadata store. A major design goal is to keep it extensible, so that all kinds of metadata may be stored.

1.4 Environment

This thesis is written in cooperation with SAP. As of today, with a total revenue of €27.34 billion, SAP is the third largest software company in the world after Microsoft and Oracle [7]. While also having gained some attention with *Business to Customer (B2C)* products like the *Corona-Warn-App*, its core products are *Business to Business (B2B)* enterprise software solutions. Originally, SAP grew around its *Enterprise Resource Planning (ERP)* system. Today, the company is also adopting *Internet of Things (IoT)* technologies and *Artificial Intelligence (AI)* to provide advanced analytics for its customers and to maximize the value of its software products [8]. Besides, SAP is also investing heavily to push its products and customers to the

cloud. Therefore, the company maintains partnerships with Amazon, Microsoft, Google and Alibaba as infrastructure providers.

The department this thesis is written with is developing and maintaining the *SAP Gardener*. SAP Gardener is SAP's own managed Kubernetes service which enables SAP itself as well as SAP customers to ship their applications to all of these different infrastructures using a unified deployment underlay. Since SAP Gardener offers the possibility to configure practically every detail, neither SAP nor its customers need to rely on the managed Kubernetes service of each hyperscaler with their individual perks and restrictions, but can leverage the functionality of a fully configurable Kubernetes cluster without actually having to fully configure it [9].

1.5 Structure of the Thesis

In order to achieve the goals laid out before, the thesis builds upon following structure:

Foundation: In this chapter, basic terminology and existing concepts in software metadata management are established.

State of the Art: In the previous motivation and goal sections, it is already mentioned that existing measures for monitoring software component metadata such as vulnerabilities are insufficient in some cases. Therefore, to further motivate the thesis, this chapter will briefly introduce the state of the art approach and point out its limitations. The identification of these problems is crucial for the successful design of the *Security and Compliance Data Lake*.

Reference Architecture: As mentioned, this thesis is written with the SAP Gardener team. It is greatly inspired by the problems they are facing and by the approaches they already implemented to overcome limitations of the state of the art. Thus, this chapter introduces their proposed standard, the *Open Component Model (OCM)*, and its capabilities. Based on that, it develops an abstract design for the *Security and Compliance Data Lake*, which is then used to derive an entire holistic solution architecture for modern development and deployment landscapes. This architecture thereby presents a way to overcome all those limitations.

System Design: The system design chapter covers the conception of the data model, the selection of a database, the overall application architecture and the design of the API. Thereby, it especially focuses on giving detailed information about the

ideas and motives leading to specific design decisions. Besides, alternatives that have been considered and the respective reasons to not follow through with them are an essential part of this section.

Results: Finally, the actual functionality of the application is showcased and the design decisions, the implementation as well as the knowledge gained through this work is evaluated and the research question is revisited.

Chapter 2

Foundations

This chapter provides the basic knowledge necessary to understand the following chapters as well as other research and literature regarding metadata management. It therefore explains several important concepts and abbreviations around *Software Identification*, *Open Source*, *Open Source Licensing* and *Vulnerability Management*. Furthermore, *Software Bill of Materials* and existing *Software Bill of Material Standards* are introduced. This is important to understand the requirements of the Executive Order mentioned in the introduction and to put everything into context.

2.1 Software Identification

Uniquely identifying a piece of software is a frequent concern when managing applications on enterprise scale. Thus, a standardized naming convention would be required. Unfortunately, almost every packet manager and tool uses its own convention.

Package URL (purl)

Package URL, or *purl* for short, is one of the most widely adopted attempts to standardize those existing approaches. As described in the projects GitHub Repository, it therefore specifies an URL string, a purl, which should be able "to identify and locate a software package in a mostly universal and uniform way across programming languages, package managers, packaging conventions, tools, APIs and databases" [10]. This Package URL consists of seven components [10]:

```
1 scheme:type/namespace/name@version?qualifiers#subpath
```

Listing 2.1: Package URL Components

```
1 pkg:docker/cassandra@sha256:244fd47e07d1004f0aed9c
2 pkg:github/package-url/purl-spec@244fd47e07d1004f0aed9c
3 pkg:deb/debian/curl@7.50.3-1?arch=i386&distro=jessie
4 pkg:npm/foobar@12.3.1
```

Listing 2.2: Package URL Examples

Each component is separated by a different specific character to allow for unambiguous parsing. The `scheme` (Required) is the constant value "pkg" which may be officially registered as an URL scheme in the future. `type` (Required) refers to a package protocol such as maven or npm. `namespace` (Optional) may be some type specific prefix such as a Maven groupid, a Docker image owner or a Github user or organization. The `name` (Required) and `version` (Optional) are the name and version of the software. `qualifiers` (Optional) is also type specific and may be used to provide extra qualifying data such as an OS, architecture or a distribution. With the `subpath` (Optional) one may specify a subpath within a package, relative to the package root [10].

Another convention to uniquely identify packages is the *Common Platform Enumeration (CPE)* format which is primarily used in the context of vulnerability management ecosystem. Therefore, the standard is also discussed in this chapter.

2.2 Open Source Software and Licensing

Open Source Software is widely used in modern software development. As to what Open Source Software actually is, the *Open Source Initiative* defined a set of rules, the *Open Source Definition*, specifying distribution terms that the license of a software must comply with. Without going into detail and examining all of these terms, this definition generally ensures that such software "can be freely accessed, used, changed, and shared (in modified or unmodified form) by anyone" [11].

This description on its own may give the impression that there is no need to keep track of what *OSS Licenses* are used within an enterprise application. But there are still some limitation that *OSS Licenses* can put to the usage, especially the distribution of the software. Specifically, the so called *copyleft* principle and corresponding licenses might pose a challenge to some companies. While OSS generally may be used for commercial purposes, this principle dictates that if a company or individual distributes the software or a derivative, it has to be done under the same license it has been received under [11].

In the context of this work, OSS is used within other OSS components, so *copyleft* is not an issue. Nevertheless, there is still a risk involved in using OSS. A company might distribute initial software versions under an OSS License until there

is some kind of customer lock-in and then switch to another more restrictive license for further versions. Also, since precedent cases may not exist regarding the legal interpretation of some specific OSS licenses, there is a risk of expensive law suits.

2.3 Vulnerability Management

Due to the widespread use of OSS, vulnerability management is an increasingly important topic. Since many organizations use the same software components, vulnerabilities often become publicly known. Also, the open source nature allows attackers to get precise information about the vulnerability itself. Thus, tracking publicly known vulnerabilities in an organization's software products is a crucial capability to keep them secure.

Vulnerability

According to the *National Institute of Standards and Technology (NIST)* a *vulnerability* is “A weakness in the computational logic (e.g., code) found in software and hardware components that, when exploited, results in a negative impact to confidentiality, integrity, or availability. Mitigation of the vulnerabilities in this context typically involves coding changes, but could also include specification changes or even specification deprecations (e.g., removal of affected protocols or functionality in their entirety).”[12]

Common Vulnerabilities and Exposures (CVE)

CVE was created by MITRE in 1999. Initially, the acronym was meant to stand for *Common Vulnerability Enumeration*. As described by the authors of the original whitepaper, many security tools and advisories used their own vulnerability identifiers. In order to integrate several of these, one had to manually compare and eventually relate the vulnerabilities to each other. Thus, a common naming convention and common enumeration of vulnerabilities was needed [13].

To solve this issue, the *CVE Program* assigns unique identifiers, so called *CVE IDs*, to each vulnerability. All the identified vulnerabilities are then maintained as *CVE Records* in the *CVE List*. A minimal *CVE Record* consists of a *CVE ID*, a description of the vulnerability and at least one public reference. It does not include technical data, information about risks, impacts or fixes. An example for such a *CVE Record* is shown in figure 2.1 below.

CVE-ID	CVE-2022-29615 Learn more at National Vulnerability Database (NVD)
	<ul style="list-style-type: none"> • CVSS Severity Rating • Fix Information • Vulnerable Software Versions • SCAP Mappings • CPE Information
Description	SAP NetWeaver Developer Studio (NWDS) - version 7.50, is based on Eclipse, which contains the logging framework log4j in version 1.x. The application's confidentiality and integrity could have a low impact due to the vulnerabilities associated with version 1.x.
References	<p>Note: References are provided for the convenience of the reader to help distinguish between vulnerabilities. The list is not intended to be complete.</p> <ul style="list-style-type: none"> • MISC:https://launchpad.support.sap.com/#/notes/3202846 • URL:https://launchpad.support.sap.com/#/notes/3202846 • MISC:https://www.sap.com/documents/2022/02/fa865ea4-167e-0010-bca6-c68f7e60039b.html • URL:https://www.sap.com/documents/2022/02/fa865ea4-167e-0010-bca6-c68f7e60039b.html

Figure 2.1: CVE Record

Source: Based on [14]

In order to make sure that all vulnerabilities listed in the *CVE List* are unique and maintained properly, *CVE IDs* can only be assigned and *CVE Records* can only be published by MITRE and several partner organizations, so called *CVE Numbering Authorities (CNA)*. Thus, to add a new vulnerability to the *CVE List*, the discoverer has to report it to a *CNA* [15].

Common Vulnerability Scoring System (CVSS)

CVSS is an open framework for rating vulnerabilities. It is owned and managed by the non-profit organization *Forum of Incident Response and Security Teams (FIRST)*. The framework captures the main characteristics of a vulnerability to produce a numerical score between 0.0 and 10.0 reflecting its severity [16].

Therefore *CVSS* is composed of three metric groups: Base, Temporal, and Environmental. The *Base Score* considers the intrinsic characteristics of a vulnerability that are constant over time and assumes the worst case impact across different environments. These characteristics take into account exploitability metrics like attack complexity, but also impact metrics like confidentiality impact. The *Base Score* is typically calculated by the organization maintaining the vulnerable product or by third party analysts on its behalf. The *Temporal Score* considers characteristics that may change over time but not across environments. Such a characteristic would be the maturity of exploit code. The *Environmental Score* considers the relevance of

a vulnerability in a specific environment. *Temporal* and *Environmental Scores* are typically calculated by the consumers of the components to adjust the vulnerability rating to their organization's use case and environment. These Scores can then be used for internal risk management [16].

Common Weakness Enumeration (CWE)

CWE is a community-developed list of software and hardware weaknesses maintained by MITRE. Each weakness in the list is assigned a *CWE ID*. The list represents weaknesses on different levels of abstraction. This is conceptually shown by figure 2.2.

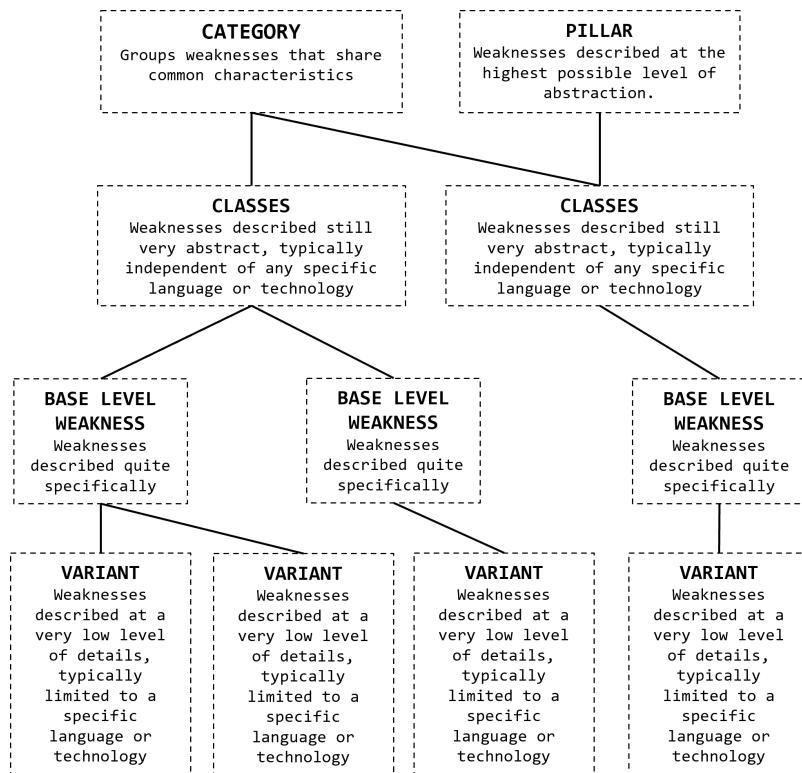


Figure 2.2: CWE Hierarchy

Source: Based on [17]

Relationships exist between elements on different levels of abstraction. As shown, *Classes* are usually member of a *Category* and might also be the child, hence a more concrete description, of a *Pillar*. Additionally to what is shown in the figure, these relationships may also skip levels of hierarchy. Thus, a *Base Level Weakness* may be a direct member of a *Category* or a child of a *Pillar*. An example for a Base Level Weakness is *CWE-478: Missing Default Case in Switch Statement* which is a member of *Bad Coding Practices* and a child of the class *Incomplete Comparison with Missing Factors*, which is a child of the pillar *Incorrect Comparison* [17].

Common Platform Enumeration (CPE)

The *CPE* specification originally created by MITRE and now maintained by NIST provides a naming scheme for IT assets such as software. It may be used to uniquely determine a specific software and its version. This way, a CPE enables cross referencing to other sources of information. The commonly used CPE naming scheme is structured as follows:

```
1 cpe:2.3: part : vendor : product : version : update : edition :  
language : sw_edition : target_sw : target_hw : other
```

Listing 2.3: CPE Formatted String Binding

Thereby `part` may be *a* for applications, *o* for operating systems, and *h* for hardware devices. `edition` is a legacy attribute in the current version of the specification and may be omitted where not required for backward compatibility. The attributes after `edition` were newly introduced in this version and are referred to as *extended attributes*. `sw_edition` should characterize a particular market or class of users a product is tailored to (e.g. online), `target_sw` a software computing environment (e.g. linux), `target_hw` the instruction set architecture (e.g. x86), and `language` the language supported in the user interface [18].

National Vulnerability Database (NVD)

The *NVD* is a database of vulnerabilities owned and maintained by NIST. In the paragraph about CVE, it was mentioned that the CVE Records do not contain technical data, information about risks and impact, or fixes. The NVD feeds from the CVE List and uses the information provided in the CVE Records to perform further analysis. As a result, a NVD entry exists for each CVE ID and provides a CVSS Base Score, a CWE ID and a CPE ID [12].

Thus, NVD combines all the aforementioned standards and concepts to provide thorough and concise human and machine-readable information about vulnerabilities. The *CPE IDs*, identifying a particular software version in use, may be queried against a NVD API to automatically check for known vulnerabilities. The *CVSS Base Score* is a valuable foundation for internal risk assessment and the *CWE ID* helps to quickly understand the type of a vulnerability.

2.4 Software Bill of Materials

A *Software Bill of Materials* is an inventory of the components used in a software. It ideally contains all direct and transitive components and their dependencies, which makes it virtually equivalent to the dependency graph of a software [19, 4].

As consequence to an *Executive Order on Improving the Nation's Cybersecurity*, the *National Telecommunications and Information Administration (NTIA)* published a document describing the minimum requirements for SBOMs [1, 4]). According to this document, these requirements are:

Data Fields (Metadata)	Baseline information about each component: Supplier, Component Name, Version of the Component, Other Unique Identifiers, Dependency Relationship, Author of SBOM Data, Timestamp of SBOM creation
Automation Support	Automatic generation and machine-readability to allow for scaling across the software ecosystem.
Practices and Processes	Implementation of policies, contracts and arrangements to maintain SBOMs.

Table 2.1: Minimum Elements of a SBOM

Source: [4]

The goal of the *Data Fields* is to sufficiently identify the components to track them through the supply chain and map them to other data sources, such as vulnerability and license databases. The *Automation Support* provides the ability to scale across the software ecosystem. The *Practices and Processes* ensure the maintenance by integration into the software development life cycle. SBOMs thereby increase software transparency, providing those who produce, purchase and operate software the means to perform proper risk assessments [4].

Due to this Executive Order, SBOMs are now required for all U.S. federal software procurements. This does not only affect direct software vendors of the U.S. government [1]. As a consequence to this Executive Order, every organization that is downstream from the U.S. government in the supply chain may be required to provide SBOMs for its products. Thus, this will be a crucial capability for most software vendors.

There are three data formats mentioned in the minimum elements document which are interoperable, able to fulfill the requirements and either human- and

machine-readable. Those are the *Software Package Data eXchange (SPDX)*, *CycloneDX* and *Software Identification (SWID)* tags [4]. SAP uses yet another SBOM format, called *Open Component Model (OCM)*, which does not fulfill the minimum requirements. The OCM is discussed further in a later chapter. SPDX is the most mature standard. It has laid out a lot of groundwork for the more recent CycloneDX. Thus, to get an better and more concrete understanding of SBOMs, SPDX will be examined in more detail.

SPDX License List and License Identifiers

SPDX is an initiative founded in 2010 and hosted at *The Linux Foundation*. In 2021 the SPDX specification even became an ISO standard [20]. The initiative focuses on solving challenges regarding the licenses and copyrights associated with software packages. SPDX therefore assembles licenses and exceptions commonly found in OSS in the *SPDX License List*. More precisely, this list includes a standardized short identifier, the full name, the license text, and a canonical permanent URL for each license and exception. By incorporating this *SPDX License Identifiers* in the sources on file level, one enables automation of concise license detection, even if just parts of an OSS project are used. Furthermore, SPDX provides *Matching Guidelines* to ensure that e.g. a “BSD 3-clause” license in a LICENSE file of an OSS project with different capitalization or usage of white space than the master license text included in the *SPDX License List* is still identified as “BSD 3-clause” license.

SPDX Documents

At the heart of the SPDX initiative are the *SPDX Documents* which leverage the *SPDX License List* and *SPDX License Identifiers* to describe the licensing of a set of associated files, referred to as *Package* in the context of SPDX. A *SPDX Document* provides the means to describe information about the document creation, the package as a whole, individual files, snippets of code within an individual file and other licenses that are not contained in the *SPDX License List* but are still relevant for the package, relationships between *SPDX Documents*, and annotations, which in a way are comments within a *SPDX Document*. The concept of relationships is a rather new addition to the specification. It is particularly useful if one has a SPDX Document describing a binary. Explicitly capturing relationships like “generated from” these source files and “dynamically linking” these libraries allows for a complete licensing picture.

These *SPDX Documents* may be represented in one of the following five file format: tag/value (.spdx), JSON (.spdx.json), YAML (.spdx.yaml), RDF/xml

(.spdx.rdf), and spreadsheets (.xls) [21, 22].

To give a more concrete idea of the basic concepts of *SPDX Documents*, an example from the SPDX GitHub repository will be briefly examined [23]. Therefore figure 2.3 below shows the directory structure of a “Hello World” project in C.



Figure 2.3: C Project Directory Structure

Source: [23]

Listing 2.4 shows a corresponding *SPDX Document*. Some tag:value pairs, which are less relevant for the overall understanding, are deliberately omitted to contain the length of the example.

```
1 SPDXVersion: SPDX-2.2
2 DataLicense: CC0-1.0
3 SPDXID: SPDXRef-DOCUMENT
4 DocumentName: hello
5 DocumentNamespace: https://swinslow.net/spdx-examples/example1/
    hello-v3
6 Creator: Person: Steve Winslow (steve@swinslow.net)
7 Created: 2021-08-26T01:46:00Z
8
9 ##### Package: hello
10 PackageName: hello
11 SPDXID: SPDXRef-Package-hello
12 PackageDownloadLocation: git+https://github.com/swinslow/spdx-
    examples.git#example1/content
13 PackageLicenseConcluded: GPL-3.0-or-later
14 PackageLicenseInfoFromFiles: GPL-3.0-or-later
15 PackageLicenseDeclared: GPL-3.0-or-later
16 PackageCopyrightText: NOASSERTION
17
18 FileName: /build/hello
19 SPDXID: SPDXRef-hello-binary
20 FileType: BINARY
```

```

21 LicenseConcluded: GPL-3.0-or-later
22 LicenseInfoInFile: NOASSERTION
23 FileCopyrightText: NOASSERTION
24
25 FileName: /src/Makefile
26 SPDXID: SPDXRef-Makefile
27 FileType: SOURCE
28 LicenseConcluded: GPL-3.0-or-later
29 LicenseInfoInFile: GPL-3.0-or-later
30 FileCopyrightText: NOASSERTION
31
32 FileName: /src/hello.c
33 SPDXID: SPDXRef-hello-src
34 FileType: SOURCE
35 LicenseConcluded: GPL-3.0-or-later
36 LicenseInfoInFile: GPL-3.0-or-later
37 FileCopyrightText: Copyright Contributors to the spdx-examples
                     project.
38
39 Relationship: SPDXRef-hello-binary GENERATED_FROM SPDXRef-hello-
                 src
40 Relationship: SPDXRef-hello-binary GENERATED_FROM SPDXRef-
                 Makefile
41 Relationship: SPDXRef-Makefile BUILD_TOOL_OF SPDXRef-Package-
                 hello

```

Listing 2.4: SPDX Document

Most of the tag:value pairs are self-explanatory, but some might require further explanation. The *Concluded License* is the license the SPDX file creator has concluded as the governing license of a package or a file. *License Information from Files* contains a list of all licenses found in a package and the *Declared License* is the license declared by the authors of the package [22]. Additionally, listing 2.4 illustrates how the concept of relationships may be used.

It is also worth mentioning that the concept of *Packages* in SPDX as a set of associated files is really rather loose. Thus, describing the project in figure 2.3 as two separate packages, one for source and one for binary, optionally in the same or also in two separate *SPDX Documents*, would be completely conform with the specification as well.

Since SPDX has been around for so long and is an accepted ISO standard, there exists a lot of useful tooling. It is therefore quite easy to automate tasks like producing, consuming, transforming and validating *SPDX Documents* [21].

2.5 Regulations

Vulnerabilities and licenses are not the only risks associated with enterprise software. Companies also need to consider governmental regulations since violations may lead to fines that have critical impact on the business.

Export Control Classification Number (ECCN)

The ECCN is a 5 character alpha-numeric designation used to determine whether a so called dual-use item needs an export license from the U.S. Department of Commerce in order to legally export it. Dual-use items are items that may be used for civil as well as military purposes. The ECCN gives some information about the product, as shown in figure 2.4 below.

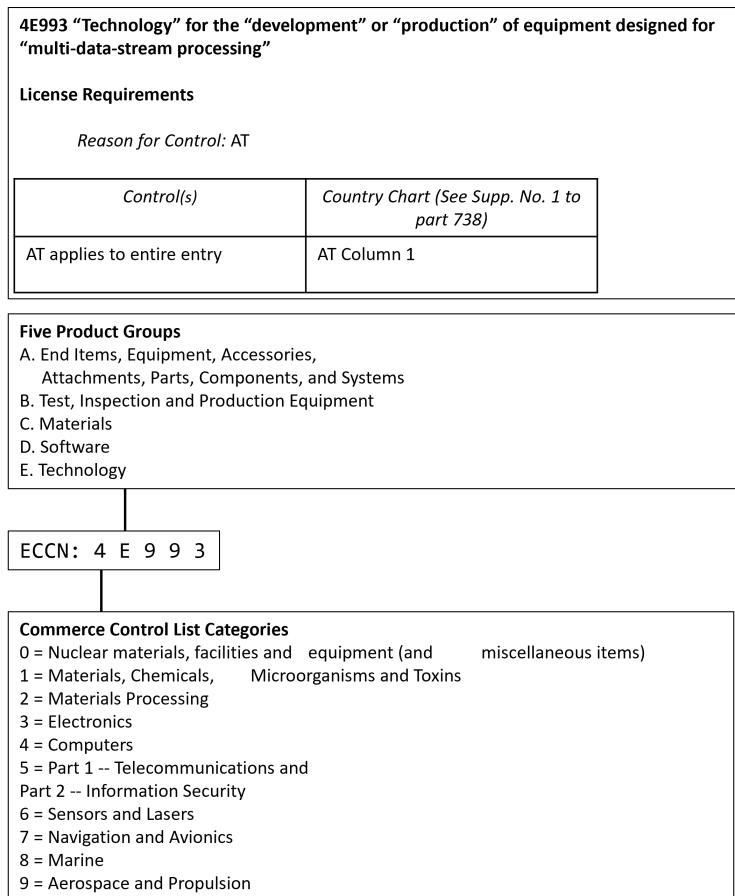


Figure 2.4: ECCN Structure
Source: Based on [24]

All ECCNs are listed in the *Commerce Control List (CCL)*. The entries in the list contain information about why an item might be under export control regulations. The top of figure 2.4 displays a snippet of such an entry. The reason for the export

2.5. Regulations

regulations of products with the classification number 4E993 is AT, which is the abbreviation for Anti-Terrorism [24].

Chapter 3

State of the Art

This chapter gives an overview of the state of the art approach of integrating security and compliance measures into the *software development life cycle (SDLC)*. That includes a brief introduction of practices such as *Continuous Integration and Continuous Delivery (CI/CD)* and *DevOps*. In the end, the limitations of this current approach are discussed, thereby further motivating this thesis.

3.1 Software Development Life Cycle

Since the emerging of cloud technologies, the software industry, especially the major software companies like Microsoft, Amazon and SAP, shifted their business model from *software as a product (SaaS)* to *software as a service (SaaS)*. This gives companies the opportunity to frequently release updates without adhering to a rigid distribution cycle. In consequence, they can react to feedback much faster and improve their software continuously [25, 26].

To keep up with this new pace, the SDLC had to be adjusted as well. Developers started to *continuously integrate (CI)* their code after they conducted some changes, while at the same time automatically testing the software. Thus, they were always keeping their code up to date and in a shippable state.

As an extension of that idea, the software is also automatically and *continuously delivered (CD)* to testing or even production environments, accelerating the process even more [27].

A quick side note – also after reviewing some literature, there seems to be some disagreement whether to distinguish continuous delivery and continuous deployment [27, 25]. The articles and papers that do differentiate between these terms, define continuous delivery as automatically delivering to production-like environments for evaluation and testing. But there are still manual steps necessary to actually deploy to production or customer environments [27, 28, 29].

DevOps, a combination of "development" and "operations", is a practice which also resulted from this change of the SDLC. Since the concept of CI/CD merges areas

of development and operations, the corresponding roles were combined to reduce communication overhead and misunderstandings [29]. DevOps is usually reliable for the setup and maintenance of the *CI/CD pipeline*. Hence, they automate the steps required for the continuous integration and delivery, such as building and testing the software after a developer integrated his code changes and moving it to the next step after the previous tests and checks are passed.

3.2 Integrating Security and Compliance Measures

In practice, in order to conduct CI/CD, DevOps usually sets up a so called *CI/CD pipeline* for the project. This includes a set of tools to automate build, test and deployment of the software. The most popular and widely used example is Jenkins [30]. But essentially, CI/CD pipelines are just simple stage-gate systems. Thus, the process is divided into a number of stages. Between each stage is a quality gate, such as a set of automated tests. The software has to pass this quality gate in order to being able to move to the next stage [31].

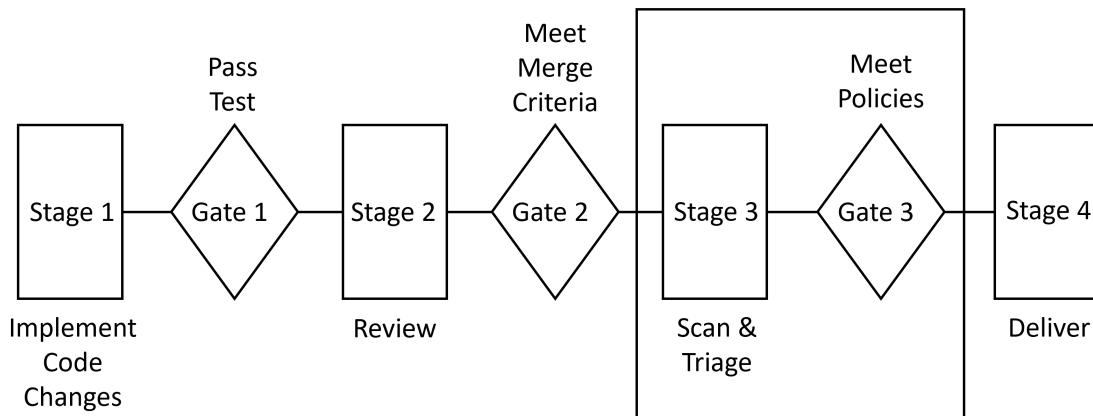


Figure 3.1: CI/CD Pipeline as Stage-Gate System

Source: Based on [31]

Figure 3.1 is an abstract technology-agnostic illustration of a CI/CD pipeline as a stage-gate system. The concept represented in this illustration also shows the currently established state of the art approach to reduce software supply chain risks such as problematic licenses or vulnerabilities in dependencies. Security and compliance scans are conducted as part of the CI/CD pipeline. The results of these scans may also be triaged before they are finally checked against a defined set of policies, forming another quality gate *before delivery*.

Triage is a term derived from medicine where it describes the classification and prioritization of patients if there are not enough resources to treat them all

immediately. In software testing or scanning, it refers correspondingly to the process of (re-)rating and (re-)classifying issues found during the scan in the context of its occurrence within the particular project. Common issues include problematic licenses and vulnerabilities. For the latter, the triage may be done based on the CVSS introduced in the previous chapter. Policies may then prevent shipping of software that contains vulnerabilities with a CVSS that exceeds a certain threshold or specific blacklisted licenses.

Unfortunately, there is not a lot of academic literature on this specific topic to back up the claim that this is the state of the art approach. Therefore, for the purpose of this work, the offerings of the major vendors of the security and compliance scanning tools were deemed a suitable source.

The security and compliance scanning tools used to retrieve this kind of information are called *Software Composition Analysis (SCA)* tools. Those tools typically scan entire software *artifacts*, thus, source code projects or binaries, and provide information about dependencies to other software projects, in other words, about libraries and frameworks used within the artifact. These software projects composing an artifact are generally referred to as *packages*. The packages are then matched against different databases such as the NVD to check for and retrieve known vulnerabilities and licenses.

According to a market research of Forrester from 2021, based on a combination of market presence, current offering and strategy, those major vendors are Mend (formerly known as WhiteSource), Synopsys, Sonatype and Snyk [32]. All of these tools advertise their smooth integration into the CI/CD pipeline, to be more precise, into the source code repositories and build tools, as another quality gate [33, 34, 35, 36]. This should generally be a suitable representation of what the industry is requesting, although it is to mention that most of these vendors also offer other integration options. Another popular one for the SCA tools that scan binaries are the artifact repositories.

After all, the reason for integrating these scans as another quality gate as part of the CI/CD pipeline into the repositories is evident. There, the source files or alternatively the binary files are conveniently accessible for these tools.

3.3 Producing Software Bill of Materials

As a solution to increasing the supply chain transparency, the US government decided that, analogous to other industries, their software vendors have to provide SBOMs. Subsequently, the NTIA conducted research and published several documents on the topic. One of those is the minimum elements for SBOMs discussed in the foundations chapter. As mentioned there, one of the minimum requirements is automation support. Therefore, the NTIA conducted a survey of existing SBOM formats and standards considering how to integrate them into the SDLC [37].

As part of this survey, the NTIA suggests to leverage existing tools such as version control systems, build tools, code scanners and binary analysis tools – thus, the tools already integrated into the CI/CD pipelines – to generate SBOMs [37].

The previous section explained that SCA tools retrieve the decomposition information about artifacts, in other words, the packages composing an artifact, as well as their respective vulnerabilities and licenses. Since this is the most important SBOM information, the SCA tools can also directly provide these results in a specific standardized SBOM format. Therefore, all of the mentioned SCA tool vendors quickly adjusted and provide such options now.

3.4 Limitations

The CI/CD integration is intuitive and might be sufficient for the scope of a single development team which deploys its software to a single infrastructure. But if this scope is exceeded, there are several serious limitations regarding the following aspects:

Synchronous Scans: Once a software artifact has passed the corresponding quality gate within the CI/CD pipeline, it will not be scanned again. But new vulnerabilities in one of the dependencies may still be discovered after this point in time (e.g. after the software is delivered). Consequently, these vulnerabilities remain undetected until the respective quality gate is reached again with the next version of the artifact.

Third Party Artifacts: The software system that is being developed might depend on additional third party artifacts at runtime. These third party artifacts have to be delivered alongside the internally developed artifacts. Since the third party artifacts are not developed internally, these may never pass one of the companies' internal CI/CD pipelines and, thus, never pass the quality gates. Therefore, these artifacts may never be scanned and their vulnerabilities may again remain undetected.

Quality Gate Mentality: With the quality gate approach, the scans are frequently only treated as gates. Thus, developers hope for their code or binary to pass the scan and discard the results if it does.

Point to Point Integration: With this approach, each development team has to configure, integrate and maintain the scanning tools on their own in each of their pipelines. In larger companies, this often leads to multiple opinionated CI/CD pipelines with additional point to point integration for such SCA tools and to individual corresponding reporting dashboards.

Distribution of Scan Results over Multiple Technologies: Development teams that use several different SCA tools, for example, Synopsys' *Black Duck Binary Analysis (BDBA)* for binary scanning and Mend's tool for source code scanning, the results are distributed over multiple different technologies with *individual formats and identifiers*.

There are popular alternatives to decouple the scans from the CI/CD pipeline. Thereby, these alternatives solve especially the first three limitations, but also introduce other problems and are still insufficient:

Scan Integration into Artifact Repositories: Through integrating the scanning tools into the artifact repositories, even vulnerabilities in artifacts that do not pass an internal CI/CD pipeline are detected. But this *only works with the binary scanners* which usually provide less precise results. Besides, this might get tedious and end up in multiple *point to point solutions* as well, in cases where the artifacts are distributed over several different artifact repositories.

Scanning based on SBOMs: Another possible option is based on the production of SBOMs during the compliance scans. When initially passing the quality gate, a SBOM can be generated, which may be used as access point to conduct further compliance scans. The tools of Mend and Synopsys already provide such an option. The problem here is that the components in the generated SBOMs usually have *tool specific identifiers*. Therefore, each tool can only conduct scans based on their own SBOM.

So both of these approaches are not ideal. The following example makes the practical

implications of these limitations more tangible.

Assuming there is a situation where the artifacts, which compose the final software system and are built from different internal as well as third party source code repositories, are stored in multiple different artifact repositories. From these artifact repositories the artifacts are deployed into multiple different environments (naturally, different environments might be comprised of different versions of the artifacts).

Given a vulnerability in specific versions of a popular OSS package is discovered, how could a question about which environments contain these particular vulnerable versions of the package be answered?

Assuming regular scans of the artifacts are conducted through one of the alternative approaches allowing scanning asynchronously from the CI/CD pipeline, the result of the scanning tools or the respective generated SBOMs of the artifact versions contain this information. Furthermore, the information to which environments the artifacts are deployed is usually stored in respective DevOps repositories. But ultimately, there is no easy way to answer this question without tediously searching through all these different data sources.

So the above identified limitation of these approaches, that the scan results are distributed over multiple technologies, may be extended and generalized to:

Software metadata is distributed over the entire software development life cycle.

To solve the described problems, it is necessary to find a technology-agnostic way to decouple the scans from the CI/CD pipeline and to bridge the gap between dedicated software installations over the effectively used software artifacts down to the various packages these artifacts are composed of and for which, finally, the desired metadata (e.g vulnerabilities or licenses) is provided. To achieve this, the structural information together with the finally desired metadata and all the relations in between can be stored in a central structured data store – the *Security and Compliance Data Lake*.

This requires a uniform identification and access scheme, used along the complete development life cycle to bind all this information together. These schemes can then be used by all kinds of tools around such an ecosystem to feed data into the *Data Lake*, as well as to provide a consolidated view on the complete set of covered software, regardless of their technology or build pipeline environment.

Chapter 4

Reference Architecture

As mentioned in the introduction, this work is written in cooperation with SAP.

In fact, the limitations mentioned in the previous chapter are deducted from the limitations the SAP Gardener team struggles with themselves. Therefore, this chapter introduces the *Open Component Model (OCM)*, SAP Gardener's proposed standard which enables to automate software delivery, and therefore decouple the compliance scans from the CI/CD pipeline, in a technology-agnostic way.

Thereon, the abstract design for the *Security and Compliance Data Lake* is developed. This abstract design is then used to derive an entire holistic solution architecture for modern development and deployment landscapes to overcome all limitations of the current state of the art approach.

4.1 Open Component Model

The OCM is a SBOM format created and used by SAP Gardener. It does not fulfill the minimum requirements as defined by the NTIA. But this is due to the fact that the OCM has a different focus than SPDX or CycloneDX. While those two were deliberately designed to be a bill of materials, thoroughly listing the decomposition of a software, the OCM was specifically developed to decouple CI from CD and thereby overcome related limitations such as the ones mentioned in the previous chapter.

So, SPDX effectively describes *arbitrary metadata* (e.g. licenses and file type) about logical units of software. SPDX refers to this logical units of software as *Packages*. As described in section 2.4 "Software Bill of Materials", a *Package* may be anything, ranging from a snippet of code over single file over an artifact to an entire software product version. Thus, SBOMs are *containers for arbitrary software metadata*.

Opposed to that, the OCM describes the *delivery relevant information about and the access to artifacts* composing a logical unit of software. OCM refers to these logical units of software as *Components*, or rather *Component Versions*. Thus, *Component Versions* are *containers for a structured set of delivery relevant artifacts*.

For this reason, SAP Gardener stopped referring to OCM as an SBOM format and instead invented the term *Software Bill of Delivery (SBOD)*.

Hereafter is a technical explanation of the OCM deducted from the specification [38], an internal presentation [39] and interviews with responsible developers. Even though this explanation focuses on understanding the rationale behind the design decisions of the proposed standard rather than technical completeness, it may still be hard to grasp. This is due to the abstract nature of the OCM. Therefore, emphasis and examples are used where possible. Also, while the textual description explains the OCM as the abstract model that it is, listing 4.1 below shows a *Component Descriptor*, the serialization format of a *Component Version*. This aims to support in understanding the abstract concepts.

4.1.1 Model Elements

Component

”A *Component* is an abstract entity describing a dedicated usage context or meaning for provided software. It is technically defined by a *globally unique identifier*”[38] (l. 2). With this globally unique identifier, a *Component* acts as a namespace for multiple *Component Versions*. So the *Component* identified by the name `github.com/gardener/etcdruid` contains software versions of a product, a druid, that helps configuring *etcd*, a central key-value store, in Kubernetes clusters deployed by SAP Gardener.

Component Version

As already stated, a *Component Version* describes a structured set of concrete *Artifacts*. It is technically defined by the globally unique identifier of the *Component* and a version (ll. 2-3). Thus, a *Component Version* is an instance of a *Component* adhering to the semantic, or in other words, the usage context or meaning, given by the name. A *Component Version* leverages essentially two mechanisms to describe this structured set of *Artifacts*.

1. A *Component Version* may describe the delivery relevant information about and the access to an artifact directly (ll. 9-22, ll. 23-39).
2. A *Component Version* may describe a reference to another *Component Version* to include its described set of artifacts indirectly (ll. 40-56).

Artifacts and Component References

The two model elements enabling those artifact composing mechanisms are *Artifacts*

(ll. 9-22, ll. 23-39) and *Component References* (ll. 40-56). Thereby, drawing an analogy to file systems, *Artifacts* and their identities correspond to files and their names and *Component References* and their identities are analogous to directories and their names. Thus, *Component References* also impose a structure on the artifact set, like a directory tree that provides access to a structured set of named data content (the files).

Both those model elements share common sets of attributes [38]:

Identity: The *Identity* attribute set composes a *Component Version-Local Identity*. Thus, it uniquely identifies the elements in the context of a *Component Version*. The *Identity* at least consists of a *name*.

The `name` (*string*) (l. 10, l. 24, l. 41) is *required* and should be chosen to be unique within the *Component Version*. But it also carries semantic information, the meaning or purpose of the element. So, a *Component Version* may describe a REST application. This REST application may use a nginx version as a HTTP server and a different nginx version as a load balancer. In order to have the flexibility to name them both "nginx" the *Identity* has two additional optional attributes, *version* and *extra identity*.

The `version` (*string*) (l. 16, l. 25, l. 43) is *optional* and describes the version of the respective model elements.

The `extraIdentity` (*map[string]string*) (l. 17, l. 32, ll. 45-46) is *optional* and allows adding arbitrary identity attributes.

As this is a *Component-Version-Local Identity*, *Artifacts* and *Component References* may be identified by the triple (*Component Name*, *Component Version*, *Local Identity*).

Labels: The *Labels* attribute set is represented by a single `labels` (*[]any*) attribute. With this, *Labels* allows to attach additional arbitrarily nested metadata to such an element, which is not directly described by the existing model elements. This enables the use of application specific attributes without the need to extend the model for new arising use cases [38]. When a *Component Version* is used by a tool to retrieve the artifacts of the respective product version and put them into a scanning tool, a tag within the labels attribute could be used to mark exceptions, e.g. an artifact within that product version that never have to be scanned.

```

1 component:
2   name: github.com/gardener/etcd-druid
3   version: v0.15.3

```

```

4 repositoryContexts:
5   - type: ociRegistry
6     baseUrl: eu.gcr.io/sap-se-gcr-k8s-private/cnudie/gardener/
7       development
8       subPath: null
9       provider: internal
10      sources:
11        - name: github_com_gardener_etcd-druid
12          access:
13            type: github
14            repoUrl: github.com/gardener/etcd-druid
15            ref: refs/tags/v0.15.3
16            commit: a6112534eb79021fa191edb8efd6512760ae050f
17            version: v0.15.3
18            extraIdentity: {}
19            type: git
20            labels:
21              - name: cloud.gardener/cicd/source
22                value:
23                  repository-classification: main
24 resources:
25   - name: etcd-druid
26     version: v0.15.3
27     type: ociImage
28     access:
29       type: ociRegistry
30       imageReference: >
31         eu.gcr.io/sap-se-gcr-k8s-public/eu_gcr_io/gardener-project/
32           gardener/etcd-druid:v0.15.3-mod1
33       digest: null
34       extraIdentity: {}
35       relation: local
36       labels:
37         - name: cloud.gardener.cnudie/migration/original_ref
38           value: eu.gcr.io/gardener-project/gardener/etcd-druid:v0.15.3
39         - name: cloud.gardener.cnudie/sdo/lssd
40           value:
41             processingRules:
42               - purge_berkeleydb_to_public
43             srcRef: []
44 componentReferences:
45   - name: etcd
46     componentName: github.com/gardener/etcd-custom-image
47     version: v3.4.13-bootstrap-8
48     digest: null

```

```

47 extraIdentity:
48     imagevector-gardener-cloud+tag: v3.4.13-bootstrap-8
49 labels:
50 - name: imagevector.gardener.cloud/images
51   value:
52     images:
53 - name: etcd
54   repository: eu.gcr.io/gardener-project/gardener/etcd
55   resourceId:
56     name: etcd
57   sourceRepository: github.com/gardener/etcd-custom-image
58   tag: v3.4.13-bootstrap-8
59 labels: []
60 signatures: []

```

Listing 4.1: Component Descriptor

Artifacts

Generally, an *Artifact* is a blob containing some data in some technical format [38]. Besides *Identity* and *Labels*, *Artifacts* have the following further attributes [38]:

- A dedicated globally unique `type` (l. 18, l. 26) representing the kind of content and how it can be used.
- A formal description of the `access` (ll. 11-15, ll. 27-30). This description can be used to technically access the content of the artifact in form of a blob with a format defined by the `type` of the artifact (l. 18, l. 26) (the `type` shown within the access in listing 4.1 specifies the access method and thereby determines the formal procedure of how to access the content and the fields of the access specification required to follow this procedure).

The OCM distinguishes two kinds of *Artifacts*, *Sources* and *Resources*.

Sources (ll. 9-22) describe artifacts that are sources for the delivery relevant artifacts (e.g. source code).

Resources (ll. 23-41) describe delivery artifacts, intended for deployment into a runtime environment (e.g. executables or OCI Images) or additional content relevant for deployment mechanisms (e.g. helm charts). Thus, these are the artifacts built from *Sources*. The connection may be expressed through the `srcRef` attribute which is an additional attribute only existing for *Resources* [38].

Component References

A *Component Version* may refer to other *Component Versions* by adding a *Component*

Reference (ll. 42-58). As already mentioned, through this mechanism, the referring *Component Version* includes the artifact set described by the referred *Component Version*.

It is important to highlight, that *Component References* do not only specify the *Identity* of the referenced *Component Version* but do also have the *Component Version-Local Identity*. Thus, the reference itself has an *Identity* within the *Component Version* just like *Artifacts*. Since the name of this *Local Identity* carries the semantic information about the meaning or purpose of the identified piece of software in the context of the *Component Version*, this allows referencing the same *Component Version* twice with different meaning. These different purposes expressed through different names can then be used by tools evaluating a *Component Version* to find the correct *Artifact* or *Component Reference* for the specific use case. Therefore, local identities should be kept stable among successive versions of a *Component*.

But since the referenced *Component Version* still has to be specified, the *Component Reference* has an additional attribute, `componentName` (l. 44), which, together with `version` (l. 45), uniquely identifies the *Component Version* [38].

Additional Information

There are some attributes in listing 4.1 that are less important, but shall still be mentioned for completeness of the description [38]:

- A *Repository Context* (ll. 4-7) describes the access to an *OCM Repository*. An *OCM Repository* thereby is a repository providing technical access to the *Component Version*, or rather the *Component Descriptor*.
- The *Provider* (l. 8) specifies the company or organization providing the *Component Version*.
- A *Component Version* may have *Labels* (l. 59) itself, to attach additional information.
- A *Component Version* may be signed by some authority to ensure that it can be trusted. The *signatures* (ll. 60) attribute may contain multiple signatures, where each entry specifies a name for the signature, the digest value for the signature alongside its respective hash- and normalization algorithm and the signature itself. The digest of the *Component Version* thereby includes the digests of the delivery relevant artifacts, thus, the *Resources*, which may be stored alongside their respective descriptions (l. 31).

For further or more detailed information about specific elements or the Open Component Model and its ecosystem in general, refer to the official documentation [38].

4.1.2 Capabilities of the Open Component Model

After the abstract description, this brief section revises the major concepts of and points out the most important capabilities enabled by the *Open Component Model*.

Aggregation Mechanism

Components and *Component Versions* respectively are an abstract concept, defined in a way that allows to describe the artifact set and its delivery relevant information for arbitrary aggregations ranging from a specific *software product* over *artificial aggregations* like all artifacts used by the SAP Gardener CI/CD team to an entire *application landscape*.

So, a *Component Version* named "SAP Gardener CI/CD" could have *Component References* to all *Component Versions* describing a software product used within the SAP Gardener CI/CD team. The *Component Version* hierarchy created through these *Component References* would thereby indirectly specify their complete set of *Artifacts*.

Access Automation

As each of these *Artifacts* describes the access to the technical artifact in a machine-readable format, a tool may automatically find all the relevant technical artifacts and deploy them into a specific environment or provide them to scanning tools.

With that, the Open Component Model enables a *holistic deployment automation* and a *decoupling of compliance scans from the CI/CD pipeline* based on *Component Versions*.

Technology-Agnostic and Access-Independent Identification Scheme

The OCM introduces a *technology-agnostic identification scheme* based on the component namespace. Thus, an artifact is addressed through its semantic within the *Component* – to be concrete the triple (*Component Name*, *Component Version*, *Local Artifact Identity*) – independent of its technology.

A *Component Version* might describe a software product. This software product includes a Java library as an *Artifact* with a *Local Artifact Identity*, e.g. "JavaLib". This Java library may be stored in a Maven repository where it is identified by its group id (uniquely identifies a project across all maven projects), artifact id (uniquely identifies a artifact within a project) and version (e.g. org.example:javalib:1.0.0).

Additionally, it may be stored in a local artifact repository where it is identified by a digest of its contents.

To switch the location where the artifact will be accessed by tools evaluating the *Component Version* describing the artifact, only the *access* has to be exchanged without affecting the identity through which the *Artifact* is addressed.

This is also important in complex development and deployment landscapes as it provides additional flexibility. Due to national restrictions, one might be obligated to deploy artifacts from artifact repositories located in the respective country. Hence, different *Component Repositories* may store *Component Descriptors* for a specific *Component Version* with different *access* properties. Thus, the *access* properties may be exchanged transparently without affecting the deployment automation.

As *Component Versions*, or rather their serialization format *Component Descriptors*, are also stored in repositories, the statement in section 3.2 "Integrating Security and Compliance Measures", that it is convenient to feed respective tools from repositories, still fits.

But opposed to this, *the OCM is technology-agnostic and thereby forms a transparency layer, linking different storage locations and content technologies (e.g. artifact repositories or source repositories) and tools (e.g. scanning tools)*.

To conclude this, the Open Component Model solves several limitations of the state of the art approach. It provides an *aggregation for artifacts and thereby bridges the gap between a specific software installation to the therein used artifacts*. With this, it also introduces a *uniform identification and access scheme* for these artifacts, making them conveniently accessible for SCA tools independent of their technology or CI/CD pipeline.

The data lake still has to use this structural information about the artifact composition of software products and connect it to the information about contained packages and further metadata (e.g. vulnerabilities or licenses) provided by arbitrary data sources (e.g. scanning tools or build tools) with different individual identification schemes.

4.2 Abstract Solution Idea

The *Open Component Model* with its abstract properties, or rather capabilities, explained in the previous section, builds the foundation for the *Security and Compliance Data Lake*. Thereby, the *Open Component Model* could also be replaced by other standards providing similar capabilities. From hereon, such standards are

generally referred to as *Component Models* for the further context of the thesis.

Based on such a *Component Model*, this section deduces and outlines an abstract data lake design. Building upon this, it describes a *reference architecture* for integrating the *Data Lake* into a modern development and deployment landscape and a *holistic solution architecture* to overcome all the limitations of the state of the art approach.

4.2.1 Abstract Data Lake Design

Figure 4.1 outlines the abstract data lake design based on a *Component Model*.

As shown, the *Component Model* is at the center of the *Data Lake*. The information about *Components* and the therein effectively used artifacts is provided by the *Component Model* directly.

Generally, the *Components* themselves are generated during the build process alongside the respective described and aggregated artifacts. Thereby, the build tool may also supply the corresponding *Build Information* to the *Data Lake*.

Through the access automation, the artifacts may be fed into scanning tools to obtain the decomposition information about the *Packages* contained in *Artifacts* and their respective attached metadata, here *Vulnerabilities* and *Licenses*. The *Component Model* enables this process for any kind of tool which may serve as a data source for the *Data Lake*.

The information about *Packages*, *Licenses*, *Vulnerabilities* and other metadata may be provided by different data sources with multiple different individual identifiers. The *Data Lake* combines all this information and establishes a uniform identification. With this, the *Data Lake* provides a consolidated view on the complete set of software and its metadata, regardless of their technology or the data sources. Thus, it may serve as the central source for *Reporting* tools and eliminates the need for multiple individual reporting dashboards.

Most prominently, through modeling the structural information, thus, connecting everything from *Component Version* over *Artifacts* and *Packages* to *Vulnerabilities* and *Licenses*, the *Data Lake* may conveniently answer which *Component Version* contains a certain *Vulnerability*.

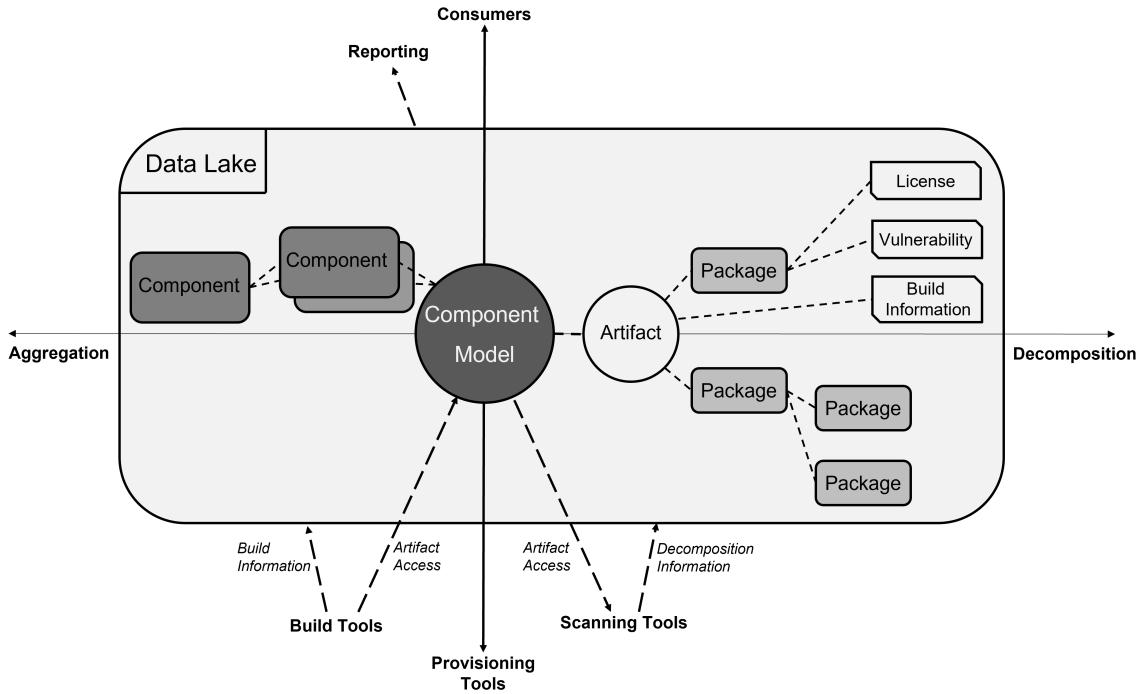


Figure 4.1: Abstract Data Lake Design

Source: Own Representation

4.2.2 Reference Integration Architecture

As already indicated in the previous chapters, although the *Security and Compliance Data Lake* is an application for storing all kinds of metadata, the initial primary use case is storing the data resulting from scanning tools, thus dependencies, vulnerabilities and licenses. Figure 4.2 shows the architecture designed for the integration into a modern development and deployment landscape based on the *Open Component Model*.

The *Data Collector Service* is the central component of this architecture. Its job is to orchestrate the process flow to derive metadata from information provided by the OCM and ingest it into the *Data Lake*. In practice, this would be done based on policies, requiring to trigger the data collection once a day or based on some kind of event.

In order to actually start the process, the Data Collector Service sends a request to the *Access Service* (1). The Access Service is a transparency layer already built into the *Component Repositories* as part of the OCM. Thus, if the Data Collector Service requests a set of Component Versions and their referenced Artifacts, the Access Service first fetches the corresponding Component Descriptors from the Component Repository (2). After receiving this information, the Access Service evaluates the access attribute in the referenced Sources and Resources and sends respective requests

to the *Source* and *Resource Repositories* to fetch the Artifacts (3). Finally, the Access Service may return all the requested information to the Data Collector Service.

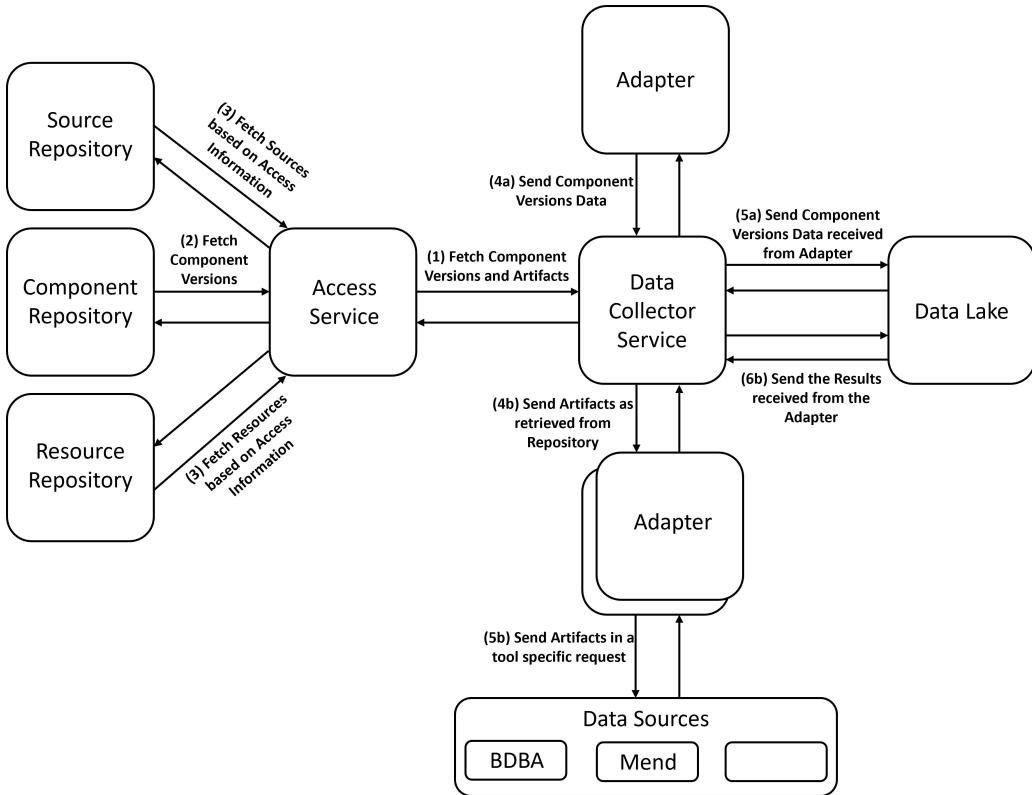


Figure 4.2: Data Lake Integration
Source: Own Representation

As shown in the figure, based on this information, there are two requests flows that may be triggered by the Data Collector Service, (4a) and (4b). The request flow illustrated above the Data Collector Service (4a) sends the Component Versions data, so the Component Descriptors, to an *Adapter* which adjusts and returns this data in the format required for the consumption by the Data Lake.

The request flow below the Data Collector Service (4b) distributes the Artifacts to a set of Adapters. Here, the Adapters initially wrap the Artifacts into the format required by a specific Data Source. In this case, there would be an Adapter for each scanning tool. Upon receipt of the scan results, these Adapters also adjust the data to the format required for consumption by the Data Lake, before returning it to the Data Collector Service.

Finally, the Data Collector Service may send all the information, thus, the information acquired from the Component Versions about Components, Resources, Sources and their relationships (5a) as well as the information from the scanning tools about dependencies, vulnerabilities and licenses within these Artifacts (6b) to

the Data Lake.

There were no communication protocols mentioned so far. This is due to the facts, that it does not really matter on the conceptual level and that the final implementation of this integration architecture is out of the scope of this thesis. But in practice, the communication between the Data Collector Service and the Access Service, between the Adapters and the Data Sources and between the Data Collector Service and the Data Lake is being handled using HTTP. The Adapters however will most likely not be implemented as discrete services but as components within Data Collector Service and therefore communicate over shared memory. Since it is very likely that additional Data Sources shall be added later on, the important part on the conceptual level is to keep components separated on a logical level. Thus, there should be well defined interfaces between the Adapters and the Data Collector Service.

Since the *Component Model* enables automated access to all artifacts of any software product (or other kinds of artifact aggregation), this integration architecture is completely decoupled from the rest of the development and deployment landscape. Thus, even in very complex development and deployment landscapes with several CI/CD pipelines based on different technologies (e.g. Jenkins or GitHub Actions) and including multiple different repositories, as long as it uses a *Component Model*, neither the reference architecture nor the implementation of the non-adapter parts has to be adjusted.

Even though the implementation of the integration architecture is out of scope of this thesis and is therefore not further discussed, it is still developed on a proof of concept basis as part of this work to feed actual real world data of the SAP Gardener into the *Security and Compliance Data Lake*.

4.2.3 Holistic Solution Architecture

To conclude this chapter, this section provides a complete picture of a development and deployment landscape overcoming the limitations identified in the previous chapter for the state of the art approach through leveraging the *Open Component Model* and the *Security and Compliance Data Lake*.

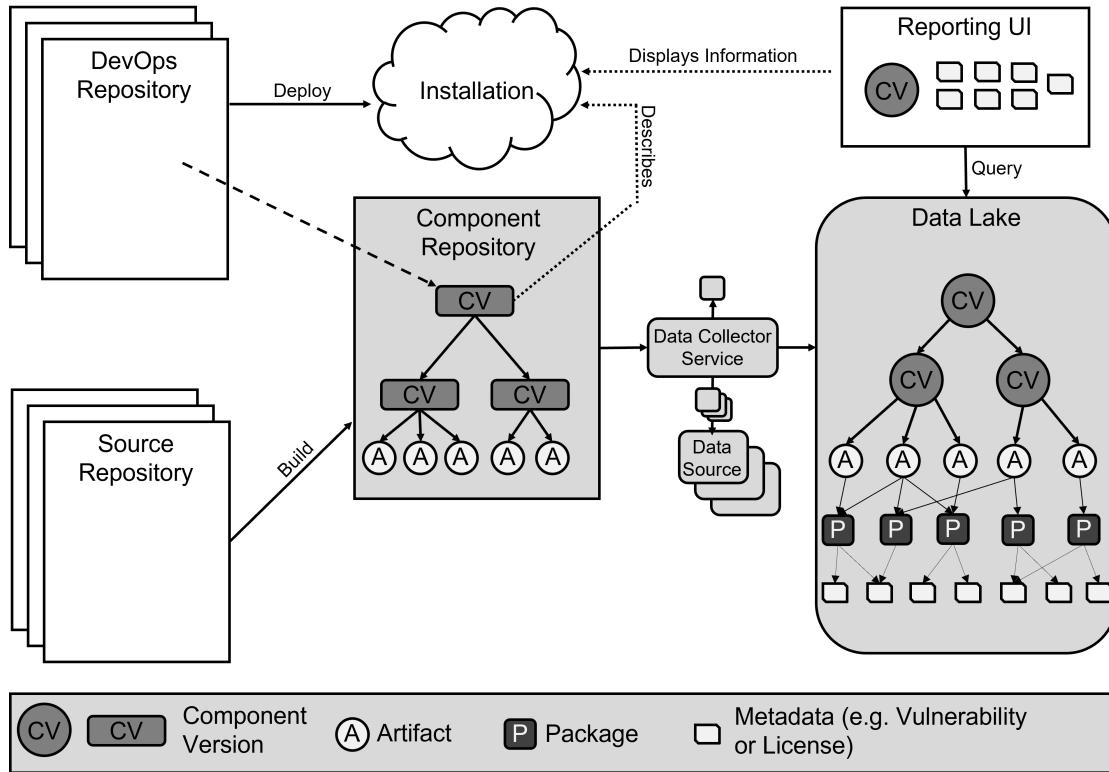


Figure 4.3: Holistic Solution Architecture

Source: Own Representation

On the left side, figure 4.3 shows the whole process from build to deployment, or in other words from the *Source Repository* over the *Component Repository* and the *DevOps Repository* to the *Deployment* of a concrete *Software Installation*.

On the right side, it shows the integration architecture from figure 4.2. This visualizes the previously mentioned decoupling of the integration architecture. Everything is linked only through the *Component Model*.

In a *Build* process, *Artifacts* may be build from *Sources* stored in a *Source Repository*. The build descriptions within the *Sources* contain information to generate a corresponding *Component Version* describing the *Sources* and *Resources*, their *Identities* and their relationship. Thus, the OCM assigns the uniform technology-agnostic identification, which is then used throughout the entire software life cycle, at the very beginning of the life cycle.

Thereby, it is irrelevant where the technical artifacts are actually stored. They may be stored in the *Component Repository* alongside the *Component Versions* as indicated in the figure or they may be stored in specific artifact repositories. This information is provided by the *Access Specification*. Since other tools may address the artifacts with their *Component Model* identity and the *Access Specification* may

be resolved by a built in *Access Service*, as introduced in the previous section, the location of the artifacts is completely transparent.

Consequently, the *DevOps Repository* may merely contain a file specifying the *Component Version* describing the software product version that shall be deployed to the *Installation* environment together with configuration information required for the deployment. The rest may be automated with respective deployment tools.

The *Data Lake* collects and combines metadata about the *Component Versions* in the *Component Repository* from different *Data Sources* providing a consolidated view onto the complete set of software.

Thereon, a *Reporting UI* may query information from the *Data Lake* about the *Component Version* describing a specific installation. So, an administrator or user may ask, which vulnerabilities and licenses are contained in this installation and the *Data Lake* may conveniently traverse the relationships to retrieve this information. The development of a respective *Reporting UI* is out of the scope of this thesis.

Ultimately, this holistic solution architecture enables *answering* the question posed in section 3.4 "Limitations" of the state of the art approach, about *which product versions might be affected by a specific vulnerability*.

Chapter 5

System Design

This chapter describes the technical design of the *Security and Compliance Data Lake*. It covers the conception of the data model, the selection of a database, the design of the API and the complete application architecture. Thereby, it especially discusses alternatives and focuses on giving detailed information about the ideas and motives that led to specific design decisions.

5.1 Requirements

Before describing the details of the system design, the requirements are specified, since they are at the core of every design decision.

5.1.1 Functional Requirements

Table 5.1 provides a condensed list of the functional requirements for the *Security and Compliance Data Lake*. Within that list, every requirement is described by the following elements: a short and precise but abstract statement of what functionality the system must have, an additional explanation which includes an example, and a column categorizing the requirements as priority 1 or 2.

Priority 1 is functionality deemed necessary for the *Security and Compliance Data Lake* in order to solve the limitations and problems of the state of the art approach. Furthermore, priority 1 functionality has to be considered in the design process and otherwise cannot be easily added without foundational remodeling.

Priority 2 functionality describes convenience features which are less urgent and may easily be added later on.

5.1. Requirements

Requirements		
Ref.#	Functionality	Prio.
R.1	<p>The SCDL shall be able to consume and to store metadata from multiple different data sources.</p> <p>The SCDL shall be able to work with any kind of metadata about software components. Therefore, it has to be able to handle the results of multiple different scanning tools as well as other kinds of data sources like build tools. As an example, it may have to consume data provided by BDBA, Mend but also Jenkins. Thus, it has to be considered that, besides vulnerabilities and licenses, a variety of other metadata types may need to be added in the future.</p>	1
R.2	<p>The SCDL shall store the metadata from different data sources without aggregation¹.</p> <p>Different tools that generally serve the same purpose may provide similar information. As an example, BDBA and Mend are both SCA tools and therefore provide overlapping results. To ensure that no data is lost, this information shall not be combined and aggregated before storing.</p>	1

Continued on next page

¹*aggregation* in this context means to merge the data about a package of e.g. a BDBA scan and a Mend scan to a single package entity instance

5.1. Requirements

Requirements		
Ref.#	Functionality	Prio.
R.3	<p>The SCDL shall provide the metadata from different data sources with aggregation¹.</p> <p>As mentioned before, to ensure no data is lost, the data from different data sources shall be stored without aggregation. Anyway, to be consumed by a user, this data shall be aggregated. As an example, when querying all packages contained in a specific resource, the result returned by the SCDL shall not contain the same package twice in different representations, if it was identified by BDBA and by Mend. Instead, the result shall contain an aggregated representation of the package. Thus, some kind of aggregation layer is needed which provides transparency regarding the data sources.</p>	1
R.4	<p>The SCDL shall provide a level of aggregation² to group sources and resources.</p> <p>As pointed out before, one problem with SBOMs is the disconnection of the artifact metadata and the software installations' information. To bridge this gap, an additional aggregation level for grouping artifacts is necessary. As an example, this additional aggregation level shall enable to group all resources contained in a specific product version.</p>	1
Continued on next page		

²aggregation in this context refers to the "whole/part" semantic of the word [40]. Thus, since resources and sources are comprised of packages, they are both aggregations of packages. On a model level, the same applies for the relationships between packages and vulnerabilities or licenses as well as between software installations and the contained resources.

5.1. Requirements

Requirements		
Ref.#	Functionality	Prio.
R.5	<p>The SCDL shall enable users to query the metadata on different levels of aggregation².</p> <p>As an example, a user shall be able to query for all vulnerabilities in a specific resource, thus, query on the aggregate level of resources. But a user shall also be able to query for all vulnerabilities in an entire specific product version, thus, querying on the aggregate level of product versions (querying on this level of aggregation enables to answer which environments are affected by the Log4j vulnerability).</p>	1
R.6	<p>The SCDL shall enable users to perform assessments.</p> <p>The relevance of specific pieces of information such as vulnerabilities or licenses depends on the usage context. As an example, while the internal usage of an altered OSS with a copyleft license is lawful, the distribution is not. Therefore, a possibility has to be provided to assess such pieces of information in the context of their occurrence.</p>	2
R.7	<p>The SCDL shall provide common data aggregation and filter functions for the queries.</p> <p>As an example, a user shall be able to filter for the vulnerability with the highest CVSS within a resource or shall be able to get the count of vulnerabilities within a resource.</p>	2
R.8	<p>The SCDL shall enable users to query the metadata in the common SBOM formats.</p> <p>In order to be able to fulfill the governmental requirements of the executive order mentioned in 2.4 "Software Bill of Materials", the SCDL has to provide a way to query the metadata in the common SBOM formats. As an example, a user shall be able to query the SPDX document for a specific resource.</p>	2

Table 5.1: Application Requirements

So, by fulfilling this functional requirements, the *Security and Compliance Data Lake* may serve as a central application for storing and querying structured software metadata. Thereby, solving the problem of metadata being distributed throughout the development life cycle and bridging the gap between dedicated software installations and the effectively used artifacts.

5.1.2 Non-functional Requirements

Since this shall be a prototypical implementation, there is a strong focus on fulfilling the functional requirements. Thus, no concrete limits regarding performance or scalability such as a maximum response time of 5 seconds or the support for up to 1000 concurrent users are set at this point. Considering the novelty of the topic, there is very few reference data and therefore, such specifications would be premature. Nevertheless, for a central metadata store which may prospectively power dashboard web applications for monitoring purposes, scalability and performance definitely have to be considered in the design decisions.

5.2 Data Model

Here is a brief revision of the relevant terms and entities, before the data model is discussed.

Artifact is an umbrella term for sources and resources. The definition of sources and resources is the same as introduced by the OCM.

Resources describe delivery artifacts, intended for deployment into a runtime environment (e.g. executables or OCI Images) or additional content relevant for deployment mechanisms (e.g. helm charts).

Sources describe the artifacts that are used to generate the delivery relevant artifacts (e.g. source code).

Compliance scanners usually scan entire source code repositories or binaries. Through different methodologies, these tools detect the packages contained in these scanned artifacts on a best effort basis. In the context of this work, *packages* are defined as functional units composing artifacts, whereby it is usually a collection of files forming a library which is imported in the source code.

By subsequently matching these packages against different databases such as the NVD, introduced in section 2.3 "Vulnerability Management", known *vulnerabilities* and *licenses* are identified. To give a better idea of these results, figure A.1 in the appendix shows a snippet returned from the API of BDBA. The results on their own are useful and provide interesting data about the above mentioned entities. But it is

still loose metadata that lacks context information such as which product versions contain the described artifacts with the corresponding decomposition information.

Therefore, an additional entity type to conduct further aggregation is required. The OCM already introduced such an entity type, the *component*.

To conclude this, from a high level perspective, the important entities are *components*, *sources*, *resources*, *packages* and the metadata attached to these entities such as vulnerabilities and licenses. To generalize this and abstract away from specific data sources, the entities representing a kind of metadata are called *info snippets*.

These entities are the starting point for the data model. From here on, it gets rather complex and abstract. To still keep the explanations tangible, figure 5.1 shows an *entity-relationship model (ERM)* describing the data model. This may be used as a reference point throughout the following paragraphs, discussing the design decisions leading up to the specific entity types, relationships, and cardinalities.

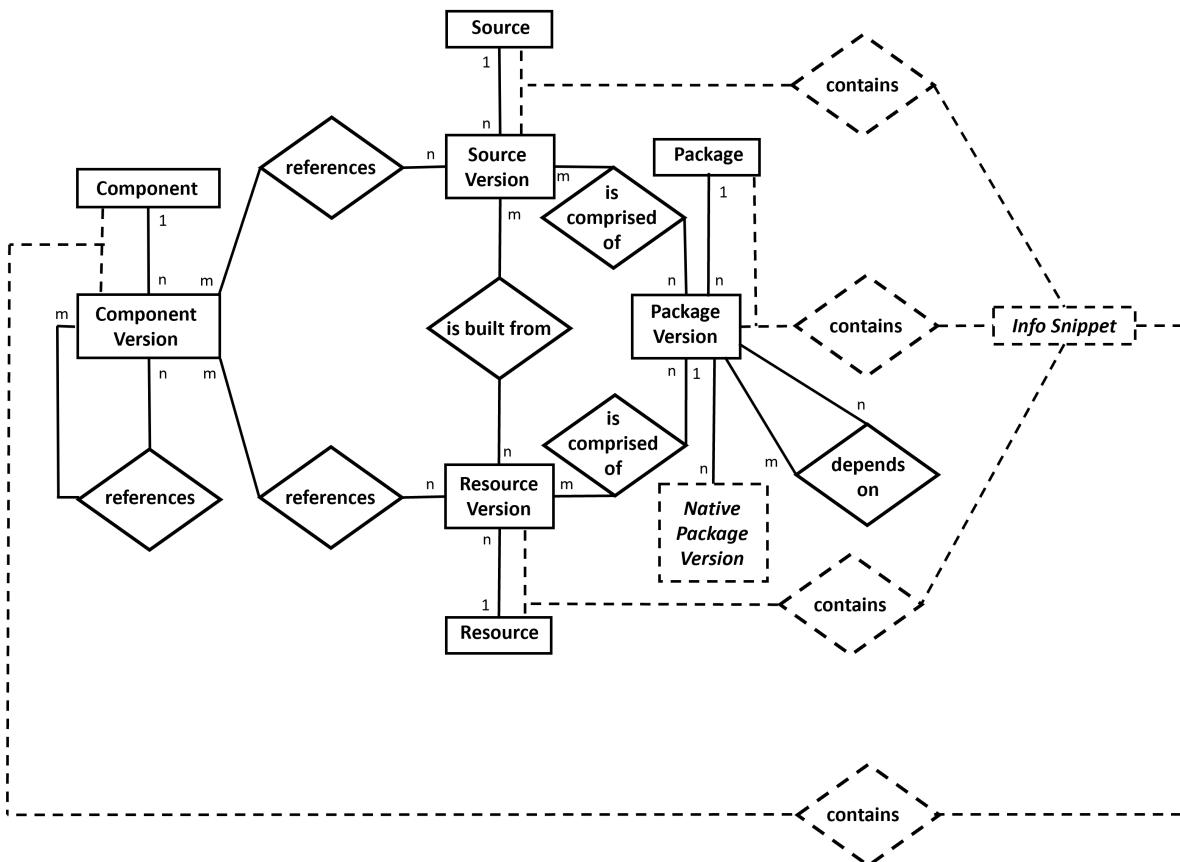


Figure 5.1: Meta Data Model
Source: Own Representation

Although the motivation might not be obvious for all elements immediately, the model as a whole should feel quite familiar from the abstract data lake design.

The previous chapter shows that the *Data Lake* is built around and based on the OCM. Consequently, the data model is also inspired by the OCM. As mentioned above, especially the entity type *component* is lend from its model elements. Since this thesis is written with SAP Gardener, a seamless integration of the *Security and Compliance Data Lake* such as described in the previous chapter is also a requirement, even though it is not explicitly listed.

But to stress this again, the *Data Lake* is nonetheless designed to be independent of the OCM. Thus, it is entirely possible to use a different kind of *Component Model*.

As an example, if one is able to provide the capabilities of a *Component Model* and express the concept of *components* and *artifacts* with the means of the SPDX standard, one could use SPDX instead of OCM to provide this structural information. Or, since SPDX is not optimal for this purpose, one could create and use an own *Component Model*.

An additional notion, there is generally no necessity to distinct between *sources* and *resources*. *Sources* could be treated as *resources*, at the only cost of losing the connecting "is built from"-information between the two entity types.

5.2.1 Meta Data Model

Contrary to common ERMs, the ERM in figure 5.1 does not have any properties. There are two major reasons for this. Firstly, the just mentioned independence of a specific *Component Model* would hardly be possible if the data model would define fixed predefined properties for each entity type. Secondly, the different scanning tools provide a wide range of information about packages, and other data sources apart from scanning tools may also be added. Consequently, it is practically impossible to predict what properties may be needed. Besides, these may vary depending on the user of the *Data Lake*.

Another special feature of the ERM are the entity types and relationships illustrated with dashed lines. These represent *classes of entity types* and their potential relationships. Since the whole set of data sources cannot be known upfront, the whole set of potentially required *Info Snippets* cannot be known upfront either. As already mentioned in several examples before, when adding a build tool as data source, an entity type *Build Information* may be needed. Also, the relationships of different *Info Snippet* entity types may vary. While a *Vulnerability* or a *License* is usually *contained* in multiple *Package Versions* leading to a (n:m)-relationship, a *Build Information* is usually associated to one *Resource Version* leading to a (1:n)-relationship. But generally, *Info Snippets* could be associated to any other entity type in the data model with any cardinality.

In conclusion, the data model allows to configure arbitrary properties for each entity type. Moreover, it allows to instantiate multiple entire entity types of these classes of entity types, with corresponding relationships and also with arbitrarily configurable properties. Therefore, it is actually a meta data model.

This kind of flexibility is necessary to enable R.1 (consume and store metadata from multiple different data sources).

The *Native Package Version* correspondingly illustrates the representation of a package, native to a concrete data source. Thus, entity type instances of *Native Package Versions* may be *BDBA Package*, *Mend Package* or even *Jenkins Package*. So, if all three data sources provide information about the exact same *Package Version*, each representation may be stored without a need to merge their properties before storing. Thereby, this enables R.2 (store metadata from different data sources without aggregation).

A set of properties commonly provided by all of the data sources may be aggregated on *Package Version* level, thereby also enabling R.5 (provide the metadata from different data sources with aggregation). To elaborate on this, all the *Native Package Versions* representing the exact same package are related to the same abstract *Package Version* on the model level. As the different data sources may use different identifiers for the packages, the merging process cannot be triggered automatically. Hence, until a human defines that the *BDBA Package*, the *Mend Package* and potentially also the *Jenkins Package* are representations of the same package, no merging is done and each of these *Native Package Versions* is related to an individual abstract *Package Version*.

The previous chapter frequently mentioned the need for a technology-agnostic uniform identification scheme. The aggregation of *Native Package Versions* through *Package Versions* allows to configure such a technology-agnostic uniform identification scheme for packages within the context of the *Data Lake* (thereby, the uniform identification scheme may just be a simple ID or UUID).

So, after explaining the special features of above ERM, the general model may be discussed.

Component, Component Version and Relationships

The basic entity type *component* is broken down into two distinct entity types, *Component* and *Component Version*. As immediately noticeable, this distinction is done for each of the basic entity types. *Component* is a purely abstract entity type.

It merely groups all the versions of the same component together. The *Component* may provide information about the semantics of this grouping such as whether this *Component* describes a specific software product or whether it describes all software products used by a department. Thus, information that is identical for all versions of this component and would have to be stored redundantly for each *Component Version* otherwise. Naturally, there are multiple *Component Versions* of each *Component*. Therefore, the (1:n)-cardinality here is self-explaining.

As established by the previously described grouping semantic of *components*, a *Component Version* may reference multiple other *Component Versions*. For example, a *Component Version* describing a specific version of a software product may reference multiple other *Component Versions* such as *Component Versions* describing specific versions of a web server, a service and a database. Reciprocal, a *Component Version* may of course be referenced by multiple *Component Versions*. For example, a *Component Version* describing a web server may be referenced by several *Component Versions* describing different versions of the same software product or entirely different software products. Thus, this is a recursive (n:m)-relationship. There may also be a need to store additional occurrence specific metadata as properties of the *references*. To continue with the web server example, such occurrence specific metadata may provide information about the usage of the web server within the software product, hence whether it is used as a HTTP server or as a load balancer. Together, these model elements fulfill requirement R.4 (provide an aggregation level to group sources and resources).

Furthermore, a *Component Version* may also reference multiple *Source Versions* and *Resource Versions*. As an example, the *Resource Versions* composing the web server and the *Source Versions* from which the respective *Resource Versions* were built. As before, with the recursive relationship of *Component Versions*, the *Source Versions* and *Resource Versions* may also be referenced by multiple *Component Versions*, resulting in a (n:m)-relationship. Again, there may be a need to store additional occurrence specific metadata as properties of the *references*. Specifically, these *references* may be used to store *triage* information. As this *reference* describes the usage context of the *Artifact*, one may, for example, decide whether a copyleft license is or is not acceptable in this case.

Artifact, Artifact Versions and Relationships

The relationship between *Source* and *Source Version* as well as between *Resource* and *Resource Version* is similar to the relationship between *Component* and *Component Versions*. But the abstract *Source* and *Resource* entity types actually do have a

concrete purpose apart from preventing redundant storage of certain properties. In this case, these abstract entities may have properties to store *triage policies*. As an example, one may store that a specific vulnerability may be ignored for the usage of *Resource Versions* v1.0.0 to v1.2.3 of a respective *Resource* within *Component Versions* v1.4.2 to v1.4.12 of a specific *Component*. The *references* and the abstract *Source* and *Resource* entity types thereby enable the fulfillment of requirement R.6 (enable users to perform assessments).

Resource Versions may also reference the *Source Versions* they are built from. A *Resource Version* may be built from multiple *Source Versions* and a *Source Version* may be used to build multiple *Resource Versions*. This also results in a (n:m)-relationship.

Since *Artifacts* are comprised of *Package Versions* and the same *Package Version* may occur in multiple *Artifacts*, both *Source Version* as well as *Resource Version* have a (n:m)-relationship to *Package Version*. Again, there may be a need to store additional occurrence specific metadata as properties of the *is comprised of* relationship.

Package, Package Version, Native Package Version and Relationships

The relationship between *Package* and *Package Version* is again similar to the relationship between *Component* and *Component Versions*. But as already explained, *packages* are broken down even further into three different entity types and thereby three aggregate levels. Since several data sources may provide different representations, thus, different *Native Package Versions*, of the same *Package Version*, the cardinality of this relationship is (1:n).

Package Versions frequently *depend on* other *Package Versions* and so on., which may lead to long chains of dependencies. This has to be kept in mind as these transitive dependencies are also relevant when trying to answer the question, whether a certain *Artifact* or *Component* contains a specific *Package* such as Log4j, and its corresponding vulnerabilities.

Info Snippet and Relationships

Finally, there is the *Info Snippet* class of entity types. As explained above, different *Info Snippet* entity types may have relationships to different entity types with different cardinalities.

5.2.2 Insights into the Development Process

At this point, some insight into the development process may be beneficial to understand the design decisions. The scope of this central data store was initially

much narrower. The first PoC was strictly bound to the OCM. The data model predefined the properties of *component* and *artifact* entity types. As it was bound to the OCM, there was no issue in doing so.

But the data model also predefined the properties of the *package* entity types. In fact, the third aggregate level for *packages*, *Native Package Version*, did not exist at all. Instead, the *Package Version* entity type had a set of properties that was hoped to be common and harmonizable throughout all prospective data sources. At this time, the application was tailored towards only having scanning tools as data sources. To define this set of properties, the API documentations of different scanning tools – especially the ones of BDBA and Mend [41] – were analyzed for the common and most important properties. Additionally, to get a better understanding, both scanners were used on some artifacts to get sample data. After that, the provided attributes were narrowed down and some interviews with developers were conducted. A huge effort was made here, as this was such an important decision. Also, instead of having the *Info Snippet* class of entity types, there was only a *Vulnerability* and a *License* entity type whose sets of properties were defined in the same process. Apart from the fact that there was already a substantial amount of disagreement between different developers about which properties are actually required, by the time the PoC was finished, several new use cases were discovered that required additional properties and even entire additional entity types to provide other metadata than vulnerabilities or licenses.

This led to a change in perspective, interpreting the task of defining the right entity types with the right properties rather as a task to make the entire application extensible regarding the respective entity types and properties. But this flexibility and extensibility comes at the cost of a highly increased overall complexity. Apart from remodeling the data model, it also required a completely new architecture and a completely different implementation. Thus, only after that, the scope widened drastically, also allowing other component models. Consequently, as a kind of disclaimer, the design process presented here is not exactly as chronological as it may initially seem, since it hides a complete development iteration leading up to the final concepts.

5.2.3 Application of the Meta Data Model

Although a great effort was made to make the meta data model as tangible as possible, it might still be hard to grasp due to its abstract nature, omitting properties and introducing classes of entity types. Therefore, this section discusses how this meta data model could be applied. As reference component model, the OCM is

5.2. Data Model

used, and as reference data sources, only BDBA is used. This thereby also reveals a problem that has to be faced when applying this data model to a real world use case. The figure 5.2 shows the corresponding data model instance.

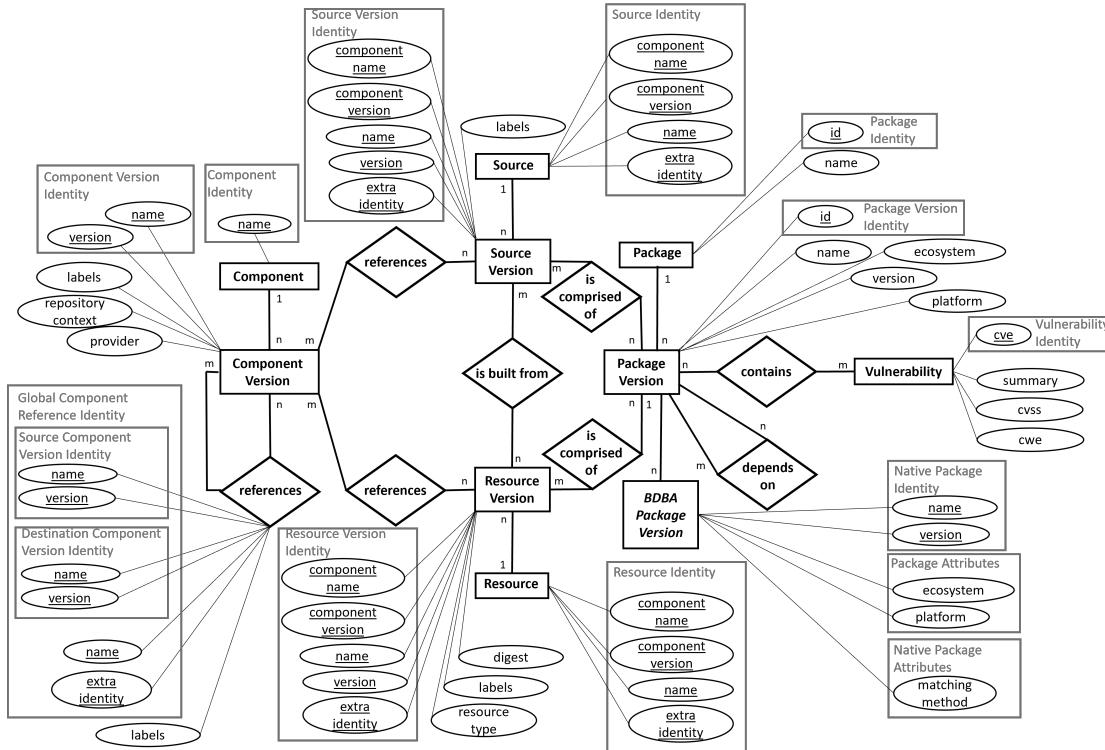


Figure 5.2: Example Data Model Instance

Source: Own Representation

The figure is very crowded. But this is necessary as the purpose of it is to provide a concrete and thorough illustration of how the meta data model may be applied to specific component models and data sources. Thus, it shows a concrete example of how to instantiate the meta model's entity types for a concrete environment with appropriate attributes.

The important aspect to point out here is the relationship between *Component Version* and *Source Version* and between *Component Version* and *Resource Version*. Although depicted as a relationship with (n:m)-cardinality in compliance with the meta data model, due to the *Component Version-Local Identity* of *Source Version* and *Resource Version*, these can actually only be (1:n)-relationships. There were detailed explanations about this in section 4.1 "Open Component Model". The SAP Gardener team chose this initially rather non-intuitive specification of *Source Versions* and *Resource Versions* with *Local Identities* since in practice, it is difficult to reliably determine whether two referenced technical artifacts are actually the same technical artifact.

The following section further explains the issue and discusses different approaches of dealing with this *artifact identity problem*. It thereby points out problems of each approach.

As the common software developer immediately attempts to translate this data model into tables, a corresponding exemplary relational model is shown in the appendix A.2.

5.2.4 Artifact Identity Problem

Essentially, there exist 3 options to identify artifacts: by arbitrarily-defined names, by location and by content.

Arbitrarily-Defined Name

The initial idea and most obvious approach is to treat technical artifacts just as components. Thus, assign an arbitrarily-defined globally unique name and a version to each technical artifact.

This works for components since these are *technology-agnostic* aggregations of artifacts *representing a specific semantic*. As components are *technology-agnostic*, they can all conveniently be stored in a central component repository (as shown in the architecture in 4.3). There, the uniqueness of the name and accessibility through the name can be ensured. As the arbitrarily-defined name of a component identifies a dedicated usage context or meaning for the provided software, the name of two components with the exact same technical content may be different.

Opposed to that, for unification, or rather de-duplication, in the *Data Lake*, artifacts are intuitively expected to *represent a technology specific content*. As artifacts are technology specific (OCI Images, Maven Package, ...), they are frequently stored in multiple different technology specific repositories. There, the name has to abide to technology specific naming schemes and its uniqueness is only ensured within the specific technology or repository. As the arbitrarily-defined name of an artifact would identify a specific content, the name of two artifacts with the exact same technical content should be the same.

Thus, to be able to identify artifacts through a arbitrarily-defined globally unique name within the *Data Lake*, the artifacts from multiple different technologies and location with their technology or location specific names would have to be mapped to their globally unique arbitrarily-defined name during the storage process. But thereby, it is difficult to determine whether two artifacts stored in different technology specific repositories are the same technical artifact (an approach for this and its corresponding problems are presented further down under *Content*).

Location

To avoid such a mapping and its implications while still being able to uniquely identify and access artifacts independent of the technology specific repository they are stored in, the name has to be directly coupled to the location. Thus, something like a URL pointing to the artifact has to be used.

But this approach is not flexible as an artifact can only be stored in one location. Otherwise, artifacts with the exact same technical content could have different names.

Content

Another option to identify artifacts would be through their content. Thus, to determine whether two artifacts are technically the same, even though they are located in different repositories is to calculate a suitable *normalized digest* of its content.

A *digest* refers to a short, fixed-length string calculated through a hash function [42]. There are several features in which artifacts may vary but that are irrelevant for their comparison. For example, the exact same source code file may provide different digests when hashed with the exact same hash method, depending on whether it is hashed on a windows or a linux machine. This is due to the different line endings, thus "Carriage Return and Line Feed" on windows and "Line Feed" on linux. In order to calculate a normalized digest, the input data for the hash function is prepared so that such irrelevant features do not have an impact on the digest. So, the input data is normalized.

This approach is generally reliable and clean, but there are again several problems. When looking at figure 5.2 above, there is already a digest property available for *Resource Versions*. As described in section 4.1 "Open Component Model", this digest property specifies a hash function, a normalization function and the corresponding digest value. Although the primary reason for this is that different kinds of resources, for example executables and OCI Images, may need different kinds of normalization functions, this also means that the same artifact may have different digest values in different *Component Versions*.

Consequently, the digests need to be calculated in a standardized way, independent of the OCM. But even then, there is still another issue regarding the data model. There is no way to decide whether artifacts with different normalized digests may be different versions of the same artifact. Thus, the relationship to a corresponding abstract entity, *Source* and *Resource*, required for triage policies is practically impossible to maintain.

So none of these approaches is practical. But there are 2 variations of the basic options that might be leveraged for the *Data Lake*:

Central Origin Location

All artifacts could be replicated into a central origin repository. The URL pointing to the artifacts within this central origin repository could be used as their identifier.

Additionally, the artifacts could have an access attribute similar to the definition of the OCM. When artifacts are replicated into other repositories, they only change their access attribute but keep the original URL of the central repository as name. This ensures assures uniqueness of the name and technical identity per construction while allowing multiple different technology specific repositories.

SAP Gardener did not consider this approach anyway, as it is more complex and globally unique identifiers for artifacts were not relevant for the original scope of the OCM, which was deployment automation.

Component Local Arbitrarily-Defined Names

This is the approach used in the OCM. Compared to what was explained in the context of *Arbitrarily-Defined Names*, with this approach, the arbitrarily-defined name of an artifact also describes its semantic instead of representing its content. Thus, it waives the necessity to globally uniquely identify the technical artifacts in terms of their contents. This concept was explained in detail in section 4.1 "Open Component Model". Although it is not intuitive, only having *local identifiers for artifacts is still functional regarding the goals which shall be achieved by the data model*.

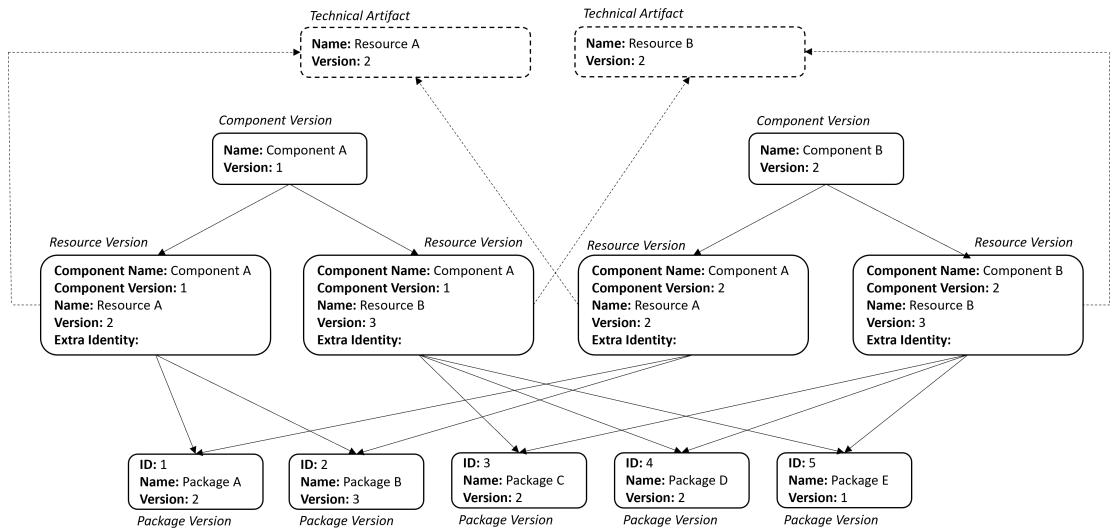


Figure 5.3: Local Artifact Identities

Source: Own Representation

Figure 5.3 shows two components that reference the exact same technical artifacts. But since the *Identity of Resource Versions* is local, the component name and version are part of the *Resource Versions' Identity*, as also shown in figure 5.2. Thus, it appears as they are referencing different artifacts. But as the packages comprising the *Resource Versions* are identified by scanning the technical artifact accessed through the access property, they are the same for the *Resource Versions* which practically reference the same technical artifact.

So, for answering questions such as: "Which software installation contains Log4j and its corresponding vulnerabilities?", the local artifact identities do not make a difference.

In practice, it may also be assumed that within a *Component*, *Resource Versions* with the same name but different version numbers are in fact pointing to different versions of the same technical artifact. *Thus, local identities are usually kept stable among successive version of a Component.*

Therefore, the *Resource* aggregate level is also possible. But its scope and the scope of *triage policies* within these entities respectively is of course limited to a certain *Component*.

5.3 Database

After the application context and data model are defined, a suitable database has to be selected. Therefore, this section analyzes the most relevant database technologies regarding their applicability as a central data store based on the previously defined

data model.

5.3.1 Requirements for the Database

The most important aspect regarding the suitability of a database technology is the purpose of the data store and the respective kind of usage. Several relevant factors depend on this information. Is read or write performance more important? What may be acceptable delays when reading or writing data? What kind of queries may be used most frequently? Will the data be used for statistical analysis and require a lot of aggregation over dimensions such as time periods or locations? Are the entities highly inter-connected and require traversing relationships efficiently? How many users may need to access the data concurrently? May the database need to be distributed?

Write Performance: The data may be provided by all sorts of data sources. As already described in the reference architecture in section 4.2.2 "Reference Integration Architecture", the data may be fetched based on policies. Considering scanning tools, such policies may require to scan all Component Versions or respectively the referenced technical artifacts once a day or based on particular events, such as upon creation of a new Component Version. Depending on the number and size of artifacts, the scans themselves take something between several minutes up to several hours. As the information in the database may be updated concurrently with the scanning process, so for example every time the scanning tool has finished scanning a particular artifact, the results may be fetched, the *write performance may be quite low*. Theoretically, the database could correspondingly to the scanning tools take up to several minutes to write or update the information about an artifact.

Read Performance: As already pointed out in section 5.1.2 "Non-functional Requirements", the Security and Compliance Data Lake may prospectively be used as backend for dashboard web applications. Of course, these dashboards are for technical professionals. Thus, the response time does not necessarily need to abide to common UX design requirements. But still, common queries should ideally not take more than several seconds. So *read performance should be high*.

Queries: The most popular example for a query is "Which software installations contain Log4j?" or more precisely in terms of the previously introduced data model "Which Component Versions contain a vulnerable Log4j Package Version?" or even "Which Component Versions contain the Log4j Vulnerability?". Considering the exemplary data model in figure 5.2, to resolve this query, one has to find the *Vulnerability* with *CVE CVE-2021-44228* [2], find all *Package Versions* that contain this

Vulnerability and also all *Package Versions* that directly or transitively depend on one of those *Package Versions*, find all *Source* and *Resource Versions* that contain one of the respective *Package Versions* and finally find all *Component Versions* that reference one of the respective *Artifacts*. Thus, resolving this query requires traversing a lot of relationships. This is prospectively also the most popular kind of query in general. Queries including aggregations may be something like "What Vulnerability contained in a Deployment X has the highest CVSS?". Therefore, even the queries including aggregation require a traversal of relationships to retrieve the subset of data the aggregation function has to be applied to.

Distribution: The database is the central metadata store of a company, used to monitor the application landscape and increase the transparency of the infrastructure. Therefore, the number of concurrent users will not be exceedingly high. Besides, the whole application and thus, also the underlying database is not business critical for a company. So, database distribution to increase scalability in terms of parallel queries or to increase availability and fault tolerance are not a primary concern for the Security and Compliance Data Lake.

From here on, this rough overview provides a good idea of the most important properties to consider when selecting a database technology for the Security and Compliance Data Lake.

5.3.2 Relational Databases

The first database technology analyzed are *relational databases*. They are by far the most popular and widely used database technology. This is represented in Stack Overflow's 2021 developer survey [43]. There, the top 3 most used databases are MySQL, PostgreSQL and SQLite, which are all relational databases. Among the technologies discussed here, it also has been around for the longest as the original paper introducing the relational model for databases by Edgar Codd was published in 1970 [44]. Therefore, the technology itself is very mature and there is a lot of know-how in the industry. Thus, relational databases are worth considering for every enterprise project.

Theoretical Foundations

A quick repetition of the foundations of relational databases. The name "relational database" stems from the mathematical definition of *relation*:

Given sets S₁, S₂,..., S_n (not necessarily distinct), R is a relation on

these n sets if it is a set of n -tuples, the first component of which is drawn from S_1 , the second component from S_2 , and so on. More concisely, R is a subset of the Cartesian product $S_1 \times S_2 \times \dots \times S_n$. [45]

Thus, a mathematical relation is an unordered set of tuples of the same type. Mapping this to the *relational model*, a table represents a relation, each row represents a tuple of this relation, the ordering of the rows is irrelevant and as per definition of sets, all rows are distinct from one another [45].

These relations are then used to represent objects of a certain problem domain, therefore *entities* and their *relationships*. Each object type has a distinct type identifier, which becomes the name of the relation. Every instance of an object type must have an instance identifier, which uniquely identifies the entity among all the other entities of the same type. This identifier is commonly referred to as *primary key* [45].

Relationships between entities are mapped to relations by using the primary key of a related entity as a reference. In the scope of the referencing relation, the primary key of the referenced relation is commonly referred to as *foreign key* [45].

Furthermore, the relations are usually normalized, which leads to a low level of redundancy and reduces the risk for anomalies during data manipulation, thereby enabling the adherence to the ACID criteria [45].

In order to being able to consider and compare performance, a basic understanding of how the data is stored and accessed with each database technology is required. For relational databases, these basics are explained excellently and with attention to detail by Ramez Elmasri and Shamkant Navathe in "Fundamentals of Database Systems" [46]. Based on this book, the most important principles which are also necessary for further understanding are introduced.

The data stored on disk is organized as *files of records*. Generally, a record is a collection of data values that can be interpreted as facts about entities, their attributes, and their relationships. In the context of relational databases, a record usually refers to a tuple of a relation, or respectively to a row in a table. Thus, the data is stored in a 1-dimensional format, listing all rows of a table.

There are several techniques, determining how the file records are physically placed on the disk, and hence how the records can be accessed. These techniques are commonly referred to as *file organizations* or *primary file organizations*. A *heap file* or *unordered file*, as the name suggests, places the records on disk in no particular order by appending new records at the end of the file. Opposed to that, a *sorted file* or *sequential file* keeps the records ordered by the value of a particular field called the

sort key. Correspondingly, a *hashed file* uses a hash function applied to a particular field called the *hash key* to determine a record's placement on disk. *B/B⁺-trees* use tree structures. *Secondary file organizations* enable efficient access to file records based on *alternate fields* than those that have been used for primary file organization.

The implications are quite intuitive. In the case of heap files, without additional indexing, thus secondary file organizations, a *linear search* looking through each record has to be conducted when searching based on a specific condition. Even with file organizations and indexing, this is still the case for conditions only considering non-key fields [46].

As common relational database technologies such as InnoDB, the storage engine behind MySQL and MariaDB, use B⁺-trees for both, primary as well as secondary file organization, this is discussed in further detail [47]. Theoretically, hash key file organizations, or rather hash indexes, are capable of providing even more efficient look ups, but due to certain limitations and implications in the context of relational databases, these are not suitable for such an application and therefore not further considered here.

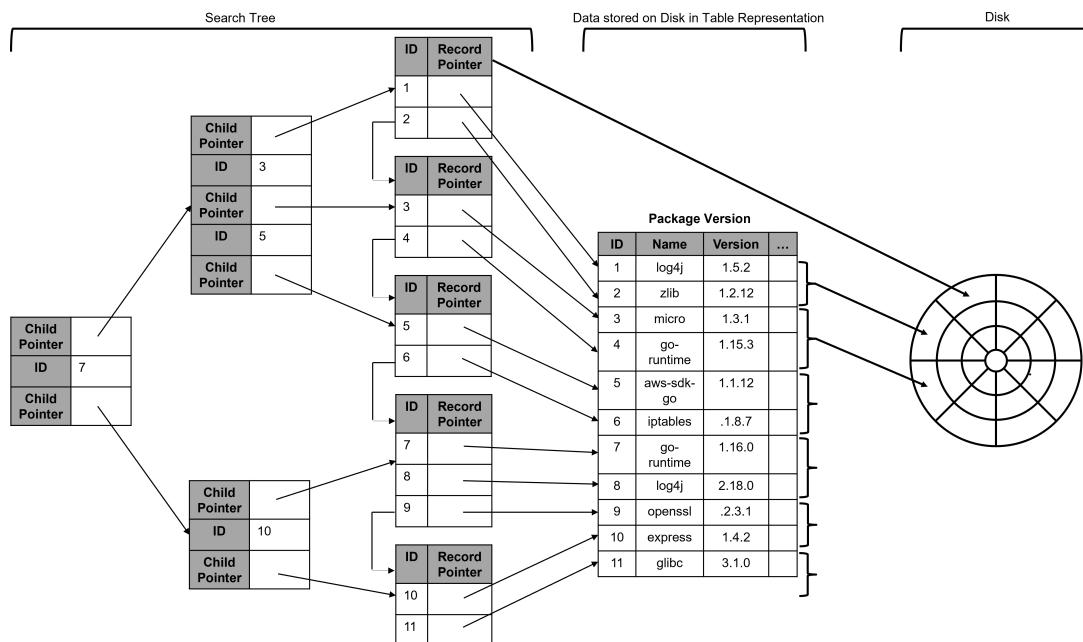


Figure 5.4: B⁺-tree as 4-way Search Tree

Source: Based on [46]

So, figure 5.4 provides a holistic view onto the search trees and how they are used in database systems with some simplifications and assumptions.

On the right of the figure is a representation of a *magnetic disk*, which is still most commonly used as storage for large databases. The concentric circles are called tracks. These tracks are further divided into *sectors*, represented by the intersections

of the concentric and straight lines. This division leads to equal-sized *disk blocks*³ which are also commonly called *pages*. The pages are the areas the arrows are pointing to in the figure. With InnoDB, the block size may be configured between 4KB and 64KB.

Data in secondary storage such as a disk cannot be processed directly by the *central processing unit (CPU)*. First, it must be copied into a primary storage such as main memory which is usually *dynamic random access memory (DRAM)*. The units in which this data is transferred between disk and main memory are the just introduced pages. Thus, to read data contained in a certain page, the whole page is copied into the main memory. And to change that data, it is edited and then the whole page is written back, or in other words copied, to the disk. As accessing pages on the disk is in the order of milliseconds while accessing RAM is in the order of nanoseconds, this is a major bottleneck and therefore, the goal is to *minimize the number of block transfers*. So this enables to understand the necessity of *multilevel indexes* or *search trees*.

By only looking at figure 5.4 without the background knowledge just provided, there would be no apparent reason to construct and store such a sophisticated search tree. It would be more efficient to just store the index as a sorted table. Then, when accessing the database with said index, load the whole table into main memory and conduct a binary search. Thus, the only reason to use such search trees is when the index itself outgrows the page size and therefore has to be stored in multiple pages. So, in practice, each node of a search tree is stored in a separate page as indicated by the arrow from the upper index leaf node. By traversing this search tree, the number of necessary page accesses may be reduced significantly. The index size with a search tree of order 4, so a node can at most have 4 children, is probably the most obvious simplification in above figure. The node and leaf indexes are so small that they could easily be stored together in a single index table.

The example shows the search tree as used for secondary file organization. So the records could generally be stored as a heap file, with no ordering at all. If used for primary file organization, instead of having pointers to the actual records, the leaves would directly store the records enforcing a respective order.

To traverse such a search tree, for example to find the record with ID 4, one starts at the root node. This is the node illustrated on the far left. As $4 < 7$, one follows the pointer above the 7 to the corresponding child node. In practice, this is a

³One may wonder that these blocks are not equal sized in the figure. The outer ones are actually much larger than the inner ones. In practice, there are different types of sector organization. The type shown in the picture maintains a fixed angle. To provide equal-sized disk blocks nonetheless, the outer sectors usually have a lower record density [46]

pointer to the page of the child node and consequently leads to copying this page into main memory. Then, this process repeats, as $3 < 4 < 5$, one follows the pointer between 3 and 5 to the next child node. In this example, this is already a leaf node, containing the pointer to the actual record with ID 4. So here, 4 pages have to be loaded, 3 index pages and the page containing the actual record.

So up until now, instead of B-tree or B⁺-tree, the term search tree was used. This is because B/B⁺-trees are essentially search trees which have to abide to an extra set of constraints. To be more concrete, a *B-tree of order m* is a tree which satisfies the following properties [48].

1. *Every node has at most m children.*
2. *Every node, except for the root and the leaves, has at least $\lceil m/2 \rceil$ children*
3. *The root has at least 2 children (unless it is a leaf).*
4. *All leaves appear on the same level*
5. *A nonleaf node with k children contains k-1 keys*

The goals of *balancing*, which is another term for all leaf nodes being on the same level, is to make the search speed uniform. Thus, the average time to find any random key is roughly the same. Furthermore, these constraints ensure that the nodes stay relatively full and do not end up empty if there are many deletions, thereby preventing the waste of storage space and an unnecessary high number of levels.

Figure 5.4 shows a B⁺-tree. A B⁺-tree is actually a variation of a B-tree which stores record pointers only at the leaf nodes. Consequently, the leaf nodes have an entry for every ID leading to some IDs – specifically 7, 3, 5 and 10 – being contained twice in the search tree. On the contrary, in a B-tree every ID is only present once at some level in the tree and the record pointer is stored directly alongside the child pointers. Thus, in above figure, the root node would additionally have a pointer to the record with ID 7. But of course, the whole tree structure would be different, as several leaf nodes could be omitted. Furthermore, B⁺-trees usually link the leaf nodes to provide ordered access. This is indicated by the arrows between the leaf nodes. In practice, this linking is also done with an additional pointer [46].

Finally, the *upper bound on running time* for searches with B-trees is examined. Initially, it is important to understand, the leaves carry essentially no information searching wise. Thus, for these considerations, leaves may just be regarded as

terminal nodes. As explained in Donald Knuth's "Art of Computer Programming - Sorting and Searching" [48], suppose there are N keys, and the $N+1$ leaves appear on level l (In the book the relation that a B-tree with N keys always has $N+1$ leaves is just given. For a further explanation on why this is always the case, refer to appendix A.3). Then, as per constraint, the number of nodes on levels 1,2,3, ... is at least $2, 2\lceil m/2 \rceil, 2\lceil m/2 \rceil^2, \dots$, hence

$$N + 1 \geq 2\lceil m/2 \rceil^l - 1$$

Solving the equation for l

$$l \leq 1 + \log_{\lceil m/2 \rceil}(\frac{N+1}{2})$$

Furthermore, on every level at most m keys have to be searched. As the elements in the nodes of a B-tree are sorted, binary search may be used. Therefore, the maximum number of look ups s per level is $\log_2(m)$. Consequently, the total number of look ups is

$$s \leq \lceil (1 + \log_{\lceil m/2 \rceil}(\frac{N+1}{2})) \rceil * \log_2(m)$$

Considering Big O notation, constants may be ignored. Also m is definitely smaller than and independent of N . This leads to a worst case time complexity for searching a B-tree of

$$O(\log N)$$

The estimation is also true for B⁺-trees.⁴

Suitability for the Central Metadata Store

To evaluate and properly illustrate the suitability of the database technology regarding the aforementioned relevant properties, a representative example based on the established data model is used. The complete data model is quite complex and its actual instantiated form, thus the final entity types, relationships and especially the properties, depends on the use case. Therefore, to keep the considerations general

⁴Although this section covers several details, it is still just a very brief introduction into relational databases, the internals of database systems and B-trees, only covering the topics in the depth necessary for understanding the further sections. Aspects such as transactions and ACID criteria, how the records are actually placed on the disk, how variable-length fields are handled or what happens when records exceed page boundaries have been omitted. Regarding B/B⁺-trees, the algorithms used for insertion or deletion and their respective bounds on running time are not discussed either. For further and more detailed information on these topics, refer to the respective sources "The Relational Model for Database Management" [45], "Fundamentals of Database Systems" [46] and "The Art of Computer Programming - Searching and Sorting" [48].

and clear, the representative example only considers the *Package Version* entity type. Since *Package Versions* may be contained in other *Package Versions*, this example allows for considering arbitrarily complex relationship chains.

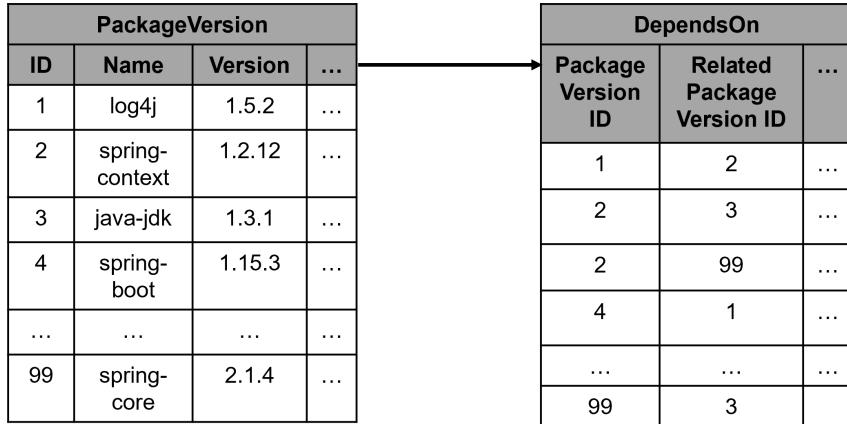


Figure 5.5: Package Versions and Dependencies
Source: Based on [49]

Figure 5.5 shows the *Package Version* entity type and the corresponding *depends on* relationship type mapped to tables and filled with some example entities. As indicated by the ordering of the tables, the example also assumes there are indexes on *ID* and *PackageVersionID*. Furthermore, it may also be assumed, that there is a composite index on *Name* and *Version*.

For simplicity reasons, initially pretend there are no transitive dependencies. In this case, to answer the question "Which Package Versions does the *spring-context:1.2.12* Package Version depend on?" in a relational database, the following SQL query would have to be used:

```

1  SELECT p1.Name, p1.Version
2    FROM PackageVersion p1
3    JOIN DependsOn
4      ON DependsOn.RelatedPackageVersionID = p1.ID
5    JOIN PackageVersion p2
6      ON DependsOn.PackageVersionID = p2.ID
7    WHERE p2.Name = 'spring-context' AND p2.Version = '1.2.12'

```

Listing 5.1: SQL Query – Package Version Dependencies

Based on the sample data in figure 5.5, this returns *java-jdk 1.3.1* and *spring-core 2.1.4* [49]. Although slightly hard to read, the query is still relatively simple and not particularly computationally expensive, because it constraints the number of rows under consideration by applying the filter `WHERE p2.Name = 'spring-context' AND`

`p2.Version = '1.2.12'`. Based on the indexes, the initial look up of this record has a time complexity of $O(\log(n))$ with n being the number of rows in the *PackageVersion* table. The consecutive join then has a time complexity of $O(m * \log(n))$ with m being the number of records found during the previous look up and n being the number of rows in the *DependsOn* table. Each additional join in the query adds a $O(m * \log(n))$. This assumes the so called nested loop join algorithm is used. This is usually the most efficient join algorithm in such scenarios, where indexes on the join condition exist and it is a small number of rows that have to be joined. Relational databases automatically estimate and choose the best strategy [50]. Thus, depending on the number of joins, the general time complexity for such a relationship traversal query may be estimated with

$$O(\log(n)) + O(m * \log(n)) + \dots + O(m * \log(n))$$

Now, to answer the reciprocal question "Which Package Versions depend on the `spring-context:1.2.12` Package Version?" in a relational database, the following SQL query would have to be used:

```

1  SELECT p1.Name, p1.Version
2   FROM PackageVersion p1
3   JOIN DependsOn
4     ON DependsOn.PackageVersionID = p1.ID
5   JOIN PackageVersion p2
6     ON DependsOn.RelatedPackageVersionID = p2.ID
7 WHERE p2.Name = 'spring-context' AND p2.Version = '1.2.12'

```

Listing 5.2: SQL Query – Package Version Reciprocal Dependencies

Based on the sample data in figure 5.5, this returns `log4j 1.5.2` [49]. Although this query looks very similar to the one before, it is computationally more expensive, because there is no index on *RelatedPackageVersionID* and consequently, the whole table has to be considered for the join with a complexity of $O(n)$. In practice, an additional index could be added without problems, as the write performance which would suffer from maintaining an additional index, does not matter too much in this use case. Then joining each record would also have a complexity of $O(\log(n))$, leading to the same overall complexity as the previous query.

Finally, the transitive dependencies have to be considered. To traverse these transitive dependencies in SQL, recursive joins have to be used. Therefore, joining a table with itself. But as the dependency chains may be of arbitrary length, a

simple recursive query is not sufficient. Nowadays, most of the popular relational databases provide a SQL feature, the WITH clause, or rather *Common Table Expressions (CTE)*, to support recursive queries [51, 52, 53]. CTEs are auxiliary statements which can be thought of as defining temporary tables existing for just one query [52]. So, to actually answer the question "Which Package Versions does the `spring-context:1.2.12` Package Version depend on?" correctly, thus also including the transitive dependencies, the following SQL query would have to be used:

```

1  WITH RECURSIVE cte(PackageVersionID) AS (
2      -- Anchor member.
3      SELECT d.RelatedPackageVersionID
4          FROM DependsOn d
5              JOIN PackageVersion p
6                  ON d.PackageVersionID = p.ID
7                  WHERE p.Name = 'spring-context' AND p.Version = '1.2.12'
8
9      UNION ALL
10         -- Recursive member.
11         SELECT d.RelatedPackageVersionID
12             FROM DependsOn d
13                 JOIN cte
14                     ON d.PackageVersionID = cte.PackageVersionID
15
16     SELECT p.Name, p.Version
17         FROM PackageVersion p
18             JOIN cte
19                 ON cte.PackageVersionID = p.ID

```

Listing 5.3: SQL Query – Package Version Dependencies (transitive)

Based on the sample data in figure 5.5, this returns `java-jdk 1.3.1` twice (as union all does not remove duplicates) and `spring-core 2.1.4`. And correspondingly to answer reciprocal question "Which Package Versions depend on the `spring-context:1.2.12` Package Version?", the following SQL query would have to be used:

```

1  WITH RECURSIVE cte(PackageVersionID) AS (
2      -- Anchor member.
3      SELECT d.PackageVersionID
4          FROM DependsOn d
5              JOIN PackageVersion p
6                  ON d.RelatedPackageVersionID = p.ID
7                  WHERE p.Name = 'spring-context' AND p.Version = '1.2.12'
8
9      UNION ALL

```

```
9  -- Recursive member.
10 SELECT d.PackageVersionID
11   FROM DependsOn d
12   JOIN cte
13     ON d.RelatedPackageVersionID = cte.PackageVersionID
14 )
15
16 SELECT p.Name, p.Version
17   FROM PackageVersion p
18   JOIN cte
19     ON cte.PackageVersionID = p.ID
```

Listing 5.4: SQL Query – Package Version Reciprocal Dependencies (transitive)

Based on the sample data in figure 5.5, this returns *log4j 1.5.2* and *spring-boot 1.15.3*. The queries are tested with a PostgreSQL 15 database. The respective DDL statements to set up the sample data are included in appendix B.1.

Besides the computational complexity and the general difficulty of the queries, another aspect to consider are the *schemas* required by relational databases. This may lead to a lot of null values in a table in cases where data sources do not provide the same amount of information for entities of the same type. Depending on how the data is stored on disk, this may lead to a waste of storage space. But most importantly, by *enforcing a schema on database level*, the flexibility to store further raw data provided by the data sources is lost. Furthermore, this makes it difficult to create or adjust schemas at application runtime. This is an issue for a central metadata store, considering that the properties and even several entities depend on the specific component model, data sources and also the use case.

So, as a conclusion for this section, a quick summary of the most important aspects when considering a relational database for the use case. (n:m)-relationships lead to large tables in relational databases. To be able to efficiently traverse these relationships and run queries with a reasonable performance, it is important to define *indexes* on the columns used in the respective join conditions. Relational databases that support recursive CTEs are able to *traverse hierarchies of arbitrary depth* without additional logic in the client side code. The *queries to traverse several relationships become quite large and may be hard to read*. The necessity for a schema on database level leads to a considerable loss of flexibility.

5.3.3 Document Databases

Document Databases have grown tremendously in popularity throughout the last several years. Especially MongoDB which is ranked the top 4 most used database in Stack Overflow's 2021 developer survey has found widespread use as a primary database for applications and therefore, as a substitute for relational databases [43].

Theoretical Foundation

Document databases are a type of NoSQL database which is commonly interpreted as an acronym for "Not only SQL". It is thereby somewhat of an extension of the most simplistic kind of databases, the *key-value databases*. These enable very fast look ups at the cost of only allowing users to address records by their key. The value is usually *opaque* to the database. Correspondingly, document databases also do not require a schema for the values they store, but they store the values in a standard format such as XML, PDF or most frequently JSON [54]. Hence, the values are not opaque and the database may still parse the data and create indexes on non-key fields even though being *schema-less*.

Generally, the primary reason for using document databases, or rather any NoSQL databases, are the horizontal scalability issues of relational databases [54]. While vertical scalability refers to providing the database machine with more compute resources, thus, CPU and RAM, horizontal scalability refers to distributing the database over multiple machines. The relevant questions at this point are: Where do the scalability issues come from? How are they overcome by document databases?

De-normalization is one of the key features of document databases. Thus, while relational databases map entity types and relationships to multiple tables connected by foreign keys, document databases allow storing entities and relationships as nested structures through their schema-less nature. These nested structures are often referred to as *documents* or also more generally *aggregates*, a term from *Domain-Driven Design* describing a collection of related objects that is treated as a unit. Documents are grouped to *collections*, which is kind of the equivalent to a table in relational databases [55]. A common example is shown by figure 5.6 and listing 5.5.

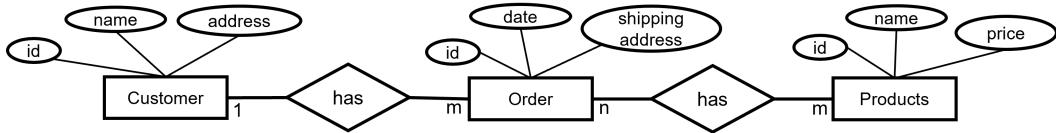


Figure 5.6: Sales Data Model

Source: Based on [55]

```

1  {
2      'id': 1,
3      'name': 'Martin',
4      'address':
5          {'city': 'Berlin'}
6      'orders': [
7          {'id': 99,
8              'date': '04.01.2023',
9              'shippingAddress': {'city': 'Berlin'}
10             'products': [
11                 {'id': 42,
12                     'name': 'The Hitchhiker's Guide to the Galaxy',
13                     'price': 7.50 },
14                 {'id': 12,
15                     'name': 'The Art of Computer Programming',
16                     'price': 120.00 }
17             ]
18         }
19     ]
20 }
21 ...

```

Listing 5.5: Document Database – Data as Nested JSON Document

This solves the so called *impedance mismatch*, the difference between the relational model and in-memory data structures. Hence, with relational databases developers frequently have to translate the nested in-memory data structures to a relational representation to persist them. Of course, there are object-relational mapping frameworks such as Hibernate, but these often lead to performance issues [55]. By allowing to store objects as JSON objects, document databases such as MongoDB enable developers to store and retrieve their *JavaScript* objects or *Python* dictionaries as they are, without any further conversion. This is definitely one of the major reasons for their success, also in use cases where scalability is not a primary concern.

Furthermore, through this de-normalization, document databases *remove the necessity of joining a lot of data*. On one hand, this makes querying objects such

as customer with all its orders easier for developers, as join queries are comparably difficult to write and understand. On the other hand, this increases the performance and scalability of these queries, as joins add quite some computational complexity which even depends on the number of stored objects. Thus, the performance of joins becomes worse as the database grows.

Besides, this kind of data organization enables efficient *sharding*. Sharding is a technique where different parts of the data are put onto different servers. Tables, or in the context of document databases rather collections, are split up into smaller shards which may be placed on multiple servers. A common sharding technique is to hash a key, also called *shard key* in MongoDB [56] or *partition key* in DynamoDB [57]. The respective hash value determines where, so on which server, to place the shard. This hashing provides an even distribution. But there are also other techniques such as ranged sharding. Thereby, the shards are distributed by the key value itself. This may be useful, when the records have some sort of order and it is common to access multiple sequential records. Or the key may be a post code and the distribution is based on the sever location to minimize latency [55]. This whole technique is enabled through the nesting, as documents are collection of related objects that are treated as units. Therefore sharding, or generally any form of partitioning, is rather difficult with relational databases. Since the data is scattered over multiple tables and the query language SQL puts no restrictions on how to query and connect this data, such partitioning would regularly require to send requests over the network to collect all the information [55].

So de-normalization obviously has numerous advantages. And this is also where most getting started guides of document databases such as MongoDB stop [58]. But naturally, it introduces the issues which were solved by normalization in the first place. Considering the example above, there is a (n:m)-relationship between orders and products. As the product is nested in the order which is again nested within the customer, updating a particular product while maintaining data integrity requires searching through all customers and through all orders within those customers, to update every single occurrence. Of course, in some cases write performance may not be relevant. But the nesting also affects queries. Imagine the above example describes the database of a retailer and this retailer wants to analyze its "Hitchhiker's Guide to the Galaxy" product sales over the last year. Again, to do this, the database has to search through all orders within all customers. There is the option to create indexes on nested fields which works pretty similar to relational database indexes to increase this performance. The problem that remains in this case is that document databases still have to retrieve the entire documents that contain the respective product from

disk, thus, every customer record with all orders and all products. This leads to enormous RAM usage and may also slow down the performance significantly [59].

Generally, there are several restrictions which may serve as an orientation whether to nest or embed documents or whether to normalize and link them.

If the application frequently has to *access the embedded document independently of the embedding document(s)* and the *documents are quite large*, it may be best to normalize and link the document rather than embedding due to the otherwise excessive RAM usage [59]. In above example, orders may be queried independently of customers. So it may be better to normalize this and make the orders reference the respective customer.

If documents grow, the database may have to relocate it on disk to an area with more space. In the above example, as a customer may regularly order products from the retailer, the order array within the customer document would grow and potentially require relocation. Such a relocation decreases performance significantly [59]. Moreover, document databases may even have a *hard limit for the document size*. For MongoDB, this limit is 16MB [60]. So, this is another argument to normalize the orders.

Considering all these limitations, many-to-many relationships are especially problematic to de-normalize. Therefore, MongoDB even suggests to *model complex many-to-many relationships in normalized form* [61]. For simple many-to-many relationships, where the references may rarely be updated, it is also a possibility to embed the references instead of creating the equivalent of a join table.

As especially MongoDB and its query language provide a lot of the functionality also provided by SQL, the final question may be, why not store the relational model in a document database? The primary concern therefore may be data integrity. Although MongoDB provides ACID conform transactions, there is no way to enforce referential integrity on database level. Without a database schema, there is no way to tell the database about foreign keys and corresponding constraints.

Furthermore, several operations such as joining are usually more efficient with a relational database. The syntax to perform the join is also simpler with SQL. To make this more tangible, below listing shows the syntax to perform the logical equivalent to a left outer join on orders with products. The example thereby assumes the data is stored according to a relational model as shown in appendix B.2.

```
1 db.orders.aggregate([
2   { '$match': { '_id': 1 } },
3   { '$lookup': {
```

```

4   'from': 'relations',
5   'let': { 'orderId': '$_id' },
6   'pipeline': [ { '$match': { '$expr': { '$eq': [ '$order', '$$orderId' ] } }
7     },
8     { '$lookup': {
9       'from': 'products',
10      'let': { 'productId': '$product' },
11      'pipeline': [ { '$match': { '$expr': { '$eq': [ '$_id', '$$productId' ] } } },
12        ],
13        'as': 'products' }
14      },
15      { '$project': {
16        '_id': 1,
17        'date': 1,
18        'shippingAddress': 1,
19        'products': '$order_products.products' }
20      }
21    ]

```

Listing 5.6: Document Database – MongoDB Join Example

Without going into too much detail about this query, it is obvious that it is rather complex and resembles a query execution plan. Thereby, the output of each stage (\$match, \$lookup, \$lookup, \$project) is the input of the consecutive stage [62]. This drastically limits the automatic optimization capabilities. So MongoDB will generally always perform a nested loop join even without indexes when a relational database would probably perform a merge or hash join.

So document databases and de-normalization work especially well, in cases where the application only requires a limited set of specific repetitive access patterns. But in order to actually leverage the advantages, these access patterns have to be identified and the de-normalization has to be done accordingly. This requires a very good understanding of the application already in the data modeling phase.⁵

⁵Again, although this section covers several details about document databases, it is still just a very brief introduction to provide the foundations necessary for understanding the further sections. Especially the topics of transactions and ACID criteria as well as the CAP theorem and the general treatment of consistency in document databases have been omitted.

Suitability for the Central Metadata Store

In contrast to the corresponding section about the suitability of relational databases, this one is kept short. Based on the background knowledge provided in the theoretical foundations, it is apparent that the advantages of document databases may barely be leveraged for this use case. Extensive examples on how to implement the central metadata store with the capabilities of a document database are therefore omitted as well.

The data model is based on (n:m)-relationships. Also, the database's primary purpose is to perform analysis about different entities. Thus, "Which Component Versions contain a vulnerable Log4j Package Version?" with Package Version being the access point. Or reciprocal, "Which Package Versions are contained in a particular Component Version?" with Component Version being the access point. So, there is not a limited set of specific repetitive access patterns as every entity may need be used as access point to answer common questions. Therefore, the *central metadata store cannot really be optimized through de-normalization*.

Consequently, at least concerning the major entities, the data model would have to be stored as a relational model. In this case, the document database would have the same scalability issues as relational databases and thereby mitigate a lot of the general advantages of document databases. Relational databases are practically designed with the primary purpose of handling relational models. Consequently, the query language, SQL, as well as the general data processing is better adjusted to such use cases.

But, being *schema-less* on database level may be an advantage. Although there is no way to enforce referential integrity on database level, enforcing a schema on application level allows for easier creation and adjustments at application runtime.

5.3.4 Graph Databases

The last database technology analyzed are *graph databases*. They are also a type of NoSQL database. Although especially Neo4j is doing a lot of advertising, graph databases are still the least known database technology, as also represented by their absence in Stack Overflow's 2021 developer survey [43].

Theoretical Foundation

Again, first a quick repetition of the foundations. The mathematical definition of a graph is the following:

A graph is a pair of sets $G = (V, E)$. In an undirected graph, the elements of E are 2-element subsets $\{v, w\}$ of V . In a directed graph, the elements of E are 2-element tuples (v, w) , thus ordered pairs, of V . [63]

In this definition, V refers to *vertices*, or rather *nodes*, and E refers to *edges*. Quite frequently, two nodes v and w connected through an edge (v, w) are elements of different sets, thus, $v \in S_1$ and $w \in S_2$ with $V = S_1 \cup S_2$. So, E is a subset of the Cartesian product of $S_1 \times S_2$. Looking back at the definition of relations in section 5.3.2 "Theoretical Foundations" of relational databases, this is also the definition of a relation.

Thus, E is a *relation* and V is the union of the *source set* and *destination set*. Now, this is generally well known from functions, which are defined by a *domain* S_1 , a *codomain* S_2 and a *mapping* f , which assigns elements of S_1 to elements of S_2 . Hence, drawing a function in a coordinate system is essentially the visualization of a graph, where commonly the *x-axis* denotes the *source set* and the *y-axis* denotes the *destination set* and the dots, or rather line, for most common functions denotes the *relation*.

In practice with business data, the source set and the destination set are usually sets of discrete and finite elements and the relationship between them is rather artificial and application specific. A primitive example, S_1 may be the set of all package names, S_2 may be the set of all version numbers and the relation E (v, w) may be package versions. Thereby, the source set and destination set are rather irrelevant. Besides, one could argue that for the scope of an application, the union of all v in E is equal to the source set and the union of all w in E is equal to the destination set.

log4j	1.5.2
log4j	1.3.1
spring-context	1.2.12
java-jdk	1.3.1

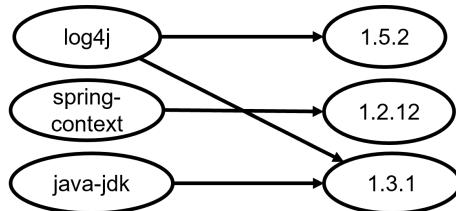


Figure 5.7: Graph Representations
Source: Own Representation

So, for the scope of an application, the table and the vertices and edges in figure 5.7 may be interpreted as two different representation of the same graph. This comparison of relations and graphs may be a slightly abstract and the example is oversimplified, but its purpose is to show how these two concepts and their common representations are tightly coupled.

The first part approached the topic from a mathematical perspective. From here on, this shifts to the actual technologies. Although the previous example may indicate this, in graph databases the graph is not used as a different representation for classical relations, hence to represent the relation of the properties of an entity. It is rather used to represent the relationships between the different entities. Thus, even though figure 5.7 is theoretically accurate, the following figure 5.8 is more suitable from a database technology perspective.

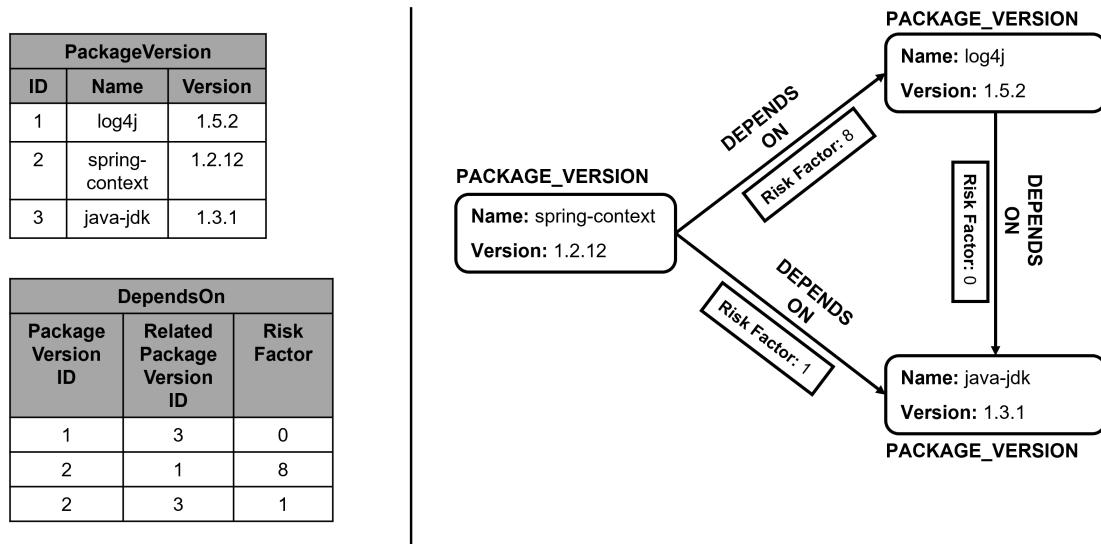


Figure 5.8: Graph Database Representation
Source: Own Representation

This graph model is called *labeled property graph*. As defined in the "Graph Databases - New Opportunities for Connected Data" written by engineers of the Neo4j database, a labeled property graph has the following characteristics [49]:

1. *It contains nodes and relationships.*
2. *Nodes contain properties (key-value-pairs).*
3. *Nodes can be labeled with one or more labels.*
4. *Relationships are named and directed, and always have a start and end node.*
5. *Relationships can also contain properties.*

This model is generally quite intuitive. But obviously, even in this more complex example, with the suitable interpretation of the table (thus, instead of interpreting the relation itself, the foreign key relationships are considered), it still expresses the

same graph as the nodes and relationships. And, as discussed previously, document databases are able to express this graph as well, either also through a normalized model or through nested documents.

So the key difference of graph databases is not that they store graphs in general. The key difference is that their file organization and data processing is optimized to work with graph data.

As shown in figure 5.8 on the bottom left, table representations and consequently relational databases use an index to link entities, or in this context rather nodes. To increase the performance of reciprocal queries, an additional index may be created on *Related Package Version ID*. These *indexes are global*, so as the graph grows, so do the indexes.

In contrast, in most graph databases with native processing capabilities each node maintains direct references to link nodes. Thus, each node maintains kind of a *local index* whose size and consequently performance do not depend on the total size of the graph. This concept of linking adjacent nodes is commonly referred to as *index-free adjacency*.⁶ For further assumptions and explanations, the architecture of the Neo4j database is considered [49].

So, when querying "Which Package Versions does the spring-context:1.2.12 Package Version depend on?", the initial look up works quite similar as in relational databases and therefore still has a complexity of $O(\log(n))$. But finding the related packages is $O(1)$ instead of $O(\log(n))$. Since nodes maintain references to nodes with incoming as well as to nodes with outgoing relationships, the reciprocal query "Which Package Versions depend on the spring-context:1.2.12 Package Version?" may be answered with the exact same complexity [49].

So given these complexities, in theory graph databases scale and perform better than other database technologies when the application involves a lot of graph traversals.⁷

⁶There has been a lot of discussion whether index-free adjacency is a requirement for graph databases [64]. The current consensus seems to be that it is not, since especially ArangoDB developed another approach to the problem which allows similar complexities as index-free adjacency for traversing graphs [65].

⁷As before, this theoretical foundation chapter is obviously not complete. Especially the file organization which enables the fast traversals are an interesting topic. For a good and thorough explanation of a implementation, refer to "Graph Databases - New Opportunities for Connected Data" [49].

Suitability for the Central Metadata Store

This suitability section takes a similar approach as the first one, evaluating the suitability of the relational database. Thus, examining the same representative *Package Version* example. This allows for a direct comparison afterwards. Below figure 5.9 therefore illustrates the equivalent example from 5.3.2 "Theoretical Foundations" of relational databases as a classical graph.

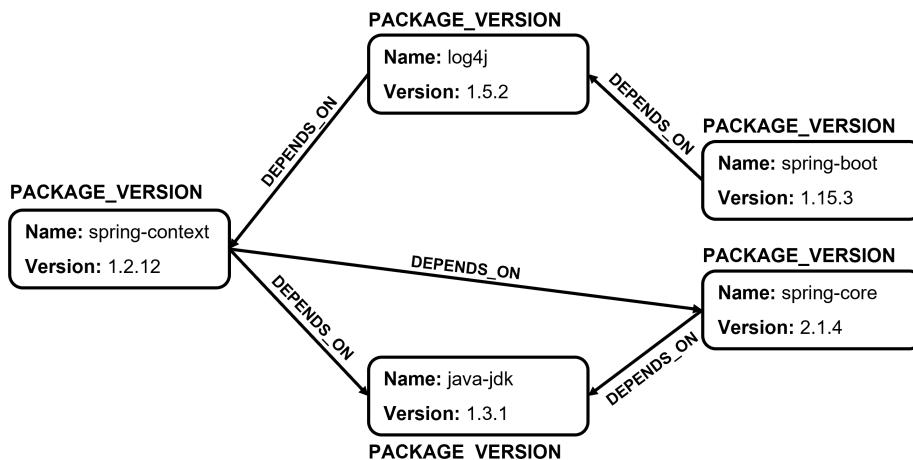


Figure 5.9: Package Versions and Dependencies
Source: Own Representation

Neo4j defines its own SQL inspired language called *cypher*. To answer the respective question "Which Package Version does the `spring-context:1.2.12` Package Version depend on?" in Neo4j, the following cypher query would have to be used:

```

1 MATCH (p1:PACKAGE_VERSION {Name: 'spring-context', Version
  : '1.2.12'})
2 MATCH (p2:PACKAGE_VERSION)
3 WHERE (p1)-[:DEPENDS_ON*1..]->(p2)
4 RETURN (p2)
  
```

Listing 5.7: Cypher Query – Package Version Dependencies (transitive)

The first line in listing 5.7 matches the node with the *label* `PACKAGE_VERSION` and the key-value-pairs `Name: 'spring-context'` and `Version: '1.2.12'` and assigns the resulting nodes, or in this example rather node as there is just a single one that fulfills this pattern, to the variable `p1`. The second and third line match all nodes with the *label* `PACKAGE_VERSION` that `p1` depends on. Thus, the `WHERE` uses a path as a predicate. Thereby, `DEPENDS_ON` specifies the *label* of the relationship and the part after the asterisk specifies the minimum and maximum amounts of hops between the

nodes. In this case, the minimum is 1 because with the default 0 it would consider itself and as there is no explicit maximum given, it defaults to infinite, to consider all transitive dependencies. The nodes that fulfill these criteria are assigned to `p2` and returned, in this case, *java-jdk 1.3.1* and *spring-core 2.1.4*.

Assuming there is a composite index on *Name* and *Version*, the initial look up of the node has a time complexity of $O(\log(n))$ with n being the number of all nodes with the label *PACKAGE_VERSION* [49]. The consecutive look ups then have a time complexity of $O(m)$ with m being the number of nodes found during the previous look up. Thus, depending on the number of hops, the general time complexity for such a relationship traversal query may be estimated as

$$O(\log(n)) + O(m) + \dots + O(m)$$

So theoretically, the graph database performs and scales significantly better than the relational database for this kind of queries. Besides, the query is definitely simpler to read and understand than the SQL query involving Common Table Expressions. Although, it is arguably still easier to learn the Common Table Expression than learning an entire new query language.

Now just for completion, to answer the reciprocal question "Which Package Versions depend on the *spring-context:1.2.12* Package Version?", the following cypher query would have to be used:

```

1 MATCH (p1:PACKAGE_VERSION {Name:'spring-context', Version
   : '1.2.12'})
2 MATCH (p2:PACKAGE_VERSION)
3 WHERE (p1)-[:DEPENDS_ON*1..]-(p2)
4 RETURN (p2)

```

Listing 5.8: Cypher Query – Package Version Reciprocal Dependencies (transitive)

The only difference in this query is the reversed direction of the relationship in line 3. In this case, it return *log4j 1.5.2* and *spring-boot 1.15.3*. The respective cypher statements to set up the sample data are included in appendix B.3.

Besides the efficient traversal of relationship with rather simple queries, there are some other aspects to consider with graph databases. Generally, most graph databases support transactions with *ACID compliance* [66, 67]. Also, there are some mechanisms to support *referential integrity*, but those do not work exactly as in relational databases. So, in graph databases there is no way to create a relationship to or from a node that does not exist. Furthermore, nodes cannot be deleted as long as relationships to or from that node exist. But compared to relational databases,

there is no option to cascade a delete [68].

Regarding *schemas*, or the general storage model, different graph databases vary. Neo4j, for example, only supports key-value properties on nodes or relationships. This is considered as *single model graph database*. In contrast, in ArangoDB every node is a JSON document, which is considered as *multi model graph database* [69]. Therefore, Neo4j does not really have a classical schema. But it still offers several types of constraints such as *unique node property constraints*, *node property existence constraints* and *relationship property existence constraints* to enforce a specific data model [70]. Similar to MongoDB, for ArangoDB a schema may be enforced on application level [71].

5.3.5 Database Selection

In the following, a conclusion of this analysis of the suitability of database technologies for a central metadata store will be drawn.

Relational databases are a solid choice. As shown in the corresponding suitability section, modern SQL provides all features necessary to conveniently implement the required functionality. But the computational complexity of joins is

$$O(\log(n)) + O(m * \log(n)) + \dots + O(m * \log(n))$$

As traversing relationships is involved in the prospectively most frequently used queries, this is a quite important measure.

Document databases are not the technology for this application. As it does not allow for significant de-normalization since the access patterns are versatile, the document database does not have an advantage over a relational database.

Graph databases are optimized for relationship traversal. The previous section has also shown that the cypher queries are simple in comparison to SQL. Due to its purpose, its computational complexity of traversing relationships is better than with relational databases

$$O(\log(n)) + O(m) + \dots + O(m)$$

So, the graph database seems to be the theoretically best database technology for this application. In practice, it also has to be considered that the application has to be supported and that there is less experience around graph databases in the industry. This also increases the chances of unforeseen queries or other situations, where the graph database may be a problem.

For the scope of this work, the graph database is chosen for the implementation to further evaluate its suitability and practicality.

The concrete graph database solutions considered were primarily Neo4j and ArangoDB. Neo4j because it has an open source community edition, good documentation and a supportive community. ArangoDB mainly because it is multi model and therefore enables storing JSON, which offers greater flexibility.

Finally, Neo4j was chosen as it seemed to be the better fit over all. Even in the situations where JSON has to be stored, it is rarely necessary to traverse into the document structure. So it is usually sufficient to access the JSON document by its key which may also be done in Neo4j.

5.4 Programming Language

The *programming language* used to implement the services comprising the Security and Compliance Data Lake is *Go* [72]. Generally, the services could probably be implemented with any common programming language. Therefore, the following application architecture design explanations and discussions are kept *programming language-agnostic*.

Since SAP Gardener is a managed Kubernetes service and Kubernetes is primarily written in Go, Go is the main programming language of the team. Therefore, Go was primarily chosen to fit in into the SAP Gardener ecosystem. This increases interoperability with other services of the team. But most importantly, this enables other members of the team to further maintain the application.

5.5 Application Architecture

To further discuss the application and the other important design decisions involved, especially concerning the API, a rough understanding of how the application is composed may be helpful. Thus, the below figure 5.10 illustrates the application's architecture.

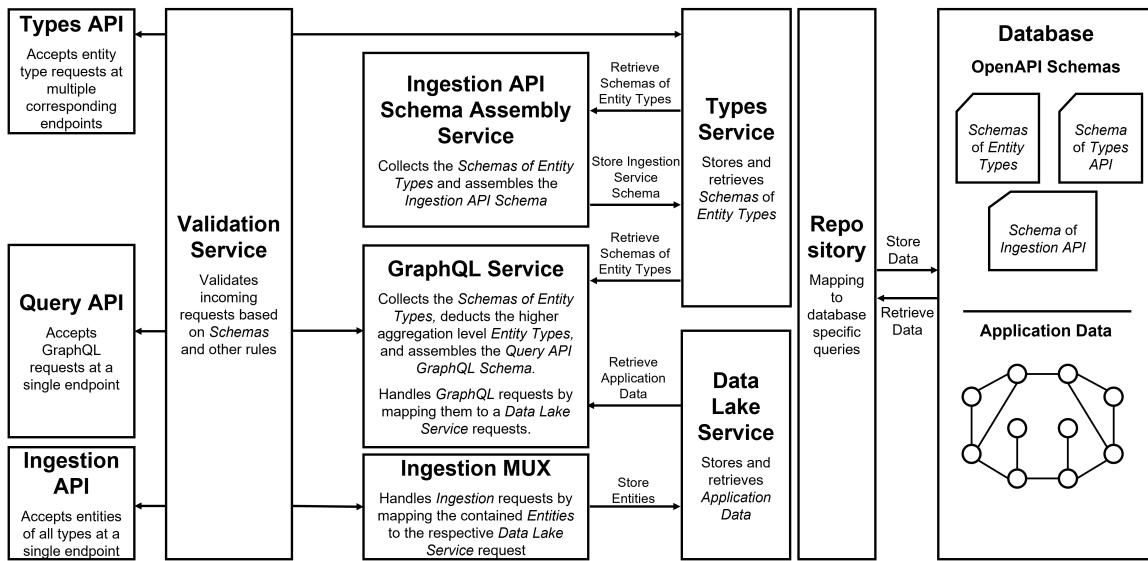


Figure 5.10: Application Architecture

Source: Own Representation

This section only provides an overview of the application. It thereby refers to standards and technologies, especially *OpenAPI* and *GraphQL*, that have not been properly introduced yet. This is done in the following *API & Services* sections along detailed explanations of the respective design decisions.

The application shall be independent of a specific component model and also open to a variety of data sources. To provide this flexibility and extensibility, the application's data model has to be configurable at runtime. This configurability, or in other words ability to instantiate the meta data model as shown by the example in section 5.2.3, is offered through the *Types API*, which allows the application administrator to define and adjust the data model to a specific use case. Therefore, *OpenAPI* schemas of entity types have to be sent to an entity type specific endpoint of the *Types API*. These *OpenAPI* schemas are then validated by the *Validation Service*. Finally, the *Types Service* stores the schemas in the database. Through this *Types API*, the administrator is also able to retrieve existing schemas.

As shown in the figure, the *Types Service* is the central service to store and retrieve schemas from the database. Correspondingly, the *Data Lake Service* is the central service to store and retrieve application data. To make the application rather persistence agnostic, thus, to make it easy to exchange the underlying database, the services' interaction with the database is performed through an abstraction layer referred to as *Repository*.

To enforce the data model configured by the administrator, the *Ingestion API Schema Assembly Service* retrieves the entity type schemas through the *Types Service*

and assembles them to a complete *Ingestion API Schema*. This *Ingestion API Schema* is then also stored in the database.

The entire *Ingestion API* is based on this dynamically, at runtime, constructed schema. Thus, the *Validation Service* validates every request with all entities to be stored against this schema, before it is passed on for further processing. The *Ingestion API* exposes only a single endpoint instead of entity type specific endpoints. This makes the upload easier and more flexible from the adapter's perspective, as requests do not have to be routed to a specific endpoint. Furthermore, considering future performance optimizations, this allows sending larger chunks of data at once.

As the entities of different entity types still have to be processed differently, the *Ingestion MUX* identifies the type of each entity contained in a request and routes it to the respective *Data Lake Service* function.

Finally, in order to being able to query the stored data based on the data model, the *GraphQL Service* constructs a *GraphQL* schema based on the *OpenAPI* schemas of the entity types. The *Query API* therefore also only exposes a single *GraphQL* endpoint. The queries against this endpoint, thus, against the dynamically constructed *GraphQL* schema, are resolved by corresponding functions of the *Data Lake Service*.

The application architecture is essentially a *layered architecture*. In the traditional sense of a the three principle layers, the *APIs* represent the *presentation logic layer* in above figure. Thus, the *APIs* receive HTTP requests, interpret those requests into actions upon the *domain layer* and return information to the user in form of a HTTP responses [73].

The *Validation Service*, *Types Service*, *Ingestion API Schema Assembly Service*, *GraphQL Service*, *Ingestion MUX* and *Data Lake Service* are on the *domain logic layer*. Thus, these services validate data that comes in from the *presentation layer* and perform tasks based on inputs and stored data [73].

Finally, the *Repository* is on the *data source logic layer* as it is responsible for communicating with the database to carry out tasks on behalf of the application [73].

Beyond that, these services are layered even further, as the services further left, thus, service on a higher layer, use functionality provided by services further right (relative to themselves), thus, services on a lower layer. But the lower layer services are unaware of the higher level services [73].

5.5.1 Types API & Services

As explained in the overview, every service of the application builds on the schemas which may be created through the *Types API* and respective *Services*. Thus, the rationale behind this service as well as the design decisions involved, are discussed first.

Schema

As the whole complexity of the architecture revolves around the dynamically configurable schemas, a common question especially considering the *schema-less* database may be why the application even requires such rigid complete schemas?

The term *schema-less* often leads to the misconception that the data somehow does not need to abide to any schema. Technically, this is true for most databases attributed with the term schema-less, like Neo4j. So the database will allow storing any key-value pairs. The problem is, without prior knowledge about the keys, the data cannot be retrieved. And without prior knowledge about the data types, the data may hardly be aggregated. Thus, even with schema-less databases, *partial schemas* are enforced in most use cases.

In this context, partial schema usually means that the application requires a specific predefined set of properties with particular data types. But the application is not limited to those properties. Therefore, a set of unspecified properties may be stored alongside the data defined by a schema.

Such partial schemas may be enforced by the database through constraints or by the application through frameworks, libraries and standards like *JSON Schema* [74] and *OpenAPI*.

The *OpenAPI Specification (OAS)* is a technical specification defining a standard, programming language-agnostic interface description for HTTP APIs. The specification is maintained by the *OpenAPI Initiative (OAI)* and is based on the *Swagger Specification*, donated by *SmartBear* who are also among the founding members [75].

An OpenAPI document may be represented either in JSON or YAML format. To make this tangible, the below listing 5.9 shows an example OpenAPI Document describing an exemplary HTTP API by leveraging the core concepts of the specification.

```
1 openapi: 3.0.0
2   info:
3     title: Sample API
```

```
4     description: Serve as a Sample to explain OpenAPI
5     version: 0.1.9
6   servers:
7     - url: http://api.example.com/v1
8       description: Production Server
9     - url: http://dev-api.example.com
10      description: Development Server
11
12 paths:
13   /customers:
14     get:
15       summary: Returns a list of customers.
16       responses:
17         200:
18           description: A JSON array of users.
19           content:
20             application/json:
21               schema:
22                 $ref: '#/components/schemas/Customers'
23     put:
24       summary: Merge a list of customers.
25       requestBody:
26         content:
27           application/json:
28             schema:
29               $ref: '#/components/schemas/Customers'
30             required: true
31       responses:
32         200:
33           description: Customers merged successfully.
34         400:
35           description: Bad Request.
36           content:
37             application/json:
38               schema:
39                 $ref: '#/components/schemas/Error'
40   /customers/{id}:
41     get:
42       summary: Return a specific customer.
43       parameters:
44         - name: id
45           in: path
46           description: ID of a customer
47           required: true
48           schema:
```

```
49         type: string
50     ...
51 ...
52 components:
53   schemas:
54     Customers:
55       type: array
56       items:
57         $ref: '#/components/schemas/Customer'
58     Customer:
59       type: object
60       required:
61         - id
62         - name
63       properties:
64         id:
65           type: integer
66         name:
67           type: string
68         address:
69           type: object
70           properties:
71             city:
72               type: string
73             country:
74               type: string
75     Error:
76       type: object
77       properties:
78         title:
79           type: string
80         status:
81           type: string
82         detail:
83           type: string
```

Listing 5.9: Example OpenAPI Document

The *openapi* and the *info* fields are both required. The *openapi* field describes the version number of the OpenAPI Specification used. The *info* field provides metadata about the API. The other fields are optional, although it is required that an OpenAPI document contains at least one *paths*, *components* or *webhooks* field. As listing 5.9 has a *paths* as well as a *components* field, this requirement is fulfilled.

The *paths* describe the available endpoints of the API and their respective

operations. Thus, `/customers` and `/customers/{id}` and the operations exposed through the HTTP verbs such as GET and PUT. Thereby, the formats of the *content* in the *requestBody* or in the *responses* may be described, in this case *application/json*. And most importantly, even a *schema* may be provided for each respective content. Thus, the input and output data types may be specified. The specification, or rather schema, for these schemas is a superset of the JSON Schema Specification [76].

The *components* field, or rather object, is a container for defining reusable objects which may be referenced. Most commonly, it is used to define reusable schemas as shown in listing 5.9 for *Customers*, *Customer* and also *Error*. But it is also possible to specify entire *responses*, *requestBodies* and several other *components* to make them reusable. If the *component* is only used once, it may also just be defined in place, as shown for the schema of the *id* parameter of the `/customers/{id}` path's GET operation.

The schemas themselves are quite intuitive to read, write and understand. *Customers* is an *array* where each item is a *Customer*. Meanwhile, *Customer* is a *object*, thus, an abstract data type, with two required properties, *id* and *name*, both of type *string* and an optional property *address* which itself is an object.

As shown by listing 5.9, an OpenAPI document specifies an HTTP API in a format which allows both humans and computers to understand the capabilities of a service without requiring access to source code or additional documentation.

There is a ecosystem of tooling around OpenAPI documents. Most prominently, *Swagger UI*, a tool which automatically generates a Web UI from OpenAPI documents for convenient visualization and interaction with the API [77]. Furthermore, there are tools for server and client code generation, for request and response validation, for document serialization and deserialization and for testing in various programming languages [76].

Several such tools are leveraged in the Security and Compliance Data Lake.

There is a rather simple reason, why this partial schemas approach to flexibility and extensibility is not sufficient for this application. The data, thus, the entities, are provided by other systems based on policies. Looking back at section 4.2.2 "Reference Integration Architecture", the component model data may be fetched from a component repository once a day. Correspondingly, the artifact composition data, thus the packages contained in an artifact, may also be fetched from scanning tools once a day. To determine whether an entity such as a component or package already exists in the database, the application has to know the identity of the respective entity. But different component models may identify components by different sets

of properties. Different scanning tools may identify packages by different sets of properties. Furthermore, different *Info Snippets* such as vulnerabilities and license are unlikely to be identified by the same set of properties. Therefore, it is also impossible to predefine even a partial schema.

The dynamically configurable schemas of this application are even implemented to be complete. So, the application does not really leverage the flexibility of using a schema-less database. This is because allowing properties that are not defined by a schema often leads to a decreased data quality as it entices to just store everything. Since properties not included in the schema are also rarely used in queries, they frequently do not add value.

Besides, the real value provided by the application is the queryability of software composition. In the rare occasions where some information may be actually missing, it is easy to manually look it up in the original data source.

Meta Data Model

Providing simple OpenAPI schemas would not be sufficient. A common OpenAPI component schema enables to define objects with a set of properties. Thereby, it may also be specified what data type these properties have and whether they are required [75]. While such a specification is enough for validation purposes, it lacks the information which properties are part of the entity type's identity.

To solve this problem, the *Types API* enforces a *Meta Data Model*, which defines schemas for the entity type schemas. Below listing 5.10 shows the YAML representation of the *Component Type* schema schema.

```
1 ComponentType:
2   type: object
3   required:
4     - kind
5     - version
6     - name
7     - identity_schema
8     - attributes_schema
9   properties:
10     kind:
11       type: string
12     version:
13       type: string
14     name:
15       type: string
16     identity_schema:
```

```
17     type: object
18     attributes_schema:
19         type: object
```

Listing 5.10: Entity Type Schema Schema

Thus, the entity type schemas have to have a *kind*, *version*, *name* and especially an *identity_schema* and an *attributes_schema* property. As the names suggest, the *identity_schema* and the *attributes_schema* of type object are themselves again schemas. The *identity_schema* thereby defines all the properties comprising the identity of the entity type and the *attributes_schema* specifies all other properties. As this is very abstract and meta level, below listing shows an exemplary JSON representation of a request body to create a *Component Version* schema.

```
1 {
2     'kind': 'ComponentType',
3     'version': 'v1alpha1',
4     'name': 'Component',
5     'identity_schema': {
6         'type': 'object',
7         'required': [
8             'name',
9             'version'
10        ],
11         'properties': {
12             'name': {
13                 'type': 'string',
14                 'maxLength': 255
15             },
16             'version': {
17                 'type': 'string'
18             }
19         }
20     },
21     'attributes_schema': {
22         'type': 'object',
23         'properties': {
24             'labels': {
25                 'type': 'array',
26                 'items': {
27                     'type': 'object',
28                     'required': [
29                         'name',
30                         'value'
31                     ]
32                 }
33             }
34         }
35     }
36 }
```

```
31     ],
32     'properties': {
33       'name': {
34         'type': 'string',
35       },
36       'value': {}
37     }
38   }
39 }
40 }
41 }
42 }
```

Listing 5.11: Request Body for Component Version Schema Creation

So the identity of this *Component Version* entity type is comprised of a *name* property of type string and a maximum length of 255 characters and *version* property also of type string with no further restriction. Both of these properties are required. The only further property in this example is *labels*, which is an array of objects with a required *name* property of type string and a *value* property of any type. Thus, this array may hold arbitrarily nested objects.

The other properties provided in this request, so *kind*, *version* and *name*, are not part of the entity types schema. These properties are used for processing logic.

kind describes the *class of entity types*.

version describes the version of the particular entity type schema. This property theoretically provides the possibility to update the entity type schema in case an additional property is required. But currently, the application does not provide any functionality to support such version upgrades by automatically updating the existing entities of that entity type.

name is the actual name of the entity type. Thereby, every node of that entity type will be labeled with this name (correspondingly, in a relational database, this would make up the name of the table). Thus, the *name* is required when querying for nodes of that entity type in the database.

Based on the data model defined in the section 5.2, the combination of *kind* and *name* would only be necessary for *Native Package Version* and *Info Snippet* where multiple entity types of the same class of entity types, or rather *kind*, may be created. In order to have uniform schemas but also especially for extensibility reasons, the application assumes all entity types in the data model may be classes of entity types. Therefore, the combination of the *kind* and *name* property is required to be included with all entity type schemas. This way, the application could be easily extended to store multiple component models.

Currently, these schemas cannot be provided for the abstract entity types, thus, *Component*, *Resource*, *Source*, *Package*, *Package Version*, but only for their respective concrete entity types *Component Version*, thus, *Resource Version*, *Source Version*, *Native Package Version* and *Info Snippets*.

This is because these are the types of entities that are provided by the component model and data sources, thus, the types of entities that are actually ingested. Their abstract entities are rather groupings of those entities, thus, aggregations. As discussed in the data model section, these entities therefore mostly store information that has to be manually created by the user such as *triage policies* or information that is merged from the concrete entities such as the common properties of *Native Package Versions* in *Package Versions*. Furthermore, especially the identities of the abstract entities depend on the respective concrete entity's identity.

For example, the identity of a *Component* is equal to the identity of a *Component Version* without considering the version property. Similarly for all other entities with a corresponding relationship. Therefore, although this unfortunately is not and cannot be represented in the OpenAPI schema of the entity types schema, the application enforces all the entity type's *identity_schemas* except the ones' of *Info Snippets* to provide a *version* property. While consuming and storing entities through the *Ingestion API*, the higher aggregation level entities are then created implicitly with the respective identity.

The *Ingestion API Schema Assembly Service* which builds the whole *Ingestion API schema* retrieves the created schemas of the concrete entity types and assembles *identity_schema* and *attributes_schema* to a single component schema for the entity type as shown in the listing below.

```
1 'Component.v1alpha1': {
2   'type': 'object',
3   'properties': {
4     'metadata': {
5       'type': 'object',
6       'required': ['kind', 'type', 'version'],
7       'properties': {
8         'kind': {
9           'type': 'string',
10          'pattern': '^ComponentType$'
11        },
12        'type': {
13          'type': 'string',
14          'pattern': '^Component$'
```

```
15      },
16      'version': {
17        'type': 'string'
18      }
19    },
20  },
21  'identity': {
22    'type': 'object',
23    'required': ['name', 'version'],
24    'properties': {
25      'name': {
26        'type': 'string',
27        'maxLength': 255
28      },
29      'version': {
30        'type': 'string'
31      }
32    }
33  },
34  'attributes': {
35    'type': 'object',
36    'properties': {
37      'labels': {
38        'type': 'array',
39        'items': {
40          'type': 'object',
41          'required': ['name', 'value'],
42          'properties': {
43            'name': {
44              'type': 'string',
45            },
46            'value': {}
47          }
48        }
49      }
50    }
51  }
52}
53}
```

Listing 5.12: Component Version Schema

So, this is the schema against which entities of type *Component Version* that shall be created through the *Ingestion API* are validated. The *metadata* properties are needed by the *Ingestion MUX* in order to route the entity to the appropriate

Data Lake Service function.

This schema creation process is almost identical for all entity types, or rather classes of entity types, except for *Native Package Version*. As common properties of different *Native Package Version*, for example of BDBA and Mend, shall be merged on the *Package Version* aggregation level, all *Native Package Versions* have to share this set of common properties. Therefore, this creation process is slightly more modular, as will be shown in the next section discussing the endpoints.

Furthermore, there are currently predefined schemas for the *Relationships* between the entity types. These schemas are all quite similar. They all require a *source*, which refers to the identity of the entity type for which the relationship is outgoing, and a *destination*, which refers to the identity of the entity type for which the relationship is incoming. Therefore, the schema for a *Component Reference*, thus, for a *Component Version to Component Version* relationship requires the identity of a *Component Version* as a *source* and as a *destination*. Correspondingly, a schema for a *Component Version to Resource Version* relationship requires the identity of a *Component Version* as a *source* and the identity of a *Resource Version* as a *destination*. This is necessary in order to create the relationships independently of the entities which makes the creation process even more flexible.

Imagine each *Component Version* schema also includes references to all related *Component Versions*. As entities, or rather nodes, have to exist before relationships including them may be created, this would require to resolve these dependencies and create a corresponding execution order. While the same may be done with the above approach, it also offers the possibility of creating all entities first and then all relationships. Of course, this comes at the cost of having to send the identities redundantly.

Apart from that, all *Relationship* schemas only have two further properties, the *labels* property, which allows arbitrarily nested objects and a *name* property of type string. As mentioned before, Neo4j can store JSON documents. But as each node and relationship only stores key-value pairs, nested properties are stored as string and cannot be leveraged in queries (except for using string comparisons, which is quite inefficient).

Endpoints

The *Types API* exposes the following HTTP endpoints:

```
1 /types/component
```

```
2 /types/resource
3
4 /types/source
5
6 /types/info-snippets
7
8 /types/packages
9 /types/package-identity
10 /types/package-attributes
11 /types/package-native-attributes
12
```

Listing 5.13: Types API Endpoints

Each endpoint accepts HTTP POST, PUT, GET and DELETE requests⁸ for managing one *kind* of schemas, the `/types/component` endpoint for *Component Version* schemas, the `/types/resource` endpoint for *Resource Version* schemas, the `/types/source` for *Source Version* schemas, the `/types/info-snippets` for *Info Snippets* schemas and the other four endpoints for *Native Package Version* schemas.

As indicated in the previous section, the creation process of the *Native Package Version* schemas is similar to the *Component Version* example, but more modular. Therefore, the *identity_schema* and the *attributes_schema* have to be created in separate requests to separate endpoints. Moreover, the *attributes_schema* is further broken down into *attributes* and *native attributes*, where *attributes* is the set of properties that may be merged in the corresponding *Package Version* entity. The individual schemas are then tied together by a request to the `/types/packages` endpoint, referencing the respective *identity*, *attributes* and *native attributes schemas*.

Finally, an open question may be, why even expose multiple endpoints? Especially considering that the *kind* is present in the request anyway, as shown in listing 5.11. Thus, the further routing could also be handled by a custom multiplexer evaluating that property.

The primary reason that such a single endpoint approach is taken for the *Ingestion API* but not for this *Types API* is that they are used completely differently.

The *Ingestion API* only has to *consume* huge amounts of data that are provided by other systems. The *Types API* is used by a human administrator to configure the application by *creating*, *reading* and in edge cases also *updating* and *deleting*

⁸The current version of the application at the time of writing the thesis actually just supports POST requests to those endpoints to create the schemas. But it is planned to extend this to also support GET request to retrieve the respective schemas and also PUT and DELETE requests to edit accidental mistakes in newly created schemas that are not used yet.

schemas. Also, this is probably done only a few times during the entire application life time. Therefore, considering the entity type schemas as REST resources and exposing corresponding endpoints allows to benefit from the intuitiveness of a REST API leveraging the common HTTP verbs.

Besides, there is another inconvenience with APIs that expose only a single endpoint which enables to submit multiple different entity types in a single request.

OpenAPI has a construct called "oneOf" which allows to specify that all entities and relationships in the request have to fulfill exactly *one of* a specified set of schemas. Thus, generally OpenAPI supports specifying that kind of endpoint. The problem is, validation functions provided by OpenAPI libraries in any programming language cannot produce helpful error messages for that kind of endpoint when trying to validate the whole request.

Imagine a request contains multiple entities. One of them is faulty. Now, the general validation function does not know what kind of entity this is supposed to be, thus, against which specific schema it has to be validated. It only knows that this faulty entity does not fulfill any of the provided schemas.

Therefore, in order to get helpful error messages, the validation would have to be implemented at a later point, evaluating the *kind* property and calling the validation function individually for each entity in the request with the corresponding schema. But this is more of an implementation detail.

5.5.2 Ingestion API & Services

The *Ingestion API* is fairly simple, once the *Types API & Services* are understood. It exposes a single HTTP endpoint:

```
1 /datalake
```

Listing 5.14: Ingestion API Endpoint

This endpoint only accepts HTTP PUT and DELETE requests⁹. The operation performed by the application when receiving a PUT request is *merging*, or in the SQL context rather known as *upsert*, of the respective entities and relationships. The request is validated against the *Ingestion API Schema* assembled by the *Ingestion API Schema Assembly Service*.

⁹The current version of the application at the time of writing the thesis actually just support PUT requests. But as in production there will likely be cases where some entity or relationship has to be removed, it is planned to extend this to also support DELETE requests

A REST API approach similar to the *Types API* would not add much value here. The API has to consume huge amounts of data provided by other systems. Multiple endpoints would result in additional unnecessary requests. A single endpoint, on the other hand, makes the implementation of the adapters simpler and entities in a request may be bundled more flexibly.

The downside of such an endpoint, as previously mentioned, is having to deal with the obscure error messages or alternatively implementing a more complex validation. But this is more of an implementation detail and the effort to implement the more complex validation is comparably low considering the rest of the application.

Furthermore, such an API also requires additional efforts, as a custom multiplexer has to be implemented to handle the routing which is commonly done by libraries and frameworks based on the URL. Popular examples therefore are *Java's Spring* [78] or *Javascript's Express* [79]. In the context of the *Ingestion API*, this routing is handled by the *Ingestion MUX* individually for each entity in the request based on its *kind* property.

5.5.3 Query API & Services

While the *Types API* and *Ingestion API* may be interesting to administrators and adapter developers, the *Query API* is the part of the application which actually provides value. Therefore, the *Query API* and respective *Services* have to enable users to conveniently perform queries and analysis by traversing relationships to answer questions such as "Which Component Versions contain the Log4J Vulnerability?".

A brief summary of the API development process may help in understanding the design decisions that led to the current *GraphQL API*.

REST API

REST APIs are currently probably the most popular types of APIs. This is due to their simplicity and intuitiveness leveraging the built in semantics of the web, the HTTP verbs, and mapping the main entity types of an application to resources and exposing them through individual endpoints.

Therefore, the initial idea was to design a single REST API for *querying* as well as *consuming* the data. Thus, at that point a separation of *Ingestion API* and *Query API* was not planned yet.

There were several different approaches for endpoint design. There was the common mapping of entity types and relationships to resources and corresponding endpoints. And there were several different deviations of the common REST approach,

trying to deal with the same *three main problems* the classic REST API is facing in this use case.

The below listings outlines the endpoints of a common mapping of the entity types and relationships to resources for the example of the component entity types:

```

1 /components
2 /components/{component_name}
3 /components/{component_name}/versions
4 /components/{component_name}/versions/{version}
5 /components/{component_name}/versions/{version}/
   component_version_references
6 /components/{component_name}/versions/{version}/
   source_version_references
7 /components/{component_name}/versions/{version}/
   resource_version_references

```

Listing 5.15: REST API Endpoints

There is a problem arising from such common REST API design regarding the *data ingestion*. Theoretically, to create *Component Versions*, the adapter would have to send a *huge number of separate requests* and route them to the correct endpoints. Thus, to merge a version v0.15.0 of the component named *github.com/gardener/etcddruid*, a PUT request would have to be sent to `/components/github.com/gardener/etcddruid/`. To merge the relationships, further PUT requests would have to be sent to the respective endpoints. This would result in a lot of separate requests. The relationships could not be created before the *Component Version* exists. Thus, even the sequence of the requests is important. Therefore, this approach is impractical.

But, considering that the data sources only provide the non abstract entities, thus, *Component Versions*, and the respective abstract entities, thus, *Components*, usually are created implicitly anyway, a deviation from the REST standard for data ingestion was deemed acceptable. Therefore, to solve this problem, all *Component Version* entities could just be sent to the `/components` endpoint with the references as nested objects.

There is another problem concerning the *queryability* of this REST API. Designing endpoints like this does enable to retrieve all information about the individual entities, but it does not provide a possibility to traverse the relationships with a single request. As the relationship itself is not really interesting for most queries, this problem could be solved by adding following endpoints:

```
1 /components/{component_name}/versions/{version}/
   component_versions
2 /components/{component_name}/versions/{version}/source_versions
3 /components/{component_name}/versions/{version}/
   resource_versions
4 /components/{component_name}/versions/{version}/package_versions
5 /components/{component_name}/versions/{version}/vulnerabilities
6 /components/{component_name}/versions/{version}/licenses
7 ...
8 /components/{component_name}/components
9 /components/{component_name}/sources
10 /components/{component_name}/resources
11 /components/{component_name}/packages
12 /components/{component_name}/vulnerabilities
13 /components/{component_name}/licenses
14 ...
```

Listing 5.16: REST API Query Endpoints

So, the REST API would also enable to traverse the graph and intuitively answer the most common questions. Through URL parameters, further filtering of the corresponding result set would also be possible.

This deviation of a REST API seemed simple and intuitive and still powerful. But there is a third problem, *parsing*. In the end, this also led to the final reason for not pursuing this approach further.

Part of the problem is already visible in the data ingestion example. Common routing frameworks and libraries are not able to handle URLs such as `/components/github.com/gardener/etcdruid/` since the slashes in the name would be interpreted as a path separator.

Intuitively, this could be solved by replacing the slash by a different character such as an underscore. But for this to work properly, the respective character cannot be allowed to be used in the name. This would definitely cause compatibility issues with OCM and other component models.

Another approach is to implement a custom parser. But in order for this to work, the URLs have to be adjusted to always end on a keyword. For example, the previous URL to retrieve a component with a specific name would have to be changed to something like `/components/github.com/gardener/etcdruid/list` with *list* being the keyword. Then a custom parser could parse the URL backwards.

But even with URLs always ending on keywords, the parsing may not be unambiguous, as shown by the listing below.

```
1 /components/github.com/gardener/etcdruid/versions/1.2.3/list
```

Listing 5.17: Ambiguous REST API URL

This request could be meant to list all *Component Versions* with *component_name* = *github.com/gardener/etcdruid*. But it could also be meant to list all *Components* with *component_name* = *github.com/gardener/etcdruid/versions/1.2.3*.

The problems become even more difficult to deal with when considering corresponding endpoints for *Sources* and *Resources*. Using the identity properties of the OCM, respective endpoints would look like this:

```
1 /sources/{component_name:component_version:source_name:extra_id}
```

Listing 5.18: REST API Source Endpoint

While this is not really clean, it still seems as this may work as long as the character separating the individual properties of the sources' identity, in this example the colon, is guaranteed to not occur within one of the properties. But as previously established, such a separator would lead to compatibility issues. Besides, per OCM specification, the *extra_id* is a flat map containing an arbitrary number of key-value pairs. Thus, the key-value pairs would have to be included in the URL.

Therefore, it is apparent that there is no clean solution to identifying entities by integrating their dynamically configurable identities into the URL.

To solve this parsing issues, assigning additional *surrogate keys*, thus *UUIDs* or *auto incremented IDs*, as it is done in common REST API use cases was considered. But this would require additional requests to find the surrogate key of the respective entity. Thus, the API would have to support requests using query parameters as indicated by the listing below.

```
1 /sources?component_name=github.com/gardener/etcdruid&version
      =0.15.0&source_name=etcd&extra_identifier=etcd-druid-source
2
3 /sources/{id}/versions
```

Listing 5.19: Querying for Entity ID

The user would first send the first request to retrieve the ID of the respective *Source*. Once the ID is acquired, it might be used to navigate further through the API and, for example, retrieve all versions of the *Source* with that ID.

Although every API that exposes its operations through individual resource endpoints is called REST API today [80], this form of usage is already much closer to the original REST architectural style defined by Roy Fielding. It is derived from the web and promotes hypermedia as the engine of application's state [81]. Thus, the original idea is that RESTful architectures and consequently also REST APIs may be navigated just as websites. So a user may know the host address of the API. Sending a GET request to that host address will then return a list of hyperlinks of further resources. For this API this may be `/components`, `/sources`, Again, following one of those hyperlinks by sending a GET request will return another list of hyperlinks of further resources available under the currently selected resource, in the case of `/components` a list of hyperlinks to all available individual *Components*. This is great for APIs that have a potentially huge number of unknown users, as it is very intuitive and requires almost no documentation to navigate through, just as a website. But it is not suitable and creates a lot of overhead for this kind of application which requires running specific potentially expensive queries based on already known entities.

Finally, hashing the identity properties and using the hash as ID was considered. While this would allow to issue the queries directly without having to retrieve the ID first, it would require to calculate the hash. Due to implications regarding normalization, this would require an additional small program on client side in order to be able to use the API. Therefore, this approach also was not pursued further.

Remote Procedure Calls

As the REST API approach was deemed unsuitable after thorough analysis, the next best idea was to implement the most common queries as *Remote Procedure Calls (RPC)*.

This approach is simple and clean. The application specifies a set of functions, or rather procedures, that may be called through the API. Furthermore, it defines which properties have to be provided for each procedure and the format in which these properties have to be provided. In the context of web services, those procedures are typically called by HTTP POST requests, specifying the procedure to be called as well as the corresponding properties in the corresponding format in the request body. Thus, there would be no issues regarding embedding properties in the URL and parsing.

While implementing the most common queries is convenient, it is also quite inflexible. An intuitive, more flexible approach which gives users the full capability

to run precise and complex queries through the API and to answer questions such as "Which Package Versions depend on the spring-context:1.2.12 Package Version?" is indicated by the listing below.

```

1 POST /query
2
3 processCypher(
4     'MATCH (p1:PACKAGE_VERSION {Name: "spring-context", Version
5         : "1.2.12"})'
6     MATCH (p2:PACKAGE_VERSION)
7     WHERE (p1)-[:DEPENDS_ON*1..]-(p2)
8     RETURN (p2)'
9 )

```

Listing 5.20: Exposing Cypher through API

Thus, exposing the entire database query language through the API. But compared to the REST API, where answering such a question did not require any knowledge about graph databases, for this approach, the user has to know the entire database specific query language. And even then, the queries may be quite error prone. Besides, there are several security implications which make this approach impractical.

The basic idea of exposing a query language through the API is still powerful. But the set of available operation has to be smaller and the permissions stricter than the ones' of the database's query language. Thus, in the past there were a few attempts to design such API query languages - for example SPARQL, FIQL and a lesser-known one by Meta, FQL [82]. While these API query languages may be unsuitable for graph traversal, designing a language tailored to this application would be an option. But this would be quite time intensive and complex. Besides, the general problems with such API query languages are similar to exposing the database query language. Users have to be skilled in those languages in order to being able to use the API properly.

Thus, these RPC approaches were not satisfying either.

GraphQL

On a basic level, a *GraphQL API* is also essentially RPC, as *GraphQL* is also a query language for APIs [83] designed by Meta. Thereby, a GraphQL API exposes a single endpoint which is typically called `/graphql`.

But GraphQL has a special approach to specifying APIs and defining queries, which makes it very intuitive and easy to use. A GraphQL API is defined by *types*

which correspond to resources in a REST API, and *fields* on this *types*. To specify those *types* and *fields*, GraphQL provides its own *GraphQL schema language*. Listing 5.21 shows an exemplary GraphQL schema language representation of a *Component Version* based on the OCM data model instance in section 5.2.3 "Application of the Data Model".

```
1 type Query {
2   Component(Identity: ComponentIdentityArgs): [Component]
3   ComponentVersion(Identity: ComponentVersionIdentityArgs,
4     Attributes: ComponentVersionAttributesArgs): [
5     ComponentVersion]
6 ...
7 }
8
9 type ComponentVersion {
10   Identity: ComponentVersionIdentity
11   Attributes: ComponentVersionAttributes
12   IncomingReferences(type: String, name: String, labels: String)
13     : [Reference]
14   OutgoingReferences(type: String, name: String, labels: String)
15     : [Reference]
16   ReferencedComponents(Identity: ComponentIdentityArgs): [
17     Component]
18   ReferencedComponentVersions(Identity:
19     ComponentVersionIdentityArgs, Attributes:
20     ComponentVersionAttributesArgs): [ComponentVersion]
21   ReferencedResources(Identity: ResourceIdentityArgs): [Resource]
22
23   ReferencedResourceVersions(Identity:
24     ResourceVersionIdentityArgs, Attributes:
25     ResourceVersionAttributesArgs): [ResourceVersion]
26   ReferencedSources(Identity: SourceIdentityArgs): [Source]
27   ReferencedSourceVersions(Identity: SourceVersionIdentityArgs,
28     Attributes: SourceVersionAttributesArgs): [SourceVersions]
29   ContainedPackages(Identity: PackageIdentityArgs, Attributes:
30     PackageAttributesArgs): [Package]
31   ContainedPackageVersions(Identity: PackageVersionIdentityArgs,
32     Attributes: PackageVersionAttributesArgs): [PackageVersion]
33   ContainedBDBAPackageVersions(Identity: BDPAIdentityArgs,
34     Attributes: BDPAAttributesArgs): [BDBAPackageVersion]
35   ContainedVulnerabilities(Identity: VulnerabilityIdentityArgs,
36     Attributes: VulnerabilityAttributesArgs): [Vulnerability]
37 }
38
39 type ComponentVersionIdentity {
40   name: String
```

```
24     version: String
25 }
26 type ComponentVersionAttributes {
27     labels: String
28 }
29 ...
30
31 input ComponentVersionIdentityArgs {
32     name: String
33     version: String
34 }
35 input ComponentVersionAttributesArgs {
36     labels: String
37 }
38 ...
```

Listing 5.21: GraphQL Component Version Schema

The GraphQL schema language is quite intuitive. The *Query type* is a special type which defines the entry point for GraphQL queries, listing all *types* a user may query for. These available *types* further specify their fields and therefore declare what may be queried on them, and so on. Thereby, the GraphQL schema resembles a REST API which properly implements hypermedia as the engine of application state. But instead of having to send multiple queries to navigate through the API step by step and actually retrieve data, with GraphQL the entire schema may be retrieved at once. This reduces the overhead in requests and data transfer significantly.

The *Component Version object type* is a field within the *Query type* and has two *input objects* as arguments, the *ComponentVersionIdentityArgs* object and the *ComponentVersionAttributesArgs* object, defined at the bottom of the listing. *[ComponentVersion]* thereby represents an array of *Component Version* objects.

The basic structure of the previously described REST API and of the GraphQL API defined by this schema is quite similar. Both hide the relationship traversals necessary to answer questions such as "Which Vulnerabilities are contained in Component Version "github.com/gardener/etcd-druid:0.15.0"?". Instead, for the user of the API, it looks and feels like *Component Versions* are stored like documents in a document database with all related entities directly nested inside of them. Thus, a GraphQL query to answer the respective question based on the above GraphQL schema looks like this:

```
1 {
2     ComponentVersion(Identity: {name:'github.com/gardener/etcd-
```

```
druid', version: '0.15.0'}) {
  ContainedVulnerabilities {
    Identity {
      cve
    }
    Attributes {
      summary
      cvss
    }
  }
}
```

Listing 5.22: GraphQL Query

In reality, the application assigns functions to the respective fields on a GraphQL schema. These functions generate a corresponding *cypher* query to retrieve the data from the database and resolve the request.

As shown in listing 5.22, the query language even allows to specify which properties of the *Vulnerabilities* shall be retrieved. While this is not really significant for this application, compared to a REST API this may reduce the overhead even further by exactly describing what shall be queried and retrieved.

More importantly, the GraphQL query language gets rid of parsing issues as the data model and queries do not have to be mapped to URL paths. Additionally, it is also more powerful than the REST API, as it supports arbitrary nesting of the queried objects. Thus, considering above query, it is possible to specify to retrieve all packages the respective *Vulnerability* occurs in. Depending on the exact GraphQL schema of the *PackageVersion* type the respective query would look something like this:

```
{
  ComponentVersion(Identity: {name:'github.com/gardener/etcdruid', version: '0.15.0'}) {
    ContainedVulnerabilities {
      Identity {
        cve
      }
      Attributes {
        summary
        cvss
      }
    }
  }
}
```

```
11     PackageVersions{  
12         Attributes {  
13             name  
14         }  
15     }  
16 }  
17 }  
18 }
```

Listing 5.23: Nested GraphQL Query

Thus, as GraphQL provides a clean, intuitive and powerful API, this is the approach that was finally selected for the Query API. GraphQL also has capabilities to manipulate data, but for this application, these are also less suitable than the *Ingestion API*.

As already mentioned, the application cannot define the GraphQL schema statically. Instead, the schema has to be created and adjusted dynamically by the application based on the OpenAPI schemas created through the *TypesAPI*. This is the biggest downside of this approach, as this is cumbersome to implement and adds a lot of complexity to the application.

5.5.4 Limitations and Problems

The current application design as presented in the previous sections is not flawless and has limitations.

The main limitation of the current design is the *lack of explicit schemas for higher aggregation level entity types* and the *implicit creation of the respective entities*. Generally, this is rather convenient for administrators, adapter developers as well as users of the API, as the ability to query the software composition on multiple aggregation levels comes seemingly for free. The administrator does not have to configure respective schemas for *Component*, *Source*, *Resource*, *Package*, *Package Version* and *Relationships*. Instead, only the schemas of entity types he is actually dealing with have to be configured, as the respective entities are provided by the component model and other data sources. And the adapter developer does not have to deduct these entities and send them to the API as their creation is handled implicitly by the application itself. So, for most use cases, this *hides complexity* and may even be seen as an *advantage*.

But answering questions about the software composition is not the sole purpose of the application. Looking back the requirements in sections 5.1, the application shall also enable assessments (R.6). The data model section 5.2 therefore mentioned

that *triage policies* may be defined on *Sources* and *Resources* and the actual *triage* information shall be stored on the relationships such as between *Component Versions* and *Source Versions* or *Resource Versions*. Furthermore, these higher aggregate level entities shall be able to store semantic information about the grouping. These capabilities are not enabled by the current application design.

There are essentially two approaches to adjusting the application to support this functionality. One approach would be to explicitly deduct the *identity_schemas* for higher level entity types such as *Component*, thus, the *Component Versions* *identity_schema* without the *version* property, and expose additional *Types API* endpoints to add *attribute_schemas*. Possibly, to provide even more flexibility, the *identity_schema* of the higher level entity types could even be configured themselves as any subset of the lower level entity type instead of being predefined as the set of the lower entity types identity properties excluding *version*. Correspondingly, *Types API* endpoints may be added to define the schemas of *Relationships*. This approach matches the general extensibility and flexibility of the application, but also requires the most efforts to implement and adds a substantial amount of complexity for the administrators.

The other option is to statically predefine the properties available for *triage policies* and *triaging*. To add general semantic information, a *labels* property which may store arbitrary key-value pairs may be added. This approach is quite low effort and does not add any complexity, but it is significantly more rigid. Also, the *labels* property does not enable proper filtering upon the nested properties.

As the *Ingestion API* is tailored towards handling huge amounts of data provided by other systems, an additional *Triage API* may be added to manually store this information.

Furthermore, the flexibility and extensibility of the application comes at the cost of *complexity*. This is partially represented in the previous sections, especially the *Types API & Services* with its *Meta Data Model*. But it is also especially represented in the application's code. The code dealing with schema creation and translation, thus, the code making up the *Ingestion API Schema Assembly Service* and *GraphQL Service*, is quite lengthy. As the data types are not predefined, the whole application code has to work extensively with abstract data types such as maps¹⁰, while leveraging reflection and casting. Furthermore, the whole processing logic had to be designed to be able to handle this level of flexibility. Even data

¹⁰The name of this abstract data type varies between programming languages. Other common terms are associative array or dictionary.

base queries are constructed dynamically based on the respective configured schemas. This makes *application maintenance* more difficult and error prone. Moreover, the implications of this amount of flexibility and extensibility for the *application's security*, especially also concerning attacks like "SQL injection", have not been sufficiently analyzed and estimated so far.

Chapter 6

Result

After the technical and architectural design decisions are discussed in detail, this chapter finally presents the results. This includes a presentation of the actual functionality of the *Security and Compliance Data Lake* and a evaluation of the design decisions, the implementation and thereby, the actual value of this research based practical contribution to increase the trust and transparency placed in our digital infrastructure. To conclude this work, the research question is revisited, summarizing the thesis and highlighting the greatest challenges throughout the application design and development.

6.1 Presentation

This presentation section briefly introduces the infrastructure, the data model configuration and the data set of the application instance used to demonstrate the functionality, before showing queries answering interesting questions about software composition through the GraphQL API. After all, the *Security and Compliance Data Lake* is a backend application providing its functionality through a standard interface, a GraphQL API, and returning its results in a machine readable format, JSON, which limits the possibilities for impressive visual showcasing.

6.1.1 Application Infrastructure

Originally, the application was developed and tested on a *Windows Subsystem for Linux 2 (WSL2)*. Thus, the database as well as the application were running there.

After the first PoC of the application was finished, SAP Gardener was used to set up a Kubernetes cluster to host presentable versions of the application. Therefore, a simple development pipeline was configured with *Github Actions* to automatically containerize the application and push the respective images into a *DockerHub* repository. Furthermore, *Kubernetes deployment files* were created, to deploy the application itself, as well as other Kubernetes resources such as an *ingress*

to expose the API. To deploy the Neo4j database, the corresponding *helm chart* provided by Neo4j was leveraged. Additionally, in order to be able to perform all operations to deploy and undeploy all Kubernetes resources at once, small bash scripts were set up.

In this Kubernetes cluster, all resources are deployed on a single *worker node* with 4 CPUs and 16Gi of RAM. Furthermore, the cluster has 50Gi of persistent memory.

6.1.2 Application Configuration & Data

The *data model* configured through the *Types API* for this application instance which is used for presentation purposes corresponds to SAP Gardeners *Open Component Model* and the usage of *Black Duck Binary Analysis* as scanning tool and sole *data source*.

Thus, the *data* created through the *Consumption API* is provided by a *Component Descriptor Repository* and *Black Duck Binary Analysis*. The presentation application instance contains 22.204 *nodes* and 97.335 *relationships* between those *nodes*. These nodes and relationships are effectively a version of the SAP Gardener, thus, a *Component Version* and its complete decomposition. Figure 6.1 shows a visual representation of a part of the graph¹ rendered by Neo4j to showcase the database's functionality and to give a better feeling for the data.

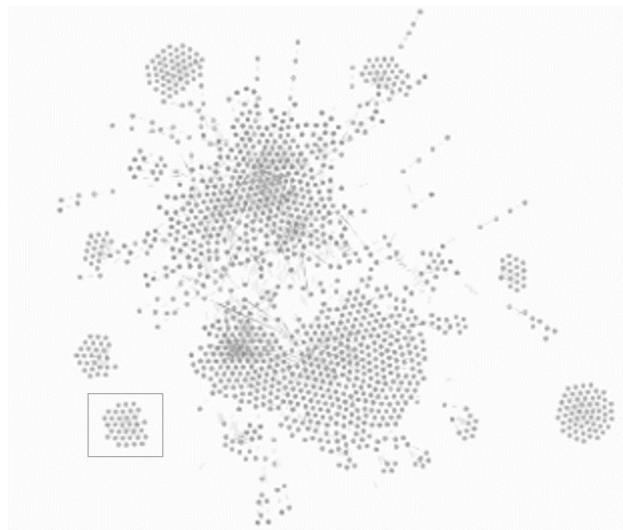


Figure 6.1: Neo4j Graph Representation
Source: Own Representation

¹The Neo4j web interface becomes unresponsive and aborts the request if the amount of nodes and relationships to be rendered becomes too high. Therefore, it is not possible to show a visual representation of the entire graph here.

6.1. Presentation

The lines connecting different nodes and thereby representing the relationships are barely visible and often not explicitly rendered, but Neo4j visualizes nodes that are highly interrelated close together leading to these visible clusters. The main clusters in the middle revolve around *Component Versions* and contain all kinds of nodes. The larger clusters on the side are *Package Versions* with an exceptionally high number of *Vulnerabilities*. For example, the highlighted cluster on the bottom left revolves around a version of *HashiCorp's Vault*, a software package for managing secrets and protecting sensitive data [84].

Although this is difficult to arrange clearly, to give an even better understanding of how the data model manifests itself in form of nodes and relationships, figure 6.2 zooms in on a specific component.

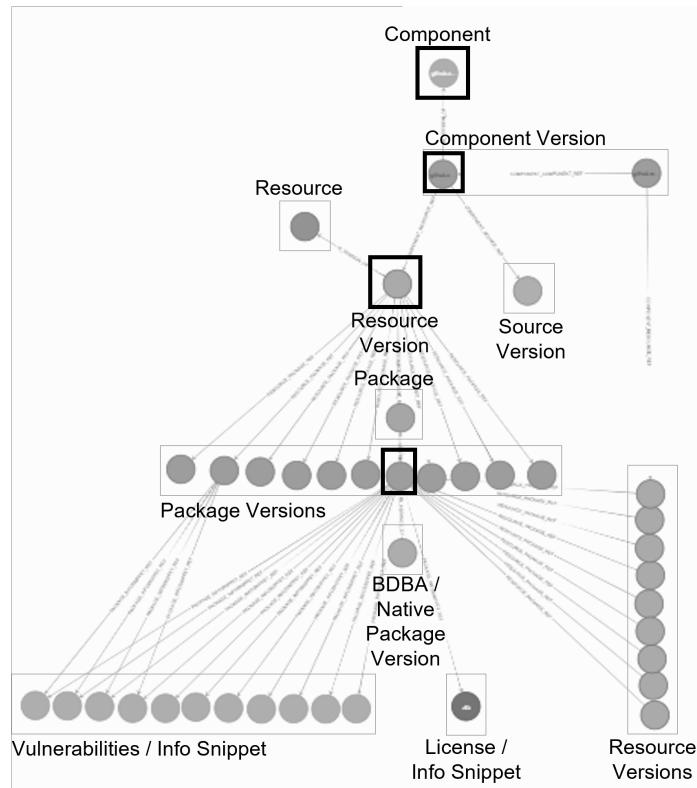


Figure 6.2: Neo4j Component Tree Representation
Source: Own Representation

The figure shows a part of the tree formed by the relationships of the *Component* at the top.² To be more precise, the nodes that are highlighted by a thicker black

²This visualization rendered by Neo4j also shows the names of relationships along the lines and the names of some nodes within the circle. These are not intended to be readable in this figure. They are not important for the explanations and there is no way to zoom in close enough to keep them readable while still showing the whole tree.

border are the ones that are unfolded. Thus, for these nodes all related nodes are visualized.

The *Component* at the top has only one relationship to a *Component Version*. Thus, there is only one version of this *Component* stored in the database. As the highlighted *Component Version* straight under the *Component* is unfolded, the figure shows that it references one other *Component Version*. Moreover, it also references a *Resource Version* and a *Source Version*. Since the *Source Version* is not unfolded, its related nodes are not shown. If it was unfolded, it would most likely open up a subtree similar to the one of the highlighted, and therefore unfolded, *Resource Version*. The same is true for the *Resource*. If it was unfolded, it would possibly show several more *Resource Versions*. Besides the *Resource*, the unfolded *Resource Version* is especially related to several *Package Versions*. The unfolded *Package Version* is related to its other two aggregate levels, *Package*, shown directly above itself and below the *Resource Version*, and *BDBA*, or rather, *Native Package Version*. The relationships to the *Resource Versions* on the right show that this particular *Package Version* occurs in several more *Resource Versions*. The topmost *Resource Version* is even referenced by the *Component Version* which is referenced by the unfolded *Component Version*. Thus, there are two different paths through which this particular *Package Version* is contained in the *Component*. Furthermore, the *Package Version* is related to, or rather has, a number of *Vulnerabilities* and a *License* which are both *Info Snippets*. The *Vulnerabilities* on the left even occur within another *Package Version*.

So, although this has been theoretically discussed before, this conveniently visualizes the paths that have to be traversed to answer questions related to software composition.

6.1.3 Demonstration

The *data model configuration* is not really interesting, as it is just sending several schemas to the respective endpoints. There is not much more to show than what was already presented in section 5.5.1 "Types API & Services". Moreover, these schemas are SAP Gardener OCM specific and are even provisional for the PoC. For the productive use of the application, these will have to be adjusted in close cooperation with the corresponding stake holders within the SAP Gardener team.

The *data upload* is handled by an *Adapter* written in *Python*. This *Adapter* is tailored for the PoC use case. Therefore, the upload may be triggered manually by specifying a specific *Component Version*. Consequently, all the *References* to other *Component Versions*, to *Resource Versions* and *Source Versions* are resolved

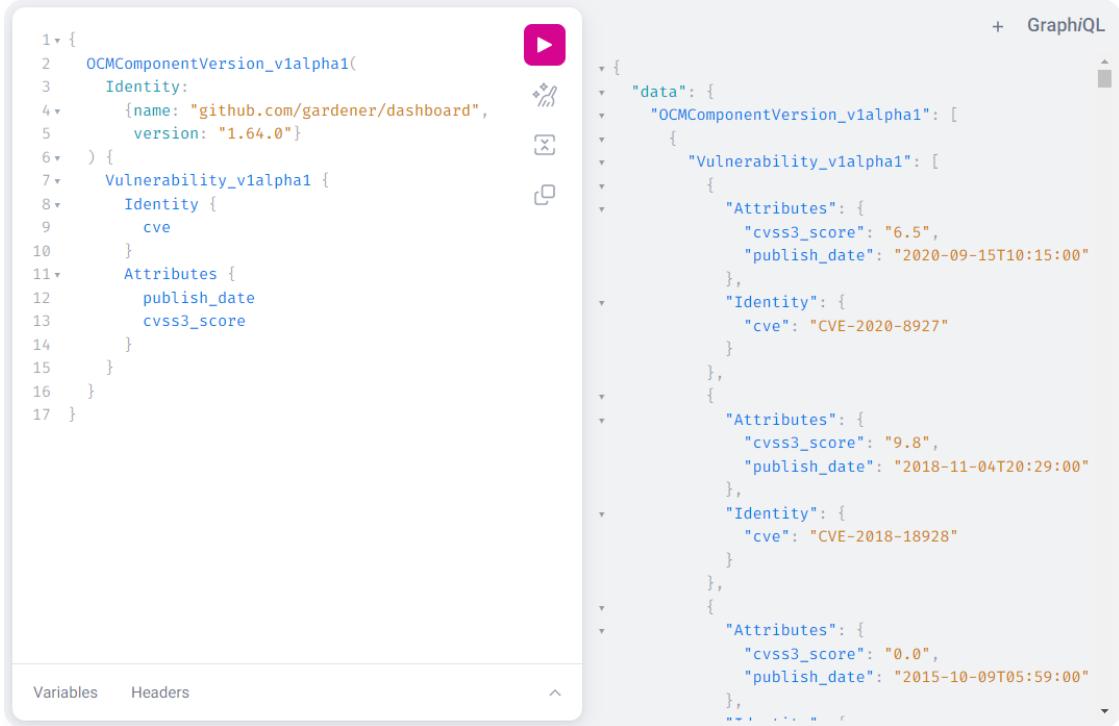
recursively and uploaded to the *Security and Compliance Data Lake* in the correct order, so that the application does not attempt to create a relationship to an entity that might not exist yet. Furthermore, the *Adapter* retrieves the scan results for respective *Resource Versions* from the BDBA API and uploads them as well. The actual scans are triggered independently. This way, the results may be uploaded without performing a scan each time which is especially useful during development.

Finally, the *data analysis* may be performed through the GraphQL API. To demonstrate the functionalities and capabilities of the *Security and Compliance Data Lake*, the following paragraphs show how it may be used to answer questions about software composition. The application is currently not connected to a front end. Therefore, provisionally, the only way to retrieve specific data is to actually type the respective request body. To make this more comfortable, the *Security and Compliance Data Lake* also hosts *GraphiQL*, a Web IDE for GraphQL, offering syntax highlighting, auto-completion and automatic documentation [85]. This Web IDE is also used to run the following queries.

Throughout the thesis, the two most common questions are "*Which Vulnerabilities are contained in a specific Component Version?*" and the reciprocal "*Which Component Version contains a specific Vulnerability?*".

So, to address the first question, a *Component Version* has to be picked. For this presentation, the github.com/gardener/dashboard version *1.64.0* is used. This *Component* represents an app to manage Kubernetes clusters deployed with SAP Gardener.

6.1. Presentation



The screenshot shows the GraphQL Web IDE interface. On the left, a code editor displays a GraphQL query:

```
1  {
2    OCMComponentVersion_v1alpha1(
3      Identity: {
4        name: "github.com/gardener/dashboard",
5        version: "1.64.0"
6      }
7    ) {
8      Vulnerability_v1alpha1 {
9        Identity {
10          cve
11        }
12        Attributes {
13          publish_date
14          cvss3_score
15        }
16      }
17    }
}
```

On the right, the results of the query are shown as a JSON document:

```
{
  "data": {
    "OCMComponentVersion_v1alpha1": [
      {
        "Vulnerability_v1alpha1": [
          {
            "Attributes": {
              "cvss3_score": "6.5",
              "publish_date": "2020-09-15T10:15:00"
            },
            "Identity": {
              "cve": "CVE-2020-8927"
            }
          },
          {
            "Attributes": {
              "cvss3_score": "9.8",
              "publish_date": "2018-11-04T20:29:00"
            },
            "Identity": {
              "cve": "CVE-2018-18928"
            }
          },
          {
            "Attributes": {
              "cvss3_score": "0.0",
              "publish_date": "2015-10-09T05:59:00"
            }
          }
        ]
      }
    ]
  }
}
```

Figure 6.3: Vulnerabilities in Component
Source: Own Representation

The screenshot from the GraphQL Web IDE shows the GraphQL query on the left and the corresponding response on the right. The query specifies to return the *Identity* property *cve* and the *Attribute* properties *publish date* and *cvss3_score* for each *Vulnerability* in the *Component Version*. The right side shows that the query returns a JSON document containing an array of *Vulnerabilities*. Thereby, the scrollbar indicates that the array contains several more *Vulnerabilities* than visible in the screenshot. As briefly mentioned in section 5.5.1 "Meta Data Model", the entity type schemas themselves may be versioned. Since schema names in GraphQL have to be unique, a suffix, in above figure *v1alpha1*, accounts for this schema versioning.

To demonstrate how to answer the reciprocal query, the upper most *Vulnerability* with the *cve* of *CVE-2020-8927* in above figure is used.

6.1. Presentation

The screenshot shows a GraphQL interface with a query editor on the left and a results panel on the right. The query is:

```
1 ▶ {  
2   Vulnerability_v1alpha1(  
3     Identity:  
4       {cve: "CVE-2020-8927"}) {  
5     OCMComponentVersion_v1alpha1  
6       Identity {  
7         name  
8         version  
9       }  
10    }  
11  }  
12 }
```

The results panel displays a list of component identities, each associated with a specific GitHub repository and version. A scrollbar is visible on the right side of the results panel.

Identity	Name	Version
1	github.com/gardener/dashboard	1.64.0
2	github.com/gardener/gardener-extension-networking-cilium	v1.18.0
3	github.com/gardener/gardener-extension-networking-calico	v1.27.1
4	github.com/gardener/vpn	0.20.0
5	github.com/gardener/ingress-default-backend	0.1.0

Figure 6.4: Components with Vulnerability
Source: Own Representation

The result in figure 6.4 is scrolled down a bit, as indicated by the scrollbar, to show that the array of *Component Versions* naturally also contains the *Component Version* where the *Vulnerability* was found in, in the example before. The popular question "*Which Deployments contain the Log4j vulnerability?*" could be answered correspondingly. But as SAP Gardener does not contain any Java coding and therefore also no Log4j, this cannot be explicitly shown.

To further stress the capabilities of the GraphQL API, the following query shows how to request all *Package Versions* contained in the *github.com/gardener/dashboard Component Version* and return the respective *Licenses* at the same time.

The screenshot shows a GraphQL interface with a query editor on the left and a results viewer on the right. The query is:

```

1+ {
2   OCMComponentVersion_v1alpha1(
3     Identity: {
4       name: "github.com/gardener/dashboard",
5       version: "1.64.0"
6     }) {
7     PackageVersion {
8       Attributes {
9         name
10        version
11      }
12      License_v1alpha1 {
13        Identity {
14          name
15        }
16        Attributes {
17          type
18          url
19        }
20      }
21    }
22 }
  
```

The results pane displays the JSON output of the query. It includes a 'data' field containing an array of 'OCMComponentVersion_v1alpha1' objects. Each object has a 'PackageVersion' field, which contains an array of 'Attributes' objects. One attribute is 'zlib' with version '1.2.11'. Another attribute is 'brotli' with an empty version string.

Figure 6.5: Packages in Component

Source: Own Representation

Figure 6.5 shows that the result is a JSON document containing an array of *Package Versions* and again, nested in each *Package Version*, an array of *Licenses*. Thus, the queries may even be nested.

There is a lot more functionality that could be showcased. But as the above examples are sufficient to convey a basic understanding of the functionality and each further example takes up a lot of space, the following section transitions to an evaluation.

6.2 Evaluation

To provide an evaluation of the PoC implementation of the *Security and Compliance Data Lake* developed during this scientific work, the currently available functionality is compared to the requirements specified in section 5.1. Thereby, this section provides a holistic view, summing up how each requirement is fulfilled by the application or what has to be done in order to satisfy a currently unfulfilled requirement.

Requirement R.1 specifies that the application shall be able to consume and store metadata from multiple different data sources. This was a crucial aspect during data model and application design as also stressed in the corresponding chapters.

The data model introduced *classes of entity types* to express this extensibility and the application architecture exposes the *Types API* to dynamically create and configure the actual data model tailored to the respective use case. As the application thereby provides the possibility to not only configure different data sources for package information such as BDBA and Mend, but to even configure the *component model* and theoretically supports the option to use multiple *component models* in parallel, requirement R.1 is exceeded.

Requirement R.2 is to store the metadata from different data sources without aggregation. The data model splits the package type into three aggregation levels, *Package*, *Package Version* and *Native Package Version*. On the *Native Package Version* aggregation level, the metadata for a package may be stored as provided by the respective data source without the necessity to merge, or rather aggregate, the data. Thus, requirement R.2 is fulfilled.

The *Package Version* aggregation level is intended to store the aggregated data. As described in the meta data model, the *Native Package Version* schema therefore has to specify the attributes specific to the data source, the *native attributes* and the attributes that shall be aggregated, the *attributes*. Although the current implementation creates this aggregate level implicitly, the merging has to be done manually. The API currently does not expose an operation to do this. Therefore, while the application architecture and design are prepared to fulfill this functionality, it currently is not available and is yet to be added. So requirement R.3 is not completely fulfilled. As the SAP Gardener team for whom the application was primarily developed, currently only actively uses a single data source for this kind of information, BDBA, the actual implementation of the requirement had a lower priority for the PoC.

Requirement R.4, providing an aggregation level to group sources and resources, is implemented through the component entity types. Thereby, the application bridges the gap between artifact metadata and deployment information.

Requirement R.5, enabling to query data on different levels of aggregation, is probably the most visible requirement and is based on several of the previously mentioned requirements. As discussed in section 5.5.3 "Query API & Services" and also shown in the previous presentations, this is enabled by the GraphQL API. Thus, questions such as which vulnerabilities are contained in a specific resource or in a specific deployment may be answered effortlessly. The underlying complex graph traversals necessary to query this information are thereby completely transparent to the user.

As already discussed in section 5.5.4 "Limitations and Problems" of the applica-

tion's architecture, *Requirement R.6* currently is not fulfilled by the PoC. Either, an extension of the *Types API* and respective services or a static extension of the schemas will have to be added to enable assessments.

Requirements R.7, to provide aggregation and filter functions, and *Requirement R.8*, to enable querying the metadata in common SBOM formats, are not implemented either. But based on the application's architecture, these functionalities may easily be added.

Regarding *non-functional requirements* no specific boundaries were defined. Still, it is mentioned several times that performance, especially of read operations, is relevant. Neither the amount of data in the presentation application instance, nor the compute resources of the cluster correspond exactly to the conditions of a productive application instance.

Anyway, as a point of orientation, the queries performed in section 6.1 usually return in approximately 1-5 seconds, also depending on the internet connection and the amount of data requested.

Theoretically, especially the performance of "top-down queries", thus, queries answering questions like which *Vulnerabilities* are contained in a *Component Version*, should stay relatively constant even with far more nodes in the graph over all. "Top-down" as *Component Versions* build up a tree that is self-contained and aside of some discovered vulnerabilities (and possibly other *Info Snippets*) usually stays relatively constant over time. Thus, corresponding to the explanations in section 5.3.4 "Theoretical Foundations" of the graph database, the look up of the initial node is the only part of the query which may decrease significantly in performance. An important side note at this point, this is only true for "top down" queries starting on *Component Version* level and below, as the number of *Component Versions*, and thus, the number of entire trees under a *Component* will naturally grow significantly over time.

Generally, a similar concept applies to most of the "bottom up queries". So, answering questions like which *Component Versions* contain a specific *Package Version* naturally spans several of the *Component Version* trees. But as new *Component Versions* will generally switch to newer *Package Versions*, this growth is also limited in most cases.

Thus, theoretically, the query performance should scale quite good, as most of them are likely limited to self-contained subgraphs.

To summarize this, the current PoC of the *Security and Compliance Data*

Lake implements almost all of the *priority 1* requirements R.1 - R.5. Thereby, the application is already able to perform the most relevant queries and analysis about software composition and answer important questions such as "*Which deployments contain the Log4j vulnerability?*" within seconds. As shown by the system design section, the decisions which lead the application's design and the application's architecture are based on thorough research. Thus, the application definitely poses a scientific and practical contribution to increase the transparency of our digital infrastructure.

Furthermore, the reasons for not fulfilling the other requirements are not actual hard limitations of the application's architecture, but merely the lack of time. So, the currently missing functionality may be added in future.

6.3 Research Question and Retrospective

As common in software development, or rather in science in general, the complexity of a problem only really unfolds, once one actually starts trying to solve it. Correspondingly, the complexity to design and develop a sustainable solution for such a central software metadata store became much bigger than originally anticipated. This section sums up the thesis, highlighting the greatest challenges and concluding with a retrospective onto implementation and decisions.

First of all, the *requirement elicitation* was an important part of this work. Due to the versatile ecosystem concerning software supply chain security and software metadata, this task proved challenging. Therefore, *chapters 2 and 3* not only serve as foundation for the reader, but serve as valuable preparation for the entire work. *Chapter 2* provides an overview and basic understanding of the subject area and its existing standards and technologies. Thereon, *chapter 3* introduces the current state of the art approach, pointing out its limitations, such as the coupling of CI and CD, the distribution of the software metadata over the software development life cycle and the resulting informational gap between artifacts and deployments.

Chapter 4 is focused on developing a holistic solution architecture. Therefore, it introduces SAP Gardener's Open Component Model. Based on its capabilities to aggregate sets of artifacts under a technology-agnostic identification and access scheme which allows to decouple compliance scans from the CI/CD pipeline, the chapter develops an abstract design for the *Security and Compliance Data Lake*. It then uses this abstract design to derive a reference architecture for modern develop-

ment and deployment landscapes. Through utilizing the *Open Component Model* in combination with the *Data Lake*, this architecture proposes a practical solution for modern development and deployment landscapes to overcome the limitations of the current state of the art approach. As the Open Component Model and its tooling is itself relatively new, the documentation still lacks depth, explanations and extensive examples. Therefore, collecting the respective information required a lot of interviews and discussions with the responsible developers which were particularly difficult due to the abstract nature and the ambiguous terms of the Open Component Model. Consequently, this supposedly simple task of *familiarizing with and explaining of existing concepts and tools*, turned out to be a real challenge.

Chapter 5 concludes the efforts of the previous chapters, leveraging the preparations to specify requirements and to design an entire software system. Thereby, the greatest challenges were the *inconsistency of identities* and the *volatility, or rather unpredictability, of the data model*.

The problem around artifact identities, thus, the ambiguity of artifact references regarding the actual technical artifacts, has been discussed in detail in the context of component models. But even more generally, in order to be able to reliably identify data delivered from different data sources as data about the same entity, this entity needs a *globally unique identifier*. Although there are efforts to establish standards for identifying different entities, such as *purl* for packages or the *SPDX License Identifier* for licenses, these are rarely used by scanning tools at all or at least not consistently.

Considering the unpredictability of the data model, extending the entity relationship model with classes of entity types added some complexity to the explanations but was not particularly difficult itself, but dealing with such a generic data model during system design and implementation was. Thus, it was hard to come up with an architecture which enables enforcing the data model while maintaining the extensibility and flexibility to add new entity types and even configure the properties of entity types dynamically at runtime. Thereby, the greatest challenges include coming up with the *Meta Data Model*, implementing code on the data source logic layer that *constructs the statements in the database's query language, cypher, based on the concrete data model configured by the user* and developing *APIs that dynamically adjust to this concrete data model*. So, if the user creates a new *Info Snippet Type*, for example *Build Information*, the *Consumption API* subsequently also considers this entity type during *validation*. Simultaneously, the *Query API* automatically adds an option to query entities of type *Build Information* based on their respective

properties.

Besides, the research to determine the best database technology for the application was time intensive. But, as the decision to use a graph database was frequently questioned by colleagues and peers, the research has been proven to be worthwhile. Most of the skepticism may be traced back to unfamiliarity and has been mostly eliminated by the successful PoC of the *Security and Compliance Data Lake*.

Other time and research intensive tasks that are not represented within the thesis are the proper containerization of the application and the configuration of the SAP Gardener Kubernetes cluster to run the presentation instance of the application, including authentication and encryption. Furthermore, implementing the adapter for the data upload based on existing coding of the SAP Gardener team was cumbersome. As the underlying code has undergone multiple incompatible changes during this work, the adapter had to be adjusted multiple times.

As the application utilizes several of the standards and their respective Go libraries, especially OpenAPI and GraphQL, way beyond their common, well-documented use cases, it was difficult to find suitable examples. Therefore, the development involved a lot of time intensive experimenting.

This and several other factors which significantly increased the cost of implementation may be traced back to the design decisions concerning the data model. But, while the flexibility and extensibility of the data model and the consecutive introduction of the *Meta Data Model* lead to an uncomfortable amount of complexity, it still seems necessary, even in retrospective. As an effort to reduce the complexity, it may be worthwhile to explore the possibilities to replace the *Open API* schemas and tooling entirely by means of *GraphQL*. The briefly mentioned initial prototype which had a rigid data model already leveraged *Open API* which is why the application evolved in a way where such options were not explored thoroughly.

Still, to conclude this scientific work: Compared to the state of the art approach, leveraging some kind of component model, such as the *Open Component Model*, in combination with a central application for storing and querying software metadata, such as the *Security and Compliance Data Lake*, as proposed in this thesis, is shown to undoubtedly be a significant practical advancement which can help companies in fulfilling the requirements resulting from the recently announced Executive Order.

6.4 Outlook

Regarding the *Security and Compliance Data Lake*, the outlook includes the completion of the application by implementing the missing features mentioned during the evaluation. Thus, the aggregation, or rather merging, of data from different data sources, the assessments or triaging functionality and aggregation and filter functions will be added. Specifically within SAP Gardener, Mend shall be integrated as an additional data source, especially since it also provides information about package dependencies. Generally, there are a lot of data sources that may be added but could not all be covered in this work - vulnerability databases such as the NVD, the *Open Source Vulnerability (OSV)* database or the *Global Security Database (GSD)*, or even more scanning tools such as Sonatype or Snyk, to just name a few.

But to loosen the focus from the application and widen the scope to the general problem domain, the research area and the demand for solutions such as the *Security and Compliance Data Lake* is rapidly growing. Even only during this scientific work, there were two publications which have the potential to change the ecosystem around software supply chain security and the general transparency of our digital infrastructure.

The first is the "Securing Open Source Act of 2022" released in September 2022 by the U.S. government [86]. This primarily discusses the implementation of a framework "for assessing the risk of open source software components" [86]. As this framework shall publicly be published within 1 year, it may be interesting to integrate respective assessments into the *Security and Compliance Data Lake*. Due to its extensibility, this may easily be done as an additional type of *Info Snippet*.

The second is the announcement of "GUAC" released in October 2022 by Google [87]. GUAC stands for *Graph for Understanding Artifact Composition*. The article describes it as an application which "aggregates software security metadata into a high fidelity graph database" [87]. Goals include the answering of questions such as "Which parts of my organization's inventory is affected by new vulnerability X?" or "Where am I exposed to risky dependencies?" [87]. Thus, the foundational idea and underlying technology choice is actually quite similar to the *Security and Compliance Data Lake*. Overall, the project is still in its early stages, but it might be worth watching.

Generally, developments such as the *Security and Compliance Data Lake* at SAP and the *Graph for Understanding Artifact Composition* at Google indicate

that the message supplied by the introductory quote, that the trust we place in our digital infrastructure should be matched by how trustworthy and transparent that infrastructure is, seems to have successfully reached the consciousness of the industry leaders.

Bibliography

- [1] THE WHITE HOUSE. “Executive Order on Improving the Nation’s Cybersecurity”. In: *The White House* (2021-05-13). URL: <https://www.whitehouse.gov/briefing-room/presidential-actions/2021/05/12/executive-order-on-improving-the-nations-cybersecurity/> (visited on 08/05/2022).
- [2] NATIONAL VULNERABILITY DATABASE. *Log4j-Vulnerability*, 2022-11-29. URL: <https://nvd.nist.gov/vuln/detail/CVE-2021-44228> (visited on 11/29/2022).
- [3] MARC ANDREESSEN. “Marc Andreessen on Why Software Is Eating the World”. In: *The Wall Street Journal* (2011-20-08). URL: <https://www.wsj.com/articles/SB1000142405311903480904576512250915629460> (visited on 11/29/2022).
- [4] THE UNITED STATES DEPARTMENT OF COMMERCE. *The Minimum Elements for a Software Bill of Materials (SBOM)*, 2021-07-12. URL: https://www.ntia.gov/files/ntia/publications/sbom_minimum_elements_report.pdf (visited on 08/05/2022).
- [5] INA SCHIEFERDECKER, JUERGEN GROSSMANN, and MARTIN SCHNEIDER. “Model-Based Security Testing”. In: *Electronic Proceedings in Theoretical Computer Science* 80 (2012), pp. 1–12. DOI: [10.4204/EPTCS.80.1](https://doi.org/10.4204/EPTCS.80.1).
- [6] MICHAEL FELDERER et al. “Security Testing”. In: *Advances in computers vol. 101*. Ed. by ATIF MEMON. Vol. 101. Advances in Computers. Amsterdam: Elsevier, 2016, pp. 1–51. ISBN: 9780128051580. DOI: [10.1016/bs.adcom.2015.11.003](https://doi.org/10.1016/bs.adcom.2015.11.003).
- [7] STATISTA. *Top software & programming companies by revenue 2017-2022 / Statista: Global 2000*, ed. by FORBES. 2022. URL: <https://www.statista.com/statistics/790179/worldwide-largest-software-programming-companies-by-sales/> (visited on 09/07/2022).
- [8] SAP. *About SAP*, 2022-09-07. URL: <https://www.sap.com/about.html> (visited on 09/07/2022).

- [9] VEDRAN LERENC et al. *SAP Gardener: Value Proposition*, 2022-09-06. URL: <https://www.mend.io/wp-content/media/2021/03/The-state-of-open-source-vulnerabilities-2021-annual-report.pdf> (visited on 09/07/2022).
- [10] GITHUB. *Package URL (purl)*, 2022-09-05. URL: <https://github.com/package-url/purl-spec> (visited on 09/05/2022).
- [11] OPEN SOURCE INITIATIVE. *Open Source Initiative*, 2022-09-02. URL: <https://opensource.org/> (visited on 09/02/2022).
- [12] NATIONAL INSTITUTE OF STANDARDS AND TECHNOLOGIES. *National Vulnerability Database (NVD)*, 2022-08-15. URL: <https://nvd.nist.gov/vuln> (visited on 08/15/2022).
- [13] DAVID E. MANN, STEVEN M. CHRISTEY. *Towards a Common Enumeration of Vulnerabilities*, 1999. URL: <https://www.cve.org/Resources/General/Towards-a-Common-Enumeration-of-Vulnerabilities.pdf> (visited on 08/15/2022).
- [14] MITRE. *CVE*, 2022-08-15. URL: <https://cve.mitre.org/index.html> (visited on 08/15/2022).
- [15] MITRE. *CVE*, 2022-08-09. URL: <https://www.cve.org/> (visited on 08/15/2022).
- [16] FIRST — FORUM OF INCIDENT RESPONSE AND SECURITY TEAMS. *CVSS v3.1 Specification Document*, 2021-01-19. URL: <https://www.first.org/cvss/specification-document#1-2-Scoring> (visited on 08/15/2022).
- [17] MITRE. *CWE*, 2022-08-15. URL: <https://cwe.mitre.org/> (visited on 08/15/2022).
- [18] NATIONAL INSTITUTE OF STANDARDS AND TECHNOLOGY. *Common platform enumeration: naming specification version 2.3*, 2011-08-19. URL: <https://nvlpubs.nist.gov/nistpubs/Legacy/IR/nistir7695.pdf> (visited on 08/11/2022).
- [19] OPEN WEB SECURITY PROJECT FOUNDATION. *OWASP CycloneDX Software Bill of Materials (SBOM) Standard*, 2022-08-05. URL: <https://cyclone-dx.org/> (visited on 08/05/2022).
- [20] INTERNATIONAL ORGANISATION FOR STANDARDIZATION. *Information technology — SPDX Specification V2.2.1*, 2022-08-04. URL: <https://www.iso.org/standard/81870.html> (visited on 08/10/2022).
- [21] SOFTWARE PACKAGE DATA EXCHANGE. *Software Package Data Exchange (SPDX)*, 2021-10-28. URL: <https://spdx.dev/> (visited on 08/09/2022).

- [22] THE LINUX FOUNDATION. *Software Package Data Exchange (SPDX) Specification*, URL: <https://spdx.dev/wp-content/uploads/sites/41/2020/08/SPDX-specification-2-2.pdf> (visited on 08/09/2022).
- [23] SPDX. *Examples of SPDX files for software combinations*, 2022-08-10. URL: <https://github.com/spdx/spdx-examples> (visited on 08/10/2022).
- [24] BUREAU OF INDUSTRY AND SECURITY. *Export Control Classification Number (ECCN)*, 2022-08-31. URL: <https://www.bis.doc.gov/index.php/licensing/commerce-control-list-classification/export-control-classification-number-eccn> (visited on 08/31/2022).
- [25] HÅVARD MYRBAKKEN and RICARDO COLOMO-PALACIOS. “DevSecOps: A Multivocal Literature Review”. In: *Software Process Improvement and Capability Determination*. Ed. by ANTONIA MAS et al. Cham: Springer International Publishing, 2017, pp. 17–29. ISBN: 978-3-319-67383-7.
- [26] THORSTEN RANGNAU et al. “Continuous Security Testing: A Case Study on Integrating Dynamic Security Testing Tools in CI/CD Pipelines”. In: *2020 IEEE 24th International Enterprise Distributed Object Computing Conference (EDOC)*. IEEE, 10/5/2020 - 10/8/2020, pp. 145–154. ISBN: 978-1-7281-6473-1. DOI: [10.1109/EDOC49727.2020.00026](https://doi.org/10.1109/EDOC49727.2020.00026).
- [27] MOJTABA SHAHIN, MUHAMMAD ALI BABAR, and LIMING ZHU. “Continuous Integration, Delivery and Deployment: A Systematic Review on Approaches, Tools, Challenges and Practices”. In: *IEEE Access* 5 (2017), pp. 3909–3943. DOI: [10.1109/ACCESS.2017.2685629](https://doi.org/10.1109/ACCESS.2017.2685629).
- [28] BRIAN FITZGERALD and KLAAS-JAN STOL. “Continuous software engineering: A roadmap and agenda”. In: *Journal of Systems and Software* 123 (2017), pp. 176–189. ISSN: 0164-1212. DOI: [10.1016/j.jss.2015.06.063](https://doi.org/10.1016/j.jss.2015.06.063). URL: <https://www.sciencedirect.com/science/article/pii/S0164121215001430>.
- [29] INGO WEBER, SURYA NEPAL, and LIMING ZHU. “Developing Dependable and Secure Cloud Applications”. In: *IEEE Internet Computing* 20.3 (2016), pp. 74–79. DOI: [10.1109/MIC.2016.67](https://doi.org/10.1109/MIC.2016.67).
- [30] JENKINS.IO. *Jenkins*, 2022-12-01. URL: <https://www.jenkins.io/> (visited on 12/01/2022).
- [31] ROBERT G. COOPER. “Stage-gate systems: A new tool for managing new products”. In: *Business Horizons* 33.3 (1990), pp. 44–54. ISSN: 0007-6813. DOI: [10.1016/0007-6813\(90\)90040-I](https://doi.org/10.1016/0007-6813(90)90040-I). URL: <https://www.sciencedirect.com/science/article/pii/000768139090040I>.

- [32] SANDY CARIELLI et al. *The Forrester Wave™: Software Composition Analysis, Q3 2021*, ed. by FORRESTER. 2021. URL: <https://reprints2.forrester.com/#/assets/2/679/RES176091/report> (visited on 12/02/2022).
- [33] MEND. *Mend Native Integrations for Developers' Environments*, 2022-10-04. URL: <https://www.mend.io/native-integrations-for-developers-environments/> (visited on 12/02/2022).
- [34] SYNOPSYS. *Black Duck DevOps Integrations / Synopsys*, 2022-11-28. URL: <https://www.synopsys.com/software-integrity/security-testing/software-composition-analysis/integrations.html> (visited on 12/02/2022).
- [35] SONATYPE INC. *Software-Entwickler / Software-Tools / Sonatype*, 2022-12-02. URL: <https://de.sonatype.com/solutions/software-developers?topnav=true> (visited on 12/02/2022).
- [36] SNYK. *Open Source Security Management / SCA Tool / Snyk*, 2022-11-09. URL: <https://snyk.io/product/open-source-security-management/> (visited on 12/02/2022).
- [37] NTIA. *Survey of Existing SBOM Formats and Standards*, URL: https://www.ntia.gov/files/ntia/publications/sbom_formats_survey-version-2021.pdf (visited on 12/01/2022).
- [38] SAP GARDENER. *Open Component Model*, 2022. URL: <https://github.com/mandelsoft/ocm/tree/main/docs/ocm> (visited on 12/05/2022).
- [39] CHRISTIAN CWIENK, INGO KOBER, and VASU CHANDRASEKHARA. *The Open Component Model (OCM) fka CNUDIE for Cloud-Native LCM*, 2022. URL: https://video.sap.com/media/t/1_v7rq3sma/145787301 (visited on 09/05/2022).
- [40] OBJECT MANAGEMENT GROUP. “Unified Modeling Language, v2.5.1”. In: (2017). URL: <https://www.omg.org/spec/UML/2.5.1/PDF> (visited on 12/12/2022).
- [41] MEND. *Mend API Documentation*, 2022. URL: https://docs.mend.io/bundle/api_sca/page/http_api_v1_3.html (visited on 12/15/2022).
- [42] CHRISTOF PAAR and JAN PELZL. *Understanding Cryptography*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010. ISBN: 978-3-642-44649-8. DOI: [10.1007/978-3-642-04101-3](https://doi.org/10.1007/978-3-642-04101-3).
- [43] STACK OVERFLOW. *Stack Overflow Developer Survey 2021*, 2021. URL: <https://insights.stackoverflow.com/survey/2021#most-popular-technologies-database> (visited on 12/21/2022).

- [44] E. F. CODD. “A Relational Model of Data for Large Shared Data Banks”. In: *Broy, Denert (Hg.) – Software Pioneers*, pp. 263–294. DOI: [10.1007/978-3-642-59412-0_16](https://doi.org/10.1007/978-3-642-59412-0_16). URL: <https://www.seas.upenn.edu/~zives/03f/cis550/codd.pdf> (visited on 12/21/2022).
- [45] E. F. CODD. *The relational model for database management: Version 2*. 1st ed. 1990, repr. with corr. Reading, Mass. <etc.>: Addison-Wesley, 1990. ISBN: 0201141922. URL: <https://dl.acm.org/doi/pdf/10.5555/77708>.
- [46] RAMEZ ELMASRI and SHAM NAVATHE. *Fundamentals of database systems*. Seventh edition. Boston and Munich: Pearson, 2016. ISBN: 9780133970777. URL: <https://docs.ccsu.edu/curriculumsheets/ChadTest.pdf> (visited on 12/27/2022).
- [47] MYSQL. *The Physical Structure of an InnoDB Index*, 2022-12-27. URL: <https://dev.mysql.com/doc/refman/8.0/en/innodb-physical-structure.html> (visited on 12/27/2022).
- [48] DONALD ERVIN KNUTH. *The Art of Computer Programming - Sorting and searching*. Second edition, forty-second printing. Vol. volume 3. The art of computer programming / Donald E. Knuth. Boston: Addison-Wesley, 2021. ISBN: 9780201896855.
- [49] IAN ROBINSON, JAMES WEBBER, and EMIL EIFREM. *Graph databases: New Opportunities for Connected Data*. Second edition. Beijing: O'Reilly, 2015. ISBN: 9781491930861.
- [50] POSTGRESQL DOCUMENTATION. *PostgreSQL Join Algorithm*, 2022. URL: <https://www.postgresql.org/docs/current/planner-optimizer.html> (visited on 01/09/2023).
- [51] MYSQL. *WITH (Common Table Expressions)*, 2023-01-02. URL: <https://dev.mysql.com/doc/refman/8.0/en/with.html> (visited on 01/02/2023).
- [52] POSTGRESQL DOCUMENTATION. *WITH Queries (Common Table Expressions)*, 2022. URL: <https://www.postgresql.org/docs/current/queries-with.html> (visited on 01/02/2023).
- [53] SQLITE. *The WITH Clause*, 2023-01-02. URL: https://www.sqlite.org/lang_with.html (visited on 01/02/2023).
- [54] AMEYA NAYAK, ANIL PORIYA, and DIKSHAY POOJARY. *Type of NOSQL databases and its comparison with relational databases*. 2013. URL: <https://research.ijais.org/volume5/number4/ijais12-450888.pdf>.

- [55] PRAMOD J. SADALAGE and MARTIN FOWLER. *NoSQL distilled: A brief guide to the emerging world of polyglot persistence*. Always learning. Upper Saddle River, NJ and Munich: Addison-Wesley, 2013. ISBN: 0321826620. URL: <https://bigdata-ir.com/wp-content/uploads/2017/04/NoSQL-Distilled.pdf>.
- [56] MONGODB. *Shard Keys — MongoDB Manual*, 2022-12-15. URL: <https://www.mongodb.com/docs/manual/core/sharding-shard-key/> (visited on 01/04/2023).
- [57] AMAZON WEB SERVICES. *Choosing the Right DynamoDB Partition Key / Amazon Web Services*, 2017. URL: <https://aws.amazon.com/blogs/database/choosing-the-right-dynamodb-partition-key/> (visited on 04/01/2023).
- [58] MONGODB. *Getting Started — MongoDB Manual*, 2022-12-15. URL: <https://www.mongodb.com/docs/manual/tutorial/getting-started/> (visited on 01/04/2023).
- [59] RICK COPELAND. *MongoDB Applied Design Patterns*, 2023-01-05. URL: <https://www.oreilly.com/library/view/mongodb-applied-design/9781449340056/ch01.html> (visited on 01/05/2023).
- [60] MONGODB. *Documents — MongoDB Manual*, 2022-12-15. URL: <https://www.mongodb.com/docs/manual/core/document/> (visited on 01/05/2023).
- [61] MONGODB. *Data Model Design — MongoDB Manual*, 2022-12-15. URL: <https://www.mongodb.com/docs/manual/core/data-model-design/> (visited on 01/05/2023).
- [62] MONGODB. *Aggregation Pipeline — MongoDB Manual*, 2022-12-15. URL: <https://www.mongodb.com/docs/manual/core/aggregation-pipeline/> (visited on 01/09/2023).
- [63] THOMAS H. CORMEN and RONALD L. RIVEST. *Introduction to algorithms*. Third edition. Cambridge: MIT Press, 2014. ISBN: 9780262270830. URL: <http://site.ebrary.com/lib/alltitles/docDetail.action?docID=10397652>.
- [64] WIKIPEDIA. *Graph Database Discussion*, URL: https://en.wikipedia.org/wiki/Talk:Graph_database#Changed_opening_paragraph.
- [65] CLAUDIO WEINBERGER. “Index Free Adjacency or Hybrid Indexes for Graph Databases”. In: *ArangoDB* (2016-04-18). URL: <https://www.arangodb.com/2016/04/index-free-adjacency-hybrid-indexes-graph-databases/> (visited on 01/11/2023).

- [66] NEO4J GRAPH DATA PLATFORM. *Transaction management*, 2023-01-10. URL: <https://neo4j.com/docs/java-reference/current/transaction-management/> (visited on 01/11/2023).
- [67] ARANGODB. *Transactions / Manual / ArangoDB Documentation*, 11.01.2023. URL: <https://www.arangodb.com/docs/stable/transactions.html>.
- [68] NEO4J GRAPH DATA PLATFORM. *DELETE - Cypher Manual*, 2023-01-09. URL: <https://neo4j.com/docs/cypher-manual/current/clauses/delete/> (visited on 01/11/2023).
- [69] ARANGODB. *ArangoDB as Document Store*, 2022-09-29. URL: <https://www.arangodb.com/document-store/> (visited on 11/01/2023).
- [70] NEO4J GRAPH DATA PLATFORM. *Constraints - Cypher Manual*, 2023-01-11. URL: <https://neo4j.com/docs/cypher-manual/current/constraints/> (visited on 01/11/2023).
- [71] ARANGODB. *Schema Validation for Documents / ArangoDB Documentation*, 2023-10-01. URL: <https://www.arangodb.com/docs/stable/data-modeling-documents-schema-validation.html> (visited on 11/01/2023).
- [72] GO. *The Go Programming Language*, 2023-01-20. URL: <https://go.dev> (visited on 01/20/2023).
- [73] MARTIN FOWLER. *Patterns of enterprise application architecture*. Nineteenth printing. The Addison-Wesley Signature Series. Boston et al.: Addison-Wesley, 2013. ISBN: 9780321127426.
- [74] JSON SCHEMA. *JSON Schema*, 2022-11-28. URL: <https://json-schema.org/> (visited on 01/16/2023).
- [75] SWAGGER. *OpenAPI Specification - Version 3.0.3 / Swagger*, 2023-01-12. URL: <https://swagger.io/specification/> (visited on 01/12/2023).
- [76] IETF DATATRACKER. *JSON Schema: A Media Type for Describing JSON Documents*, 2023-01-23. URL: <https://datatracker.ietf.org/doc/html/draft-bhutton-json-schema-00> (visited on 01/23/2023).
- [77] SWAGGER. *API Documentation Tool / Swagger UI*, 2023-01-23. URL: <https://swagger.io/tools/swagger-ui/> (visited on 01/23/2023).
- [78] SPRING. *Spring*, 2023-01-18. URL: <https://spring.io/> (visited on 01/18/2023).
- [79] EXPRESS. *Express - Node.js web application framework*, 2022-09-13. URL: <https://expressjs.com/> (visited on 01/18/2023).

- [80] MARTIN FOWLER. *Richardson Maturity Model*, 2023-01-05. URL: <https://martinfowler.com/articles/richardsonMaturityModel.html> (visited on 01/18/2023).
- [81] ROY FIELDING. “Architectural Styles and the Design of Network-based Software Architectures”. In: (2000). URL: https://www.ics.uci.edu/~fielding/pubs/dissertation/fielding_dissertation.pdf (visited on 01/12/2023).
- [82] PHILIP STURGEON. *Build APIs You Won't Hate 2: This Time It's Serious*, 2023-01-12. URL: <https://github.com/apisyouwonthate/book-build-apis-2> (visited on 01/12/2023).
- [83] GRAPHQL. *GraphQL / A query language for your API*, 2023-01-12. URL: <https://graphql.org/> (visited on 01/12/2023).
- [84] VAULT BY HASHICORP. *Vault by HashiCorp*, 2023-01-26. URL: <https://www.vaultproject.io/> (visited on 01/26/2023).
- [85] GITHUB. *GraphiQL*, 2023-01-26. URL: <https://github.com/graphql/graphiql> (visited on 01/26/2023).
- [86] U.S. GOVERNMENT. *Securing Open Source Act of 2022*, 2023-01-29. URL: <https://www.govinfo.gov/content/pkg/BILLS-117s4913is/pdf/BILLS-117s4913is.pdf> (visited on 01/29/2023).
- [87] BRANDON LUM et al. *Announcing GUAC, a great pairing with SLSA (and SBOM)!*, 2023-01-28. URL: <https://security.googleblog.com/2022/10/announcing-guac-great-pairing-with-slsa.html> (visited on 01/28/2023).

Appendix A

Additional Information

A.1 Black Duck Binary Analysis Result

```
725 ~ | { "lib": "glob",  
726 ~ | "objects": [  
727 | | "etcd-druid"  
728 | ],  
729 | "version": "v0.2.3",  
730 | "vendor": "gobwas",  
731 | "extended-objects": [  
732 | | {  
733 | | | "name": "etcd-druid",  
734 | | | "fullpath": [  
735 | | | | "etcd-druid_v0.13.0_github.com_gardener_etcd-druid",  
736 | | | | "sha256:3db0932b4abd14ebbb6b385a931c3b1ea74d2408124739a97265b8887f182ec5.tar",  
737 | | | | "etcd-druid"  
738 | | | ],  
739 | | | "timestamp": 1662401610,  
740 | | | "size": 62745750,  
741 | | | "sha1": "a86cccd14767d5adf117f9daf20586fadde8c2fdc",  
742 | | | "confidence": 1,  
743 | | | "matching-method": "go-mod-package",  
744 | | | "binary-type": "elf-executable-x86_64",  
745 | | | "package-type": "go-mod",  
746 | | | "type": "go",  
747 | | | "bd-id": "github.com/gobwas/glob",  
748 | | | "detected_version": "v0.2.3"  
749 | | }  
750 | | },  
751 | | "vulns": [],  
752 | | "tags": [  
753 | | | "glob"  
754 | | ],  
755 | | "homepage": "https://github.com/gobwas/glob",  
756 | | "latest-version": "1.0.4",  
757 | | "codetype": "go",  
758 | | "coverity_scan": null,  
759 | | "license": {  
760 | | | "name": "MIT",  
761 | | | "url": "https://opensource.org/licenses/MIT",  
762 | | | "type": "permissive"  
763 | | },  
764 | | "vuln-count": {  
765 | | | "total": 0,  
766 | | | "exact": 0,  
767 | | | "historical": 0  
768 | | },  
769 | | },  
770 | },
```

Figure A.1: Snippet of BDBA Analysis Result

Source: [34]

A.2 Relational Model of the Example Data Model Instance

This list shows a exemplary relational model corresponding to figure 5.2 in section 5.2.3. The properties are deducted from the OCM, but as stressed several times, they are just one possible example. As this relational model is mainly supposed to support the general understanding of the problem domain, it treats each attribute as atomic. Thus, complex attributes such as repository context are not further normalized for the scope of this relational model.

Primary key attributes are bold and underlined and *foreign key attributes* are italic and underlined. Naturally, in this data model which rarely uses surrogate keys, but still has to represent several hierarchical relationships, attributes are frequently part of both, the primary and the foreign key.

- **Component**(**name**)
- **ComponentVersion**(**name**, **version**, labels, repository-context, provider)
- **ComponentReference**(**source-component-name**,
source-component-version, **destination-component-name**,
destination-component-version, **name**, **extra-identity**, labels)
- **Resource**(**name**, **extra-identity**, **component-name**,
component-version)
- **ResourceVersion**(**name**, **version**, **extra-identity**, **component-name**,
component-version, digest, labels, resource-type)
- **Source**(**name**, **extra-identity**, **component-name**, **component-version**)
- **SourceVersion**(**name**, **version**, **extra-identity**, **component-name**,
component-version, labels)
- **IsBuiltFrom**(**component-name**, **component-version**, **resource-name**,
resource-version, **resource-extra-identity**, **source-name**,
source-version, **source-extra-identity**)
- **Package**(**id**, name)
- **PackageVersion**(**id**, **package-id**, name, version, ecosystem, platform)

- **BDBAPackageVersion**(name, version, package-version-id, ecosystem, platform, matching-method)
- **SourceVersionIsComprisedOf**(component-name, component-version, source-name, source-version, extra-identity, package-version-id)
- **ResourceVersionIsComprisedOf**(component-name, component-version, resource-name, resource-version, extra-identity, package-version-id)
- **DependsOn**(package-version-id, package-version-id)
- **Vulnerability**(cve, summary, cvss, cwe)
- **Contains**(package-version-id, cve)

The IsBuiltFrom-relationship contains a particularly interesting detail about this mapping. As explained in the respective sections, due to the *Component Version-Local Identities of Artifacts* in the OCM, *Resource Versions* may only be built from *Source Versions* within the same *Component Version*. Thus, even though the component name and component version are part of the identities of both *Source Version* and *Resource Version*, it is only contained once in the IsBuiltFrom relation manifesting this (n:m)-relationship.

Another detail, of which the purpose may not be immediately obvious, are the surrogate keys, id, of *Package Version* and *Package*. The scanning tools and other data sources may use different identifiers for the same technical package or even identify different technical packages with the same identifier. Thus, the decision, that for example a *BDBA Package Version* and a *Mend Package Version* actually represent the same technical package has to be done by humans. Until then, the *BDBA Package Version* and the *Mend Package Version* have to be stored without merging and therefore relating to two different *Package Versions* and *Packages*. With natural keys, this could therefore lead to duplicate *Package Versions* and *Packages*.

A.3 Relation of Number of Keys and Number of Leaves in a B-tree

Given $N + 1$ is the number of leaf nodes of a B-tree. $n_l \in \mathbb{N}$ is the number of nodes on a particular level $l \in \mathbb{N}$. Thereby, $l = 1$ is the lowest level, thus the leaf node level, and consequently n_1 is the number of leaf nodes, thus $n_1 = N + 1$. As a B-tree is a tree, it always has a root node. So $l = \text{root}$ is the highest level with $n_l = 1$.

As per constraint, the number of keys in a node is $k - 1$. Since all the nodes on $l = 1$ are children of some node on $l = 2$ and every node on $l = 2$ has one key less than it has children, the total number of keys on l_2 is equal to $K_2 = n_1 - n_2$. This correlation is generally true, the number of keys on a level is equal to the number of nodes on the level below minus the number of nodes on the respective level (with the only exception of $l = 1$).

$$K_l = n_{l-1} - n_l$$

The resulting sequence $(K_l)_{l=2}^{root}$ can be used to form a series calculating the cumulative number of keys up to a particular level.

$$\Sigma K_l = \sum_{i=2}^l n_{i-1} - n_i$$

And as every tree has a root node, n always reaches 1 at some point. This can be written as follows.

$$\Sigma K_{root} = (n_0 - n_1) + (n_1 - n_2) + (n_2 - n_3) + \dots + (n_{l-1} - 1)$$

Obviously, all elements but the first, n_0 , and the last, 1, appear twice, once with a positive and once with a negative sign and thus reduce each other.

$$\Sigma K_{root} = n_0 - 1$$

And since $n_0 = N + 1$, the number of keys in a B-tree is

$$\Sigma K_{root} = N + 1 - 1 = N$$

Appendix B

Data Definition Statements for Sample Data

B.1 Data Definition Statement for SQL Sample Data

```
1 CREATE TABLE PackageVersion(
2   ID int,
3   Name varchar(100),
4   Version varchar(100)
5 );
6
7 INSERT INTO PackageVersion VALUES
8   (1, 'log4j', '1.5.2'), -- depends on spring-context directly
9   (2, 'spring-context', '1.2.12'),
10  (3, 'java-jdk', '1.3.1'), -- direct dependency of spring context
11  (4, 'spring-boot', '1.15.3'), -- depends on spring-context
12    transitivity (spring-boot -> log4j -> spring-context)
12  (99, 'spring-core', '2.1.4'); -- direct dependency of spring
13    context
13
14 CREATE TABLE DependsOn(
15   PackageVersionID int,
16   RelatedPackageVersionID int
17 );
18
19 INSERT INTO DependsOn VALUES
20   (1, 2),
21   (2, 3),
22   (2, 99),
23   (4, 1),
24   (99, 3);
```

Listing B.1: SQL Example Data Creation Statements

B.2 Data Definition Statement for MongoDB Sample Data

```
1 db={
2   "orders": [
3     "_id": 1,
4     "date": "04.01.2023",
5     "shippingAddress": {
6       "city": "Berlin",
7       "street": "Kreuzberg" },
8     },
9     {
10     "_id": 2,
11     "date": "05.01.2023",
12     "shippingAddress": {
13       "city": "Hamburg",
14       "street": "Sesame Street" },
15     },
16     "relations": [
17       "_id": 1,
18       "order": 1,
19       "product": 1
20     },
21     {
22       "_id": 2,
23       "order": 1,
24       "product": 2
25     },
26     {
27       "_id": 3,
28       "order": 2,
29       "product": 2
30     },
31     "products": [
32       "_id": 1,
33       "name": "The Hitchhiker's Guide to the Galaxy",
34       "price": 7.50,
35     },
36     {
37       "_id": 2,
38       "name": "The Art of Computer Programming",
39       "price": 120.00,
40     }
41 }
```

Listing B.2: Document Database Example Data Creation

B.3 Data Definition Statement for Neo4J Sample Data

```
1 MERGE (p1:PACKAGE_VERSION {Name:"log4j",Version:"1.5.2"})
2 MERGE (p2:PACKAGE_VERSION {Name:"spring-context",Version:"1.2.12"
   "})
3 MERGE (p3:PACKAGE_VERSION {Name:"java-jdk",Version:"1.3.1"})
4 MERGE (p4:PACKAGE_VERSION {Name:"spring-boot",Version:"1.15.3"})
5 MERGE (p5:PACKAGE_VERSION {Name:"spring-core",Version:"2.1.4"})
6
7 MERGE (p1)-[r1:DEPENDS_ON]->(p2)
8 MERGE (p2)-[r2:DEPENDS_ON]->(p3)
9 MERGE (p2)-[r3:DEPENDS_ON]->(p5)
10 MERGE (p4)-[r4:DEPENDS_ON]->(p1)
11 MERGE (p5)-[r5:DEPENDS_ON]->(p3)
```

Listing B.3: Cypher Example Data Creation