

Ostbayerische Technische Hochschule Amberg-Weiden
Fakultät Elektrotechnik, Medien und Informatik

Studiengang Medieninformatik

Bachelorarbeit

von

Albert Hahn

**Konzeption und Implementierung einer Microservice
Architektur in einem hybriden kubernetes Cluster für
industrielle KI-Anwendungsfälle**

Conceptual Design and Implementation of a Microservice
Architecture in a Hybrid Kubernetes Cluster for Industrial
AI Use Cases

Ostbayerische Technische Hochschule Amberg-Weiden
Fakultät Elektrotechnik, Medien und Informatik

Studiengang Medieninformatik

Bachelorarbeit

von

Albert Hahn

**Konzeption und Implementierung einer Microservice
Architektur in einem hybriden kubernetes Cluster für
industrielle KI-Anwendungsfälle**

Conceptual Design and Implementation of a Microservice
Architecture in a Hybrid Kubernetes Cluster for Industrial
AI Use Cases

Bearbeitungszeitraum: von 4. Oktober 2021
 bis 3. März 2022

1. Prüfer: Prof. Dr.-Ing. Christoph Neumann

2. Prüfer: Prof. Dr. Dieter Meiller

Bestätigung gemäß § 12 APO

Name und Vorname
der Studentin/des Studenten: **Hahn, Albert**

Studiengang: **Medieninformatik**

Ich bestätige, dass ich die Bachelorarbeit mit dem Titel:

**Konzeption und Implementierung einer Microservice Architektur in einem
hybriden kubernetes Cluster für industrielle KI-Anwendungsfälle**

selbständig verfasst, noch nicht anderweitig für Prüfungszwecke vorgelegt, keine
anderen als die angegebenen Quellen oder Hilfsmittel benutzt sowie wörtliche und
sinngemäße Zitate als solche gekennzeichnet habe.

Datum: 18. Januar 2022

Unterschrift:

Bachelorarbeit Zusammenfassung

Studentin/Student (Name, Vorname):	Hahn, Albert
Studiengang:	Medieninformatik
Aufgabensteller, Professor:	Prof. Dr.-Ing. Christoph Neumann
Durchgeführt in (Firma/Behörde/Hochschule):	Krones AG, Neutraubling
Betreuer in Firma/Behörde:	Ottmar Amann
Ausgabedatum: 4. Oktober 2021	Abgabedatum: 3. März 2022

Titel:

**Konzeption und Implementierung einer Microservice Architektur in einem
hybriden kubernetes Cluster für industrielle KI-Anwendungsfälle**

Zusammenfassung:

Das Ziel dieser Bachelorarbeit ist es, eine flexible und nahtlose Lösung für ein Hybrides Cluster aus on-premise Edge Devices und Cloud Ressourcen bereitzustellen. Produktionslinienanwendungen/Microservices sollen zukünftig beliebig skalierbar und agil sein, dabei sollen für die Anwendungen generell keine Differenzierung zwischen offline und online Ressource getroffen werden. Im Zuge dessen wird die Umsetzbarkeit und Relevanz von cloudbasierten Microservices im Bereich der künstlichen Intelligenz auf einer zukünftigen Produktionsanlage untersucht.

Schlüsselwörter:

Inhaltsverzeichnis

1	Einleitung	2
1.1	Motivation	3
1.2	Zielsetzung	3
2	Theoretische Grundlagen	4
2.1	Docker	4
2.1.1	Aufbau	4
2.1.2	Architektur	4
2.2	Microservice	4
2.2.1	Aufbau	4
2.2.2	Entwicklung	4
2.2.3	Dezentrale Datenmanagement	4
2.3	Kubernetes	4
2.3.1	Architektur	4
2.3.2	Lightweight Kubernetes	4
2.3.3	Hybrid Cloud	4
2.3.4	Rancher	4
3	Microservice Architektur	5
3.1	Anforderungen	5
3.2	Konzept	5
3.3	Herausforderungen	5
3.4	KI-Anwendungsfall	5
4	Implementierung der Architektur	6
4.1	Frameworks und Bibliotheken für Microservices	7
4.1.1	Flask	7
4.1.2	Gunicorn	7
4.1.3	SocketIO	7
4.1.4	OpenCV	7
4.1.5	MongoDB	7
4.2	Containerisierung	7
4.2.1	Volumes	7
4.2.2	Netzwerk	7
4.2.3	CUDA Toolkit	7

4.2.4	DockerHub	7
4.3	Orchestrierung	7
4.3.1	Deployment	7
4.3.2	Ingress	7
4.3.3	Loadbalancer	7
4.3.4	Taints and Tolerations	7
4.3.5	Node Affinity	7
4.3.6	Helm	7
4.4	Konfiguration und Einrichtung	7
4.4.1	Virtueller Privater Server	7
4.4.2	Domain	7
4.4.3	SSL-Verschlüsselung	7
4.5	Aufbau der Implementierung	7
4.5.1	Software-Layer	7
5	Ergebnisse	8
5.1	Microservice	8
5.1.1	Frontend-Service	8
5.1.2	Backend-Service	8
5.1.3	Loadbalancer	8
5.1.4	Kubernetes Cluster	8
6	Diskussion und Ausblick	9
	Literaturverzeichnis	10
	Abbildungsverzeichnis	11
	Tabellenverzeichnis	12

Kapitel 1

Einleitung

Die Krones AG bietet Anlagen für die Getränkeindustrie als auch Nahrungsmittelhersteller, von der Prozesstechnik bis hin zur IT-Lösung. Die Komplettlinie beinhaltet auch das bereitstellen von Software auf den einzelnen Produktionsanlagen. Hierfür werden eine Vielzahl von Produktionslinienanwendungen auf den Anlagen installiert, gewartet und verwaltet. Ein riesiger Aufwand der Fehleranfälligkeiten wie fehlende Frameworks, Bibliotheken und anderer Abhängigkeiten mit sich bringt. Eigene Server müssen für die Kommunikation der Anlagen verbaut und gewartet werden, was zusätzlich Ressourcen beansprucht und automatisch die Kosten für die Inbetriebnahme einer solchen Linien erhöhen. Die Weiterentwicklung der zukünftigen Bereitstellung von Produktionsanlagensoftware erfolgt mithilfe eines Proof of Concept (PoC), welcher die Möglichkeiten einer wartungsfreien Infrastruktur durch ein continuous delivery System evaluiert. Dies verläuft in Zusammenarbeit mit dem Kooperationspartner und Softwareunternehmen SUSE GmbH, welches das wartungsfreie Betriebssystem SUSE Linux Enterprise Micro und die multi-cluster Orchestrierungsplattform Rancher anbietet.

Als Grundlage hierfür dient das Open-Source-System Kubernetes, welches zur Automatisierung, Skalierung und Verwaltung von containerisierten Anwendungen bestimmt ist. Künftige Produktionsanlagen sollen mittels zusätzlichen Edge Devices als Knotenpunkte in einem Kubernetes Cluster fungieren, Ressourcen teilen, untereinander kommunizieren und Softwarepakete unkompliziert bereitstellen. Die Integration der kompakten Linux Rechner ermöglichen den Variablen Einsatz von Hardwareressourcen beim Kunden, der je nach Leistungsanspruch Knotenpunkte erweitern kann. Dabei soll es für die einzelnen Anwendungen möglich sein, sowohl auf cloudbasierten als auch auf on-premise Hardware zur Verfügung gestellt zu werden. Ein hybrides Kubernetes Cluster ermöglicht es somit lokale Rechenleistung oder öffentliche Cloudressourcen in der selben Softwareumgebung zu nutzen.

1.1 Motivation

Die Vorteile von Kubernetes und dem stetigen Paradigmenwechsel der Softwarelandschaft im Cloudbereich, welcher den Wechsel von monolithischen Architektur zu einer mehr flexibleren microservice Architektur bevorzugt, sind das Hauptmotiv der Auswertung neuer agiler Distributionsmöglichkeiten. Die Containerisierung von Anwendungen ermöglichen erst die Aufteilung großer Projekte in kleine unabhängige Services die mittels Orchestrierungsplattformen sinnvoll gebündelt werden können. Namenhafte Unternehmen wie Netflix, Amazon und Uber entwickeln und verwenden bereits robuste und komplexe Microservices die containerisiert auf Plattformen verwaltet werden [1].

Durch die Flexibilität einer solchen Infrastruktur ist es möglich Anwendungsfälle im Bereich der künstlichen Intelligenz für die Industrie zu testen. Die Anlage Linatronic AI der Krones AG nutzt bereits Deep-Learning-Technologie, um in der Linie mittels Vollinspektion Schäden, Dichtflächen oder Seitenwanddicken zu erkennen und Prozesse zu optimieren [2]. Allgemein sind Anwendungen mit künstlicher Intelligenz durch ihre Komplexität und Vielzahl an Abhängigkeiten schwierig zu entwickeln und bereitzustellen. Eine passende Plattform für Anwendungsfälle mit Bezug zur künstlichen Intelligenz muss eine Vielzahl an Services anbieten. Verwaltung von Ressourcen wie Speicher, Rechenleistung und Verbindungsgeschwindigkeit für die Datenübertragung, bei der Ausführung einzelner Phasen von der Datenverarbeitung bis hin zur Evaluierung und Entwicklung [3].

1.2 Zielsetzung

Ziel dieser Arbeit ist die Entwicklung einer Microservice Architektur in einem hybriden Kubernetes Cluster. Das Endresultat soll eine Anwendung werden die mittels einer Weboberfläche, welche über eine Domain erreichbar ist, ein Login-Verfahren mittels einem backend Service ermöglichen der ein Authentifizierungsverfahren per Gesichtserkennung verwendet. Diese Daten sollen schließlich verarbeitet und persistent gespeichert werden, um bei erneuten Aufruf der Website bestehen zu bleiben. Die Konzeption der Anwendung findet containerisiert auf mehreren Software und Hardware Layern statt. Beispielhaft an der Auswahl von Ressourcen zur Auswertung von Gesichtsdaten, dabei steht die Option zwischen einem „Fast“ oder „Performance“ Modus zur Verfügung, welcher entweder einen Knotenpunkt mit Grafikprozessor auswählt oder einer herkömmlichen Prozessoreinheit. Das ganze System wird auf einem Kubernetes Cluster bereitgestellt und verwaltet. Ein Ingress Controller dient dabei als Loadbalancer und verteilt die Last beim Aufrufen der Website und der Kommunikation zwischen backend Service.

Kapitel 2

Theoretische Grundlagen

Dieses Kapitel erläutert die Grundlagen die zum Verständnis dieser Arbeit notwendig sind.

2.1 Docker

2.1.1 Aufbau

2.1.2 Architektur

2.2 Microservice

2.2.1 Aufbau

2.2.2 Entwicklung

2.2.3 Dezentrale Datenmanagement

2.3 Kubernetes

2.3.1 Architektur

2.3.2 Lightweight Kubernetes

2.3.3 Hybrid Cloud

2.3.4 Rancher

Kapitel 3

Microservice Architektur

3.1 Anforderungen

3.2 Konzept

3.3 Herausforderungen

3.4 KI-Anwendungsfall

Kapitel 4

Implementierung der Architektur

4.1 Frameworks und Bibliotheken für Microservices

4.1.1 Flask

4.1.2 Gunicorn

4.1.3 SocketIO

4.1.4 OpenCV

4.1.5 MongoDB

4.2 Containerisierung

4.2.1 Volumes

4.2.2 Netzwerk

4.2.3 CUDA Toolkit

4.2.4 DockerHub

4.3 Orchestrierung

4.3.1 Deployment

4.3.2 Ingress

Nginx-Ingress

4.3.3 Loadbalancer

4.3.4 Taints and Tolerations

4.3.5 Node Affinity

Kapitel 4. Implementierung der Architektur

4.3.6 Helm

4.4 Konfiguration und Einrichtung

Kapitel 5

Ergebnisse

5.1 Microservice

5.1.1 Frontend-Serivce

5.1.2 Backend-Serivce

5.1.3 Loadbalancer

5.1.4 Kubernetes Cluster

Kapitel 6

Diskussion und Ausblick

Literaturverzeichnis

- [1] M. Villamizar, O. Garces, H. Castro, M. Verano, L. Salamanca, R. Casal- ´ las, and S. Gil, "Evaluating the monolithic and the microservice architecture pattern to deploy web applications in the cloud," in 2015 10th
- [2] Krones Linatronic 735," Krones.com. [Online]. Available: <https://www.krones.com/de/produkte/maschinen/leerflaschen-inspektionsmaschine-linatronic-735.php>. [Accessed: 16-Jan-2022].
- [3] Y. Zhou, Y. Yu and B. Ding, "Towards MLOps: A Case Study of ML Pipeline Platform,"2020 International Conference on Artificial Intelligence and Computer Engineering (ICAICE), 2020, pp. 494-500, doi: 10.1109/ICAICE51518.2020.00102.

Abbildungsverzeichnis

Tabellenverzeichnis