

Erstellung und Validierung eines Leitfadens zur
Migration einer Microservice-Architektur nach
Function as a Service

Frankfurt University of Applied Sciences

Gianni Pasqual

18.03.2020

Abstract

Mit der Technologie Function as a Service ist nach der Microservice-Architektur eine noch feinere Modularisierungsstufe von Software erreicht worden. Mit dieser ist es möglich Software noch schneller bereitzustellen und im gleichen Zuge Einsparungen bei dem Betrieb der IT-Infrastruktur zu erzielen. Dies erscheint auf den ersten Blick natürlich sehr lokrativ für viele Unternehmen, jedoch sollte die Entscheidung für den teilweisen oder gesamten Umzug der Service-Landschaft wohl überlegt sein. Wie bei jeder anderen Technologie hat auch diese ihre Potentiale und Herausforderungen die gemeistert werden müssen. Für Unternehmenn, die sich dazu entschieden haben Function as a Service (FaaS) in ihre Infrastruktur einzubinden, soll diese Arbeit als Orientierung bei der Migration dienen. Da viele Unternehmen bereits von den monolithischen Anwendungen auf kleinere modularere Services umgestiegen sind, geht diese Arbeit von einer bereits vorliegenden Microserive-Architektur aus. Es gilt herauszufinden in welchem Maße und ob überhaupt strukturelle Anpassungen vorgenommen werden müssen, als auch der Frage nach "Best-Practices", nach nun knapp vier Jahren des Bestehens, nachzugehen.

Zudem soll sich vor allem mit dem von den verschiedenen Anbietern ausgehenden Vendor-Lock-In auseinadner gesetzt werden und erarbeitet werden, welche Möglichkeiten bestehen, diesen zu umgehen bzw. zu mildern. Um mögliche Schwächen des Leitfadens aufzuzeigen, soll dieser beispielhaft an einem Service erprobt werden und auftretenden Fehler dokumentiert und behandelt werden.

Abkürzungsverzeichnis

ADF	Azure Durable Functions
ASF	Amazon Step Functions
AWS	Amazon Web Services
BaaS	Backend as a Service
BDD	Behaviour Driven Development
CaaS	Container as a Service
CD	Continuous Delivery
CI	Continuous Integration
CNCF	Cloud Native Computing Foundatino
DevOps	Development and Operations
DEV	Development
DSL	Domain Specific Language
FaaS	Function as a Service
IAM	Identity Access Management
IaaS	Infrastructure as a Service
INT	Integration (Testing)
IT	Information Technology
NIST	National Institute of Standards and Technology
PaaS	Platform as a Service
PROD	Production
RPC	Remote Procedure Call
SAM	Serverless Application Model
SDK	Software Development Kit
UAT	User-Acceptance (Testing)
VCS	Version Controle System
XaaS	Anything as a Service

Inhaltsverzeichnis

1	Einleitung	5
1.1	Forschungsfrage und Ziele der Arbeit	5
1.2	Erläuterung der Problemstellung	5
1.3	Aufbau der Thesis	5
2	Aufarbeitung des Themengebietes	5
2.1	Definition des Begriffes Cloud Computing	5
2.2	Function as a Service	8
2.3	Abgrenzung von FaaS zu Serverless	12
2.4	Potentiale und Herausforderungen	14
3	Guidline on Migration	20
3.1	Zu grunde liegenden Architektur	20
3.2	Challenging the Decision	21
3.3	Choosing a suitable service	23
3.4	Auswirkungen auf die Team-/ Entwicklungsstruktur	23
3.5	Abwägung Open-Source und Cloud-Vendor	25
3.5.1	Vendoren-Analyse	26
3.5.2	Tools zur Entwicklungsunterstützung	29
3.5.3	Cross-Functions in Multi-Cloud Lösung	29
3.6	Auswirkungen auf die Testumgebungen	29
3.7	Konsequenzen für das Monitoring	31
4	Anwendung des Leitfadens	32
4.1	Auswertung der Ergebnisse	32
4.1.1	Beurteilung der Kollaborationsauswirkungen	32
4.1.2	Beurteilung der Stabilität	32
4.1.3	Beurteilung der Skalierbarkeit	32
4.1.4	Beurteilung des Monitorings	32
4.1.5	Beurteilung der Testmöglichkeiten	32
4.2	Korrektur und Anpassungen	32
5	Abschließende Betrachtung	32
5.1	Absehbare Entwicklungen	32
5.2	Zusammenfassung	32
5.3	Weiterführende Forschung	32

1 Einleitung

1.1 Forschungsfrage und Ziele der Arbeit

1.2 Erläuterung der Problemstellung

1.3 Aufbau der Thesis

2 Aufarbeitung des Themengebietes

Ist eine Technologie, wie Function as a Service, noch relativ jung, so sind die Schritte der Findung einer in sich schlüssigen und allgemein vertretenen Definition oftmals noch nicht abgeschlossen. Auch bei der Definition von FaaS, Serverless und der Einordnung dieser beiden Konzepte in die Infrastruktur des Cloud Computings⁸ [MG+11], ist dieser Prozess noch im Gange, wobei mittlerweile die unterschiedlichen Definitionen der öffentlichen FaaS- und Serverless-Anbieter einige Gemeinsamkeiten aufweisen. Trotz alledem bestehen weiterhin Ungenauigkeiten, die es in den nächsten Jahren noch zu beseitigen gilt. Hierzu später mehr.

2.1 Definition des Begriffes Cloud Computing

Der Begriff Cloud Computing ist laut eines Technology Report bis auf das Jahr 1996 zurückzuführen, wo er in einem Business-Plan des Unternehmens Compaq von ein paar Entwicklern genutzt worden sein sollen, um über die zukünftigen Entwicklungen des Internet-Businesses zu diskutieren [Reg11].

Auch wenn das Konzept des Cloud-Computings, die dynamische, skalierbare, zuverlässige und unbegrenzte Bereitstellung von Ressourcen als Dienst über das Internet, immer mal wieder in der Literatur auftauchte [Fox+09], so erlangte es erst ab 2006 einen immer größer werdenden Bekanntheitsgrad. 2006 veröffentlichte AWS das Produkt Elastic Compute Cloud (EC2), gefolgt von Goolges App Engine 2008. Bereits in diesem früher Stadium wurde mit App Engine das Prinzip einer sog. „stateful“-Ebene, für das Speichern von des Aktuellen Anwendungsstatus und einer „stateless“-Ebene für die Ausführung des eigentlichen Programmes unterschieden [Fox+09]. Des Weiteren wurde 2009 in einer Berkeley-Studie die sechs größten Potentiale des Cloud-Computings herausgearbeitet welche bis heute in den unterschiedlichen XaaS-Konzepten umgesetzt wurden. Im folgenden sind diese kurz auf

gelistet.

1. Unbegrenzte Ressourcen wann immer sie benötigt werden.
2. Durch die Möglichkeit mit wenigen Ressourcen zu starten, vielen Nutzern Zugang zu der Plattform zu gewähren.
3. Die genutzten Ressourcen so genau wie möglich nach dem „Pay Per Use“-Prinzip zu bezahlen.
4. Größenvorteile (Economies of scale) nutzen, um die Kosten durch eine dauerhafte, optimale Ausnutzung riesiger Datacenter auf ein Minimum zu reduzieren.
5. Operationelle Kosten so weit es geht zu senken und die Ressourcennutzung durch Virtualisierung so weit wie möglich auszunutzen.
6. Eine hohe Hardwareausnutzung durch „Multiplexing“ der Auslastungen verschiedener Unternehmen zu erreichen.

Das National Institute of Standards and Technology (NIST) veröffentlichte 2011, nach 16 vorausgehenden Definitionen, schließlich eine finale Definition des Cloud-Computings, welche bis heute als Orientierung Anwendung findet und in einer Vielzahl an Werken aufgegriffen wird [MG+11]. Die Definition nach NIST beschreibt das Cloud Model als eines aus fünf essentiellen Charakteristiken und drei Service Modellen, IaaS, PaaS und SaaS bestehendes Modell, mit insgesamt vier Deployment-Models.

“Cloud Computing is a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, storage, applications, and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction.” [MG+11].

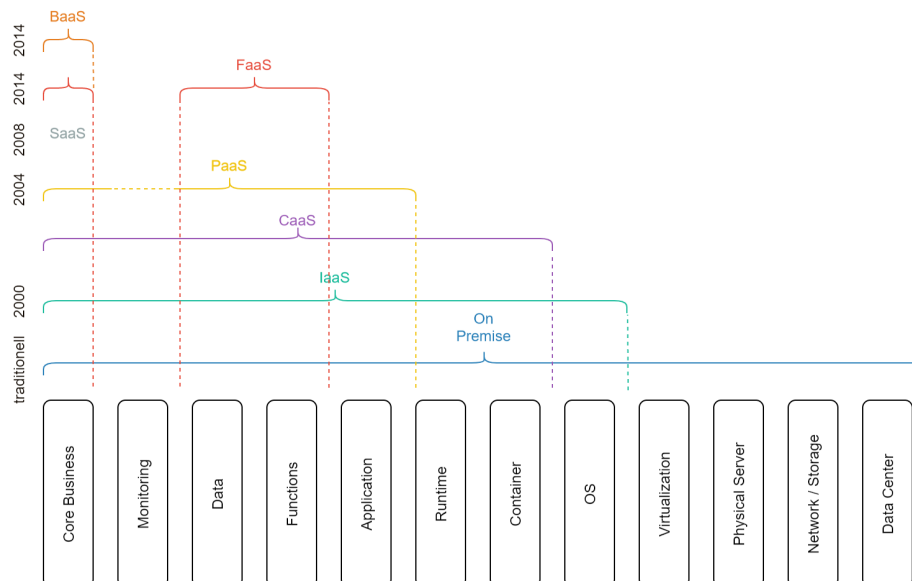
NIST bricht die Definition mit Essentiellen Charakteristiken, Service-Models und Deployment-Models in drei Unterkategorien herunter, welche erfüllt sein müssen um dem Anspruch an Cloud Computing zu genügen.

Grob gefasst sind dies bei den Charakteristiken, die automatische Bereitstellung der benötigten Ressourcen je nach Bedarf (*On-demand self-service*), die Möglichkeit von allen gängigen Geräten darauf zugreifen zu können (*Broad network access*), die optimale Verteilung von Ressourcen auf die Kunden welche sie gerade benötigen, wobei hier die exakte Lokation dieser von geringer

Relevanz ist (*Resource pooling*), die automatische Skalierung von Ressourcen (*Rapid elasticity*) und die Möglichkeit der Überwachung und Limitierung der Ressourcenausnutzung, welche für den Kunden als auch den Anbieter transparent sein muss (*Measured service*).

Bei den Service Models werden hier lediglich IaaS, PaaS und SaaS unterschieden. Laut NIST gibt es um den Code in der Cloud bereitzustellen vier verschiedene Deployment-Models, welche genutzt werden können. Mit der *Private cloud*, obliegt es dem Unternehmen seine Infrastruktur komplett selber zu betreiben, jedoch muss diese dem Unternehmen dabei nicht selber gehören, sondern kann von einer dritten Partei gehostet werden. Mit dem deployment auf eine *Community cloud* teilen sich verschiedene Unternehmen eine Cloud, welche entweder von ihnen gemeinsam, oder von einem Drittanbieter betrieben wird. Zuletzt gibt es noch die Möglichkeit die von einem Cloud-Computing-Anbieter bereitgestellte öffentliche Cloud *Public cloud* zu nutzen, bei welcher das Unternehmen die Räumlichkeiten bzw. Ressourcen von einem Drittanbieter nutzt. Das letzte Deployment-Model stellt die sog. *Hybrid cloud* dar, welche eine Kombination aus den vorherigen darstellt [MG+11].

Abbildung 1: Übersicht Service-Models Cloud-Computing



Neben den bereits genannten Service Models IaaS, PaaS und SaaS gibt es jedoch noch weitere, welche sich im Laufe der Zeit etabliert haben. Abbildung 1 gibt einen Überblick hierüber.

2.2 Function as a Service

Function as a Service (FaaS) ist ein sogenanntes „Serverless“ Cloud-Computing Konzept, welches erstmals 2014 von AWS mit Lambda Functions als Preview Release und schließlich 2015 zur kommerziellen Nutzung zur Verfügung gestellt wurde ¹. Es kann nach IaaS und PaaS als ein weiterer Schritt in der Entwicklung des Cloud Computings gesehen werden, welcher das Management von Infrastruktur, Servern und Ressourcen weg von dem Entwickler nimmt und hin zu dem Cloud-Anbieter delegiert. Ein knappes Jahr später, 2016, traten Microsoft mit Azure Functions, Google mit Google Cloud Functions und IBM mit OpenWhisk in den bis dahin von Amazon dominierten Markt ein. Mittlerweile gibt es eine Vielzahl an Open-Source sowie proprietären Anbietern, welche um die Gunst der Kunden werben und einen stetigen Wettbewerb aufrecht erhalten. Eine ausführliche Übersicht über die jeweiligen Anbieter findet sich online von der CNCF ²

Da es sich mit FaaS und Serverless um eine noch recht junge Technologie handelt, gibt es zu diesem Konzept keine offiziell dokumentierte Definition. bsp. von NIST o.ä., in der Literatur, wie es beispielsweise beim Cloud-Computing [MG+11] der Fall ist. Die Anbieter sind sich jedoch über die Funktionalität von FaaS größtenteils einig, was das Spektrum an Funktionalitäten angeht. So definiert Microsoft Azure Functions als: „[...] event-driven serverless compute platform that can also solve complex orchestration problems. Build and debug locally without additional setup, deploy and operate at scale in the cloud, and integrate services using triggers and bindings.“ ³. Hiermit beschreibt der Anbieter bereits die Kernfunktionalitäten von FaaS und drückt den Kerngedanken hinter FaaS aus. Das Konzept soll genutzt werden können, um bei dem Auftreten eines zuvor definierten sog. Triggers auf ein Event aus der eigenen Applikation oder der des Anbieters reagieren zu können. Dabei stellen die meisten Anbieter wie Google, Azure oder AWS neben einem einfachen event-basierten HTTP-Trigger oder wiederkehrenden zeitbasierten Triggern noch weitere, an die jeweilige Infrastruktur des Anbieters angepasste, Trigger zur Verfügung. Diese ergeben sich meist

¹<https://docs.aws.amazon.com/lambda/latest/dg/lambda-releases.html>

²<https://github.com/cncf/wg-serverless>

³<https://azure.microsoft.com/en-us/services/functions/>

aus den BaaS Angeboten, hierzu mehr in *Abgrenzung von FaaS zu Serverless*, welche in Form von Datenbank Triggern bei einer CRUD-Operation, dem Anlegen eines Nutzers oder der Integration mit einem Push-Notification Service [AWS SSN, Google Pub/Sub etc.] auftreten können.

Abbildung 2: ⁴

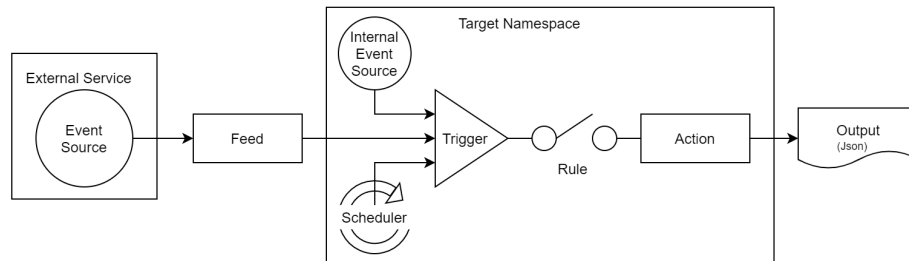


Abbildung 2 gibt hierbei eine Übersicht über die verschiedenen Event-Trigger, welche derzeit bei den großen, proprietären Anbietern (AWS, Google, Microsoft) existieren. Da hinter OpenWhisk mit IBM eine großes Unternehmen steht, wird dieses Modell in der Literatur häufig als Referenz in Bezug auf die anderen Anbieter verwendet. Es kann angenommen werden, dass andere Anbieter ein ähnliches Konzept bei ihrer FaaS Architektur verfolgen, da Apache OpenWhisk in einer großen Cloud (IBM) bereitgestellt wird [VE+19]. Eine weitere fundamentale Eigenschaft ist die erwähnte Skalierbarkeit, welche nach dem sog. „Pay Per Use“-Modell abgerechnet wird. Gemeint ist hiermit, dass der Kunde nur das bezahlt was er auch wirklich verbraucht hat, wobei die Verrechnung extrem granular nach genutzten Ressourcen und gelaufener Zeit erfolgt. Der Kunde bezahlt daher nicht für die benötigten Ressourcen, wie er es bei PaaS (wobei es hier variierende und bereits granularere Modell gibt) oder IaaS der Fall ist, sondern nur für die tatsächlich genutzten Ressourcen.

Grundsätzlich ist zwar der Preis der für das Laufen einer einminütigen Funktionsausführung verglichen mit dem Laufen des selben Codes auf einem von den Ressourcen her ähnlich bestückten PaaS-Server günstiger [Jon+19], jedoch der Anwendungszweck ein ganz anderer. Während mit PaaS Anwendungen bedient werden, die eine dauerhaft hohe Nachfrage erfahren, sollen mit Funktionen primär einfache stark frequentierende Aufgaben gelöst werden, bei denen der Nutzer nur für die Ressourcen bezahlen muss die die

Funktion bei dem Aufruf nutzt und nicht für jene die oben einkommende Aufrufe provisioniert wurden.

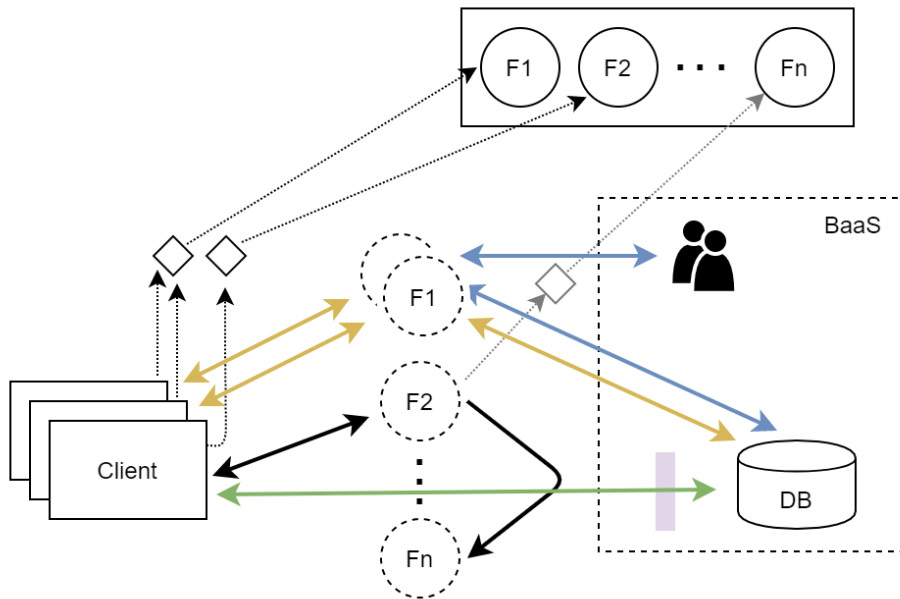
Ein weiterer wichtiger Punkt, welcher sich auch bei anderen Anbietern findet, ist in der Definition von OpenWhisk zum einen mit „[...] you can focus on building amazing and efficient applications [...]“ und zum anderen mit „developers write functional logic [...] in any supported programming language, that can be dynamically scheduled [...]“⁵. Ersteres meint die Abstraktion der Infrastruktur, welche dem Nutzer zum Bereitstellen von lauffähigem, skalierbarem, sicheren Backend-Code nicht bekannt sein muss. Die automatische Skalierung und optimale Verteilung von Ressourcen obliegt der Hoheit des Anbieters, und ist dem Nutzer nicht zugänglich, zumindest bei den proprietären Anbietern. Letzteres ist zwar keine einzigartige Eigenschaft von FaaS, da dieses Konzept bereits lange bekannt ist und auch bei der Microservice-Entwicklung häufig zum Einsatz kommt, jedoch ist die Auswahl aus vielen verschiedenen Programmiersprachen wie Java, Go, Javascript (Typescript), Python, PHP usw. (abhängig von dem jeweiligen Anbieter) unter dem Gesichtspunkt der direkten Nutzung zu sehen. Es muss keinerlei zusätzliches Setup oder andere Konfigurationen in der Umgebung vorgenommen werden, da dies der Anbieter übernimmt. Er kümmert sich um Patches und das Upgraden auf die aktuellste Version, was dem Entwickler ein großes Maß an Flexibilität bei der Entwicklung der Anwendung einräumt. Eine weitere häufig in der Literatur referenzierte Definition ist [Fow18], von Mike Roberts, welche die Beschreibung von AWS Lambda näher erläutert und im Großen die oben genannten Punkte wiedergibt. Natürlich sind diese Eigenschaften mit einem gewissen Vendor Lock-In dieser jedoch weder als schlecht noch gut zu bewerten ist und vielmehr nüchtern auf seine Vor- und Nachteile analysiert werden sollte. verbunden, da diese Entscheidungen über die Infrastruktur trifft, mehr dazu in *Potentiale und Herausforderungen*.

Abbildung 3 verdeutlicht hierbei noch einmal die Wirkungsweise von FaaS. Oben rechts ist der Funktionspool zu sehen, welcher bei dem jeweiligen Cloud Provider steht und in welchen der Entwickler seine Funktionen lädt. Wird eine Funktion von einem Client aufgerufen, so wird bei Aktivierung eines Triggers eine Kopie der jeweiligen Funktion instantiiert und alle mit dieser Funktion verbundenen Abhängigkeiten geladen (bspw. benötigte Module). Je nachdem wie frequentiert die Funktionen von den Clients angefordert

⁵<http://openwhisk.apache.org/>

werden, werden die vorhandenen Instanzen automatisch von dem Provider

Abbildung 3: Anwendung mit FaaS und (P)BaaS, angelehnt an [SKM20]



hochskaliert, um den eingehenden Anfragen gerecht zu werden. Hierbei kann die Funktion eine Aufgabe direkt ausführen und das Resultat an den Client zurückgeben, oder wie bei F1 auf die im hintergrund laufende BaaS Infrastruktur zugreifen. Diese reicht von einfachen Datenbankabfragen bis hin zu anlegen eines neuen Nutzers etc. Limitiert werden die Möglichkeiten hierbei lediglich von dem Service-Ökosystem des jeweiligen Anbieters. Genauso ist es natürlich möglich Funktionen sich gegenseitig aufrufen zu lassen oder direkt mit BaaS-Servicen zu interagieren, siehe *Abgrenzung von FaaS zu Serverless*.

Um FaaS besser und einheitlicher definieren zu könne und damit einen plattformübergreifenden Standart zu schaffen, wird in der Industrie und Forschung das Verlangen nach einer Referenzarchitektur, wie sie beispielsweise für das Cloud oder Grid Computing vorhanden ist [Liu+11], [FK03] lauter. Die Abwesenheit einer solchen Architektur behindere die Etablierung von Best-Practices, Design Pattern und einen genaueren Überblick darüber zu erhalten, wie sich das Feld rund um Function as a Service entwickele

[Lei+19]. Erste Vorschläge wie die Referenzarchitektur von FaaS aussehen könnte, haben sich mit [VE+19] durch die über die Jahre, von 2016 an, gestiegenen Popularität von FaaS in der Forschung und der Literatur bereits herausgebildet. War hier die Thematik bis 2016 nicht adressiert worden, so stieg die Präsenz von hier an kontinuierlich bis 2019, in Journals, Konferenzen und Workshops an [Yus+19].

2.3 Abgrenzung von FaaS zu Serverless

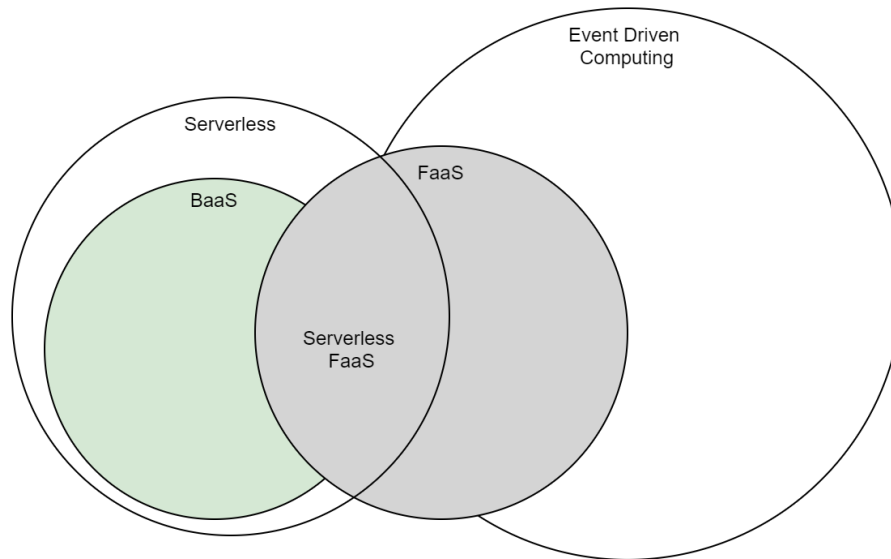
Serverless-Computing, oftmals auch als Serverless bezeichnet, ist eine Teildisziplin des Cloud-Computings, welche sich aus der Virtualisierung von Rechenleistung, Speicher und Netzwerken aus entwickelt hat [JC18]. Wie so häufig ist die Abgrenzung bzw. die Unterscheidung bei sich neu entwickelnden Technologien nicht ganz einfach. Zunächst stand Serverless für Anwendungen welche teilweise oder komplett auf Drittanbieter zurückgriffen, auf sog. Backends as a Service (BaaS) [siehe Abbildung. 4], um serverseitige Aufgaben wie Datenbankabfragen, Nutzerverwaltung o.ä. zu regeln [Fow18].

Mit FaaS wurde das Serverless-Konzept dahingehend erweitert, dass die Serverlogik nicht mehr vollständig von einer dritten Partei zur Verfügung gestellt wird, sondern vom Entwickler selber implementiert werden kann. Serverless ist dabei eines der in den letzten Jahren häufig genutzten Buzzwords in der IT, wobei der Begriff an sich etwas anderes suggeriert, als das was damit tatsächlich gemeint ist. Der Begriff impliziert die Abwesenheit von Servern, wobei damit lediglich eine Verschiebung der Zuständigkeiten einher geht. Entwickler einer „Serverlosen“-Anwendung bspw. mit FaaS, müssen sich nicht mit den operationellen Tätigkeiten wie dem *Provisioning*, Monitoring, der Wartung der Infrastruktur, der Skalierbarkeit dieser und der Robustheit des Systems befassen [Bal+17]. Jedoch geht mit der Abgabe an Zuständigkeiten auch ein gewisser Vendor Lock-In einher, wobei die Anbieter darauf achten, das mit FaaS eine möglichst große Zahl der in ihrem Service-Ökosystem vorhandene Services genutzt werden [KS18]. Es sollen möglichst viele der bereits im vorherigen Abschnitt vorgestellten Trigger bei dem Bau von Anwendungen genutzt werden.

Function as a Service ist somit zum Großteil aus dem Event-Driven Computing, welches vor allem bei der UI-Entwicklung genutzt wird, abgeleitetes Konzept, welches Serverless-Computing adaptiert hat. Nichts desto trotz

wird der Begriff Serverless in vielen Fällen, laut [Lei+19] in 58% der Fälle, mit FaaS gleichgesetzt, was eine Abgrenzung erschwert. Abbildung 4 verdeutlicht in welcher Beziehung die unterschiedlichen Konzepte stehen.

Abbildung 4: Serverless Concept, including FaaS and BaaS



Besteht eine Anwendung in der Folge lediglich aus den zwei „serverlosen“ Komponenten, FaaS und BaaS, so ist häufig die Rede von einer „pure serverless“, also einer Anwendung, deren operationeller Teil vollständig an einen Cloud-Anbieter ausgelagert wurde und einer „hybrid serverless“ Anwendung unterschieden [Lei+19]. Bei letzterem fungieren die Funktionen häufig als sog. „Glue“, der meisten zeitlich variierenden Aufgaben übernimmt, da IDLE-Zeit, wie man sie aus IaaS oder PaaS kennt, nicht berechnet wird.

Nun mag der ein oder andere argumentieren, dass ein Großteil der infrastrukturellen Aufgaben auch bei PaaS von dem Cloud-Anbieter übernommen werden, womit er natürlich recht hat. FaaS geht an dieser Stelle jedoch noch einen Schritt weiter. Bei PaaS werden vorgefertigte *Packages* einer Anwendung auf der Runtime der entsprechenden Plattform bereitgestellt, womit sich die Entwickler immer noch um die Anwendungsstruktur kümmern müssen [Kap19]. In FaaS hingegen wird dieser Teil übernommen und der Entwickler konzentriert sich lediglich auf die Business logik und die jeweiligen Event-

Trigger.

Zuletzt soll noch auf einen weiteren Punkt, „NoOps“, eingegangen werden, der oft mit Serverless gleichgestellt wird. Es wird zwar durch die Übernahme des Loadbalancing, der automatischen Skalierung, Sicherheitsaspekten und Patches ein Großteil der Wartung ausgelagert, jedoch ist dies nicht mit der Annahme des Wegfalls der Operations von DevOps gleichzusetzen, wie suggeriert wird [Fow18]. Es obliegt weiterhin den Entwicklern qualitativ hochwertigen Code zu schreiben und diesen in der Cloud-Umgebung zu testen, um dessen Performanz sicherzustellen. Die Komplexität entfällt daher nicht vollkommen, sondern verlagert sich zu einem gewissen Teil [Eiv17]. **Schaut man sich beispielsweise AWS an, so muss eine Vielzahl an Konfigurationen vorgenommen werden. Zwar kann später eine Vorhanden auf alle Funktionen angewendet werden, jedoch ist diese in der Folge von essentieller Bedeutung in einer riesigen „Multitenancy“-Umgebung.**

2.4 Potentiale und Herausforderungen

Die im Folgenden aufgelisten Potential und Herausforderungen, denen sich Literatur und Anwender von FaaS gegenübersehen, sind zu Teilen dem Konzept selber, zu Teilen aber auch der Implementation der jeweiligen Anbieter bzw. Open-Source Lösungen geschuldet. Es werden daher auf Seiten der Herausforderungen diese zunächst erläutert und falls in der Literatur bereits adressiert, vorgeschlagen Lösungen aufgezählt. Bei beiden, Potentiale und Herausforderungen, wird sowohl Bezug auf FaaS als auch auf BaaS genommen, da diese beiden Kategorien in den meisten Fällen in Kombination genutzt werden bzw. architektonisch beding in Kombination genutzt werden müssen [siehe Abbildung 4 und „statelessness“ Abschnitt *Function as a Service*].

Potentiale

Kosteneffizienz

Time to Market / Lead Time Development

DevOps

Testing

Optimale Auslastung von Rechenzentren Greener Computing [SKM20] [Fow18]

Herausforderungen

Vendor Lock-In Auch wenn diese Eigenschaft sowohl **Physische Lokation**

Eine weiterer Punkt ist die physische Lokation der Funktionen. Dadurch, dass es dem Provider obliegt, die Ressourcen seiner Infrastruktur optimal zu nutzen, entscheidet dieser auch auf welcher Node eine Funktion ausgeführt wird. Platziert der Provider in der Folge Funktionen, die eine hohe Datenabhängigkeit haben, physische weit voneinander entfernt, so wird sich dies auf die Performance der Anwendung auswirken und letztlich als Latenz bei dem Endverbraucher zu spüren sein [SKM20]. Es gibt zwar bereits bei vielen Anbietern, wie AWS, Microsoft oder Google die Möglichkeit die Region beziehungsweise ein sog. Cluster festzulegen, jedoch lässt sich damit die Lokation lediglich eingrenzen, aber nicht rechenzentren-genau festlegen. Mit der Performance von „serverful“-Applikationen kann damit nicht gleichgezogen werden [SKM20].

„Serverless“

Es ist an dieser Stelle zwar relativ trivial, trotz alledem ein Nachteil verglichen mit Lösungen wie PaaS oder IaaS. Die Rede ist von dem Verlust der serverseitigen Optimierung. Wie bereits durch den Namen deutlich gemacht, existieren die Server für den Nutzer nicht bzw. hat er bei einer auf BaaS basierenden Softwarelösung keine Kontrolle über das Backend. Er kann lediglich die Datenbanken, Objektspeicher oder Authentifizierung an die Struktur seiner Daten anpassen und über „Regeln“ Zugriffsbeschränkungen umsetzen. Die Services sind von dem jeweiligen Vendor vorgegeben und können bspw. nicht nach dem „Backend For Frontend“-Muster ⁶ auf die entsprechenden Client-Typen, wie Tablet, Mobile-Phone oder Desktop angepasst werden.

⁶<https://samnewman.io/patterns/architectural/bff/>

Es existieren von jeder Sorte (Datenbanken, Objektspeicher, Authentifizierung usw.) nur eine Version auf welche man beschränkt ist. Jegliche benutzerdefinierte Logik muss daher in den Client ausgelagert werden, da es nicht auf backendseitig implementiert werden kann [Fow18]. Mit FaaS kann dieser Effekt gemindert werden, indem ein leichtgewichtige Logik in form von Funktionen serverseitig umgesetzt wird. Hierbei erfolgt die Interaktion mit den unterschiedlichen BaaS-Servicen über die vom Anbieter beschränkten Trigger [siehe Abschnitt *Function as a Service*].

„Statelessness“

Mit der „statelessness“, also der Zustandslosigkeit von Funktionen, geht eine weitere Problematik von FaaS einher. Der Kern von Anwendungen ist es Aufgaben zu lösen, wofür viele verschiedene Schritte von Nöten sind, welche sinnvoll miteinander verbunden werden müssen. So ist es auch bei dem funktionsbasierten Aufbau von Applikationen wichtig, dass diese untereinander kommunizieren können und den Programmstatus voneinander abfragen können, ohne dass es zu Inkonsistenzen kommt. Bedingt durch die kurze Laufzeit der Funktionen ist es essentiell, dass auch der Austausch in entsprechender Geschwindigkeit erfolgt. Wie in einem Report von Berkeley [Jon+19] herausgearbeitet, erweist sich das schnelle und exakte „State-Sharing“ jedoch immer noch als problematisch dar, betrachtet man die Geschwindigkeit von „Serverless“-Anwendungen im Vergleich mit „Serverful“-Anwendungen.

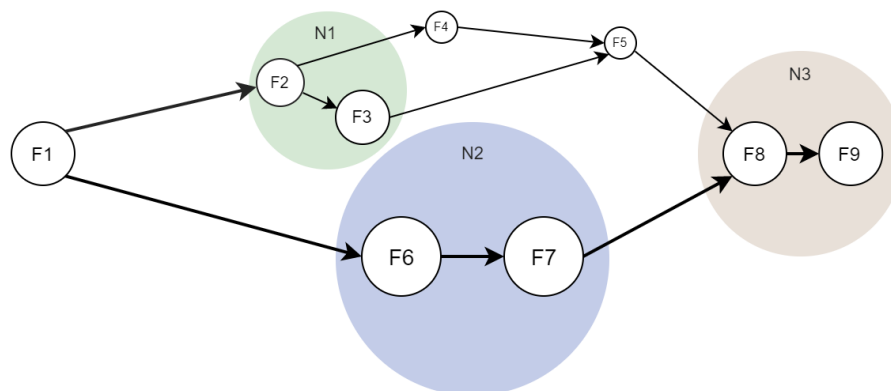
Grund hierfür sind die von den Anbietern zur Verfügung gestellten BaaS-Lösungen für das persistente Speichern des Programmstatus. Objektspeicher der verschiedenen Anbieter (AWS S3, Google Cloud Storage, Azure Blob Storage) sind zwar nicht teuer sehr teuer wenn es um das Speichern mehrerer GB geht, jedoch sind die Kosten beim Zugreifen auf die Speicher hoch und die Latenzzeiten von bis zu 20ms zu hoch [Jon+19]. Die Key-Value-Speicher der Anbieter sind in diesem Falle die bessere Wahl, da ihre Ansprechzeiten mit 8 - 15ms geringer sind, jedoch sind sie, bezogen auf ihre Input/Output Operationen pro Sekunde, deutlich teurer als die der Objektspeicher.

Neben dem schnellen Informationsaustausch der Funktionen mit einem externen Speichermedium, kann es bei einer Kopplung von Funktionen aus performancetechnischen Gründen sinnvoller sein, den aktuellen Programmstatus direkt an die nächste Funktion weiter zu geben. Wie zuvor bei der *physischen Lokation* der Funktionen erwähnt, ist es wahrscheinlich, dass

Funktionen nicht auf der selben Node gestartet werden. Der Loadbalancer des jeweiligen Providers entscheidet je nach Auslastung der Infrastruktur, auf welcher Node eine Kopie der Funktion ausgeführt wird. Auch wenn im Idealfall die Weitergabe des „application states“ schneller ist als der Zugriff auf ein externes Speichermedium, so wird dies durch „Startup-Latencies“ und physischer Distanz wieder relativiert. Um dem entgegenzuwirken schlagen Shafiei et. al. [SKM20] vor, mit einem *Dependency-Graphen* miteinander verwobene Funktionen zu kennzeichnen und direkt bei dem initialen Aufruf einer Funktion mögliche Folgefunktionen zu starten.

Würde man dem Nutzer die Möglichkeit geben zu bestimmen auf welcher Instanz (Node) eine Funktion laufen soll, so das Konzept von FaaS damit untergraben werden (siehe *Optimale Auslastung*). Es würde dazu führen, dass erneut Kapazität zurückgehalten wird, über welche der Provider nicht mehr frei verfügen könnte.

Abbildung 5: Dependency Graph nach [SKM20]



Wenn jeder Nutzer bestimmen könnten auf welcher Node, bspw. Node-X, seine Funktionen laufen sollen, um den Programmstatus mit so wenig physisch bedingter Latenz wie möglich weiterzugeben, so wäre eine optimale Auslastung nicht mehr garantiert. Der Provider müsste stets Kapazität von Node-X zurückhalten und könnte sie nicht für andere Nutzerfunktionen freigeben. Er muss stets einen Teil der Kapazität von Node-X zurückhalten, für den fall dass inaktive, aber der Node zugeteilte Funktionen aufgerufen werden. Mit einem *Dependency-Graphen* wie in Abbildung 5 zu sehen würde dieses Problem umgangen werden. Der *Load-Balancer* des Anbieters könnte

unbeeinträchtigt die Funktionen wie gewohnt je nach Auslastung auf freie Nodes verteilen, da er die korrelierenden Funktionen „gleichzeitig“ starten kann. Da die Latenzen beim Starten eines Services erheblich höher sind, wobei die Programmiersprache dies zusätzlich beeinflussen kann [Man+18], als jene durch die Lokation bedingt [Adi+19] [JC18], würde deren Eliminierung die Weitergabe des *States* erheblich beschleunigen.

Generell betrachtet könnte die Geschwindigkeit der Kommunikation von Funktionen untereinander beschleunigt werden. Verbindet man das ganze noch mit einem einfachen neuronalen Netz, um beispielsweise vorherzusagen mit welcher Wahrscheinlichkeit Funktion zwei (F2) im Vergleich zu Funktion sechs (F6) angesprochen wird, so könnte die Performance sogar noch ressourcenschonender auf Seiten des Cloud-Providers gestaltet werden. Zunächst käme es entweder zu einer höheren Belastung, sollte man zu Beginn alle Funktionen starten oder zu einer höheren Latenz, sollte man die Funktionen einzeln starten und aus dem sich ergebenden Verlauf lernen. Nach einer kurzen Lernphase würde die Auslastung der Infrastruktur jedoch erhöht werden.

„Multitenancy“

Obgleich das Konzept von FaaS darauf beruht, dass der Anbieter eine Vielzahl von Funktionen verschiedener Kunden steuert und pflegt, so soll sich jeder Kunde fühlen, als sei der der Einzige. Allerdings stellt dies die Betreiber der Plattformen vor weitere Herausforderungen, da sie die Ressourcenisolation sowie Zugriffsbeschränkungen zuverlässig gewährleisten müssen. Dies ist jedoch leichter gesagt als getan, da es teilweise nicht nur von dem Provider selber, sondern auch von der richtigen Konfiguration der Sicherheitseinstellungen der Nutzer ⁷, ⁸ abhängt, ob die Daten unzugänglich für Dritte sind.

Obwohl Anbieter wie AWS mittlerweile erfahren genug sein sollten, dass solche Probleme nicht mehr zu erwarten sind [Fow18], ist weiterhin Vorsicht geboten. Es nichts an der Tatsache, dass auch sie sich mit Sicherheitsaspekten, der Robustheit und Performance ihrer Infrastruktur ständig weiter auseinandersetzen müssen. Kein Kunde darf die Daten eines anderen sehen, kein Fehler die Stabilität anderer Funktionen gefährden und kein plötzlich auftretender *Spike* eines Kunden die Performance der Funktionen eines anderen beeinträchtigen. Das Zusammenspiel von RPCs und der Con-

⁷<https://awsinsider.net/articles/2017/06/20/voter-data-leak.aspx>

⁸<https://searchsecurity.techtarget.com/news/450422962/Another-AWS-cloud-data-leakage-due-to-misconfiguration>

ainersicherheit muss dauerhaft gegeben sein und durch sorgfältiges *Security Managament* sichergestellt werden [MB17].

Testing

Cold-Starts

Cold Starts beziehen sich auf das Starten eines Containers, in welchem eine Kopie der benötigten Funktionen zur Verfügung gestellt wird. Ist eine Funktion lange nicht genutzt worden oder wird zum ersten Mal angesprochen, so muss für die Ausführung zunächst alles vorbereitet werden. Vor allem bei sehr kurzen Funktionen stellen sie ein erhebliches Problem dar, da ihre *Startup*-Zeit bis zu ein Zehnfaches der eigentlichen Ausführungszeit einnehmen kann [SBW19]. Ein Container muss auf einer Node gestartet werden und die Funktion, sowie deren Abhängigkeiten geladen werden. Die Zeit die es dauert, bis die Instanz bereit ist hängt dabei von verschiedenen Faktoren, wie der Programmiersprache, der benötigten Bibliotheken (*Libraries*), der Größe der Funktion (menge an Code) und der Hardware [SKM20] [Jon+19] ab, auf der die Funktion instantiiert wird. Wird eine Funktion dabei sehr frequentiert aufgerufen, so halten die Plattformen diese Funktionen für eine gewisse Zeit vor („warm“, bei AWS bis zu einer Stunde [RC17], wodurch die Wahrscheinlichkeit eines erneuten *Cold Starts* gegen Null tendiert. Bei Funktionen die hingegen nur einmal pro Stunde bemüht werden, sind die Folgen von *Cold Start* deutlich häufiger zu spüren [RC17].

Um dem entgegenzuwirken schlägt [Cas+19] vor, eine Art stammzellen Container vorzuhalten. Diese würden für die unterstützten Programmiersprachen bereits instantiiert aber leere Container vorhalten. Diese Container sind nicht Kundenspezifischen und generisch nutzbar. Damit würden zwei der obene genannten Punkte, die Programmiersprache und die Hardware, welche die *Startup*-Zeit der Funktionen negativ beeinflussen, eliminiert werden. Lediglich die größe der Funktionen und deren *Dependencies* müssen noch geladen werden. [IMS18] hat zudem gezeigt, dass mit der Erhöhung der Speicherkapazität die Zeit bis eine Funktion bereit ist einen *Request* zu bearbeiten verringert werden kann. Kombiniert man dies mit dem Stammzellen-Ansatz, so wird das Laden des Programmcodes und der *Libraries* beschleunigt. Natürlich kommt der zusätzliche Speicher nicht umsonst, jedoch liegt es in dem Fall bei dem Unternehmen zwischen den Mehrkosten und dem Performancezuwachs abzuwägen.

3 Guidline on Migration

The below presented guideline for migrating parts of an existing microservice architecture to the relatively new cloud-computing concept, Function as a Service, is structured as follows. At first the decision of migrating will be questioned, by challenging the motives of adopting this new technology and giving advice on when to migrate. Afterwards criteria for selecting the right service will be stated, which gear towards the possibilities that come along with FaaS. When a reasonable service has been found a provider needs to be identified which best meets the companys requirements. There are two types of providers that primarily differ in the service-ecosystem inherent to them and the degree of vendor lock-in. Whereas proprietary platforms like AWS, Azure or Google provide a great ecosystem of feature services and backend as a service offerings, open source frameworks like OpenWhisk, Fission or Kubeless dedicate greater control to the developers. When implementing FaaS, the open source solution can be far more adjusted and configured to best suit an organisation needs. Next, impacts on the general process of software development, the „Dev“-part in *DevOps*, will be investigated. Changes on organisational structures regarding the development team, as well as effects on the CI/CD pipeline, challenges faced with local testing and the VCS will be addressed. With the development process covered, the guideline on migration will end with the operational part which will be taken into account. Alternations on the process of monitoring, maintenance and testing instances will be covered and advice on coping with them will be given.

In order to better follow along the guideline, the subsequent section will give an overview of the prevailing microservice infrastructure. The development process as well as the operational tasks will be described, along with the underlying service-landscape.

3.1 Zu grunde liegenden Architektur

Um im weiteren Verlauf dieser Arbeit Aussagen über die nötigen Anpassungen und Auswirkungen auf das Unternehmen, die Entwicklung (Dev) selber, sowie die Kultur und den operationellen Teil (Ops) treffen zu können, soll zunächst der momentan vorliegende Aufbau beschrieben werden. Hierbei handelt es sich einfach gesagt, um die Entwicklung von Microservices, wobei momentan 31 Services in unterschiedlichen Versionen und Testinstanzen vorliegen.

3.2 Challenging the Decision

Whereas migrating from a monolithic application to serverless, respectively function as a service, seems far more challenging, due to the assumed size and complexity of the application, migrating from a microservice-architecture can be challenging as well. With monolithic application, functionalities need to be identified in the first place and afterward be broken down into many small functional sections. Microservices, on the other hand, might already be small enough that they supply only certain functionalities but in some cases they do incorporate too many functionalities to be converted to functions right away. Nevertheless dealing with the issue of outsourcing the currently inherent state of services, as well as reducing dependencies and optimizing code remains.

Before starting to introduce this new technology, the current state of the architecture as well as the desired state of the system need to be untangled. When a system is already composed of many small services, it is very likely that the underlying infrastructure consists of VMs or containers. Those containers often run in an PaaS environment and it is not possible to scale them to zero to free reserved computing capacity, in favour of other services to consume it. To do so or even to outsource the expense of provisioning services, managing the underlying infrastructure and coping with operational tasks, such as monitoring, load balancing etc., Function as a Service can help to achieve this desired state. On the other hand the restrictions that come along with the ease of development have to be traded off against its benefits.

In order to profit as much as possible of introducing FaaS, the services, which should be migrated, have to be inspected regarding dependencies, size, version, language and complexity. Starting with its dependencies, it is recommended to reduce them to a minimum, due to their effect on the start-up time of the container the function runs in [Man+18]. Especially with Python, Nodejs and Java [PS17] loading all dependencies required, which might again interconnect with other dependencies, will have an impact on cold starts. The next issue which need to be addressed is state management. As mentioned in section two, *Benefits and Drawbacks*, there is no persistent state an application can rely on. There actually is persistent state in the container of a function, but whether an incoming request will hit a certain, running container, is unpredictable. When a service does not receive any request or is running for a long period of time, the provider will eventually

kill its container to free capacity. In the case of AWS Lambda the current maximum amount of time, of a frequently called function until it will be killed, is 45 minutes⁹. Later on, when comparing cloud providers with open source frameworks, there will be guidance on finding an appropriate solution to meet a company's purposes but for now the following must be considered. When using an open source framework like OpenWhisk, the external state management system can be covered with Redis or another low latency database, whereas using a cloud provider platform, one of its inherent database solutions will presumably be the best choice. If it is necessary to have full control over the performance and configuration of the database, an open source framework has to be chosen.

Latency and frequency of a service will also decide upon its aptitude of being a candidate to be migrated to FaaS. If the service is called very frequently and experiences most of the time a very high traffic, the concept of Function as a Service will not be applicable to it [Jon+19]. Due to a limit of concurrent running functions, that can vary between the different providers and frameworks, traffic can only be handled up to a certain amount. Moreover, the same application experiencing the same amount of traffic, once running on a FaaS platform and once running in a Docker container in a PaaS environment, will be more expensive implemented with FaaS, than it will be with PaaS [Jon+19]. Therefore, services with various workloads, having eventually high peaks and then some time of inactivity, are more applicable to the concept of serverless, than their counterparts. Latency should also not be a critical component due to the before mentioned cold starts. Latency sensitive applications like trading platforms, which strongly rely on real time data, are not a suitable candidate for Function as a Service.

Another pitfall is the promise of not having to maintain, provision, scale or monitor the infrastructure and thereby reducing complexity and operational tasks. On one side, the cloud provider will, to some extent, take care of load balancing, scaling functions up and down and providing monitoring solutions. On the other side, new complexity will appear in other areas. The issue of mono-repo and poly-repo arises, as discussed in a later section, as well as the necessity to learn the providers DSL, acquire the skills to correctly configure security policies, interconnect functions, define different triggers on functions, dealing with concurrency settings¹⁰ and many more

⁹<https://aws.amazon.com/de/lambda/>

¹⁰<https://docs.aws.amazon.com/lambda/latest/dg/configuration-concurrency.html>

aspects will add complexity in new fields. If the company aims at reducing infrastructural complexities, then function as a service can support at the listed tasks, but when its main purpose is to reduce operational efforts, as often suggested with the term „NoOps“ [Eiv17], FaaS, again, is not the right choice.

Lastly attention needs to be drawn to the programming language and lead time development. FaaS can provide a reduction in lead time development and time to market [SS18] [Lei+19], thanks to its small code base and rapid development. Developers are able to make changes, fix issues, test and finally deploy a new version to production in a few minutes. Prerequisite, even though this might sound trivial, is the platform's language support for the language used in the application or service. Despite the fact that many programming languages are already supported by the big cloud providers, application logic might have to be rewritten, if that specific language is not supported. The code also needs to be rewritten in case of version incompatibility, when the version of the application is ahead of the version offered by the provider.

3.3 Choosing a suitable service

3.4 Auswirkungen auf die Team-/ Entwicklungsstruktur

Agile / Dev Ops

Es gilt abzuwägen, da verschiedene Arten unterschiedliche Vor- und Nachteile haben. Review (Peerreview bleibt gleich, auch dependency management, da es von dem Kontext und der Sprache abhängt)

Deployment / CI CD:

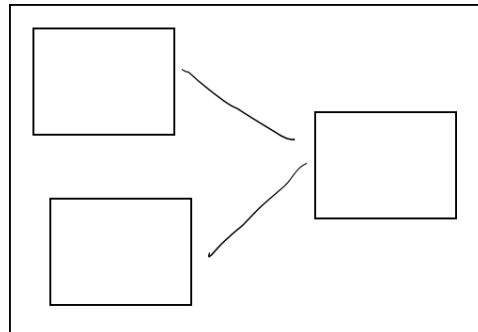
Vorsicht bei der Anzahl der Funktionen geboten, da sie sehr klein sind, Gefahr höher dass manche Funktionen dupliziert werden, daher vielleicht eher Monorepo. Zudem leidet dann die Wartbarkeit des Codes darunter, da manche Funktionen mehrfach in leicht unterschiedlicher Weise vorliegen. Falls von den ähnlichen Funktionen auch noch unterschiedliche Versionen existieren, so kann dies die Wartbarkeit zusätzlich erschweren [Rac+19].

Adopting serverless requires a different mental model, where systems are primarily constructed by composing pre-existing services, with FaaS often acting as the “glue” that brings these services together. FaaS offers technical and business-related advantages, but managing and predicting deployment costs is difficult for larger applications. Further, tooling availability and maturity, especially related to testing and deployment, remains a barrier to entry. Finally, limi-

ted support for function sharing and the absence of a service ecosystem is seen as a challenge. [Lei+19]

Um ein Aussage darüber treffen zu können, ob und in welchem Maße Anpassungen an der Team- bzw. Entwicklungsstruktur vorgenommen werden müssen, gilt es zunächst zu klären, wie die zu grunde liegende Unternehmensstruktur aufgebaut ist. Da diese Arbeit von einer bereits existierende Microservice-Architektur ausgeht, soll zunächst im Folgenden der organisatorische Aufbau dieser beschrieben werden. Im Nachhinein wird dieser noch einmal mit den entsprechenden Änderungen aufgegriffen, um eine bessere Aussage über die tatsächlich notwendigen Anpassungen treffen zu können. Durch das neue Programmier-Paradigma wird den Programmierern automatisch ein neuer Aspekt auferlegt. Dadurch, dass die Funktionen einzeln nach Laufzeit und bezahlt werden, liegt es im interesse des Programmierers diese so performant wie möglich zu gestalten, um bei jedem Funktionsaufruf Ressourcen und Laufzeit einzusparen [SKM20].

Abbildung 6: Referenzarchitektur Microservice-Entwicklung



[VE+19] [Yus+19] [LSF18] [aws2020ManagingFunctions]

Um provisioned concurrency festzulegen, sollte nachgesehen werden, wie häufig ein Funktion bzw. ein Service aufgerufen wird, und dementsprechend geschaut werden, wie viele Funktionsinstanzen immer vorrätig sein sollten. zudem sollte geprüft werden, wie lange ein Funktionsaufruf von dem Initialen Event bis zu seiner Response braucht, um zu schauen in wie weit die Performance durch die Nutzung einer alternativen Programmiersprache gesteigert werden kann. Eine weitere Möglichkeit bietet AWS mit Application Auto Scaling welches provisioned concurrency erweitert ... **Target tracking scaling policy, welches provisioned concurrency automatisch anpasst, je nach Nutzungsmetrik**

3.5 Abwägung Open-Source und Cloud-Vendor

Wie bereits im vorherigen Abschnitt beschrieben, OpenFaaS

Die Abwägung zwischen der Implementation von Function as a Service in der unternehmenseigenen Cloud bzw. Servern oder dem Outsourcen an einen proprietären Cloud-Vendor, sollte wohl überlegt sein, da hiervon in der Folge eine Vielzahl an Möglichkeiten und Restriktionen abhängt.

Entscheidet man sich für ersteres, also dem Aufbau einer privaten FaaS-Plattform so stehen hierfür, mit IBM Apache OpenWhisk, Fission, OpenFaaS oder Kubless, nahezu identisch viele Frameworks zur Verfügung wie bei den öffentlichen Vendors. [Schaubild CNCF Serverless Landscape ggf. einfügen](#). Des Weiteren bietet dieser Weg, abgesehen von der Vielfalt an Frameworks mit wiederum unterschiedlichen Eigenschaften, die Möglichkeit die Größe der Funktionen oder die maximal zulässige Laufzeit eines Containers individuell anzupassen. Dies würde einerseits eine granularer Verrechnung der in Anspruch genommenen Ressourcen der einzelnen Bereiche oder Teams zulassen, zugleich aber auch ein Operations-Team verlangen, welches sich in das jeweilige Framework einarbeitet, die Plattform zu dessen Betrieb aufsetzt und die spätere Wartung dieser übernimmt. [MPDF+18].

Eine Alternative stellen proprietäre Lösungen von öffentlichen Cloud-Vendors, wie beispielsweise Amazon Web Services (AWS) Lambda, Microsofts Azure Functions, IBM Cloud Functions oder Google Cloud Functions, bei welchen sich von Unternehmensseite aus niemand um die Wartung der Infrastruktur, die Skalierung der Services, die Behandlung von Fehlermeldungen o.ä. kümmern muss. Diese Aufgaben werden in der Folge von dem Plattformbetreiber übernommen. Natürlich sind diese Dinge in erster Linie Aufgabe des Plattformbetreibers, jedoch haben sie unmittelbare Auswirkungen auf das Unternehmen, sollten Probleme beim Skalieren von Funktionen oder dem Monitoring auftreten. Sollte dies der Fall sein, so sind die Entwickler von den bereitgestellten Debugging Möglichkeiten und Monitoring-Lösungen der Plattform abhängig um Probleme schnellstmöglich zu beheben, sollte dies die Plattform nicht tun. Daneben ist ein weiterer Punkt des Vendor Lock-Ins die von Anbieter zu Anbieter variierende Infrastruktur, welche ein einfaches Shiften von Funktionen erheblich erschwert. Zudem sind die Benutzung der Funktionen meist automatisch mit der Inanspruchnahme weiterer Services der Plattform, wie dem Message Queuing oder der Datenspeicherung, gekoppelt.

Dies hat sowohl Vor- als auch Nachteile. Ist man bei der Einrichtung der privaten FaaS-Cloud auf die unternehmensinternen Ressourcen beschränkt, so bieten die Cloud-Anbieter, neben den Restriktionen des Vendor Lock-Ins, in den meisten Fällen ein großes Ökosystem an weiteren Services, welche sich problemlos an die Funktionen anbinden lassen. Im Folgenden wird daher, auch wenn sich diese Arbeit hauptsächlich auf FaaS beschränkt, das Ökosystem der einzelnen Cloud-Anbieter kurz betrachtet, um einen besseren Überblick über deren zusätzliche Leistungen zu erhalten und eine fundierte Entscheidung treffen zu können.

3.5.1 Vendoren-Analyse

Ein weiterer Punkt bei der Migration von Teilen einer bestehenden Architektur, in diesem Fall einer Microservice Architektur, ist die von dem Cloud-Vendor zur Verfügung gestellte Möglichkeit der Orchestrierung der Funktionen. Diese entscheidet darüber wie performant die Funktionen parallel ausgeführt werden können und wie groß der Runtime-Overhead der jeweiligen Plattformen dabei ist. Zudem spielt der Preis der hierfür anfällt eine nicht unerhebliche Rolle, da die sequentielle als auch die parallele Ausführung von Funktionen bei mittleren und großen Softwareprojekten häufig auftreten.

Wie bereits in *Potentiale und Herausforderungen* unter *Statelessness* erwähnt spielt bei der Kopplung von Funktionen die Weitergabe des Anwendungs- bzw. Funktions-*States* und dessen Geschwindigkeit eine wichtige Rolle. Nur durch das richtige Zusammenspiel vieler Funktionen ist es möglich große Anwendungen zu bauen und komplexe Abläufe umzusetzen. Bei der Betrachtung der unterschiedlichen Hosting-Lösungen soll der Fokus auf die beiden am häufigsten genutzten proprietären FaaS-Lösungen [Lei+19], AWS Lambda mit Amazon Step Functions und Microsoft Azure Functions mit Azure Durable Functions, gelegt werden. Die Seite der Open-Source Lösungen wird von IBM OpenWhisk, mit IBM Composer als Orchestrierungs-Lösung, vertreten.

[Lóp+18] hat diese drei Frameworks genauer untersucht und bei der sequentiell und parallele Planung sowie der Weitergabe des *States* der drei Orchestrierungs-Tool erhebliche Unterschiede feststellen können. Vorab soll an dieser Stelle vermerkt werden, dass nur „warme“ Instanzen für jegliche nachfolgende Tests genutzt wurden, um Ungenauigkeiten, durch variierende *Cold-Start* Zeiten [siehe [Man+18] und [JC18]], vorzubeugen.

In Bezug auf die Ausführung von aufeinander folgenden Funktionen, sequentielle Verarbeitung, erwiesen sich IBM's Composer und ASF als deutlich schneller im Vergleich zu ADF. So Betrug der Overhead, also die Zeit welche nicht zur Ausführung der Funktion verwendet wurde, bei 40 hintereinander geschalteten Funktionen 1,1s für IBM Composer und 1,2s für AWS Step Functions. Azure Durable Functions brauchten hingegen für die selben 40 Funktionen ganze 8s. Bei weiteren Durchführungen mit [5, 10, 20, 40, 80] stellte sich jedoch heraus, dass IBM Composer die Orchestrierung von Funktionen nur bis zu einer Anzahl von 50 Stück unterstützt. Alles was darüber hinausgeht, müsste durch Orchestrierungs-Tool von Drittanbietern übernommen werden [Lóp+18]. ADF und ASF sind wiederum in der Lage *Workflows* festzulegen, welche über Tage und Monate lauffähig sind.

Die Evaluierung des Overheads bei parallel geschalteten Funktionen wurde für ASF und ADF durchgeführt. Dabei wurde wieder mit 5 Funktionen begonnen und sich, wie oben beschrieben, bis auf 80 hochgearbeitet. Die Ergebnisse waren eindeutig. Bei einer Anzahl von 80 Funktionen hatten Azure Durable Functions mit einem durchschnittlichen Overhead von 32.1s fast das doppelte Volumen von AWS Step Functions mit einem durchschnittlichen Overhead von 18.3s. Die Ergebnisse legten zudem nahe, dass ASF zuverlässiger bei der Vorhersage des zu erwartenden Overheads ist als ADF. Microsofts Overhead stieg nicht immer gleichbleibend exponentiell an wie der von Amazon, was eine Prognose über das Verhalten erschwert.

Bei der Evaluierung für die Eignung von parallel geschalteten Funktionen viel IBMs Composer direkt zu Beginn aus dem Testportfolio heraus, da parallele Ausführungen nicht unterstützt wurden [stand 2018 [Lóp+18]]. Mittlerweile [stand 2020] wird von Seiten IBMs die parallele Ausführung von Funktionen seitens Composer zwar unterstützt und gesagt, dass seitens Composer dies nicht auf eine bestimmte Anzahl von Funktionen beschränkt ist ¹¹. Allerdings wird explizit erwähnt, dass eine Limitierung gleichzeitig ausführbarer Funktionen Seitens OpenWhisk besteht, welche bei Überschreitung zu Fehlern führen kann: „[...] many concurrent invocations may hit OpenWhisk limits leading to failures: failure to execute a branch of a parallel composition or failure to complete the parallel composition [...]“ ¹². Die derzeitige Limitierung gleichzeitiger Funktionen in OpenWhisk liegt bei 100 pro

¹¹<https://github.com/apache/openwhisk-composer>

¹²<https://github.com/apache/openwhisk-composer>

Namespace ¹³.

Zuletzt wurden die drei Orchestrierungs-Lösungen bezüglich der Weitergabe des *Application-States* untersucht. Aufgrund der Begrenzung von ASF auf 32KB wurde bei den beiden anderen Lösungen die selbe Größe gewählt. Dieses mal wurden nur 5 sequentiell ablaufende Funktionen getestet. Die Grenze bei IBM Cloud Functions lag 2018 offiziell bei 1MB, liegt mittlerweile aber bei 5MB ¹⁴. ADF hingegen ermöglicht die Weitergabe von bis zu 60KB. Es zeigte sich, dass IBM Composer und AWS Step Functions bei der Ausführung ohne *Payload*, jeweils einen Overhead von 175.7ms und 168.0ms hatten. Mit Payload betrug der Overhead in ms für Composer 298.4 und Step Functions 287.0, was eine Zunahme von 70% darstellt [siehe Abbildung 7]. Azure Durable Functions stieß bei diesem Test deutlich hervor. Mit einem Overhead von 766.2ms ohne Payload und 859.5ms mit Payload ist der grundlegende Overhead zwar deutlich höher als bei den beiden vorherigen, steigt unter Last aber nur um 12% an [Lóp+18].

Abbildung 7: Overhead bei 5 sequentiellen Funktionen mit einem Payload von 32KB, nach [Lóp+18]

Platform	Overhead (ms)		Increase (%)
	Without payload	With payload	
<i>IBM Composer</i>	175.7	298.4	70%
<i>AWS Step Functions</i>	168.0	287.0	71%
<i>Azure DF</i>	766.2	859.5	12%

Anhand der durch Lopez et. al. gezeigten Ergebnisse lässt sich festhalten, dass AWS mit Step Function das ausgereifteste Orchestrierungs-Tool zur Verfügung stellt. Sowohl bei der sequentiellen als auch bei der parallelen Ausführung bietet AWS Lösungen für langlebige und kurzlebige Funktionskopplungen. Zudem ermöglicht die Limitierung des *States* auf 32KB eine klare Auskunft über die entstehenden Kosten zu geben. Hat man vor FaaS lediglich für leichtgewichtige Aufgaben zu nutzen, so spielt IBM Composer

¹³<https://github.com/apache/openwhisk/blob/master/docs/>

¹⁴<https://cloud.ibm.com/docs/openwhisk?topic=openwhisk-limits>

bei der Kopplung von bis zu 50 Funktionen seine Stärken aus und ist dabei geringfügig schneller als AWS. Für die Festlegung längerer Abläufe, welche auf das Laufen über Tage bis hin zu Monaten ausgelegt sind, ist Composer nicht geeignet [Lóp+18]. Steht der Austausch des *Application-States* zwischen Funktionen im Vordergrund und können höhere Latenzzeiten in Kauf genommen werden, bietet Azure mit einer Kapazität von 60KB eine Alternative zu AWS und IBM. Auch Abläufe die länger Zeit beanspruchen können mit ADF konfiguriert werden. Das Konzept von *async/await* bei sequentiellen und *fan-out/fan-in* ¹⁵ bei parallelen Abläufen bietet zusätzlich eine etwas leichtere Umsetzung als AWS und IBM [Lóp+18]. Kein Garantie, dass der folgende Funktionsaufruf auf die selbe bereits laufende Instanz einer Funktion trifft, welche dann wiederum auch Zugriff auf dem im Memory gespeicherten State hat. Es ist daher notwendig, dass der Application-State extern gespeichert wird. Auf diesen sog. "*Shared memory*" müssen in der Folge alle Funktionen drauf zugreifen, wenn Applikations-Daten benötigt werden.

3.5.2 Tools zur Entwicklungsunterstützung

Terraform

IBM Composer, Amazon Step Functions, Azure Durable Functions Serverless Framework

Chalice SAM Serverless Framework

Results

3.5.3 Cross-Functions in Multi-Cloud Lösung

CloudFormation, Terraform or Serverless Framework tool [IS18]

3.6 Auswirkungen auf die Testumgebungen

Zum Testen einer serverlosen Applikation stehen mehrere Möglichkeiten zur Verfügung, welche ihre Vor- und Nachteile haben. Unterschieden wird in lokales Unit-Testen, Canary Release Testen und A/B Testen und Integrationstests. War OpenWhisk zu Beginn einer der ersten Anbieter, welcher lokale Unit-Tests unterstützte, so zogen AWS mit SAM ¹⁶ und Azure mit seinen sog. Function Core Tools ¹⁷ nach. In diesen Unit-Tests liegt jedoch

¹⁵<https://docs.microsoft.com/en-us/azure/azure-functions/durable/durable-functions-cloud-backup>

¹⁶<https://aws.amazon.com/de/serverless/sam/?nc1=his>

¹⁷<https://docs.microsoft.com/en-us/azure/azure-functions/functions-develop-local>

bereits eine Einschränkung, die bei der Migration bedacht werden muss. Gibt es bei Anwendungen die nicht über einen *Serverless Cloud-Provider* laufen oft die Möglichkeit Teile der Anwendung durch lokale Kopien von Datenbanken oder Message-Queues, welche denjenigen in der Produktion sehr ähnlich sind, bei dem Testen zu integrieren, so ist dies bei serverlosen Funktionen schwerer [RC17]. Dadurch, dass die Funktionen in den meisten Fällen mit anderen Funktionen sowie Datenbanken und ggf. noch einigen weiteren Services des Anbieters interagieren müssen, ist es nicht möglich dies lokal zu simulieren, zumal in der Plattform gesetzte Konfigurationen lokal nicht umgesetzt werden können. Kennzahlen wie die Ausführungszeit von Funktionen, das Laden von Abhängigkeiten (Libraries) und Verzögerungen durch *Cold-Starts* können nicht akkurat wiedergegeben werden [Rac+19]. Mit der Abgabe der Hoheit über die Infrastruktur kann hinzukommend auch nicht mehr der Server bestimmt werden, auf dem die Funktion letzten Endes gestartet wurde, was im Zweifel, bei älterer Hardware, die Ausführungszeit der Funktion erhöhen kann. Es ist daher unumgänglich den Service auch auf der Plattform selber zu testen, um zuverlässige Daten über Performanz und Kompatibilität zu erhalten.

Eine Möglichkeit dafür bietet das sog. Canary-Testen, bei welchem es nicht erforderlich ist die komplette Servicelandschaft auf einen DEV-Account zu spiegeln. Stattdessen wird eine neue Version einer bereits existierenden Funktion in Produktion geladen und nur ein bestimmter Teil der Nutzer bzw. Tester darauf umgeleitet. In der Praxis werden *Canary-Testing* sowie *A/B-Testing* aber nicht sehr häufig durchgeführt [Lei+19], da dies negative Auswirkungen auf die Performance der in Produktion laufenden Anwendung haben kann. Vor allem Lasttests können dafür sorgen, dass es zu spürbaren Latenzen bis hin zu Ausfällen bei den Nutzern kommt.

Es bietet sich daher an, den Aufwand einer Spiegelung zu betreiben, da so die in Produktion laufenden Funktionen nicht von Tests beeinflusst werden. Viele Anwender haben dies bereits umgesetzt und es scheint sich als „Best Practice“ plattformunabhängig etabliert zu haben [Lei+19]. An sich ist dies jedoch nicht verwunderlich, da FaaS bzw. Serverless bei diesem Vorgehen seine Stärken ausspielen kann. So ist es preislich unabhängig, ob das Service-Ökosystem inkl. Funktionen gespiegelt wird oder die Test neben den in Produktion genutzten Funktionen durchgeführt werden.

wird daher in der Praxis häufig die Funktionen in Produktion zu Testen bzw. einen gespiegelten Testaccount der eigentlichen, in Produktion laufenden, serverlosen Anwendungslandschaft anzulegen.

So muss sich meistens auf lokale Unit test beschränkt werden. Um Komponenten bzw. Funktionen einem Integrationstest zu unterziehen, muss sich genau über die Möglichkeiten bei dem jeweiligen Cloud provider informiert werden, da Preise und infrasturkurelle Gegbenheiten variieren. AWS bietet über

SAM Für traffic shifting

Dies ist teilweise Vendorenspezifisch und bedarf der auseinandersetzung

3.7 Konsequenzen für das Monitoring

4 Anwendung des Leitfadens

4.1 Auswertung der Ergebnisse

4.1.1 Beurteilung der Kollaborationsauswirkungen

4.1.2 Beurteilung der Stabilität

4.1.3 Beurteilung der Skalierbarkeit

4.1.4 Beurteilung des Monitorings

4.1.5 Beurteilung der Testmöglichkeiten

4.2 Korrektur und Anpassungen

5 Abschließende Betrachtung

5.1 Absehbare Entwicklungen

[AAS19] Serverless immer mehr bekanntheit und mehr Papers etc. ... Später market für funktionen die optimiert sind etc. [SKM20]

Real time communication tool [SKM20]

Real time tracking gps [SKM20] da beide nicht auf dem Application state benötigen/ beruhen [SKM20]

[Hel+18]

5.2 Zusammenfassung

5.3 Weiterführende Forschung

Serverless Computing: A Survey of Opportunities, Challenges and Applications store for functions [SBW19]

Literaturverzeichnis

- [FK03] Ian Foster und Carl Kesselman. *The Grid 2: Blueprint for a new computing infrastructure*. Elsevier, 2003.
- [Fox+09] Armando Fox u. a. “Above the clouds: A berkeley view of cloud computing”. In: *Dept. Electrical Eng. and Comput. Sciences, University of California, Berkeley, Rep. UCB/EECS 28.13* (2009), S. 2009.
- [Liu+11] Fang Liu u. a. “NIST cloud computing reference architecture”. In: *NIST special publication 500.2011* (2011), S. 292.
- [MG+11] Peter Mell, Tim Grance u. a. “The NIST definition of cloud computing”. In: (2011).
- [Reg11] Antonio Regalado. “Who coined ‘cloud computing’”. In: *Technology Review* 31 (2011).
- [Bal+17] Ioana Baldini u. a. “Serverless computing: Current trends and open problems”. In: *Research Advances in Cloud Computing*. Springer, 2017, S. 1–20.
- [Eiv17] Adam Eivy. “Be wary of the economics ofSServerless” Cloud Computing”. In: *IEEE Cloud Computing* 4.2 (2017), S. 6–12.
- [MB17] Garrett McGrath und Paul R Brenner. “Serverless computing: Design, implementation, and performance”. In: *2017 IEEE 37th International Conference on Distributed Computing Systems Workshops (ICDCSW)*. IEEE. 2017, S. 405–410.
- [PS17] Hussachai Puripunpinyo und MH Samadzadeh. “Effect of optimizing Java deployment artifacts on AWS Lambda”. In: *2017 IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*. IEEE. 2017, S. 438–443.
- [RC17] Michael Roberts und John Chapin. *What Is Serverless?* O’Reilly Media, Incorporated, 2017.
- [Fow18] Martin Fowler. “Serverless Architectures”. In: *martinfowler.com* (2018). <https://www.martinfowler.com/articles/serverless.html>.
- [Hel+18] Joseph M Hellerstein u. a. “Serverless computing: One step forward, two steps back”. In: *arXiv preprint arXiv:1812.03651* (2018).

- [IMS18] Vatche Ishakian, Vinod Muthusamy und Aleksander Slominski. “Serving deep learning models in a serverless platform”. In: *2018 IEEE International Conference on Cloud Engineering (IC2E)*. IEEE. 2018, S. 257–262.
- [IS18] Vitalii Ivanov und Kari Smolander. “Implementation of a DevOps pipeline for serverless applications”. In: *International Conference on Product-Focused Software Process Improvement*. Springer. 2018, S. 48–64.
- [JC18] David Jackson und Gary Clynch. “An Investigation of the Impact of Language Runtime on the Performance and Cost of Serverless Functions”. In: *2018 IEEE/ACM International Conference on Utility and Cloud Computing Companion (UCC Companion)*. IEEE. 2018, S. 154–160.
- [KS18] Kyriakos Kritikos und Paweł Skrzypek. “A review of serverless frameworks”. In: *2018 IEEE/ACM International Conference on Utility and Cloud Computing Companion (UCC Companion)*. IEEE. 2018, S. 161–168.
- [LSF18] Hyungro Lee, Kumar Satyam und Geoffrey Fox. “Evaluation of production serverless computing environments”. In: *2018 IEEE 11th International Conference on Cloud Computing (CLOUD)*. IEEE. 2018, S. 442–450.
- [Lóp+18] Pedro García López u. a. “Comparison of faas orchestration systems”. In: *2018 IEEE/ACM International Conference on Utility and Cloud Computing Companion (UCC Companion)*. IEEE. 2018, S. 148–153.
- [Man+18] Johannes Manner u. a. “Cold start influencing factors in function as a service”. In: *2018 IEEE/ACM International Conference on Utility and Cloud Computing Companion (UCC Companion)*. IEEE. 2018, S. 181–188.
- [MPDF+18] Sunil Kumar Mohanty, Gopika Premsankar, Mario Di Francesco u. a. “An Evaluation of Open Source Serverless Computing Frameworks.” In: *CloudCom*. 2018, S. 115–120.
- [SS18] Mohit Sewak und Sachchidanand Singh. “Winning in the era of serverless computing and function as a service”. In: *2018 3rd International Conference for Convergence in Technology (I2CT)*. IEEE. 2018, S. 1–5.

- [Adi+19] Paarijaat Aditya u. a. “Will Serverless Computing Revolutionize NFV?” In: *Proceedings of the IEEE* 107.4 (2019), S. 667–678.
- [AAS19] Mohammed Al-Ameen und Josef Spillner. “A systematic and open exploration of FaaS research”. In: *ESSCA, Zurich, Switzerland, December 21, 2018*. CEUR-WS. 2019, S. 30–35.
- [Cas+19] Paul Castro u. a. “The server is dead, long live the server: Rise of Serverless Computing, Overview of Current State and Future Trends in Research and Industry”. In: *arXiv preprint arXiv:1906.02888* (2019).
- [Jon+19] Eric Jonas u. a. “Cloud programming simplified: A berkeley view on serverless computing”. In: *arXiv preprint arXiv:1902.03383* (2019).
- [Kap19] Ayhan Kaplan. “Framework for migrating deployed serverless applications”. Magisterarb. 2019.
- [Lei+19] Philipp Leitner u. a. “A mixed-method empirical study of Function-as-a-Service software development in industrial practice”. In: *Journal of Systems and Software* 149 (2019), S. 340–359.
- [Rac+19] Louis Racicot u. a. “Quality Aspects of Serverless Architecture: An Exploratory Study on Maintainability”. In: *Proceedings of the 14th International Conference on Software Technologies, ICSOFT*. 2019, S. 26–28.
- [SBW19] Mohammad Shahrads, Jonathan Balkind und David Wentzlaff. “Architectural implications of function-as-a-service computing”. In: *Proceedings of the 52nd Annual IEEE/ACM International Symposium on Microarchitecture*. 2019, S. 1063–1075.
- [VE+19] Erwin Van Eyk u. a. “The SPEC-RG Reference Architecture for FaaS: From Microservices and Containers to Serverless Platforms”. In: *IEEE Internet Computing* (2019).
- [Yus+19] Vladimir Yussupov u. a. “A Systematic Mapping Study on Engineering Function-as-a-Service Platforms and Tools”. In: *Proceedings of the 12th IEEE/ACM International Conference on Utility and Cloud Computing (UCC 2019)*. ACM, Dez. 2019, S. 229–240. DOI: 10.1145/3344341.3368803.

- [SKM20] Hossein Shafiei, Ahmad Khonsari und Payam Mousavi. “Serverless Computing: A Survey of Opportunities, Challenges and Applications”. In: (2020).