



HOCHSCHULE HEILBRONN

**Proseminar (282136)**

# **XaaS - Anything as a Service**

Suphi Pembe (207617),  
Andreas Würzer (207258),  
Christian Nguyen (207613)

Sommersemester 2022

Vorgelegt bei Claudia Pittel

# Management Summary

# Inhaltsverzeichnis

<b>Management Summary</b>	<b>ii</b>
<b>Abkürzungsverzeichnis</b>	<b>v</b>
<b>Abbildungsverzeichnis</b>	<b>vi</b>
<b>1 Einleitung</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 Ziel der Arbeit . . . . .	1
1.3 Vorgehensweise . . . . .	1
<b>2 Anything as a Service - Cloud Computing</b>	<b>2</b>
2.1 Definition . . . . .	2
2.2 Typische Servicemodelle . . . . .	2
2.2.1 IaaS: Infrastructure as a Service . . . . .	2
2.2.2 SaaS: Software as a Service . . . . .	2
2.2.3 PaaS: Plattform as a Service . . . . .	2
2.3 Vor- und Nachteile . . . . .	2
<b>3 Knappheit von Grafikkarten</b>	<b>3</b>
3.1 Preisentwicklung . . . . .	4
3.2 Ursache Halbleitermangel und Krypto-Mining . . . . .	5
<b>4 Gaming as a Service</b>	<b>6</b>
4.1 Funktionsweise . . . . .	6
4.2 Anbietervergleich . . . . .	8
4.2.1 Voraussetzung . . . . .	8
4.2.2 Angebot . . . . .	9
4.2.3 Preis . . . . .	9
4.3 Hardwarevoraussetzung um Usability zu gewährleisten . . . . .	10
<b>5 GPU as a Service</b>	<b>11</b>
5.1 Funktionsweise . . . . .	11
5.2 Einsatzgebiete . . . . .	12
5.3 Vergleich eigene GPU und GPU in der Cloud . . . . .	13
<b>6 Marktvorhersage</b>	<b>14</b>

<b>7 Fazit und Ausblick</b>	<b>15</b>
<b>Quellenverzeichnis</b>	<b>vii</b>
<b>Ehrenwörtliche Erklärung</b>	<b>ix</b>

# Abkürzungsverzeichnis

**GPU** Graphics-Processing-Unit oder Grafikkarte

**HPC** High-Performance-Computing

**RAID** redundant array of independent disk

# Abbildungsverzeichnis

3.1	M. Kords (2022)	3
3.2	Voas, Kshetri und DeFranco (2021)	4
4.1	D'Angelo, Ferretti and Marzolla (2022)	7
4.2	atene KOM GmbH (2021)	8
5.1	Wang u. a. (2017)	11

# 1 Einleitung

## 1.1 Motivation

Durch den aktuell anhaltende Halbleitermangel besteht ein Engpass an Ressourcen von dem die meisten Wirtschaftszweige betroffen sind. Einer dieser Wirtschaftszweige ist die Produktion von GPUs (graphics processing unit). Diese werden für diverse Prozesse in Computern verwendet, im betrieblichen wie auch im privaten Bereich. Primär in dieser Arbeit werden die Bereiche High-Performance-Computing (HPC) und Gaming betrachtet. Beide diese Bereiche benötigen GPU-Rechenleistung, welche konventionell von einer lokal verbauten GPU zur Verfügung gestellt wird. Als langfristige Lösung soll analysiert werden ob es möglich ist durch zentrale Services, welche GPU as a Service anbieten. Durch die zentralen Ressourcenteilung wird dem Mangel entgegengewirkt durch die Schaffung einer Alternative für den Bedarf.

## 1.2 Ziel der Arbeit

...

## 1.3 Vorgehensweise

...

## **2 Anything as a Service - Cloud Computing**

### **2.1 Definition**

### **2.2 Typische Servicemodelle**

#### **2.2.1 IaaS: Infrastructure as a Service**

#### **2.2.2 SaaS: Software as a Service**

#### **2.2.3 PaaS: Plattform as a Service**

### **2.3 Vor- und Nachteile**



### 3 Knappheit von Grafikkarten

Die Knappheit von Grafikkarten hat den aktuellen Markt durch neue Branchen, die GPU-Leistung nutzen, nachhaltig verändert. Diese Knappheit entsteht nicht nur durch den Mangel des Rohstoffes, sondern auch durch die Weiterentwicklung von verwendeten Computern in allen Einsatzgebieten.<sup>1</sup>



Abbildung 3.1: Weltweite Lieferung von Halbleiterprodukten für die Automobilindustrie von 2011 bis 2021

Im Vergleich zur Lage im Jahr 2011 wurden für die Automobilindustrie im Jahr 2021 fast drei mal so viele Halbleiter geliefert. Ebenfalls mit der Weiterentwicklung von Internet of Things Produkten wird in Zukunft der Bedarf an Halbleitern weiter steigen.<sup>234</sup>

In diesem Kapitel soll die Preisentwicklung von GPUs betrachtet werden, dabei wird ein Zusammenhang geschaffen mit den Ursachen die diese Preisentwicklung verursacht haben.

<sup>1</sup>Vgl. Voas, Kshetri und DeFranco, "Scarcity and Global Insecurity: The Semiconductor Shortage".

<sup>2</sup>Vgl. McClean, *The 2022 McClean Report*.

<sup>3</sup>Vgl. Voas, Kshetri und DeFranco, "Scarcity and Global Insecurity: The Semiconductor Shortage".

<sup>4</sup>Kords, *Weltweite Lieferung von Halbleiterprodukten (integrated circuit) für die Automobilindustrie von 2011 bis 2021*, Vgl.

## 3.1 Preisentwicklung

Die rapide steigende Preisentwicklung von GPUs ist auf zwei Kernfaktoren reduzierbar.

- Größerer Bedarf an GPUs und Halbleitern, dem Kernbestandteil von GPUs
- Mangelnde Kapazitäten zur Produktion von Halbleitern

Der Bedarf an Halbleitern und GPUs ist konstant im Anstieg. Besonders durch die Corona-Pandemie, hat sich im Vergleich zu 2019 im Jahr 2020 ein Umsatzanstieg von 5,4 % aufgezeigt.<sup>5</sup>

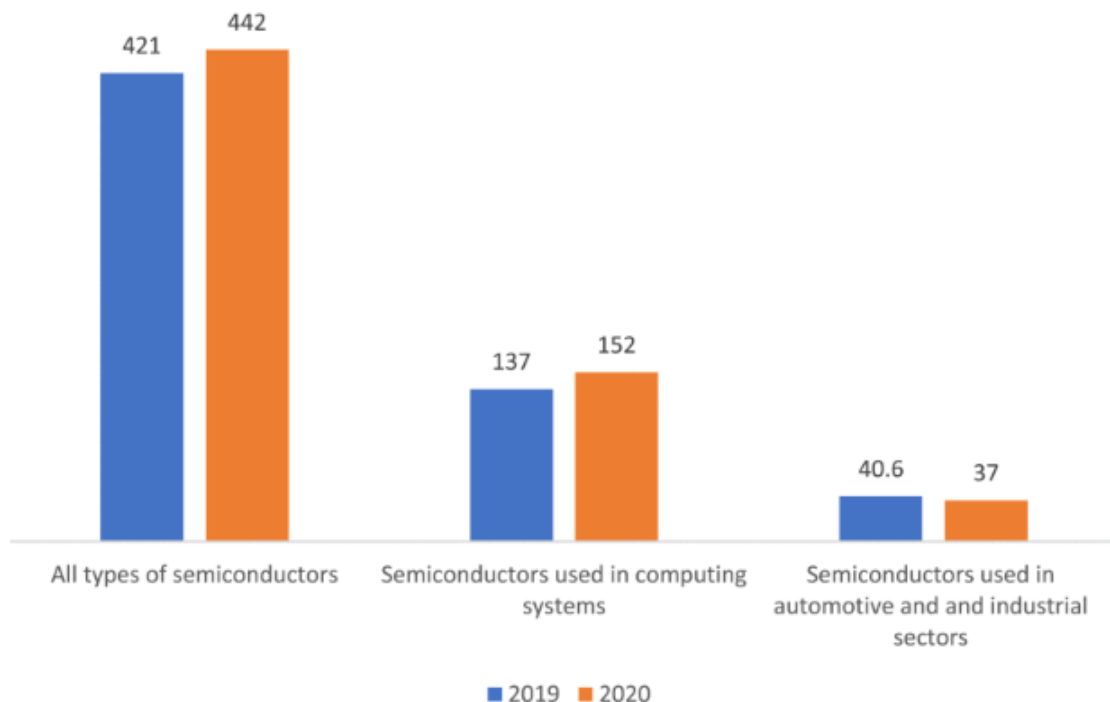


Abbildung 3.2: Worldwide semiconductor revenues in 2019 and 2020 (dollar, billions)

Wie in Abbildung 3.2 zu sehen ist, ist der Umsatzanstieg größtenteils durch Erlöse von Computersystemen entstanden. Im Vergleich dazu sind Umsätze, die durch Abnehmer in der Automobilbranche entstanden sind, gesunken. Das lässt sich auf den steigenden Bedarf an Computersystemen zurückführen. Während der Pandemie mussten sich viele Menschen an Home-Schooling und Home-Office anpassen, um weiter den Alltagsbetrieb ausführen zu können. Ein Nebenläufiger Effekt ist damit, dass durch die Digitalisierung weniger Mobilität benötigt wird. Damit lässt sich der reduzierte Bedarf an Halbleitern in der Automobilindustrie erklären. Dennoch ist damit insgesamt der Bedarf an Halbleitern gestiegen.<sup>6</sup>

---

<sup>5</sup>Vgl. Voas, Kshetri und DeFranco, "[Scarcity and Global Insecurity: The Semiconductor Shortage](#)".

<sup>6</sup>Vgl. [ebd.](#)

Der steigende Bedarf allein ist aber wie angeführt nicht der einzige Faktor. Die Produktion von Halbleitern stagniert. Das lässt sich auf verschiedene Ursachen zurückführen. Da die meisten Halbleiter in asiatischen Ländern produziert werden und diese den eigenen Bedarf zuerst decken, ist für den Export weniger verfügbar. Ebenfalls haben sich in den letzten Jahren Naturkatastrophen ereignet, die z.B. durch Dürre, die Produktion lahmgelegt.<sup>7</sup> Durch mangelnde Produktion und steigenden Bedarf hat sich nun ein sehr hoher Marktpreis entwickelt. Um diesem entgegenzuwirken ist nicht nur Ressourcenverfügbarkeit zu schaffen, sondern auch eine effizientere Methode zur Nutzung der Ressourcen.

## 3.2 Ursache Halbleitermangel und Krypto-Mining

Um den Halbleitermangel besser zu verstehen, sollte eine neue Branche, die primär GPU-Leistung nutzt, thematisiert werden. Seit dem Jahr 2021 haben Kryptowährungen ein mehr als fünffaches Investitionsvolumen im Vergleich zum Vorjahr zu verzeichnen.<sup>8</sup>

Um den Zusammenhang zu erläutern: Kryptowährungen validieren Ihre Transaktionen durch die Nutzung der Blockchain-Technologie. Bei diesem Validierungsprozess werden neue Datenblöcke in einer Datenbank gespeichert und von anderen Nutzern überprüft durch eine Prüfsumme, die dabei gebildet werden können muss. Um diese Prüfsumme berechnen zu können nutzen sogenannte "Krypto-Miner" primär GPU-Leistung.<sup>9</sup>

Da große Investitionen für das Kryptomining getätigt werden, hat sich ein neuer Markt mit einer großen Nachfrage gebildet, der GPUs benötigt.

Allerdings gibt es faktengestützte Prognosen, welche behaupten, dass der Halbleitermangel nicht zu lange anhalten wird. Nach Gartners Prognose für das zweite Quartal 2022 sollte sich der Halbleitermangel verringern, auch wenn dies nicht im vollen Umfang eingetreten ist, sind die Argumente für die weitere Zukunft nicht irrelevant.

Angeführt wird, dass durch den Halbleitermangel nun die Lieferkette enger überwacht wird. Daraus resultierend werden mehr Transparenz und Vorinvestitionen geschaffen, um die Lieferungen garantieren zu können. Ebenfalls möchte man auch die Abnahme von mehreren Lieferanten bevorzugen, anstatt von einem Lieferanten abhängig zu sein. Diese Faktoren sollen langfristig Sicherheit in der Lieferkette bieten.<sup>10</sup>

---

<sup>7</sup>Vgl. Voas, Kshetri und DeFranco, "[Scarcity and Global Insecurity: The Semiconductor Shortage](#)".

<sup>8</sup>Statista-Research-Department, [Volumen der weltweiten Investitionen in Blockchain-Technologien und Kryptowährungen von 2018 bis 2021](#).

<sup>9</sup>Vgl. Arslanian, [The Book of Crypto: The Complete Guide to Understanding Bitcoin, Cryptocurrencies and Digital Assets](#), S.259-273.

<sup>10</sup>Vgl. Rimol, [Gartner Says Global Chip Shortage Expected to Persist Until Second Quarter of 2022](#).

## 4 Gaming as a Service

Gaming-as-a-Service ist ein Zukunftstrend in der Spielindustrie. Je nach Entwicklung, steigt die benötigte Hardwareleistung wie CPU, GPU, RAM, sowie Speicherplatz zum Installieren von Spielen stetig an. Hochwertige Spiele können nicht mehr von veralteten Computern genossen werden, wodurch die Hardware regelmäßig aktualisiert werden muss, für eine bessere Umgebung.

In diesem Kapitel soll Cloud Computing mit Gaming-as-a-Service als Beispiel vertieft werden. Hierbei sollen die Funktionalität und Architektur, sowie die aktuell verschiedenen Angebote betrachtet werden.

### 4.1 Funktionsweise

Unabhängig von der Architektur wird das Gaming als Schleifenprozedur betrachtet, die eine Interaktion zwischen Endnutzern und Spiellogik ermöglicht. Hierbei stehen zwei wichtige Komponenten in Relation: Der Server, auch Cloud genannt und das Gerät des Endnutzers, auch Thin-Client genannt. Je nach GaaS-Modell, findet die Ausführung des Spiels, die Spiellogik und die Wiedergabe der Szenen innerhalb der Cloud statt. Für den Empfang der komprimierten Audio- und Videosignale, ist der Thin-Client verantwortlich.<sup>11</sup> In Betracht gezogen werden hier drei GaaS-Modelle: Remote-Rendering-GaaS, Local-Rendering-GaaS und Cognitive-Resource-Allocation-GaaS.<sup>12</sup>

Beim Remote-Rendering-GaaS-Modell (RR-GaaS) besitzt die Cloud-Infrastruktur ein Modul zum Kodieren. Dieser ist dafür verantwortlich, jeden Frame der Spielszene zu rendern und den Stream des Videos zu komprimieren, damit er an das Thin-Client übertragen werden kann. Dort wird der Stream dekodiert und angezeigt. Benutzereingaben werden vom Terminal erfasst und die Cloud an die Spiellogik zurückgesendet, die sich um die entsprechende Aktualisierung des Spielzustands kümmert.<sup>13</sup> Dies impliziert, dass die Hardwareanforderung für den Endnutzern minimiert wird, unabhängig von der Komplexität von Spielszenen, Spiellogik und Interaktionen. Folglich können hochwertige Spiele mit leistungsschwachen Geräten bedient werden.<sup>14</sup> Das RR-GaaS-Modell verbraucht jedoch eine beträchtliche

---

<sup>11</sup>Vgl. Zadtootaghaj, *Quality of Experience Modeling for Cloud Gaming Services*.

<sup>12</sup>Vgl. Wei Cai und Leung, "Toward Gaming as a Service".

<sup>13</sup>Vgl. Gabriele D'Angelo und Marzolla, *Cloud for Gaming*.

<sup>14</sup>Vgl. Wei Cai und Leung, "Toward Gaming as a Service".

Bandbreite, um den komprimierten Videostream zu übertragen, und kann besonders empfindlich auf Netzwerkverzögerungen reagieren.

Beim Local-Rendering-GaaS-Modell (LR-GaaS) wird der Stream des Videos in der Cloud als eine Folge von Rendering-Anweisungen auf hoher Ebene kodiert, die zum Thin-Client gestreamt werden. Dieser dekodiert und führt die Anweisungen aus, um jeden einzelnen Frame zu zeichnen.<sup>15</sup> Der hervorstechendste Vorteil des LR-GaaS-Modells besteht darin, dass die Cloud keine einzelnen Frames in Echtzeit mehr über das Internet an die Thin-Clients übertragen muss, was die Netzwerklast erheblich reduziert. Ansonsten bestehen die ähnlichen Vorteile wie beim RR-GaaS-Modell.<sup>16</sup>

Anders als beim RR- und LR-GaaS-Modell, ist beim Cognitive-Resource-Allocation-GaaS die Cloud logisch in eine Reihe von Modulen unterteilt. Die Module können dann wiederum auf dem Thin-Client hochgeladen und ausgeführt werden. Das CRA-GaaS-Modell verlagert die Berechnung zurück auf das Client-Terminal und reduziert so die Belastung der Cloud. Die Client-Ressourcen werden effizient genutzt, da immer nur die benötigten Komponenten lokal gespeichert werden. Dies ist ein erheblicher Vorteil, wenn man bedenkt, dass die Daten eines kompletten modernen Spiels viel Platz einnehmen.<sup>17</sup>

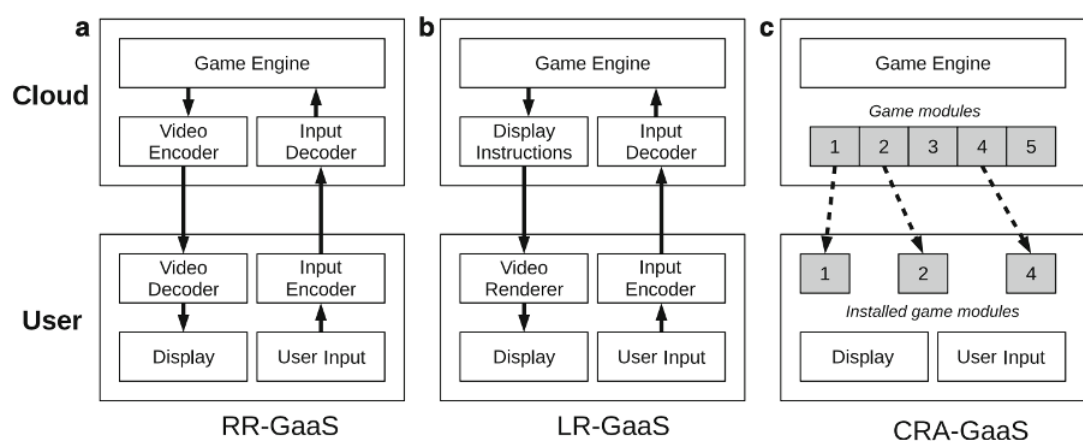


Abbildung 4.1: Gaming as a Service models

<sup>15</sup>Vgl. Gabriele D'Angelo und Marzolla, *Cloud for Gaming*.

<sup>16</sup>Vgl. Wei Cai und Leung, "Toward Gaming as a Service".

<sup>17</sup>Vgl. Gabriele D'Angelo und Marzolla, *Cloud for Gaming*.

## 4.2 Anbietervergleich

### 4.2.1 Voraussetzung

All die Prozesse und Interaktionen zwischen der Cloud und dem Thin-Client werden durch die Leistung des Netzwerks eingeschränkt. Weitere Einschränkungen wie eine begrenzte Bandbreite, würde die Erfahrung der Spieler beeinflussen. Ist es aber in Deutschland überhaupt anwendbar mit der aktuellen Breitbandverfügbarkeit? Je nach Anbieter und Qualität des Videostreams variiert sich die erforderliche Internetgeschwindigkeit. Eines der großen Cloud-Gaming-Anbieter, wie Google Stadia, erfordert eine Bandbreite bei der höchsten Auflösung eine Netzwerkgeschwindigkeit von mindestens 35 Mbit/s.<sup>18</sup> Da sich die Cloud beim Gaming-as-a-Service typischerweise im regionalen Netzwerk eines Betreibers befinden, bleibt die Ende-zu-Ende-Übertragungsverzögerung gering, da zwischen dem Client und dem Server normalerweise viel Bandbreite verfügbar ist. Falls jedoch mehrere Geräte den Breitbandanschluss zu Hause nutzen, kann die verfügbare Bandbreite beim Zugang auf der letzten Meile erheblich variieren.

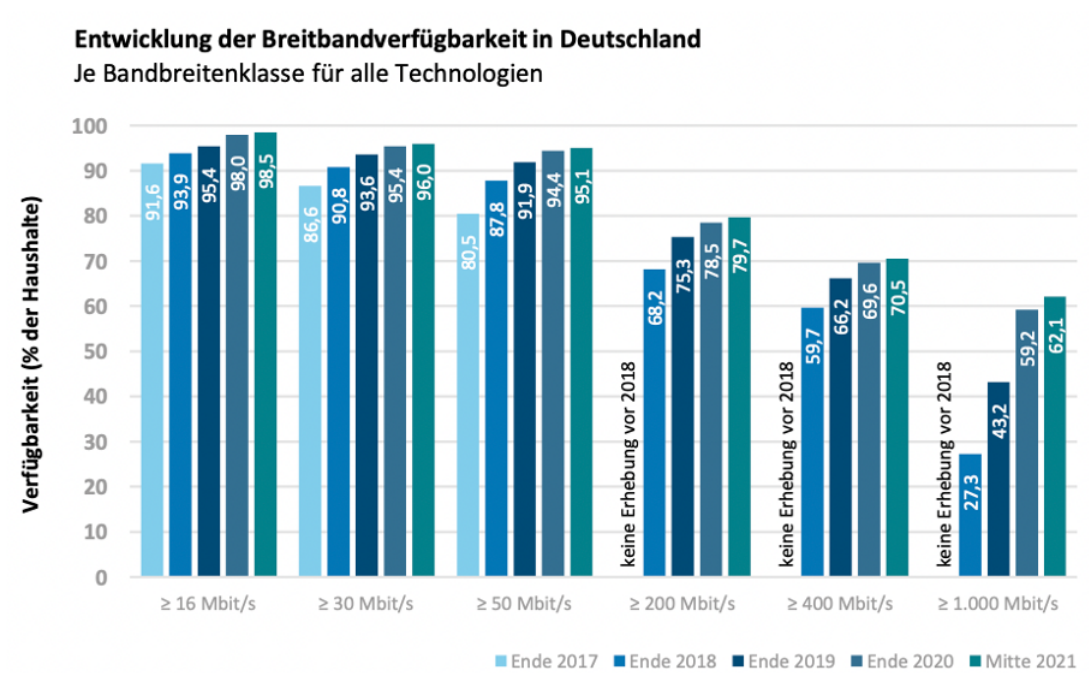


Abbildung 4.2: Entwicklung der Breitbandverfügbarkeit in Deutschland nach Bandbreitenklassen

Aus der Abbildung 4.2: „Entwicklung der Breitbandverfügbarkeit in Deutschland“ ist ein

<sup>18</sup>Vgl. Stadia, [Bandbreite, Datennutzung und Streamingqualität](#).

Zuwachs der Breitbandverfügbarkeit in den vergangenen Jahren deutlich zu sehen. Ein deutlicher Sprung ist in der höchsten Klasse zu sehen, die sich mehr als verdoppelt hat. Festzustellen ist also, dass die entnommene Anforderung von Google Stadia mit 35 Mbit/s für die höchste Auflösung von 95,1 Prozent der Haushalte in Deutschland erfüllen.

### 4.2.2 Angebot

Aktuell sind fünf große Cloud-Gaming-Anbieter im deutschen Markt vertreten: Nvidia GeForce Now, Google Stadia, Shadow, Xbox Cloud-Gaming, PS Now sowie weitere kleine Anbieter. Für die Nutzung des Angebotes der Anbieter ist ein Abonnement erforderlich und nur mit bestimmten Geräten nutzbar. Erste Unterschiede ergeben sich bei den jeweiligen Konzepten. Vergleicht man die angebotenen Abos des Anbieters Nvidia GeForce Now mit Google Stadia, bietet Nvidia GeForce Now drei mögliche Mitgliedschaften, bei der man je nach Auswahl unterschiedliche Leistungen erhält.<sup>19</sup> Im Vergleich bietet Google Stadia nur einen möglichen Abo, bei der jeder Endnutzer die gleichen Kosten und Leistungen angeboten bekommt. Aufgrund der Diversität im Angebot besteht im Konzept der einzelnen Betreiber immer noch eine Gemeinsamkeit. Und zwar schaffen die Konzepte einen potenziellen Markt für Benutzer, die nicht den Kauf von Spielsoftware, sondern den Kauf von Spielzeit für Computer- und Konsolenspiele anzustreben. Auf diese Weise können Benutzer zu geringen Kosten auf eine Vielzahl von Spielen zugreifen.

Von den großen Cloud-Gaming-Anbietern ist kein deutsches Unternehmen mitinbegriffen, zumindest aktuell nicht mehr. Auch Telekom hatte eine eigene Cloud-Gaming-Plattform, bei der keine Downloads oder Kauf teurer Hardware von Nöten ist. Nach nicht einmal zwei Jahren, wurde der Streamingdienst eingestellt, aufgrund Desinteresses der Spieler.<sup>20</sup> Laut Prognose wird der Marktwert von Cloud-Gaming weltweit bis 2024 sich um 4,27 Prozent steigern.<sup>21</sup>

### 4.2.3 Preis

Die Preisgestaltung variiert sich von Anbieter zu Anbieter, von Kostenlos bis zu 29,99 Euro pro Monat. Hier spielen zwei große Faktoren eine Rolle: Service/Leistung und Spielebibliothek. Bei der Mitgliedschaft „Kostenlos“ ist eindeutig festzustellen, dass diese Option nur für Interessenten ist, die zum Testen des Service angeregt werden. Mit einem Abo von 9,99 Euro im Monat, wird natürlich schon bei einigen Anbietern der volle Umfang geliefert.

---

<sup>19</sup>Vgl. Suznjevic, Slivar und Skorin-Kapov, [“Analysis and QoE evaluation of cloud gaming service adaptation under different network conditions: The case of NVIDIA GeForce NOW”](#).

<sup>20</sup>Vgl. Telekom, [5G Cloud Gaming](#).

<sup>21</sup>Vgl. Clement, [Cloud gaming market value worldwide from 2019 to 2024](#).

## 4.3 Hardwarevoraussetzung um Usability zu gewährleisten

Die Anforderungen an Rechenleistung und Speicherkapazität von Computer- und Videospielen werden immer höher, da die Spiele immer realistischer und komplexer werden.<sup>22</sup> Für einen einfachen Heimcomputer ist das Abspielen von hochwertigen Spielen in einer guten Umgebung nicht mehr möglich und muss aktualisiert werden.<sup>23</sup> Für die Nutzung des Cloud-Gaming reicht das jedoch. Je nach Anbieter sind die Systemanforderungen unterschiedlich, aber dennoch keine hohen Anforderungen.

Nehmen wir als Beispiel eines der größten Cloud-Gaming-Anbieter Nvidia GeForce Now (GFN). Um GFN nutzen zu können, werden GPUs, die seit 2015 veröffentlicht wurden, beim Streaming mit bis zu 3840x2160p 60 FPS und 1440p/1600p 120 FPS unterstützt. Für die CPU genügt eine Dual-Core x86-64 mit 2,0 GHz und RAM mit nur 4 GB Speicher. Mit der Mindestanforderung von GFN muss der Thin-Client keine große Verarbeitungsleistung und Speicherkapazität haben.<sup>24</sup> Mit nur niedriger Anforderung kann alte Hardware die Usability bei Spielen wie bei High-End-Computern gewährleistet werden. Es besteht keine Sorge um die Prozessorleistung, das Betriebssystem, die Grafikkarten oder andere technische Spezifikationen eines Computers.<sup>25</sup>

---

<sup>22</sup>.

<sup>23</sup>Vgl. Suznjevic, Slivar und Skorin-Kapov, "Analysis and QoE evaluation of cloud gaming service adaptation under different network conditions: The case of NVIDIA GeForce NOW".

<sup>24</sup>Vgl. Clement, *Cloud gaming market value worldwide from 2019 to 2024*.

<sup>25</sup>Vgl. Ojala und Tyrvaäinen, "Developing Cloud Business Models: A Case Study on Cloud Gaming".



## 5 GPU as a Service

Außerhalb von Gaming as a Service gibt es weitere Service-Möglichkeiten die ebenfalls GPU as a Service beanspruchen bzw. inkludieren. Wachsende Märkte dafür sind das Rendern von 3D-Modellen und animierten Videos, wie auch Deep-Learning-Model-Training für KIs.

Im folgenden Kapitel wird erläutert werden, wie GPU as a Service-Anbieter Ihren Service realisieren können und Qualität mit verschiedenen Methoden schaffen. Wie auch Einsatzgebiete für Cloud-Lösungen festgestellt und identifiziert werden.

### 5.1 Funktionsweise

Bei der Dienstleistungsinanspruchnahme werden die von einem Nutzer geforderten Prozesse, wie z.B. das Rendern von 3D-Modellen, durch die Rechenleistung des Anbieters verarbeitet. Im Gegensatz zu der Privatsnutzung bei der nur eine GPU genutzt wird, verwenden GPU as a Service-Anbieter mehrere GPUs. Allerdings, da GPUs nur als Co-Prozessoren in solchen Systemen genutzt werden, können diese nicht eigenständig betrieben werden, sondern benötigen ein zentrales Betriebssystem, welches auch als Kernel bezeichnet. Üblicherweise wird dies skalierbar angewendet mit einer Vielzahl an Kernels, welche eine Vielzahl an GPUs besitzen.<sup>26</sup>

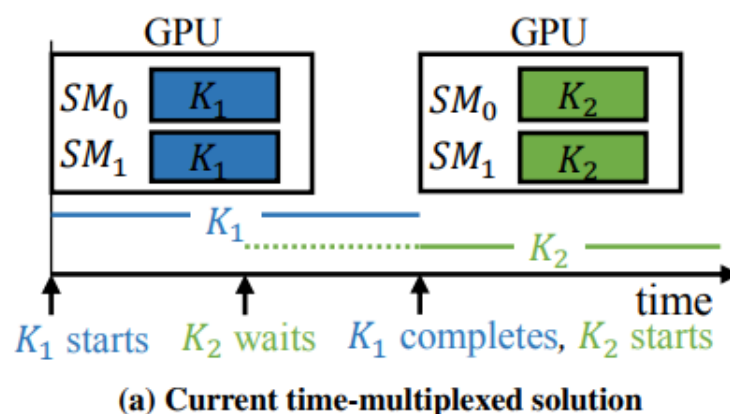


Abbildung 5.1: Current time-multiplexed solution

Es gibt verschiedene Methoden, diese Prozesse zu verarbeiten. Eine Methode davon ist

<sup>26</sup>Vgl. Wang u. a., "Quality of Service Support for Fine-Grained Sharing on GPUs".

das Zeit-Multiplexverfahren. Bei diesem Verfahren werden mehrere Prozesse auf die Kernels sequenziell aufgeteilt und verarbeitet, wie in Abbildung 5.1 zu sehen. Diese Methode hat keinen direkten Mehrwert in Bezug zur Verarbeitungseffizienz, allerdings verhindert sie, dass Prozesse Kernels unnötig blockieren können.<sup>27</sup>

Ein anderer Ansatz, um Konstanz zu schaffen, ist, dass man die Kernels über Software zu einem Kernel fusioniert. Das ist vergleichbar mit dem für Festplatten verwendeten "redundant array of independent diskSystem, auch abgekürzt RAID genannt, welches gängiger bekannt ist. Durch das Fusionieren der Kernels über die Software kann eine konstante Leistung gewährleistet werden. Ebenfalls, um eine Fairness zwischen allen Leistungsabnehmern zu schaffen, da es keine Situation geben kann, einen leistungsschwächeren Kernel zugewiesen zu bekommen.<sup>28</sup>

Eine quasi gegenteilige Methode, die Kernels in verschiedene Partitionen aufzuteilen, anstatt sie zu fusionieren. Wieder mal vergleichbar, wie man auch mit Festplatten umgehen kann. Die Leistung des Kernel wird in verschiedene Partitionen aufgeteilt. Diese können dann je nach Monetarisierungsmodell vermietet werden, für einen Zeitraum oder auf einer Bedarfsbasis. Daraus entsteht der Vorteil, dass Nutzer des Services ihren Bedarf selbst definieren können und garantiert diese Leistung in Anspruch nehmen können.<sup>29</sup>

Es gibt neben diesen Methoden weitere Maßnahmen für die Qualitätssicherung des GPU as a Service Modells. Allerdings würden diese sich zu weit vertiefen und über das Ziel dieser Arbeit hinausgehen. Für eine weitere Vertiefung ist die Arbeit von Wang u.a., "Quality of Service Support for Fine-Grained Sharing on GPUs". zu empfehlen.

## 5.2 Einsatzgebiete

Die Einsatzgebiete sind weitreichend und können auf zwei Bedarfsmethoden eingeteilt werden. Einmal in Schüben, in denen ein festes Ziel und der Aufwand definiert werden, welche der Prozess erreichen soll. Beispiele dafür sind Rendern von 3D-Modellen und animierten Videos.

Um dieses Beispiel auszuführen: Wenn ein 3D-Modell dargestellt wird, besteht es aus Matrixoperationen, die von GPUs in der Regel berechnet werden. Dies hat meist das Ziel, diese realitätsgetreu und logisch in einer Umgebung darzustellen. Dieses Beispiel lässt

---

<sup>27</sup>Vgl. Wang u. a., "Quality of Service Support for Fine-Grained Sharing on GPUs".

<sup>28</sup>Vgl. ebd.

<sup>29</sup>Vgl. ebd.

sich auf animierte Videos erweitern, bei denen dann eine veränderte Version des Vorgängermodells oder ein vollkommen anderes Modell für jedes Bild im Video dargestellt wird.<sup>30</sup>

Ebenso ist es möglich, dass mit einem optionalen Ziel, aber unbekannten Aufwand, so ein Prozess ausgeführt wird. In diesem Fall wird eine konstante Rechenleistung benötigt. Beispiele dafür sind Deep-Learning-Modelle für KIs oder Gaming as a Service.

Um auch hier ein Beispiel anzuführen, falls Gaming as a Service-Dienste in Anspruch genommen wird. Bei diesem Service ist kein Ziel festgelegt, da auch kein konkreter Prozess vorgegeben ist. Anhand der Eingaben des Nutzers variiert die benötigte Visualisierung und Spielelogik. Da die Eingaben nicht vorher konkret bekannt sind, besteht ein konstanter Bedarf an Rechenleistung um auf die Eingaben reagieren zu können. Die Rechenleistung wird benötigt bis der Nutzer das Spiel beendet.<sup>313233</sup>

## 5.3 Vergleich eigene GPU und GPU in der Cloud

Nach den angeführten Informationen lässt sich Folgendes feststellen: Der Bedarf an GPU Rechenleistung ist im konstanten Wachstum für verschiedene Märkte. Die GPU as a Service-Anbieter können Ihre Systeme entsprechend dem Monetarisierungsmodell und den Leistungsprioritäten aufstellen. Dabei handelt es sich nicht um neue Methoden, sondern um etwas mit der Massenspeicherverwaltung Vergleichbaren. Je nach geplanter Nutzung werden GPUs entweder in Schüben oder mit konstanter Rechenleistung beansprucht.

Nach unserem Ermessen spielt besonders bei der Entscheidung zwischen der eigenen GPU und GPU as a Service die benötigte Rechenleistung und die Verfügbarkeit des Services eine Rolle. Unter der Annahme, dass der Dienstleister des GPU as a Service permanent verfügbar ist. Für Prozesse, welche in Schüben erfolgen und große Rechenleistung benötigen, ist GPU as a Service ein attraktives Angebot. Diese können dadurch in kurzer Zeit abgeschlossen werden.

Eine eigene GPU ist attraktiver im Fall von konstant benötigter Rechenleistung, wenn eine verfügbare GPU in der Lage ist, die benötigte Rechenleistung zu erbringen. Andernfalls ist GPU as a Service ebenfalls eine Lösung, besonders in Zeiten, in denen die Preise von GPUs mit dem Halbleiternmangel gestiegen sind.

---

<sup>30</sup>Vgl. Loop und Blinn, "Real-time GPU rendering of piecewise algebraic surfaces".

<sup>31</sup>Vgl. Lattuada u. a., "Performance prediction of deep learning applications training in GPU as a service systems".

<sup>32</sup>Vgl. Wang u. a., "Quality of Service Support for Fine-Grained Sharing on GPUs".

<sup>33</sup>Vgl. Loop und Blinn, "Real-time GPU rendering of piecewise algebraic surfaces".

## **6 Marktvorhersage**

## **7 Fazit und Ausblick**

# Quellenverzeichnis

- Arslanian, Henri. *The Book of Crypto: The Complete Guide to Understanding Bitcoin, Cryptocurrencies and Digital Assets*. Springer eBook Collection. Cham: Springer International Publishing und Imprint Palgrave Macmillan, 2022.
- Clement, J. *Cloud gaming market value worldwide from 2019 to 2024*. URL: <https://www.statista.com/statistics/932758/cloud-gaming-market-world/> (besucht am 09.06.2022).
- Gabriele D'Angelo, Stefano Ferretti und Moreno Marzolla. *Cloud for Gaming*. Springer eBook Collection. Cham: Springer International Publishing Switzerland, 2015.
- Kords, Martin. *Weltweite Lieferung von Halbleiterprodukten (integrated circuit) für die Automobilindustrie von 2011 bis 2021*. URL: <https://de.statista.com/statistik/daten/studie/1288183/umfrage/halbleiterlieferungen-fuer-kraftfahrzeuge/> (besucht am 03.06.2022).
- Lattuada, Marco u. a. "Performance prediction of deep learning applications training in GPU as a service systems". In: *Cluster Computing* 25.2 (2022), S. 1279–1302.
- Loop, Charles und Jim Blinn. "Real-time GPU rendering of piecewise algebraic surfaces". In: *ACM SIGGRAPH 2006 Papers on - SIGGRAPH '06*. Hrsg. von John Finnegan und Julie Dorsey. New York, New York, USA: ACM Press, 2006, S. 664.
- McClean, Bill. *The 2022 McClean Report*. URL: <https://www.icinsights.com/news/bulletins/The-Real-Reason-Behind-The-Automotive-Industry-IC-ShortageA-StepFunction-Surge-In-Demand/> (besucht am 02.06.2022).
- Ojala, Arto und Pasi Tyrvaïnen. "Developing Cloud Business Models: A Case Study on Cloud Gaming". In: *IEEE Software* 28.4 (2011), S. 42–47.
- Rimol, Meghan. *Gartner Says Global Chip Shortage Expected to Persist Until Second Quarter of 2022*. URL: <https://www.gartner.com/en/newsroom/press-releases/2021-05-12-gartner-says-global-chip-shortage-expected-to-persist-until-second-quarter-of-2022> (besucht am 09.06.2022).
- Stadia, Google. *Bandbreite, Datennutzung und Streamingqualität*. URL: <https://support.google.com/stadia/answer/9607891?hl=de> (besucht am 09.06.2022).
- Statista-Research-Department. *Volumen der weltweiten Investitionen in Blockchain-Technologien und Kryptowährungen von 2018 bis 2021*. URL: <https://de.statista.com/statistik/daten/studie/1198230/umfrage/weltweite-investitionen-in-blockchain-technologien-und-kryptowaehrungen/> (besucht am 03.06.2022).
- Suznjevic, Mirko, Ivan Slivar und Lea Skorin-Kapov. "Analysis and QoE evaluation of cloud gaming service adaptation under different network conditions: The case of NVIDIA Ge-

- Force NOW". In: *2016 Eighth International Conference on Quality of Multimedia Experience (QoMEX)*. 2016, S. 1–6.
- Telekom. *5G Cloud Gaming*. URL: <https://www.telekom.de/unterwegs/was-ist-5g/5g-cloud-gaming> (besucht am 09.06.2022).
- Voas, Jeffrey, Nir Kshetri und Joanna F. DeFranco. "Scarcity and Global Insecurity: The Semiconductor Shortage". In: *IT Professional* 23.5 (2021), S. 78–82.
- Wang, Zhenning u. a. "Quality of Service Support for Fine-Grained Sharing on GPUs". In: *Proceedings of the 44th Annual International Symposium on Computer Architecture*. New York, NY, USA: ACM, 2017, S. 269–281.
- Wei Cai, Min Chen und Victor C.M. Leung. "Toward Gaming as a Service". In: *IEEE Internet Computing* 18.3 (2014), S. 12–18.
- Zadtootaghaj, Saman. *Quality of Experience Modeling for Cloud Gaming Services*. Springer eBook Collection. Cham: Springer Nature Switzerland AG, 2022.

# Ehrenwörtliche Erklärung

„Wir versichern, dass die vorliegende Arbeit von uns selbständig und ausschließlich unter Verwendung der angegebenen Quellen und Hilfsmittel angefertigt wurde. Alle Stellen, die wörtlich oder annähernd aus Veröffentlichungen entnommen sind, haben wir als solche kenntlich gemacht. Die Arbeit wurde bisher in gleicher oder ähnlicher Form, auch nicht in Teilen, keiner anderen Prüfungsbehörde vorgelegt und auch nicht veröffentlicht.“

Ort, Datum

Unterschrift

Ort, Datum

Unterschrift

Ort, Datum

Unterschrift