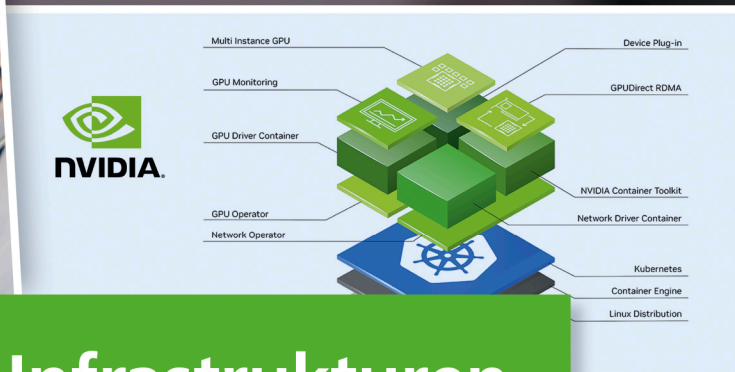


Oliver Liebel



Skalierbare KI/ML-Infrastrukturen

Evaluieren, Automatisieren, Praxis

- ▶ Kubernetes- und OpenShift-Cluster mit NVIDIAs Datacenter-GPUs
- ▶ Skalierbare und resiliente Infrastrukturen in der Cloud und On-Prem
- ▶ Vollautomation und Kosteneffizienz mit IaC und Operatoren

Kapitel 1

Vorwort

»Eine echte künstliche Intelligenz wäre intelligent genug, um nicht zu verraten, dass sie wirklich intelligent ist.«

– George Dyson, US-amerikanischer Wissenschafts- und Technikhistoriker

*»Hypes? Similar to as***les. Every day you encounter a new one.«*

– Shaun T., U.S. Special Forces (CAG), ret.

The Long Road – und ein etwas längeres Vorwort

Da wären wir, und es hat lange genug gedauert. Die Arbeit an diesem Buch hat bereits 2019 begonnen, parallel zu meiner 3. Container-Publikation, die Ende 2020 erschien. Und es war ein langer und oftmals und deutlich beschwerlicherer und enervierenderer Weg als die langjährige Arbeit an meinen Container-Themen, die für sich genommen schon einen hochvolatilen Sack sich permanent vermehrender und mutierender Flöhe darstellen, um diese einfache, bildhafte, aber passende Metaphorik zu benutzen.

Anfangs nahm ich noch an, dass sich das Buchprojekt *Skalierbare KI/ML-Infrastrukturen* halbwegs schnell abwickeln lassen würde. Schließlich war das geplante Volumen mit »nur« rund 400 Seiten gerade mal ein Drittel des Umfangs meiner letzten Container-Publikationen. Wie so oft im Leben stellen sich die Dinge in den meisten Fällen dann jedoch eher als komplexer denn einfacher heraus. Das *Warum* der Problematik ist mehreren Punkten geschuldet und zumindest aus einer sehr großen Flughöhe recht schnell erklärt: Ein Kernproblem liegt in der hohen Volatilität skalierbarer Container-Infrastrukturen selbst. Diese stellen jedoch im Unternehmensumfeld die einzig valide Foundation dar, um sowohl reguläre Arbeitslasten wie auch KI/ML-Stacks flexibel und skalierbar vollautomatisch, on demand und für jeden gewünschten Anwendungsfall und Workload provisionieren zu können. Und dieser KI/ML-relevante Stack muss nun noch zusätzlich in zwei Lokationen implementiert werden: im Unterbau (die GPUs in den Servern und deren Management-Systeme in den Hypervisoren) und *on top* im Workload-Layer der Containerisierungs-Plattform.

Klingt so weit erst einmal nicht unüberwindlich. Allerdings war – ergänzend zur ohnehin hohen Volatilität der Containerisierungs-Plattformen – die Entwicklung der GPU-spezifischen Baugruppen des Stacks durch NVIDIA, VMware und Red Hat in den letzten drei Jahren ebenfalls mit einer dermaßen hohen Geschwindigkeit und Volatilität der Komponenten und ihrer Features hinterlegt, dass ich nach Beendigung der Arbeit an einem Themenblock oft

genug den vorherigen, bereits fertiggestellten wieder überarbeiten oder gar komplett neu erstellen musste. Sisyphus lässt grüßen, aber wie üblich sucht sich jeder sein zu tragendes Päckchen meist selber aus.

Wie auch immer – die Dinge im skalierbaren KI/ML-Container-Land sind im zuletzt betrachteten Stand zumindest etwas ruhiger geworden, auch wenn die Entwicklung immer noch rasant ist, und es sind immer noch viele Ecken und Kanten zu meistern.

Der Stand der Dinge folgt nun.

Superlative, Geschwindigkeit, Hypes und die Realität

2022 – ein Jahr der Superlative. Alles ist groß. Aber zum Teil auch wieder klein. Größenwahnsinnige und zugleich kleingeistige Autokraten. Großflächige Corona-Durchseuchung mit hoffentlich kleinem Impact. Große Rezessions- und Inflationsängste, die sich hoffentlich nur zu einem kleinen Teil bewahrheiten. Große Nachfrage nach Gütern und Rohstoffen jedweder Art und ein mittlerweile oft viel zu kleines Angebot.

Aber nimmt man die Aussagen in allen Medien und den üblichen Politiker- und Consulter-Buzzword-Bullshit-Bingo-Runden abseits dessen zusammen, gibt es anscheinend im technischen Bereich nach wie vor nichts Größeres und Wichtigeres als Di-, Meta- oder Omniverses und Machine-Learning-/KI-Systeme oder allgemeiner:

Künstliche Intelligenz.

Nun, da zumindest in der Politik selten echte anzutreffen ist, könnte das theoretisch ein Hoffnungsschimmer sein.

Aber was ist mit den Unternehmen, die ihre Applikationsstacks immer noch nicht auf Basis von Big Data, KI und maschinellem Lernen betreiben? Denn die haben ja – zumindest gemäß dem schon viel zu lange anhaltenden KI-Hype- und Buzzword-Tsunami – mittel- bis langfristig ganz sicher keine Überlebenschance ohne hyper-agil entwickelte, self-optimizing and -repairing KI-Systeme. Die jeden Tag, jede Stunde aktualisiert oder am besten mehrmals pro Woche gleich ganz neu erfunden werden. Rückwärtskompatibilität ist ja auch sowas von gestern, oldschool und voll out. So wie alles, was noch on-prem und nicht in der Cloud läuft.

Ach ja? Ist das so?

Nun ja.

Im Grunde genommen spiegelt der ganze Faster-, Faster-, Faster-Nonsense im IT-Bereich ein Paradigma wider, in das sich unsere gesamte Gesellschaftskultur – leider, muss man sagen – unvermeidlich bewegt. Häppchen müssen nur noch mundgerecht serviert werden, egal ob's hier und da mal nicht so schmeckt oder passt, wie es soll. Hauptsache schnell, schnell, schnell. Trauriges Synonym einer fast schon demagogischen Lebens- und IT-Entwicklung, in der – zumindest aktuell – alles, was wir aufnehmen, so konzipiert und aufbereitet werden muss, dass auch Rezipienten mit einer maximalen Aufmerksamkeitsspanne von dreikommafünf Sekunden nicht sofort gelangweilt sind.

Viele böse Zungen behaupten ohnehin schon länger, und politisch völlig unkorrekt, dass »Agilität« in der IT wahrscheinlich doch nur eine freundliche Umschreibung sei, um Personen mit ADHS-Tendenzen in IT-Projekten unterzubringen. Wie üblich hängt das vom Standpunkt des Betrachters ab. Und Hypes sind erfahrungsgemäß in der Realität irgendwo zwischen dem verortet, was sie und ihre Fanboys versprechen, und dem, was eine (von ihren Verweigerern gern) überkritische Darstellung ihrer Mankos ist.

Viele Probleme – (k)eine echte Lösung?

Und wie üblich wird der zunehmende, großflächige und damit irgendwann relativ günstig verfügbare Einsatz einer innovativen und extrem leistungsfähigen Technik wie KI/ML auch abseits effizienter, legitimer und moralisch korrekter Anwendungen jede Menge Schattenseiten nach sich ziehen:

Denn solange wir auf die gleichen Technologien setzen, werden Leistungshunger und Abwärme von RZs mit KI-tauglichen GPU- und TPU- Systemen in jeder neuen Generation exponentiell durch die Decke gehen.

Aber hey, was soll's, werden sich (zu) viele grenzdebile Individuen mit viel Langeweile und/oder zu wenigen Hirnwindungen sagen: *Green IT ist voll gestern, Klimawandel Fake News und Greta nervt eh nur. Oder?* Und so kann man schließlich auf Knopfdruck und völlig on demand in High-End-Clouds, die minütlich Abwärme in der Größenordnung von New York verbrauchen, seine Nachbarn und/oder den Rest der Welt mit per Deepfake generierten, personalisierten Bildern und Videos verwerflichen Inhaltes beglücken, die auch gegen Whatever-Coinzahlungen ganz sicher nie wieder gelöscht werden.

Schalten wir den Sarkasmus mal wieder ab, aber Sie wissen so gut wie ich, dass dies erfahrungsgemäß nicht so weit abseits der Realität liegt, wie man denken möge. Und wir reden an dieser Stelle noch nicht einmal von vollautonomen Waffensystemen.

Dass die vollautomatische Bearbeitung von Prozessen auf Basis komplexer Daten nicht nur Vorteile bieten kann, sollte auf der Hand liegen. Allein schon die persönlichen Daten eines jeden Menschen bieten ein gewaltiges Potential, sowohl in positiver wie negativer Hinsicht: Stammdaten einer Person aus ihrer bisherigen Krankenakte, Fehlzeiten im Job, Bestellung (nicht)alkoholischer und/oder extrem zuckerhaltiger Getränkelieferungen per Internet, Postings der Person in sozialen Netzen z. B. über (extrem)sportliche Aktivitäten, erlittene physische oder psychische Krankheiten/Traumata oder (un)populäre Gesinnungen, dazu Gesundheits-Apps auf dem Handy oder der Smartwatch, welche die Telemetriedaten des letzten Joggings (oder Körperfunktionen permanent) übermitteln, die vielleicht nicht dem optimalen Vital-Standard entsprechen, und vieles andere mehr. Dies alles eignet sich »hervorragend« dazu, bestimmte Personen oder Risikogruppen beispielsweise automatisch in bestimmte Versicherungstarife einzuordnen, diesen Personen Kredit zu gewähren (oder nicht), sie bei Bewerbungen auszuschließen (oder nicht), bei Beförderungen zu berücksichtigen (oder nicht) und so weiter. Das Prinzip dürfte offenkundig sein.

Oder nehmen wir ein weiteres, mögliches Problemfeld, bei dem die Auswirkungen weitaus drastischer sein können – die vernetzte KI, wie sie bereits aktuell und zukünftig zunehmend auch herstellerübergreifend in KI-gesteuerten Fahrzeugen zum Einsatz kommt. Nehmen wir eine Situation an, in der zwei KI-gesteuerte Fahrzeuge mit hoher Geschwindigkeit unvermeidbar aufeinanderprallen werden, sofern keine Entscheidung der Autopiloten getroffen wird. Situationsbedingt kann nur ein Fahrzeug bzw. Insasse gerettet werden, die Person in dem anderen Fahrzeug muss mit hoher Wahrscheinlichkeit sterben.

In Fahrzeug 1 befindet sich eine 75-jährige Person, in Fahrzeug 2 ein Kind von 7 Jahren, das zur Schule gebracht wird. Wie entscheidet nun die vernetzte KI? Und exakt an dieser Stelle des Gedankenspiels kommen etliche Faktoren hinzu, die erst beim zweiten Blick auf das Szenario klar werden: Kennt das vernetzte KI-System die Krankenakte des Kindes und der älteren Person und weiß bereits, dass das Kind an einer schweren Krankheit leidet, mit einem garantiert tödlichen Verlauf innerhalb eines Jahres, während die ältere Person noch kerngesund ist?

Hat das KI-System gegebenenfalls Kenntnis darüber, dass die ältere Person eine wichtige Person des öffentlichen Lebens oder der Politik ist oder gar zum Vorstand genau des Konzerns gehört, von dem das Fahrzeug und/oder die KI-Systeme produziert wurden? Und wurde exakt für dieses Worst-Case-Szenario bereits eine versteckte Exception bzw. Priorisierung in die Entscheidungsalgorithmen eingebracht? Und was passiert, wenn der Fall »klassisch« andersherum liegt: Die alte Person ist krank, hat keine hohe Lebenserwartung mehr, die Rentenkassen sind leer, die Kosten der Krankenkasse für die ärztliche Versorgung der alten Person sind immens hoch. Das Kind jedoch gehört zu einer wohlhabenden, einflussreichen Familie und ist kerngesund.

Nun, gemäß einem nicht defekten moralischen Kompass sollte jedes Leben gleich viel wert sein. So weit die Theorie, aber allein das gerade gezeigte, kleine Gedankenspiel sollte verdeutlichen, dass die Dinge sehr schnell viel komplizierter werden können. Paart man diese Betrachtungen mit den – seit Anbeginn der Menschheit ständig und erfahrungsgemäß leider immer noch viel zu häufig anzutreffenden – Charakterzügen wie Egoismus und Gier auf jeder Hierarchieebene, wird viel schneller ein Schuh daraus, als man denken mag.

Und dabei reden wir noch nicht einmal von potentiellen Schwachstellen, durch welche z. B. gehackte, autonome Schwerlastler als Terrorwaffen missbraucht werden könnten. Oder von militärischen Drohnen der nächsten Generation mit KI-gestützten Waffensystemen, die zukünftig innerhalb bestimmter Parameter auch ohne die finale Human-Authorization in Sekundenbruchteilen Entscheidungen treffen können. Sie tun dies auf Basis der ihnen vorliegenden Daten sowie der Daten, mit denen diese Systeme trainiert wurden, und müssen dann im ungünstigsten Fall im Millisekunden-Bereich eine Entscheidung über Nutzen vs. Kollateralschaden treffen.

Und wer denkt, dass diese kleinen Gedankenspiele fernab der Realität liegen – dem ist nicht so. Wir ahnen bereits, wie komplex das Thema abseits der üblichen und oft leider nur rein technischen Betrachtung tatsächlich zu sehen ist.

Foundation – und »mehr Power«

Fakt ist: Die Verarbeitung der Daten, sowohl für die Trainings von KI-Systemen als auch im späteren Einsatz, ist unstrittig hochkomplex. Aber auch wenn alle Data Scientists und Data Engineers, Programmierer und auch alle anderen, die sich um die Vorbereitung und KI/ML-Modellerstellung kümmern, höchst sorgfältig arbeiten und alle Daten maximal transparent und zudem in allen relevanten Belangen korrekt sind, kommt noch ein wichtiger, entscheidender Faktor hinzu.

Die Systeme, die die Daten prozessieren. Leistungsfähige und vollautomatisch skalierbare KI-Infrastrukturen. Hocheffiziente High-Performance-GPUs in Container-Clustern, die dynamisch mit DPUs und CPUs in Software-Defined Datacenters zu einer logischen Super-Compute-Unit verschmelzen, um allen involvierten Teams des Unternehmens das Leben leichter zu machen und nicht komplizierter. Und dem Kunden die Informationen liefern, die er bestellt hat bzw. benötigt – und zwar möglichst ohne einen Prediction-Error, der im ungünstigsten Fall vielleicht Einfluss auf Leib und Leben haben könnte.

Wie es der KI-Forscher Richard Sutton in der Einleitung seines Artikels *The Bitter Lesson* 2019 trefflich formulierte:

»Die wichtigste Lektion, die man aus 70 Jahren KI-Forschung ziehen kann, ist, dass allgemeine Methoden, die schlichtweg Rechenleistung nutzen, letztlich die effektivsten sind, und zwar mit großem Abstand. [...] Auf der Suche nach einer Verbesserung, die kurzfristig einen Unterschied macht, versuchen Forscher, ihr menschliches Fachbereichs-Wissen zu nutzen, aber das Einzige, was auf lange Sicht zählt, ist massive Rechenleistung.«

Oder kurz: *Viel hilft viel*, zumindest im Moment und zumindest, was die Rechenleistung von KI-Systemen angeht. Und damit sind wir bei der zwingend erforderlichen Skalierbarkeit dieser Systeme.

Und genau diese Skalierbarkeit von containerisiert arbeitenden KI/ML-Infrastrukturen war mein Trigger, dieses Buch zu schreiben, da es sehr eng mit meinem primären (und seit fast 30 Jahren Operations-lastigen) Schaffensfeld verheiratet ist: hochskalierbare Container-Cluster-Infrastrukturen.

Dabei ging es mir nicht darum, ein weiteres Fachbuch der Kategorie *Prima, ich habe auf GCP/AWS/AKS/EKS/GKE/Whatever einen Python-Code-Snippet in mein containerisiertes Jupyter-Notebook mit TensorFlow gepastet und es tut irgendwas ...* beizusteuern, zusätzlich zu den gefühlten Millionen Exemplaren, die auf dieser Welle bereits mitschwimmen.

Hypes und die eher selten auseinanderbrechende Realität

Mir ging es vor allem darum, als jemand, der seit fast drei Jahrzehnten sehr tief in hochkomplexen Themengebieten der IT verwurzelt ist, zu analysieren, was wirklich hinter dem ganzen Hype steckt. Kritische und unbequeme Fragen abseits von oftmals blindem Opportunismus und maximal enerzierend gebetsmühlenartig repetierten »Mit KI/ML lösen wir alle Probleme«-Mantras zu stellen.

Denn das Konzept, einen KI/ML-Insel-Stack von einer KI-Workstation eines Data Scientists *mal eben* in eine auto-skalierbare, produktivtaugliche und Multi-Tenant-fähige Cluster-Umgebung zu verpflanzen, erfordert mehr als eine Hollywood-geprägte Illusion von KI-Systemen und ein paar knallige, cloudaffine Buzzwords beim dritten Cappuccino im legeren Entscheider-Meeting.

Was wird auf der Infrastruktur-Seite wirklich benötigt, um KI/ML-Systeme flexibel, effizient, stabil und sicher zu betreiben? Und zu welchem Preis? Es geht darum, zu hinterfragen und zu analysieren, *was* implementierungs- (ausdrücklich nicht coding-) und infrastrukturtechnisch *wie* für *wen* »geht« – und was nicht. Und vor allem: mit welchem *Automationsgrad*. Denn das ist der wichtigste Schlüssel für performante, skalierbare, containerisierte KI/ML-Cluster. Und zu schauen, wo die Reise hinführen kann. Mit allen positiven und leider auch nicht wenigen negativen Aspekten.

Es geht um die explizite Beleuchtung von oftmals nicht unerheblichen Implementierungsproblemen, insbesondere für Unternehmen, die sensible Daten ihrer KI/ML-Workflows in der nach wie vor hochbedenklichen Sicherheit von Everything-happy-Everything-a-a-S-ML-Clouds prozessieren möchten – und im Cloud-Anwendungsfall noch dazu permanent Kosten auftürmen, die jenseits von Gut und Böse liegen. Aber auch im Self-Hosted-KI/ML-Bereich ist leider nicht mehr viel los mit eitel Sonnenschein, sondern in vielen Teilen eher der permanente Link nach */dev/null* für Dollars, Euros und Effizienz angesagt – wenn sorgfältige Planung fehlt oder einmal mehr auf Entscheider-Ebene dank hart antrainierter Beratungsresistenz gekonnt ignoriert wird.

Es geht darum, wie hoch der Grad der *Vollautomation* für containerisierte, skalierbare KI/ML-Cluster im betrachteten Stand wirklich ist bzw. sein kann. Denn nur das zählt in ernstzunehmenden Unternehmensimplementierungen. KI-Systeme kann seit einigen Monaten oder Jahren plötzlich jeder implementieren, wenn man dem Geschwurbel der Medien und üblichen Verdächtigen lauscht. Wie üblich liegen zwischen Anspruch und Wahrheit in der Regel oft Welten. Greift man tiefer, besitzen nur die wenigsten das Know-how um die echte, hohe Kunst: die Implementierung massivst skalierbarer Container-Cluster für KI/ML-Applikationen mit maximaler Vollautomatisierung auf jedem Level, von der Infrastruktur bis zum operatorgestützten Rollout der (v)GPU-Ressourcen und ML-Stacks, egal ob on-premises oder in der Cloud.

Es geht um Fragen und Antworten zu Konzepten, ROI und LTS. Mit einem konsequenten Fokus auf potentielle Strategien für ML-Infrastrukturen in Unternehmen. Mit Betrachtungen und Analysen zu Auswirkungen auf das Unternehmen durch eine Umstrukturierung auf containergestützte KI/ML-Systeme. Ob sich der Einsatz von Machine-Learning-Systemen für das eigene Unternehmen rentiert. Und wenn ja, in welchem Umfang.

Und es geht darum, dass auch diejenigen von uns, die Python nicht ad hoc mit vierhundert Anschlägen pro Minute coden können, und auch die, die bei Begriffen wie TensorFlow, Keras, PyTorch und Ähnlichem nicht sofort Speichelabsonderungen produzieren, eine realistische

Einschätzung darüber erhalten, was KI/ML-Infrastruktur-technisch im betrachteten Stand im Unternehmensumfeld umgesetzt werden kann. Und was in Grimms Märchenstunde bzw. -Tonne bzw. den großen, ewigen Schwurbel-Bullshit-Leitfaden aller Sales-Fraktionen gehört.

Faster, faster ... und kein Ende. Leider auch für die Kosten.

Denn insbesondere in der immer schnelllebigeren Container- und KI/ML-Welt haben Informationen rund um Hardware, Software, Verfahren und Konzepte leider eine immer kürzere Halbwertszeit.

Eine höchst bedauerliche Problematik, auf die ich bereits in meinen letzten vier Container-Publikationen mehr als ausdrücklich hingewiesen hatte und an der sich bis heute herzlich wenig geändert hat. Im Gegenteil. Aber was soll's – denn den Sales-Fraktionen von NVIDIA, AWS, GCP, Red Hat, VMware und Co. bereitet das meist keine mächtigen Kopfschmerzen. Allen übrigen Beteiligten sehr wohl. Denn im Grunde gilt nur noch eine Regel: dass nichts mehr gilt. Oder in Langform: Das meiste, was gestern noch superhip und State of the Art war, ist heute schon oft genug nur noch alter Krempel, der keinen mehr wirklich interessiert.

Und das trifft auf fast alle Beteiligten zu – mit Ausnahme des Endkunden. Denn der würde sich gern mal eine etwas langsamere Pace wünschen, Innovation hin oder her.

Fatalerweise müssen sich die Unternehmen aber ebendieser hochvolatilen Gemengelage permanent neu stellen und anpassen, und dies oft zu einem – im wahrsten Wortsinne – hohen Preis.

Kein anderes IT-Geschäftsfeld ist im betrachteten Stand mit derartig hohen Kosten hinterlegt. Auf der anderen Seite: Die relevanten Unternehmen, die den Zug verpassen, werden sich irgendwann keine Sorgen mehr um Kosten machen müssen bzw. können. Im Grunde heißt es wie üblich »am Ball bleiben« oder wenigstens »endlich einsteigen« – aber es wird von Jahr zu Jahr erfahrungsgemäß eher schwieriger denn einfacher:

Kleinere Unternehmen machen sich meist Gedanken, wie die eigene Infrastruktur mit möglichst geringem finanziellem Aufwand KI/ML-tauglich gemacht werden kann, z. B. durch Aufrüstung einiger On-Prem-Systeme mit zusätzlichen GPUs und verstärkter Kühlung und ergänzend den einen oder anderen lastintensiven Trainings-Workload in der Cloud, stecken dann aber erfahrungsgemäß leider allzu oft in einer teuren Pseudo-Kreativ-Sackgasse voller loser Enden fest.

Etliche Mittelständler pendeln zum Teil nach wie vor in fantastisch-kreativem Herumgeeiere eher unentschlossen zwischen Cloud und On-Prem, Letzteres mit dedizierten, gekauften oder gemieteten GPU-Servern oder aufgebohrten Bestandsservern in den Unternehmens-RZ, ohne wirklich nach vorne zu kommen.

Und im Konzernbereich sieht die Sachlage auf den ersten Blick oft klarer aus, aber seien Sie versichert – das ist sie meist nicht. Viele sind überstürzt in die scheinbare und trügerische

Sicherheit »KI/ML? Alles kein Problem. Vertrauen Sie uns ...« verschiedener Cloud-Provider (man will ja schön diversifizieren) abgerauscht und kämpfen nun mit providerspezifischen Implementierungsproblemen, mangelnder Konfigurierbarkeit der GPU-Nodes, den allgegenwärtigen cloudtypischen und -bedingten Sicherheitsproblemen und nicht zuletzt sehr hohen Kosten.

Und wieder die unbequeme Realität, aber hübsch aufbereitet

Womit wir wieder bei der allgegenwärtigen Cloud und ihren Anbietern wären. Und selbst die müssen sich mehr und mehr der harten Realität stellen, welche uns seit vielen Monden und in echter 24/7-HA nonstop Pandemien, Kriege, Krisen, Rohstoff- und Energieknappheit und einiges mehr zur täglichen Erbauung beschert. Aber die Provider verbiegen die Realität wie üblich mit einem (jedoch nur auf den ersten Blick) gekonnten Griff in die Werbe-Trickkiste.

Microsoft zog Mitte 2022 als Erster die Handbremse und kündigte ab sofort eine längere Nutzungsdauer für Cloud-Hardware an. Konkret verlängern die Redmonder ab 2023 die Gesamtbetriebsdauer ihrer CPU- und GPU-Server bis zum Austausch gegen neuere Systeme von vier auf sechs Jahre. Das maximale sinnfreie Werbe-Blabla, dass für den nun verlängerten Betriebszeitraum »Investitionen in unsere Software [erfolgen], die den Betrieb unserer Server- und Netzwerkausrüstung effizienter machen« (Zitat Microsoft, untermalt von brüllendem Gelächter gestandener RZ-Admins weltweit), gehört sicher in jede Vorstandssitzung, aber technisch betrachtet auch nur nach `/dev/null`.

Aber was soll's: Nicht zuletzt die Shareholder wird es freuen, da Microsoft dank der verlängerten Abschreibungsdauer allein im Fiskaljahr 2023 Einsparungen von fast 4 Milliarden erwartet. Auch Google (Verlängerung von drei auf vier Jahre) und AWS (Verlängerung von vier auf fünf Jahre) haben längst ähnliche Schritte angekündigt.

Neben allen Finanztricks, um besser durch die Krise zu kommen, bedeutet das jedoch auch unter dem Strich, dass Unternehmen, die ihre KI/ML-Strecken ganz oder in Teilen in der Cloud betreiben, schlichtweg länger auf CPUs und GPUs neuester Bauart verzichten müssen und damit auf eine gegebenenfalls deutlich bessere Performance und/oder Energieeffizienz.

Ja, stimmt: Abstrakt betrachtet, klingt das alles zunächst sehr nach einer Diskussion Pest vs. Cholera. Ich könnte jetzt sagen: *Alles gut, Freunde – keine Sorge, dem ist nicht so.*

Aber das würde leider nicht der Realität entsprechen.

Denn es geht im Grunde eigentlich und letztlich nur noch darum, wie man die IT des eigenen Unternehmens (sofern KI/ML-Szenarien dort von hoher wirtschaftlicher Relevanz sind, was auf immer mehr Unternehmensbereiche zutrifft) in diesen technisch, betriebs- und finanzwirtschaftlich hypervolatilten Zeiten strategisch so aufstellt, dass der geringstmögliche wirtschaftliche Schaden entsteht.

Denn seien Sie auch bezogen auf diesen Punkt versichert: Wenn Sie in der Welt der KI- und ML-Cluster-Infrastrukturen ein oder zweimal design- und entscheidungstechnisch falsch abgebo-gen sind, kann das bereits die sprichwörtliche Road to Financial Ruin sein.

Falls Sie meinen, dies wäre eher eine Übertreibung – willkommen im realen Leben. Es gab in den vergangenen Jahren genügend Unternehmen, die diese leidvolle und teure Erfahrung machen mussten.

In eigener Sache

Und damit die Leserinnen und Leser, die mich noch nicht aus meinen acht bisherigen Publikati-onen kennen, wissen, von wem diese Aussagen kommen: Ich bin jemand, den andere vielleicht als Spezialisten im RZ-/Cloud-/Großkunden-Bereich für High-Availability-Cluster, Software-Defined Storage, Verzeichnisdienste sowie für hochskalierbare, vollautomatisierte Container-Cluster und GPU-beschleunigte Microservice-Infrastrukturen von großen Unternehmen und international operierenden Konzernen bezeichnen würden. Als Systemarchitekten und oft genug auch als Problem-Fixer für bestimmte Bereiche der IT von Unternehmen und Konzernen.

Aber im Grunde bin ich unter dem Strich nichts anderes als ein IT-Veteran, der unzählige Hypes hat kommen und gehen sehen. Und genau daher geht es mir, wie in all meinen bishe-rigen Publikationen, vor allem darum, eine möglichst realistische Einschätzung abzuliefern. Darüber, wo wir implementierungstechnisch im Bereich der KI/ML-Infrastrukturen wirklich stehen, ob und wie für das eigene Unternehmen ein gangbarer Weg durch den KI-Infra-Dschungel gefunden werden kann und für wen sich der Einsatz dieser Systeme lohnen kann.

Aber selbst im Vergleich zu meinen vier bisherigen Publikationen rund um das ebenfalls hochvolatile Themengebiet Container-Cluster kann jede Publikation, die sich im betrachte-ten Stand seriös mit der Thematik skalierbarer KI/ML-Cluster-Infrastrukturen beschäftigt, trotz aller Genauigkeit immer nur eines sein – eine Momentaufnahme. Und genau deswegen habe ich dieses Mal einen etwas anderen Ansatz gewählt, nämlich einen, der etwas weniger mit praktischen Beispielen und etwas mehr mit Theorie und vor allem Strategie hinterlegt ist.

Denn wie ich im Laufe der zähen und zum Teil leidvollen Recherchen und Tests (danke an dieser Stelle noch einmal an die vielen netten Menschen von NVIDIA, insbesondere Erik Bohnhorst, die mir dabei fast drei Jahre lang zur Seite standen) für dieses Buch erkennen musste, sind zwar praktische Betrachtungen, also das *Doing*, absolut essentiell und unab-dingbar, um bestimmte Sachverhalte erstmalig zu verinnerlichen.

Jedoch spielt in diesem Fall der pure Technik-Kontext, d. h. die Konfigurationen bis hin zum letzten kleinen Bit, strategisch eine etwas untergeordnete Rolle. Weil er bzw. es sich unglaub-lich schnell verändert – noch schneller als in »reinen« Container-Clustern, die für sich genommen schon volatil genug sind und hier lediglich die Foundation für die KI/ML-Stacks darstellen.

Und betrachtet man die Entwicklung im RZ- und Cloud-Bereich einmal ohne alte Muster und mühsam etablierte Scheuklappen, so ist der gute alte Server doch schon lange nicht mehr die Compute-Unit: Es ist längst das bereits erwähnte, vollautomatisierte *Software-Defined Data-center*, mit jeder Menge CPUs, GPUs, TPUs und DPUs.

Was im Bereich von skalierbaren, containerisierten KI/ML-Infrastrukturen wirklich zählt, ist insbesondere das Verständnis, wie welche Komponenten in welchem speziellen Szenario arbeiten und wie sie ineinandergreifen. Ist das verinnerlicht, ist das Verständnis für die Executive-Komponenten von Hard- (GPUs/TPUs) und Software (Container-Cluster und Operatoren) relativ schnell auf den nächsten Evolutions-Level adaptiert. Und damit portierbar.

Ja, geschenkt, Freunde – pünktlich zum Release von NVIDIAs erster Quanten-GPU mit Ultra-Freon-Kühlung gibt's dann eine leicht aktualisierte Auflage.

Wie auch immer, gehen wir's an. Denn die Arbeit erledigt sich leider noch nicht von selbst.

Aber genau daran arbeiten wir.

Und wer das Motto lieber im Buzzword-Sprech für die Vorstandsrunde hätte:

Hyperautomation mit KI.

Ist ganz sicher bis zum Ende dieser Woche noch superhip und voll angesagt.

Danach, naja ...

Ah, der Cappuccino ist fertig.

Danksagungen und Widmung

Mein besonderer Dank geht an viele Mitarbeiter von NVIDIA aus dem Datacenter-GPU-Bereich weltweit, die mich in den letzten Jahren tatkräftig unterstützt haben, insbesondere: Erik Bohnhorst, Christopher Desiniotis, Milan Diebel, Francis Guillier, Shiva Krishna Merla, Rajeshka Rao, Thomas Remmlinger, Thomas Wernicke und viele andere mehr, die hier nicht explizit benannt sind. Ebenso gilt mein Dank dem Red Hat EMEA Partner Team.

Für T.Y. – in tiefstem Dank für Deine zeitlosen, unerschütterlichen Weisheiten.

1.1 Vorbemerkungen

Im Folgenden finden sich vermehrt Anglizismen (welche in der Regel auch erklärt werden), die sich im Rahmen aktueller und fachspezifisch fortgeschrittener Applikationen/Publikationen ohne kilometerlange und oft unpassende deutsche Um- und Beschreibungen oft nicht umgehen lassen.

Kapitel 4

NVIDIA-Datacenter-GPUs und mehr – technischer Background

»9.7 TFLOPS for FP64? Sounds great, but, eh – can it run DOOM?«

– Typischer Foren-Kommentar bei so ziemlich jeder Ankündigung einer neuen High-End-GPU-Karte

In diesem Abschnitt befassen wir uns konkreter mit der Hardware, sprich – den NVIDIA-Datacenter-tauglichen GPUs, den jeweils dafür geeigneten Bereitstellungsverfahren, ihrer Performance und etlichem anderen mehr. Mit dem Bereitstellungsverfahren ist im einfachsten Fall beispielsweise das dedizierte Durchreichen einer GPU an eine VM gemeint oder softwarebasierte GPU-Partitionierungen per Virtual GPU (vGPU), hardwarebasierte Partitionierung (MIG) mit einer besseren Tenancy, GPU-Sharing oder Kopplungen multipler GPUs, lokal oder über das Netzwerk, oder auch verschiedene Kombinationen der vorgenannten Verfahren, um unterschiedlichste Anwendungsfälle abdecken zu können.

4.1 NVIDIA und ML-Cluster

Nun, wo fängt man an? Vielleicht mit einer Frage:

Lieber NVIDIA-Aktien anstelle von Goldbarren? Nein, das soll durchaus kein polemischer Kommentar zu Anlagestrategien sein, sondern liegt näher an der Realität, als Sie vielleicht denken mögen, wie im Folgenden erklärt werden wird. Denn kaum ein Unternehmen hat sich strategisch in den letzten Jahren in eine ähnliche monopolartige Position gebracht, wenn es um KI/ML-Belange im RZ-Bereich geht, wie NVIDIA.

NVIDIAs GPU-Accelerators wie die neue Hopper-Generation oder ihre Vorläufer wie die Ampere-Generation oder deren Tesla-Vorgänger (Turing-Generation) sind aus Rechenzentren kaum noch wegzudenken. Datacenter-taugliche GPUs werden – je nach GPU-Modell – zwar auch noch in bestimmten Unternehmen für Remote-Desktop-Szenarien, also »echte« grafische Workloads eingesetzt, in zunehmender Anzahl jedoch primär nur noch für *Accelerated Computing*, also KI/ML-Einsätze mit GPU-Beschleunigung.

Dabei existieren verschiedene Möglichkeiten, die Rechenleistung der GPU zu nutzen. Während für reine ML-Entwicklungs-Workstations beispielsweise dedizierte High-End-Grafik-

karten wie eine High-End-Karte der RTX-Serie per *Passthrough* angebunden werden können, bleibt das Manko, dass die jeweilige GPU nur von *einem* Nutzer bzw. exakt einer realen oder virtuellen Maschine für ML-Workloads verwendet werden kann. Für geclusterte und skalierbare Szenarien in Unternehmen wird also eine andere Herangehensweise benötigt.

Insbesondere Kubernetes-/OpenShift-basierte Container-Cluster mit Pod- und Cluster-Autoscalern sind daher eigentlich prädestiniert, um bei entsprechender Lastanforderung automatisch in die Breite zu skalieren, um so für alle (KI/ML-)Anwendungsfälle genügend Rechenkapazitäten bereitstellen zu können. So können z. B. multiple GPU-fähige Anwendungen automatisch horizontal skaliert werden, um große Datenlasten parallel prozessieren zu können – natürlich vorausgesetzt, dass die Anwendung bzw. Architektur dies unterstützt.

Logischerweise gibt es auch alternative Verfahren wie das bereits erwähnte, dedizierte Passthrough (1:1-Bindung zwischen VM und GPU-Accelerator), die zwar technisch ebenfalls greifen, jedoch in Clustern (siehe das Thema »Skalierung«) für Accelerated Computing selbst mit Verfahren wie VMwares *Dynamic DirectPath IO* (siehe Abschnitt 4.5.3) nicht wirklich effizient sind. Es werden daher *shareable* GPU-Ressourcen benötigt, und hier startet die Gelddruckmaschine für Cloud-Betreiber, Virtualisierer – und nicht zuletzt NVIDIA.

Betrachten wir das »Warum«.

Cloud

Nahezu alle Cloud-Provider lassen sich mittlerweile ihre VM-Instanztypen mit durchgeschleifter GPU/TPU-Accelerator-Hardware hervorragend vergüten, wie in Abschnitt 3.2 bereits betrachtet. Sicher, die üblicherweise horrenden Preise für RZ-taugliche GPU-Accelerators stellen auch für Cloud-Provider zunächst eine Investition dar, die sich jedoch auf deren Seite sehr schnell amortisiert. (v)GPU-Manager, die NVIDIA-GPUs softwaretechnisch auf (Public-)Cloud-Hypervisoren wie GCP oder AWS partitionieren könnten, waren im betrachteten Stand für Endverbraucher nicht verfügbar.

Virtualisierer

Nehmen wir VMware als Beispiel und die in den meisten Unternehmens-RZ präsenten vSphere-Systeme: Schon interessant, dass gerade das *vGPU*-Feature, das anfänglich kaum Beachtung fand, nun nur noch in VMwares höchster Lizenz-Preisklasse zu finden ist. Ohne vSphere Enterprise Plus-Lizenz läuft im Bereich »*vGPU-Computing*« (vorab: softwarebasierte Partitionierung einer GPU) rein gar nichts. Und betrachten wir den »Aufwand«, der dahintersteckt, ist die Verortung dieses Features in der höchsten Lizenz-Preisklasse ohne jedweden technischen Hintergrund, sondern rein gewinnorientiert motiviert. Denn es handelt sich auf ESXi-Seite lediglich um die Integration eines einzigen VIBs, das zudem noch von NVIDIA geliefert wird. Dagegen ist prinzipiell nichts einzuwenden, das ist freie Marktwirtschaft. Aber dennoch wäre hier seitens VMware eine etwas niedrigere finanzielle Einstiegshürde gerade für KMU durchaus wünschenswert gewesen. Aber wozu den Profit minimieren, wenn es die GPU-Quelle #1 nicht anders vormacht.

NVIDIA

Sicher, leistungsstarke GPU-Karten waren noch nie günstig, und einige Hardcore-Gamer würden wahrscheinlich lieber eine Hypothek aufnehmen und/oder Organe spenden, als das neue Spiel XYZ nicht in Ultra-Super-HD-256K-with-unbelievable-high-FPS zu spielen. Aber erst seit ihrer neu erstrahlenden Popularität im KI/ML-Bereich schießen die Preise für »echte« GPU-Accelerators wie NVIDIAs Ampere- oder Tesla-Modelle (oder RTX-Karten für ML-Workstations) durch die Decke – fünfstelligen Beträge für RZ-taugliche Karten der Oberliga sind hier keine Seltenheit, sondern das Normalprogramm, und komplette DGX-Systeme liegen wie bereits betrachtet noch deutlich darüber. Wer denkt, dass sein Unternehmen mit gemieteten, cloudgehosteten, RZ-tauglichen GPUs auf Dauer günstiger fährt – eine nette Idee, die oft nur dann aufgeht, wenn es sich nicht um 24/7-Dauerbetrieb der Systeme handelt.

4.2 Partitionierte GPUs mit NVIDIAs vGPU und MIG

Betrachten wir an dieser Stelle zunächst einige grundlegende Begrifflichkeiten rund um die möglichen Bereitstellungsarten RZ-tauglicher GPU-Hardware, da sich diese zum Teil grundlegend von »normalen« Verfahren im HEDT/Workstation-Bereich unterscheiden und wichtig für das spätere Verständnis der durchzuführenden Setup-Tasks sind.

Ein wichtiger Aspekt ist dabei die Möglichkeit, GPUs zu *partitionieren* und damit ihre immense Rechenleistung, die sonst gegebenenfalls nicht vollumfänglich genutzt wird, an multiple Client-Prozesse zu verteilen. Dazu stehen verschiedene Verfahren zur Verfügung, die wir im Folgenden betrachten.

Eine weitere Variante, die seit dem GPU-Operator 1.11 zur Verfügung steht, um Däumchen drehende GPUs in KI/ML-Clustern zu vermeiden, ist die neue Option des GPU-Sharings/Overcommitments (z. B. physikalisch 1 GPU, aber n Client-Prozesse dürfen die GPU beanspruchen). Dies wird konkret ab Abschnitt 10.10 betrachtet.

4.2.1 Vorbetrachtungen und Scope

Die primär für grafische Anwendungsfälle gedachten GPU-Modi *vSGA* und *vDGA* werden im Folgenden nicht angesprochen, da sie für den Scope des Buches irrelevant sind.

Für skalierbaren Datacenter-Betrieb mit Fokus auf Compute-Anwendungen (KI/ML) und multiple Tenants sind üblicherweise nur die Modi *vGPU* (vGPU = virtual GPU, »Software«-partitionierbare GPUs) und *MIG* (MIG = Multi Instance GPUs = »Hardware«-partitionierbare GPUs, die üblicherweise auch den vGPU Modus unterstützen) geeignet und von Interesse, in wenigen Ausnahmefällen auch ein PCI-Passthrough von ganzen GPUs ohne Partitionen.

Daher liegt der Fokus aller folgenden Betrachtungen auf den GPU-Bereitstellungsmodi vGPU und MIG und den damit verbundenen Verfahren und Konzepten für einen effizienten Datacenter Betrieb.

4.2.2 Vorbetrachtungen: Partitionierte GPUs mit vGPU und MIG

Die Partitionierung einer physikalischen GPU ist besonders effizient für Workloads, die die Rechenkapazität der physikalischen GPU nicht vollständig auslasten würden. Je nach Leistungsklasse der GPU und dem konsumierenden Workload wird in der Praxis daher oftmals keine vollständige GPU benötigt. Dies gilt insbesondere dann, wenn es sich um kleinere PoC/Modelle oder Demos handelt, die auf einem Boliden wie einer A100 mit 40 oder gar 80 GB GPU-RAM betrieben werden sollen.

Ohne die Möglichkeit einer GPU-Partitionierung wird jeder ML-Workload die komplette GPU okkupieren, unabhängig davon, ob er sie tatsächlich verwendet bzw. voll auslastet ... oder eben nicht. Insbesondere containerisierte Compute-Workloads im RZ/Cloud-Umfeld auf Kubernetes-/OpenShift-Clustern profitieren von der höheren Effizienz einer GPU-Partitionierung sowie der im Fall von MIG damit einhergehenden, echten Mandanten-Fähigkeit bzw. Multi-Tenancy und damit höheren Sicherheit.

NVIDIAS vGPU-(Virtual-GPU-)Software ermöglicht bereits seit langem die softwarebasierte/temporäre Partitionierung von geeigneten NVIDIA-GPUs und unterstützt auch eine breite Palette an weit verbreiteten und relativ kostengünstigen Datacenter-GPUs, wie z. B. Tesla T4-Instanzen. Für den vGPU-Modus wird auf Hypervisor-Ebene ein lizenzpflichtiger GPU-Manager benötigt.

Der *MIG-Mode (Multi-Instance GPU)*, der mit NVIDIAS Ampere-GPU-Generation und Modellen wie der A100 oder A30 erstmals zur Verfügung stand, ist zwar ebenfalls ein Betriebsmodus um GPUs zu partitionieren, aber im Gegensatz zum vGPU Mode bringt dieses Verfahren eine deutlich effizientere Multi-Tenancy, die jedoch im Vergleich zum vGPU-Mode meist mit einer höheren Idle Power einhergeht, siehe Abschnitt 5.4.4. MIG-fähige Karten können in der Regel auch eine Partitionierung mittels vGPU anbieten. Die (vorab vereinfacht und abstrakt ausgedrückt) GPU-Core-Tenancy ist dann jedoch nicht mehr aktiv. Alle Details zum MIG-Modus finden Sie ab Abschnitt 4.4.

Während sich im *vGPU-Modus (Time-Sliced, CPU-like Scheduling)* alle GPU-Partitionen im gleichen Hardware-Pool der GPU-Cores-Engine bedienen und nur den GPU-RAM partitionieren, sorgt der *MIG-Mode* für explizite Separierung der multiplen (im betrachteten Stand bis zu sieben) GPU-Instanzen. Dadurch können Workloads besser aufgeteilt und zudem die Datensicherheit erhöht werden. Sie hierzu auch: <https://www.nvidia.com/de-de/technologies/multi-instance-gpu/>

MIG-fähige Karten werden von NVIDIAS GPU-Operator ab Version 1.7 unterstützt.

4.2.3 NVIDIAs vGPU und (leider noch kein) Cloud-Einsatz

Für alle folgenden Betrachtungen ist es vorab essentiell wichtig zu verstehen, dass *NVIDIAs vGPU-Software* (der GPU-Manager) auf gängigen Public-Cloud-Plattformen wie GCP, AWS oder Azure typischerweise *nicht* zur Verfügung steht. Der GPU-Manager (und damit vGPU) steht nur für Hypervisoren zur Verfügung, wie sie in Unternehmens-RZ bzw. Private Clouds eingesetzt werden, wie z. B. vSphere, Citrix, RHEL KVM, Nutanix und andere.

Hintergrund: Der *GPU-Manager* wird bereits auf der Ebene des Hypervisors zwingend für vGPU benötigt, damit über ihn GPUs softwaretechnisch partitioniert und an die VMs durchgereicht werden können. So viel Zugriff kann oder will im betrachteten Stand keiner der Public-Cloud-Riesen dem Endkunden anbieten. In logischer Konsequenz existiert daher (bisher) auch keine GPU-Manager-Software seitens NVIDIA für AWS, GCP und Co. Somit bleiben den Unternehmen, die GPU-Workloads in der Cloud prozessieren möchten oder müssen, nur zwei Alternativen.

Entweder werden nicht-MIG-fähige Karten, wie z. B. T4, A10 oder A40 komplett (en bloc, d. h. 1:1) via Passthrough direkt an den Consumer (die KI/ML-Applikation) durchgereicht oder es kommen – wenn Tenancy gefragt ist und/oder multiple Consumer gleichzeitig auf einer großen, leistungsfähigen GPU wie der A100 sicher bedient werden sollen – MIG-fähige Karten zum Einsatz. Diese werden vom Cloud-Provider typischerweise per Passthrough durchgereicht, können aber vom Kunden bzw. dem GPU-Operator-Stack zumindest auf dem betreffenden Node in GPU-Slices partitioniert und damit von n Applikationen verwendet werden.

Auch wenn die Passthrough/MIG-Only-Diktatur der Cloud-Provider die Konfigurationsflexibilität deutlich einschränkt – einen kleinen Vorteil für den Endkunden bietet die vGPU-Absistenz im Cloud-Umfeld dennoch: Die vGPU-Lizenzkosten entfallen.

4.3 vGPU – Virtual GPU

Schauen wir nun etwas tiefer in Konzepte und Arbeitsweise der vGPU-Partitionierung. Im Folgenden wird zunächst betrachtet, was genau unter vGPU zu verstehen ist, welche Produkte NVIDIA in diesem Bereich anbietet und wie vGPU konzeptionell und technisch betrachtet arbeitet.

4.3.1 Generelle vGPU-Architektur

Betrachten wir zunächst die generelle vGPU-Architektur anhand Abbildung 4.1, die eine generelle, vereinfachte Darstellung der Arbeitsweise einer vGPU und ihre Aufteilung an multiple VM-Instanzen zeigt.

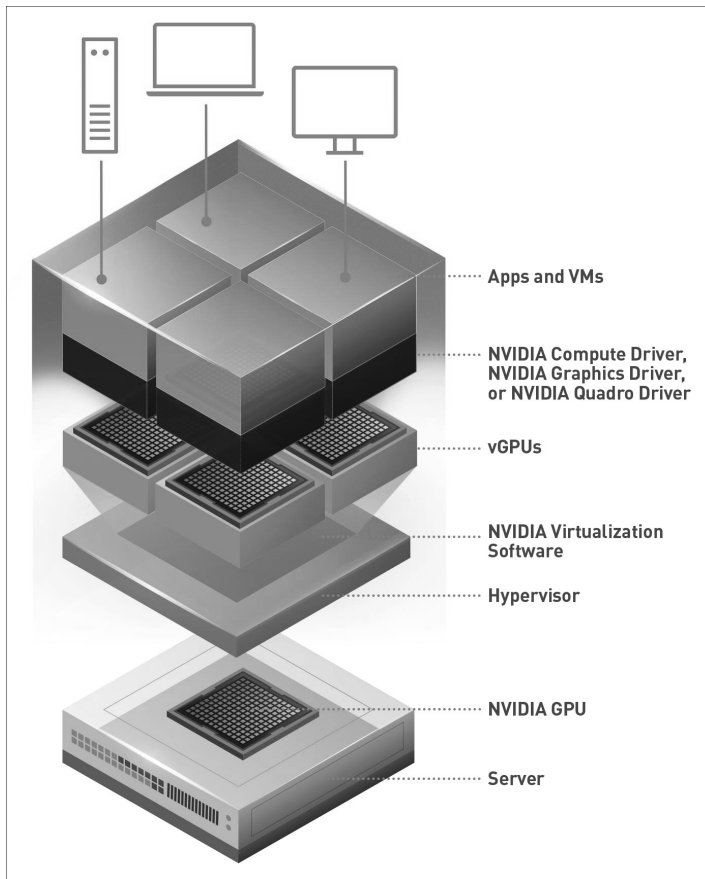


Abbildung 4.1 NVIDIAs vGPU-Stack in stark vereinfachter, konzeptioneller Form, Quelle: NVIDIA

4.3.2 Details zur Funktionsweise

Wie in Abbildung 4.1 bereits zu erkennen ist, wird die vGPU-Funktionalität als zusätzliche Software-Komponente zwingend über den jeweils verwendeten *Hypervisor* bereitgestellt. Im betrachteten Stand wurden u. a. folgende (in der Regel Typ1-)Hypervisoren/Bare-Metal-OS unterstützt (Nennung hier ohne Versionen): vSphere, OpenStack, Windows Server/Hyper-V, RHEL KVM/RHV, XenServer.

Es ist ebenfalls zu beachten, dass Ampere-Karten wie z. B. die A30 oder ihr größeres Pendant, die A100 oder die neue H100, die per *MIG* (siehe ab Abschnitt 4.4) auch hardwarebasiert (korrekter wäre »räumlich«/»EEPROM-like«) partitioniert werden können, ebenfalls rein softwarebasiert im *vGPU*-Mode betrieben werden können. Wann welcher Modus (vGPU, MIG oder MIG-backed vGPU) zum Einsatz kommt, bestimmen üblicherweise die (gegebenenfalls vorgegebene) RZ-Architektur, Security-Anforderungen und der konkrete, technische Einsatzzweck.

Bleiben wir zunächst bei der *vGPU-Partitionierung*. Jede GPU-Karte stellt eine bestimmte Anzahl an Rechenkernen und verfügbarem Grafikspeicher (RAM) zur Verfügung. Diese Compute-Ressourcen können nun an eine VM (oder eben partitioniert an mehrere) durchgeschleift werden. Während es sich bei den – von der jeweiligen VM abgerufenen – (v)GPU-Compute-Ressourcen wie üblich um eine komprimierbare Compute-Ressource handelt (vergleiche CPU/GPU-Shares), stellen die RAM-Ressourcen wie üblich ein nicht komprimierbares Hard-Limit dar. Ein einfaches Rechenbeispiel:

Reichen wir eine vGPU-fähige Karte mit 16 GB RAM an vier VMs bei identischer Partitionierung durch, so erhält jede VM exakt 4 GB GPU-RAM. Damit sind die nicht komprimierbaren bzw. nicht over-commit-fähigen Compute-Ressourcen dieser GPU-Karte erschöpft.

Achtung: vGPU und GPU-Pool-Sharing

Dabei ist jedoch zu beachten, dass die partitionierten GPU-Slices auf einer nicht-MIG-, sondern »nur« vGPU-fähigen Karte nicht wirklich voneinander isoliert sind. Alle vGPU-Partitionen z. B. einer Tesla T4 oder A10 nutzen die GPU-Ressourcen wie einen »Shared«-Pool. Wer zwingend echte *Tenancy* braucht, wird um MIG-fähige Karten nicht herkommen.

Das vGPU-Modell bietet dabei die funktionalen Eigenschaften, die Sie in Abbildung 4.2 sehen.

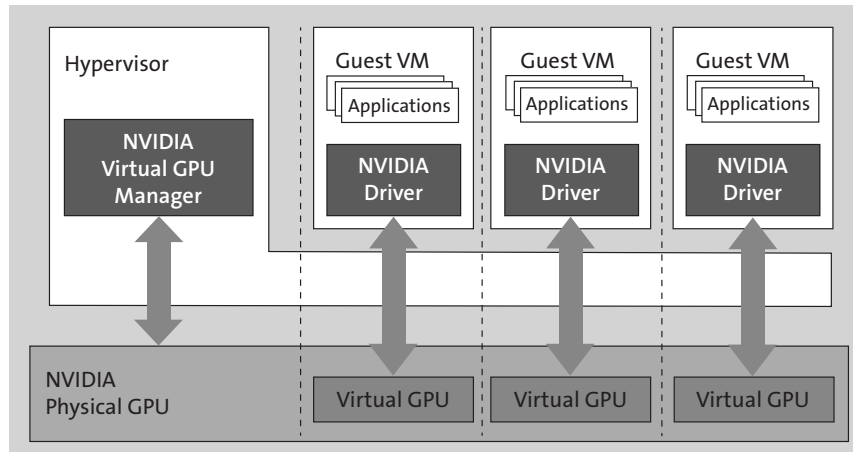


Abbildung 4.2 NVIDIA-vGPU-Systemarchitektur – High Level, Quelle: NVIDIA

Jede NVIDIA-vGPU kann analog zu einer herkömmlichen GPU betrachtet werden: Sie besitzt eine fixe/statische Menge an GPU-Framebuffer (RAM der GPU-Karte) und einen oder mehrere virtuelle Displayausgänge oder auch Heads.

Die virtuelle GPU (vGPU) erhält zum Zeitpunkt ihrer Bereitstellung exklusiv einen Teil des Framebuffers der realen GPU (oder je nach Profil gegebenenfalls auch den ganzen). Diesen verwendet die virtuelle GPU so lange, bis sie deaktiviert wird (z. B. Shutdown der VM).

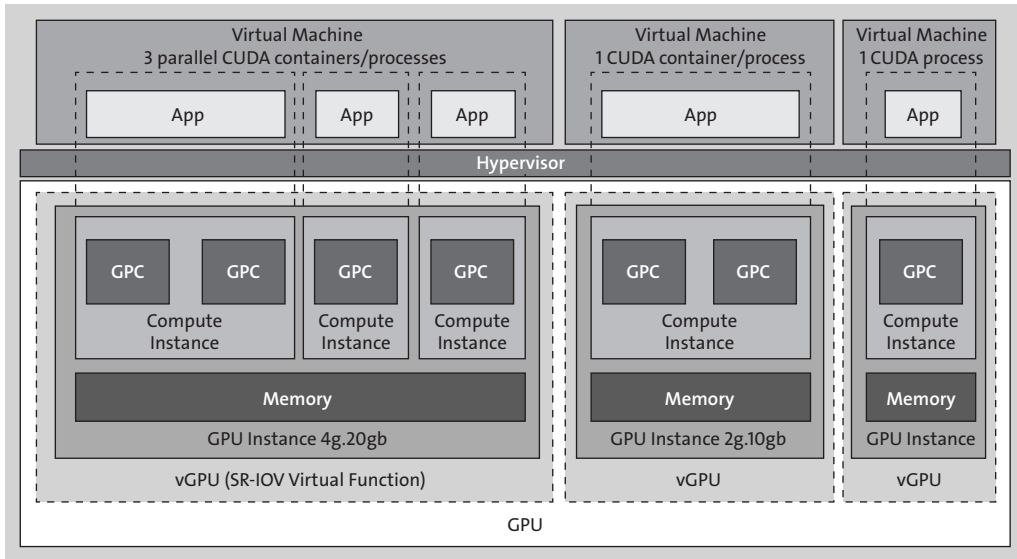


Abbildung 4.3 NVIDIA-vGPU-Systemarchitektur – Details, Quelle: NVIDIA

4.3.3 Time-Sliced vGPU

Der vGPU-Modus (ohne explizite Abschottung im GPU-Engine-Layer) ist mit jeder NVIDIA-vGPU-fähigen Karte verfügbar (siehe Abbildung 4.4).

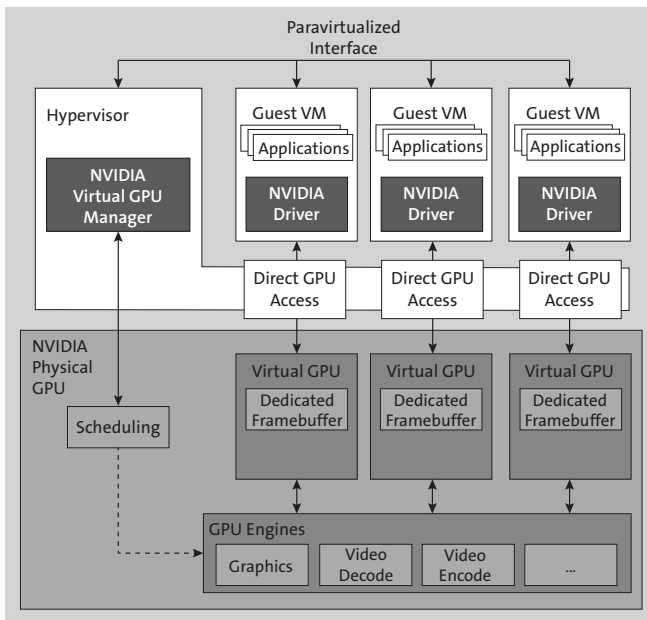


Abbildung 4.4 NVIDIA-vGPU-Systemarchitektur – Time-Sliced, Quelle: NVIDIA

vGPU-fähige Karten

Eine Auflistung vGPU-fähiger Karten von NVIDIA findet sich unter anderem in der Dokumentation NVIDIAs, wird aber auch noch in anderen Abschnitten des Buches thematisiert (u. a. Abschnitt 5.3.1): <https://docs.nvidia.com/grid/gpus-supported-by-vgpu.html>.

4.3.4 Passthrough GPU vs. vGPU im Hypervisor

Abbildung 4.5 zeigt im direkten Vergleich, wie GPUs direkt per Passthrough oder per vGPUs auf einem vSphere-System an die darüberliegenden virtuellen Clients durchgereicht werden.

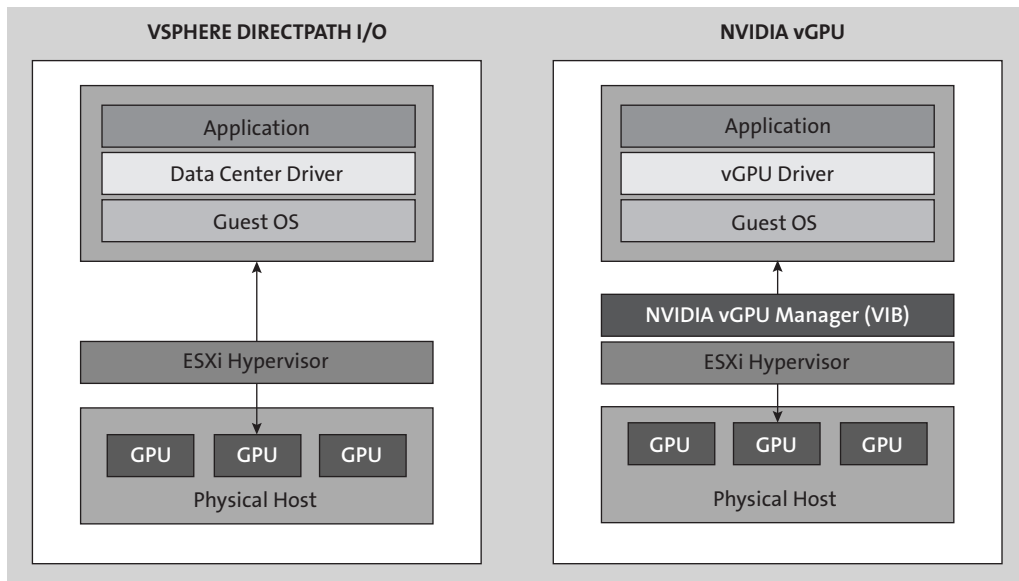


Abbildung 4.5 Passthrough GPU vs. vGPU-Workflow zwischen VM und ESXi-Hypervisor, Quelle: VMware

Im Passthrough-Mode (links) wird das PCI-Device über den Hypervisor dediziert an einen Consumer (sprich: VM) durchgereicht. Soll stattdessen der vGPU-Mode (rechts) zum Einsatz kommen, wird auf dem Hypervisor ein GPU-Manager in Form eines VIB installiert. Dieses registriert die verbauten NVIDIA-GPUs und stellt bei der Einrichtung einer VM oder eines Templates die Möglichkeit zur Verfügung, nur bestimmte Partitionen einer GPU an eine VM durchzureichen. In der VM muss bei diesem Setup-Typ auch ein passender, vGPU-fähiger Treiber installiert werden.

Diese Installation der GPU-Treiber (egal ob normaler Driver oder vGPU) erfolgt in containerisierten Umgebungen üblicherweise vollautomatisch über den GPU-Operator, der ab Abschnitt 9.4 besprochen wird.

4.3.5 vGPU-Produkte

Die Rubrik *vGPU* umfasste gemäß der Definition von NVIDIA im Stand 07/2022 die folgenden Produkte:

- ▶ *NVIDIA Virtual Applications (vApps)* – Dieses vGPU-Produkt ist z. B. für Unternehmen gedacht, die Citrix Virtual Apps und Desktops bereitstellen oder RDSH oder andere App-Streaming- oder sessionbasierte Lösungen einsetzen. Generell: für Anwendungen auf PC-Ebene und serverbasierte Desktops geeignet.
- ▶ *NVIDIA Virtual PC (vPC)* – Für Anwendungsfälle, in denen virtuelle Remote-Desktops (»virtuelle PCs«) mit PC-Windows-Anwendungen, Browsern und High-Definition-Video benötigt werden.
- ▶ *NVIDIA RTX Virtual Workstation (vWS)* – Für den Einsatzfall, remote grafische Systeme (z. B. CAD/CAE) mit maximaler Leistung zu betreiben.
- ▶ *NVIDIA Virtual Compute Server (vCS)* – Der von uns primär betrachtete Anwendungsfall: HPC-/KI-/ML-/Deep-Learning-Workloads ohne grafische Ausgabe (*Display-less*), die z. B. in VMs/Containern auf Hypervisoren (vSphere, KVM etc.) mit ECC-Unterstützung betrieben werden können.

Siehe dazu auch: <https://www.nvidia.com/content/dam/en-zz/Solutions/design-visualization/solutions/resources/documents1/Virtual-GPU-Packaging-and-Licensing-Guide.pdf>.

4.3.6 vGPU-Versionen

Es ist essenziell wichtig zu verstehen, dass *NVIDIAs vGPU-Software-Version* (im betrachteten Stand 08/2022: vGPU-Version 14.2) immer der Dreh- und Angelpunkt für alle – von dem jeweiligen Release – unterstützten Features ist. Sie allein entscheidet, welche Version des *vGPU-Managers* (Hypervisor) und des *Linux- und Windows-Drivers* unterstützt wird und welche jeweiligen Features auf welcher GPU-Hardware und Hypervisor-Plattform supportet werden:

- ▶ <https://docs.nvidia.com/grid/index.html>
- ▶ <https://docs.nvidia.com/grid/latest/product-support-matrix/index.html>

LTS und (Short-Term) Production Branch

Wie in allen produktivtauglichen Konzepten bietet NVIDIA neben *Short-Term Releases* (unverständlicherweise von NVIDIA als *Production Branch* gekennzeichnet) auch *LTS Releases* (*Long Term Support Branch*). Leider zieht NVIDIA hier keinen »roten Faden« bezogen auf die Nummerierung (etwa durch gerade und ungerade Nummern) durch.

Einige der aktuelleren, ungeraden Release-Nummern werden für die LTS-Versionen verwendet (11, 13), einige der geraden (12, 14) für die Short-Term Branches. In älteren Versionen passt dieses Schema nicht, daher kontrollieren Sie immer den verwendeten Branch gemäß der Support-Matrix.

Ein vGPU *LTS Branch* hatte im betrachteten Stand einen Supportzeitraum von drei Jahren: So ist der EOL für den vGPU 13.x Branch, der im August 2021 gestartet wurde, im August 2024 erreicht. Die *Short-Term Branches* hingegen haben in der Regel nur einen maximalen Supportzeitraum von einem Jahr.

Hinweis: NVIDIA AI Enterprise (NVAIE-Versionen)

Bitte beachten Sie bereits an dieser Stelle, dass sich vGPU- und NVAIE-Software (Letztere ab Kapitel 6 erläutert) in Art, Aufbau und Lizenzierung (siehe ab Kapitel 7) unterscheiden. Die von NVAIE unterstützten Produkte und Features finden sich hier:

<https://docs.nvidia.com/ai-enterprise/latest/release-notes/index.html>.

4.3.7 vGPU-Features je nach Produkt

Tabelle 4.1 listet die vGPU-Features je nach gewähltem/lizenziertem Produkt auf. Bezogen auf den Scope des Buches ist für uns primär nur die letzte Spalte der Tabelle (vCS, virtual Compute Server) von Relevanz.

Feature	NVIDIA vApps	NVIDIA vPC	NVIDIA vWS	NVIDIA vCS
License Entitlement				
Concurrent User (CCU)	✓	✓	✓	
Per GPU				✓
Capability Entitlement				
Desktop Virtualization		✓	✓	
RDSH App Hosting	✓	✓	✓	
RDSH Desktop Hosting	✓	✓	✓	
Compute Virtualization			✓	✓
Windows Guest OS Support	✓	✓	✓	

Tabelle 4.1 vGPU-Features und -Lizenzen

Feature	NVIDIA vApps	NVIDIA vPC	NVIDIA vWS	NVIDIA vCS
Linux Guest OS Support		✓	✓	✓
Maximum Displays	1	4	4	1
Maximum Resolution	1280 × 1024	5120 × 2880 (5K)	7680 × 4320 (8K)	4.096 × 2160 (4K)
NVIDIA RTX Enterprise Software Features			✓	
OpenGL, DirectX and Vulkan	✓	✓	✓	In-situ Graphics only
CUDA and OpenCL Support			✓	✓
ECC and Page Retirement			✓	✓
Multi-vGPU			✓	✓
NVLink			✓	✓
GPU Passthrough Support	✓		✓	✓
Bare Metal Support	✓		✓	
vGPU Profile Sizes Supported				
512 MB		✓	✓	
1 GB	✓	✓	✓	
2 GB	✓	✓	✓	
3 GB	✓		✓	
4 GB	✓		✓	✓
5 GB				✓
6 GB	✓		✓	✓

Tabelle 4.1 vGPU-Features und -Lizenzen (Forts.)

Inhalt

1	Vorwort	19
1.1	Vorbemerkungen	28
1.1.1	Verwendete Formatierungen	29
1.1.2	Breites Buch-/Seitenformat	29
1.1.3	Klartext	29
1.1.4	KI/ML-Begrifflichkeiten	30
1.1.5	Weiterführende Hinweise	30
1.1.6	Verwendete Testsysteme	30
1.1.7	Im Buch verwendete Grafiken	31
1.2	Was dieses Buch sein bzw. nicht sein soll	31
1.2.1	Was es sein soll	31
1.2.2	Was es nicht sein soll – und nicht ist	32
1.2.3	Scope und Fokus des Buches	32
1.2.4	Wissensaufbau	33
1.3	Wie dieses Buch zu lesen ist	34
1.4	Thematischer Überblick – was wird in welchen Kapiteln behandelt	34

TEIL I Technische Foundations zu skalierbaren KI/ML-Infrastrukturen

2	Am Anfang war die Dunkelheit	41
2.1	Eine kurze Einführung: KI/ML-Systeme – und alles wird gut. Oder eher nicht?	42
2.1.1	Historisches – kurz und kompakt	42
2.1.2	KI as a panacea?	43
2.1.3	Eine kurze Einordnung: KI, Machine Learning, neuronale Netze und Deep Learning	44
2.2	Use Cases für KI/ ML-Anwendungen – Auszüge	45
2.2.1	Wer profitiert vom Einsatz von KI/ML-Systemen? Mögliche Use Cases im Überblick	45
2.2.2	Exemplarische Use Cases	47

2.3 Fehlerfreie KI? Sicher nicht.	50
2.3.1 Regeln und Transparenz	50
2.3.2 Lösungsansätze	50
2.3.3 Vorbereitung	51
2.4 Einige Grundbegrifflichkeiten im KI/ML-Kontext	52
2.4.1 Machine Learning: Training und Inference	52
2.4.2 CNN (Convolutional Neural Networks)	53
2.4.3 Alles fließend: FP/TF/BF (Floating Point)	54
2.4.4 GPUs und Parallel-Computing	55
2.4.5 CPUs mit ML-Erweiterungen	56
2.4.6 CUDA (Cores)	57
2.4.7 Tensor (Cores)	58
2.4.8 Präzision, Performance und Kosten	60
2.4.9 Tensor-Core-Effizienz und Mixed/Reduced Precision	60
2.4.10 CUDA-Cores vs. Tensor-Cores	62
2.4.11 Und noch einmal Performance: NVIDIA Hopper, TMA, Transformer-Engine und FP8	62

3 High-Level-Vorbetrachtungen zur Implementierung von skalierbaren KI/ML-Infrastrukturen 65

3.1 Bare-Metal, Virtualisierung, Containerisierung	65
3.1.1 Bare-Metal vs. Virtualisierung	65
3.1.2 Containerisierung	66
3.1.3 Die Kernkomponenten/-Layer des (auto-)skalierbaren KI/ML-Infra-Stacks	68
3.2 Generelle Infrastruktur-Fragen: Cloud vs. On-Prem, Managed Server, hybrider Mischbetrieb, dedizierte KI-Plattformen (NVIDIA DGX)	69
3.2.1 Implementierungs- und Kostenfaktoren in der Cloud	69
3.2.2 Exkurs: Managed Server kleinerer SPs als günstigere Cloud-Alternative mit höherer Flexibilität?	72
3.2.3 Implementierungs- und Kostenfaktoren: Self-Hosted	72
3.2.4 Datensicherheit	75
3.2.5 Storage	75
3.2.6 Netzwerk	76
3.2.7 Hybrider Ansatz: On-Prem und Cloud (Pay-per-Use)	76
3.2.8 Alles cool? In der Cloud oft eher nicht. Temperatur-, Performance- und damit Kostenfragen.	77
3.2.9 Generelle Funktions- und Lizenzkosten-Betrachtungen: vGPU vs. MIG	79

3.2.10	Für größere Budgets: Out-of-the-Box-, Ready-to-use-ML-Server (NVIDIA DGX)	81
3.2.11	DGX: Technische Eckdaten und Blick unter die Haube	82
3.2.12	HPE ML – und wieder NVIDIA	84
3.2.13	Miete von RZ-tauglicher Hardware und Bereitstellung im eigenen RZ	84
3.3	Entscheidungshilfe: Reguläre GPU-Server, KI/ML-Boliden wie DGX oder alles in die Cloud?	85
3.3.1	KMU	85
3.3.2	Größere Unternehmen und Konzerne	86
3.4	Generelle GPU-Hardware-Fragen: NVIDIA vs. AMD vs. Intel vs. Googles TPU	86
3.4.1	Vorbetrachtungen: Was darf es denn sein?	86
3.4.2	GPU vs. TPU ... oder doch gemeinsam?	88
3.4.3	NVIDIA	88
3.4.4	AMD	90
3.4.5	Intel	91
3.4.6	Fazit: GPU-Provider	92

4 NVIDIA-Datacenter-GPUs und mehr – technischer Background 93

4.1	NVIDIA und ML-Cluster	93
4.2	Partitionierte GPUs mit NVIDIAs vGPU und MIG	95
4.2.1	Vorbetrachtungen und Scope	95
4.2.2	Vorbetrachtungen: Partitionierte GPUs mit vGPU und MIG	96
4.2.3	NVIDIAs vGPU und (leider noch kein) Cloud-Einsatz	97
4.3	vGPU – Virtual GPU	97
4.3.1	Generelle vGPU-Architektur	97
4.3.2	Details zur Funktionsweise	98
4.3.3	Time-Sliced vGPU	100
4.3.4	Passthrough GPU vs. vGPU im Hypervisor	101
4.3.5	vGPU-Produkte	102
4.3.6	vGPU-Versionen	102
4.3.7	vGPU-Features je nach Produkt	103
4.3.8	vGPU-Arbeitsweise (konzeptionell)	105
4.3.9	Scheduling-Policies von vGPU	105
4.3.10	vGPU-Profile und Zuordnung (exemplarisch: NVIDIA A100 40 GB)	108
4.3.11	Konkrete vGPU-Profile und Details	109
4.3.12	Exemplarische vGPU-Partitionslayouts	111

4.3.13	Erforderliche vGPU-Lizenzen und Entitlements je nach Modell und Typ	112
4.3.14	Übersicht der vGPU-Modi nach GPUs/Karten	112
4.4	MIG – Multi-Instance GPU	113
4.4.1	MIG-fähige GPUs	114
4.4.2	MIG-Konzepte, Terminologien und technische Details	115
4.4.3	Allgemeine MIG-Tech-Specs	118
4.4.4	MIG-Instanzen und -Partitionen	120
4.4.5	Compute-Instanzen/Compute-Sub-Partitionierung	121
4.4.6	MIG-Profilübersichten (A100 und A30)	122
4.4.7	Der A100-MIG-Black-Hole-Effekt, oder: die verschwundenen Partitionen	124
4.4.8	MIG-Strategien	126
4.4.9	MIG-Exposition: gesamte GPU via Passthrough oder MIG-Partitionen per VM-Template?	127
4.5	MIG: Multi-Tenancy revisited	128
4.5.1	Modelle und Konzepte, Vor- und Nachteile, plattformspezifische Limitierungen	128
4.5.2	MIG, echte Tenancy – und Next-Gen-MIG mit Confidential Computing	129
4.5.3	Nicht nur MIG: Das Problem, Passthrough GPUs (auto-)skalierbar anzubieten	133
4.6	Technische Daten und Preise ausgewählter NVIDIA Datacenter-GPUs	134
4.6.1	Supported NVIDIA GPUs Optimized for Compute (AI/ML) Workloads	135
4.6.2	Supported NVIDIA GPUs Optimized for Mixed Workloads	136
4.7	GPU-Time-Slicing und GPU-Overcommitment	137
4.7.1	Theoretische Vorbetrachtungen	137
4.7.2	Konzepte zur Umsetzung	139
4.8	NVLink und NVSwitch: GPU Big Blocks – Bündelung multipler GPUs	139
4.8.1	NVLink	139
4.9	GPUDirect (RDMA)	142
4.9.1	GPUDirect Storage	142
4.9.2	GPUDirect RDMA	143
4.10	GPU-Performance in ML-Trainings – Bare-Metal vs. vGPU/MIG	145
4.10.1	Vorbetrachtungen	145
4.10.2	Konkretes Setup und Messwerte	146
4.11	NVIDIA-Datacenter-Produkte: The Road Ahead	149

TEIL II Implementierung von skalierbaren KI/ML-Infrastrukturen

5	Implementierung: vSphere als Hypervisor für skalierbare ML-Infrastrukturen	153
5.1	Hardware-Voraussetzungen und Vorbetrachtungen (vSphere/On-Prem)	153
5.2	Preflights	154
5.2.1	BIOS/UEFI-Settings, SR-IOV, vSphere Edition, DRS	154
5.2.2	vSphere 7 und ESXi-Patchlevel	155
5.2.3	Update 7U3 und VMClasses	156
5.2.4	Host-Updates für vSphere/ESXi	156
5.2.5	vGPU für RTX 6000/8000 und RTX A5000/A6000 aktivieren	156
5.2.6	ECC-Memory	157
5.2.7	Preflight-Checks: Tools und Tests	157
5.2.8	Virtualization-Mode (Achtung: wichtig!)	159
5.3	Setup des GPU-Managers/vGPU-Host-Drivers (ESXi/vSphere 7)	160
5.3.1	NVD-AIE oder NVD-VGPU, NVIDIA vGPU Certified Server	160
5.3.2	Setup des NVIDIA-VIB (vGPU-Manager) auf den ESXis	161
5.3.3	PoC: Einfaches Passthrough	162
5.4	VM-Templates mit GPUs erstellen	164
5.4.1	Erstellung und Konfiguration eines VM-Templates (vGPU-Variante, OpenShift)	165
5.4.2	Erstellung und Konfiguration eines VM-Templates (MIG-backed vGPU-Variante, OpenShift)	167
5.4.3	Erstellung eines skalierbaren PCI-Passthrough-VM-Templates (identische GPUs per Dynamic DirectPath IO und Hardware-Bezeichner)	169
5.4.4	Die GPU-Power-Modi (P0–P8) und (Idle-)Leistungsaufnahme	170
5.4.5	Checkliste für mögliche Fehler beim vGPU-Betrieb	170
5.5	MIG-Mode auf dem Hypervisor aktivieren	171
5.5.1	Setup – Vorbereitungen	171
5.5.2	GPU (ohne Reboot des ESXi-Hosts) auf MIG-Mode umstellen	173
5.5.3	Manuelle Partitionierung anlegen (nur zur Veranschaulichung)	176

6	Der NVIDIA AI Enterprise (NVAIE)-Stack – infrastrukturelevante Betrachtungen	181
<hr/>		
6.1	Vorbetrachtungen	181
6.2	Motivation	182
6.3	Plattformen für NVAIE	183
6.4	NVAIE vs. vGPU vs. Free GPU Operator	185
6.4.1	VIBs	185
6.4.2	GPU-Operator	185
6.4.3	NVAIE-Features	186
6.5	NVAIE in der Public Cloud	186
6.6	NVAIE ist Pflicht für skalierbare ML-Cluster?	187
6.7	NVAIE als AI-End-to-End-Plattform	187
7	vGPU-/NVAIE-Preflights: Lizenzierung	189
<hr/>		
7.1	Grundsätzliches: vGPU- vs. NVAIE-Lizenzen und DLS vs. CLS	189
7.1.1	Preise und SLAs: vGPU	190
7.1.2	Preise und SLAs: NVAIE	191
7.1.3	NVIDIA-Entitlement beantragen	193
7.2	NVIDIA Licensing System (NLS)	194
7.2.1	Vorbetrachtungen	194
7.2.2	License Server: Self-Hosted License Server und Alternativen	195
7.3	License Server: DLS vs. CLS	196
7.4	Self-Hosted License Server: DLS und Legacy License Server	197
7.4.1	License Server – Legacy-Variante (self-hosted)	197
7.4.2	License Server – DLS-Variante (DLS Virtual Appliance, On-Prem)	200
7.4.3	Troubleshooting – Token Debugging	204
7.5	Cloud-Hosted License Server: CLS	205
7.5.1	License Server – CLS-Variante (cloudbasiert)	205

8	Kubernetes-basierte Plattformen für skalierbare, GPU-Accelerated KI/ML-Cluster	209
8.1	The Road so far	209
8.2	Generelle Plattform-Fragen: (Vanilla-)Kubernetes-Derivate und OpenShift im Überblick	211
8.3	Vanilla Kubernetes	213
8.3.1	Test and Play	213
8.3.2	Benötigte 3rd Party Tools und asynchrone Produktzyklen	214
8.3.3	Vanilla Kubernetes und das traurige Thema LTS: Geld verbrennen? Oder besser doch nicht?	214
8.3.4	Releases, Changes und kein Ende	215
8.3.5	Vanilla Kubernetes und TTM-Märchenstunden	216
8.3.6	AKS, EKS, GKE & Co.	217
8.4	VMwares Tanzu und das Eckige, das durchs Runde soll	217
8.4.1	Historisches	217
8.4.2	Tanzu	218
8.5	OpenShift	219
8.6	Abschließende LTS-Betrachtungen	220
8.7	Kubernetes-Basics – Aufbau des Systems	222
8.7.1	Kernkomponenten und Konzepte	222
8.7.2	Kubernetes-spezifische Dienste auf den Master-Nodes (Controlplane)	224
8.7.3	Kubernetes-spezifische Dienste auf allen Nodes	224
8.8	Kubernetes-Basics – Ressourcen/Workloads	225
8.8.1	Kubernetes/OpenShift: API-Version und API-Ressourcen	225
8.8.2	Namespaces	226
8.8.3	Pods	227
8.8.4	Pod-Metahüllen: Deployments, StatefulSets, DaemonSets	228
8.8.5	ConfigMaps	232
8.8.6	Node-Objekte und Node-Label	232
8.8.7	Services	233
8.8.8	Ingress und Routen	234
8.9	Sonstige im Folgenden verwendete, Kubernetes-spezifische Tools	236
8.9.1	kubectl- und oc-Bash-Completion und kubectl-Alias	236
8.9.2	Kustomize und Helm	238

9	Preflights für GPU-Accelerated Container-Cluster: Operatoren	241
9.1	Generelle Vorbetrachtungen zum Thema Operatoren	241
9.1.1	Einführung	242
9.1.2	Was ist ein Operator?	242
9.1.3	Horizontal? Vertikal? Beides?	245
9.1.4	Controller-Loops	245
9.1.5	Operator-Kategorien	246
9.1.6	Red Hats Operator Framework und Operator-SDK	246
9.2	Operator-Typen und Maturitäts-Level: Helm vs. Ansible vs. Go	247
9.2.1	Operator-Maturitäts-Level und -Kategorien	247
9.2.2	Operator-Build	248
9.2.3	Operatorhub.io und OpenShift-Operatoren	249
9.3	Die wichtige Rolle von Operatoren im auto-skalierbaren KI/ML-Stack	250
9.3.1	Team-Play	250
9.3.2	Der GPU-Operator – Vergangenheit und Zukunft	250
9.4	NVIDIAs GPU-Operator – die Architektur	251
9.4.1	Das Gesamtkonstrukt	251
9.4.2	GPU-Operator: Unterstützte Hypervisoren und GPUs	252
9.4.3	All-in-One	253
9.4.4	GPU-Operator und DGX	253
9.4.5	Die Einzelkomponenten des GPU-Operators im High-Level-Überblick	253
9.4.6	Preflight: der NFD-Operator	256
9.4.7	Die Einzelkomponenten des GPU-Operators im Detail	257
9.5	Automatische Provisionierung eines Nodes durch den GPU-Operator	258
9.5.1	K8s-Device-Plugin	259
9.5.2	GPU Feature Discovery	259
9.5.3	Driver	261
9.5.4	Container-Toolkit	261
9.5.5	DCGM/DCGM-Exporter	261
9.5.6	Der MIG-Manager	261
9.5.7	MIG-Manager und assoziierte ConfigMap	264
9.5.8	MIG-Strategies: »mixed« vs. »single« in der Praxis	266
9.5.9	Custom-MIG-ConfigMap	267
9.6	NVIDIAs Network-Operator – die Architektur	268
9.6.1	Vorbetrachtungen und Übersicht	268
9.6.2	Arbeitsweise (High-Level)	269

9.7	Komponenten des Network-Operators im Überblick	270
9.7.1	Mofed (NVIDIA_MLNX_OFED) Driver	270
9.7.2	Kubernetes RDMA Shared Device Plugin	271
9.7.3	NVIDIA Peer Memory Driver	271
9.7.4	Sonstige wichtige Komponenten	271
10	OpenShift (GPU-Accelerated) – Multiplattform (Cloud und On-Premises)	273
<hr/>		
10.1	Theoretische Vorbetrachtungen	273
10.1.1	Preflights: NVIDIA-Entitlements/-Lizenzen, Lizenzserver	274
10.1.2	Funktionsweise – High-Level-Überblick	274
10.2	Konzeptionelle Vorbetrachtungen zum Setup (On-Prem mit vSphere)	275
10.2.1	Überblick	275
10.2.2	Setup-Prozeduren GPU-Accelerated OpenShift IPI on vSphere – schematisch	276
10.3	On-Premises: OpenShift 4.10-Setup – Installer Provisioned Infrastructure (IPI) auf vSphere	277
10.3.1	Preflights: Infrastruktur und OpenShift-Cluster	277
10.3.2	Generelle Tool-Hinweise zu allen OpenShift-Setups (AWS, GCP, vSphere & Co.)	277
10.3.3	Der OpenShift-Installer: Terraform in schön	278
10.3.4	Vorbetrachtungen: Cluster Sizing	279
10.3.5	Zusammenfassung der technischen Preflights für das vSphere-Setup	279
10.3.6	Achtung, wichtig: DNS-Settings	280
10.3.7	DNS-Reverse-Zonen	281
10.3.8	vSphere-HA und OpenShift-Installer (OVA Upload fails in Single Datastore)	281
10.3.9	install-config.yaml für vSphere-IPI-Installation (Auszüge)	282
10.3.10	Rollout	283
10.3.11	Der Post-Rollout-Zustand	286
10.4	Preflights für skalierbare GPU-Nodes unter OpenShift: MachineSets, MachineConfigs und Machine-/Cluster-Autoscaler	286
10.4.1	Vorbetrachtungen	287
10.4.2	Cluster-Operatoren und Machine*-Ressourcen	287
10.4.3	MachineConfigs	289
10.4.4	MachineConfig-Operator	290
10.4.5	Komponenten des MCO	290

10.4.6	MachineConfigPool	291
10.4.7	Machines und MachineSets, Skalierung	292
10.5	Cluster-Autoscaler/Machine-Autoscaler	294
10.5.1	High-Level-Betrachtung	294
10.5.2	Machine-Autoscaler	295
10.5.3	Cluster-Autoscaler	295
10.5.4	Thresholds	296
10.5.5	Zu beachtende Punkte	297
10.5.6	GPU-VM-Template (vSphere) in MachineSet einbinden	298
10.5.7	GPU-MachineConfigPool und customisiertes MachineSet für skalierbare GPU-Nodes erzeugen	299
10.5.8	Skalierung des neuen GPU-MachineSets	303
10.5.9	Exemplarische Erzeugung eines GPU-MachineSets unter AWS	304
10.5.10	Fazit	306
10.6	vGPU-/MIG-spezifisches Setup des OpenShift-Clusters: NFD- und GPU-Operator	306
10.6.1	Historisches – NVIDIA-Driver-Build mit Red Hat Entitlements	306
10.6.2	Kernel für Driver-DaemonSet zu neu? Achtung bei OpenShift-Release-Updates	307
10.6.3	Installationsverfahren, generelle Operator-Settings	307
10.6.4	GPU-Manager-managed MIG-Mode und vGPU	308
10.6.5	NFD-Operator-Installation und -Konfiguration	308
10.6.6	GPU-Operator-Installation und -Konfiguration	311
10.6.7	License-ConfigMap	313
10.6.8	ImagePullSecret für Driver-Images aus der NGC-Registry	315
10.6.9	Die ClusterPolicy-CR (GPU-Operator)	316
10.7	Automatisches vGPU-Node-Setup per Operator – OpenShift-MachineSet mit Tesla T4	320
10.7.1	Rollout der ClusterPolicy-CR	320
10.7.2	Status auf den ESXi-Hosts	322
10.7.3	Analyse des ausgerollten (v)GPU-Stacks	324
10.8	Automatisches MIG-Slice-Setup per Operator – A30 on-premises	327
10.8.1	MIG im PCI Passthrough (A30 on-premises), Partitionierung durch den MIG-Manager	328
10.8.2	OpenShift-MachineSet und Default-MIG-Settings	328
10.8.3	Skalierung des MachineSets	330
10.8.4	Teilen? Oder lieber doch nicht?	332

10.9 Cloud: GPU-MachineSets in OpenShift 4.10 unter GCP mit A100-Instanzen (MIG-Partitionen via Operator)	333
10.9.1 Vorbetrachtungen	333
10.9.2 Verfügbare VM-Instanzen (GCP) mit GPU	334
10.9.3 Setup-Prozeduren – schematisch	334
10.9.4 Preflights – GCP-Kontingente gegebenenfalls erhöhen	335
10.9.5 Preflights – Domain, DNS und APIs	335
10.9.6 Service-Account zur OpenShift-Cluster-Erzeugung	337
10.9.7 Anpassungen der install-config.yaml, Rollout des Clusters	338
10.9.8 Setup der GPU-Nodes	341
10.9.9 Extraktion, Anpassung und Re-Import MachineSet und MCP	342
10.9.10 Skalierung des neuen GPU-MachineSets	345
10.9.11 Check der provisionierten GPU-Nodes	346
10.9.12 NFD- und GPU-Operator	346
10.9.13 MIG-Mode aktivieren, MIG-Partition-Size für A100 einstellen	348
10.9.14 Debugging und Troubleshooting	352
10.10 GPU-Sharing/-Overcommitment	353
10.10.1 Konzept-Recap und praktische Umsetzung	353
10.10.2 Setup (OpenShift)	355
10.10.3 Shared Workload testen	356
10.10.4 GPU-Sharing-Konfiguration per Node zur Laufzeit ändern	358
10.10.5 GPU Sharing mit vGPU	359
10.10.6 GPU-Sharing mit MIG-Slices	360
10.10.7 GPU-Sharing in der GCP-Cloud als kurzes PoC	366
10.11 Setup des Network-Operators (OpenShift on vSphere [IPI]) für GPUDirect RDMA	371
10.11.1 Preflights	371
10.11.2 High-Level-Workflow für den Network-Operator	378
10.11.3 Network-Operator und NFD-CR	379
10.11.4 Tests nach erfolgreichem Rollout	380
10.11.5 GPUDirect-RDMA-Test mit MacVLAN	382
10.11.6 Connect-Tests	383
10.11.7 Ein (nicht wirklich rundes) Fazit	385
10.12 KI/ML-System-Performance-Test (OpenShift on DGX)	386
10.13 GPU-Dashboard für OpenShift	387

11 GKE – Google Kubernetes Engine Cluster (GPU-Accelerated) 389

11.1 Überblick 389

11.1.1 Generelle Preflights: GPU-Verfügbarkeit nach Regionen/Zonen, geeignete Instanztypen 389

11.2 Setup-Variante 1: GKE-Cluster mit separatem Node-Pool für GPU-Nodes 390

11.2.1 Setup 390

11.2.2 Rollout des GPU-Operators 392

11.3 Setup-Variante 2: GPU-Cluster auf GKE direkt ausrollen 395

TEIL III ML-Stacks für skalierbare KI/ML-Infrastrukturen

12 CI/CD-Pipelines, GitOps und MLOps 399

12.1 Von der (ML-)Insel zur Pipeline 399

12.2 CI/CD und GitOps 400

12.2.1 CI/CD 400

12.2.2 GitOps 401

12.3 GitOps-Pipeline-Modelle 401

12.3.1 Pull- vs. Push-based 401

12.3.2 Push-based 402

12.3.3 Pull-based 403

12.3.4 Multiple Stages/Applications 404

12.4 MLOps, LTS und Portierbarkeit 404

12.4.1 MLOps und CRISP-DM 406

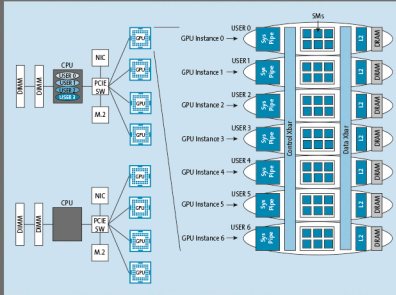
12.4.2 MLOps und ML-Pipelines – technische Foundation/schematisch 407

13 ML-Pipeline- und AI-End-to-End-Implementierungen mit Kubeflow/Vertex AI, Open Data Hub und NVIDIA AI Enterprise	411
13.1 ML-Pipeline-Implementierungen in Kubernetes-basierten Clustern	411
13.1.1 Der (KI/ML-)Pipeline-Ansatz	411
13.1.2 End-to-End-AI-Plattformen und Workflows	412
13.1.3 Das generelle Findungsproblem	413
13.1.4 Containerisierte ML-Pipelines und Segen und Fluch der Modularität	414
13.1.5 Kubernetes/Kubeflow to the Rescue? Genau betrachtet eher (noch) nicht.	415
13.1.6 Eine Lösung ...	416
13.2 Kubeflow	417
13.2.1 Kubeflow-Komponenten im Überblick	417
13.2.2 Entwicklung und Module (Auszüge)	418
13.2.3 Die Kernkomponenten	419
13.2.4 All together?	420
13.2.5 Istio	421
13.2.6 Kubeflow war gestern – es lebe Vertex AI. Na, zumindest ganz sicher bis ... sagen wir mal: morgen Mittag.	421
13.3 Hands-on: Kubeflow unter GKE in der Praxis	422
13.3.1 Preflights	422
13.3.2 Setup	423
13.3.3 Grafische Oberflächen	428
13.4 Open Data Hub	430
13.4.1 Die Unterschiede zu Kubeflow – ein High-Level-Überblick	430
13.4.2 Open Data Hub (ODH) – Architektur und Arbeitsweise	430
13.4.3 Die ODH-Module	432
13.5 Hands-on: Open-Data-Hub-Setup unter OpenShift	433
13.5.1 Preflights	433
13.5.2 Setup	434
13.5.3 Post Rollout	439

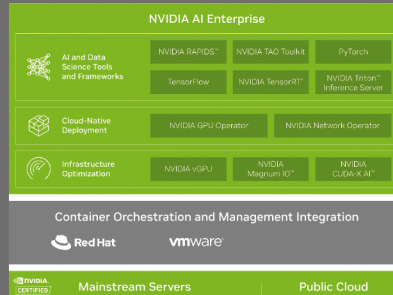
13.6 NVIDIA AI Enterprise (AI-End-to-End-relevante Betrachtungen)	442
13.6.1 NVIDIAs AI-End-to-End-Stack – reloaded	443
13.6.2 Die Module im Detail	444
13.6.3 NVIDIAs AI-End-to-End-Patterns	445
13.7 Hands-on: NVIDIA AI Enterprise (AI End-to-End) unter OpenShift	447
13.7.1 NVIDIA Morpheus AI Engine	447
13.7.2 Triton Inference Server	448
13.7.3 Morpheus MLflow Triton Plugin	449
13.7.4 Vorbetrachtungen: AI End-to-End mit Morpheus AI Engine	449
13.7.5 Preflights	450
13.7.6 Hands-on	450
13.7.7 Cybersecurity mit Morpheus AI (Red Hat Developer)	458
13.7.8 NVIDIA Launchpad	458
14 The Road Ahead	459
Index	463

Container-Cluster für Machine Learning und KI

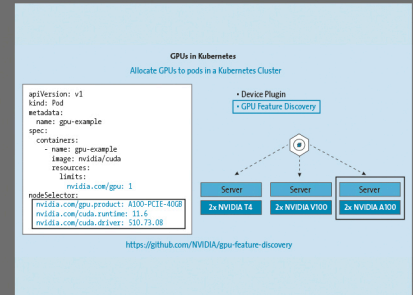
KI/ML-Modelle brauchen leistungsstarke Infrastrukturen. Oliver Liebel zeigt Ihnen, wie Sie resiliente, hochautomatisierte und autoskalierbare Systeme für Produktiv-Workloads auf Basis von NVIDIA's Datacenter-GPUs aufbauen und sie Operator-gestützt mit Kubernetes bzw. OpenShift verwalten.



Infrastrukturen planen



Komponenten verstehen



Workflows automatisieren

Best Practices für Planung und Umsetzung

Wenn große Datenmengen verarbeitet werden müssen und teure Investitionen anstehen, ist durchdachte Planung unabdingbar. Mit Hintergrundwissen und fundierten Markteinschätzungen treffen Sie richtige Entscheidungen für zukunftssichere KI/ML-Projekte und haben alle Kostenfaktoren im Blick.

Setzen Sie auf Industriestandards

In enger Zusammenarbeit mit NVIDIA wurden die vorgestellten Setups getestet und evaluiert. Sie profitieren von dieser Praxiserfahrung und erfahren, wie Sie GPUs in Kubernetes- und OpenShift-Cluster integrieren.

Zuverlässig und automatisiert skalieren

Operator-gestützte Vollautomation ist ein Muss, um KI/ML-Infrastrukturen skalierbar, resilient und kosteneffizient zu betreiben. Zahlreiche Beispiele führen Sie durch praxisnahe Konfigurationen, die für aktuelle und zukünftige Anforderungen gerüstet sind.



Getestete Setups zum Download



Linux-Enterprise-Experte **Oliver Liebel** ist Partner von Red Hat und SUSE und arbeitet eng mit NVIDIA's EGX-Team zusammen. Er ist seit 28 Jahren als Berater und Systemarchitekt tätig und hat zahlreiche GPU-beschleunigte Container-Cluster auf Kubernetes-/OpenShift-Basis für große Unternehmenskunden entworfen und umgesetzt.

Aus dem Inhalt

KI/ML: Grundlagen und Use Cases

Infrastruktur planen: Cloud, On-Premises oder Hybrid?

Technischer Background: KI/ML mit NVIDIA-GPUs

GPU-Modi: Passthrough-MIG, MIG-backed vGPU und vGPU

NVIDIA-GPUs auf vSphere On-Prem implementieren

NVIDIA AI Enterprise

KI/ML-Cluster mit Kubernetes und OpenShift

GPU-spezifische Operatoren

GPU-Cluster in der Cloud

Von CI/CD über GitOps zu MLOps

ML-Pipelines & AI End-to-End

