

KI-Hosting Made in Germany

Entscheidungsleitfaden,
Referenzarchitekturen,
Security & Compliance
für produktive KI-Systeme
im Mittelstand

PERSÖNLICH.
STARK.
SICHER.
KLUG.
RICHTIG GUT.

Inhalt

Executive Summary.....	3
Einleitung und Marktkontext.....	4
2. KI-Hosting als strategische Entscheidung.....	6
3. Hosting-Optionen im Vergleich.....	8
4. Technische Referenzarchitekturen.....	10
5. Compliance und Rechtssicherheit.....	14
6. Security und Governance für KI-Hosting.....	16
7. Betrieb, SRE und Migration.....	19
8. Wirtschaftlichkeit und Total Cost of Ownership.....	21
9. Entscheidungsleitfaden.....	23
10. Glossar und Referenzen.....	25

Executive Summary

Unternehmen in Deutschland stehen im Jahr 2026 vor einer doppelten Herausforderung: Einerseits ist Künstliche Intelligenz vom Innovationsprojekt zur geschäftskritischen Kerntechnologie geworden. Andererseits verschärfen sich regulatorische, sicherheitstechnische und geopolitische Anforderungen an Daten, Modelle und Infrastruktur mit rasanter Geschwindigkeit. Der Einsatz von KI in der Produktion, im Kundenservice, in der Dokumentenverarbeitung und in der strategischen Entscheidungsunterstützung ist für viele mittelständische Unternehmen längst keine Zukunftsvision mehr – er ist operative Realität.

Mit dem EU AI Act (Verordnung EU 2024/1689), der seit August 2026 in wesentlichen Teilen anwendbar ist, mit den verschärften Anforderungen der DSGVO, den NIS2-Umsetzungsgesetzen und den Orientierungshilfen der Datenschutzkonferenz (DSK) zu KI und Datenschutz sowie zu Retrieval-Augmented-Generation-Systemen (RAG) wird die Hosting-Entscheidung für KI-Workloads zu einer zentralen Governance-Frage.

Wo laufen GPU-Workloads? Wo liegen Trainingsdaten, Vektordatenbanken und Inferenz-Logs? Wer hat Zugriff auf Schlüssel und Metadaten? Welche Rollen entstehen nach DSGVO und AI Act?

Dieses Whitepaper bietet IT-Leitern, CISOs und Geschäftsführern im Mittelstand eine belastbare, praxisnahe Grundlage für diese Entscheidungen. Es vergleicht die relevanten Hosting-Optionen – von On-Premises über Private Cloud und Managed GPU-Hosting in Deutschland bis zu Hyper-scalern und Edge-Szenarien – und beschreibt konkrete Referenzarchitekturen, Security-Anforderungen sowie einen strukturierten Entscheidungs- und Implementierungsprozess.

Fünf zentrale Erkenntnisse auf einen Blick:

- > In Deutschland gehostet ist nicht automatisch souverän - entscheidend sind juristische Kontrolle, Subunternehmerketten, Schlüsselhöhe und Auditierbarkeit.
- > GPU-Hosting ist SRE- und Netzwerk-Engineering, nicht nur Rechenleistung.
- > S3-Objektspeicher ist der Standard-Ankerpunkt für KI-Datenpipelines.
- > Regulatorik wirkt direkt in Architekturentscheidungen hinein.
- > Für die meisten mittelständischen Unternehmen empfiehlt sich Managed GPU-Hosting in Deutschland mit strengem Vertrags- und Sicherheitsrahmen.

Einleitung und Marktkontext



1.1 Die veränderte Ausgangslage

Künstliche Intelligenz hat in den vergangenen drei Jahren den Sprung von der Forschungsumgebung in den produktiven Betrieb vollzogen. Generative KI-Modelle, Large Language Models (LLMs) und darauf aufbauende Anwendungen wie Retrieval-Augmented Generation (RAG) werden heute in nahezu allen Branchen eingesetzt: in der Fertigung für vorausschauende Wartung und Qualitätskontrolle, im Finanzsektor für Risikoanalysen und Kundenkommunikation, im Gesundheitswesen für Dokumentenanalyse und Diagnoseunterstützung, sowie in Querschnittsfunktionen wie Recht, HR, Einkauf und Kundenservice.

Für mittelständische Unternehmen bedeutet dies eine grundlegende Transformation der IT-Infrastruktur. KI-Systeme sind keine isolierten Anwendungen mehr, die auf separaten Servern laufen. Sie sind tief in Geschäftsprozesse integriert, verarbeiten vertrauliche Unternehmens- und Kundendaten, generieren Entscheidungsempfehlungen und sind damit Teil der kritischen IT-Infrastruktur. Diese

Integration schafft neue Anforderungen an Datenschutz, Informationssicherheit, regulatorische Konformität und operative Resilienz.

Gleichzeitig verändert sich der regulatorische Rahmen grundlegend. Mit dem Inkrafttreten des EU AI Acts, der Verschärfung der DSGVO-Praxis durch europäische und deutsche Aufsichtsbehörden, den NIS2-Anforderungen an Kritische Infrastrukturen sowie dem Financial Stability Board DORA für den Finanzsektor entsteht ein komplexes Gefüge von Pflichten, das direkte Auswirkungen auf Hosting-Entscheidungen hat. Die Frage, wo und unter welchen Bedingungen KI betrieben wird, ist damit zur Governance-Frage geworden.

1.2 Warum Made in Germany im KI-Hosting?

Der Marktkontext treibt Unternehmen in Richtung souveräner, deutsch- oder EU-basierter Infrastrukturen aus mehreren sich verstärkenden Gründen:

- > **Digitale Souveränität und Abhängigkeitsrisiken:**
Laut Bitkom sehen viele Unternehmen eine hohe strategische Abhängigkeit von außereuropäischen Cloud-Anbietern und wünschen sich europäische Alternativen. Der US CLOUD Act (Clarifying Lawful Overseas Use of Data Act) ermöglicht US-Behörden unter bestimmten Bedingungen den Zugriff auf Daten von US-Unternehmen, unabhängig davon, wo diese Daten physisch gespeichert sind. Dies schafft ein rechtliches Restrisiko, das für viele Unternehmen schwer kalkulierbar ist.
- > **Regulatorische Dynamisierung:**
Der EU AI Act gilt progressiv, mit wesentlichen Pflichten seit August 2025 (GPAI-Anbieter) und umfassender Anwendbarkeit seit August 2026. Gleichzeitig publizieren deutsche Aufsichtsbehörden konkrete Orientierungshilfen zu KI-Datenschutz und RAG-Systemen, die direkt auf Architekturentscheidungen einwirken.
- > **Energie- und Rechenzentrumsrealität:**
KI-Workloads treiben extrem hohe Rackdichten und steigenden Strombedarf. Die International Energy Agency (IEA) prognostiziert bis 2030 einen erheblichen Anstieg des globalen Rechenzentrumsstromverbrauchs, wobei KI-Beschleuniger der wesentliche Treiber sind. Deutsche Rechenzentren sind hier oft besser positioniert hinsichtlich Energieeffizienz und Nachhaltigkeitszertifizierungen.
- > **Vertrauen als Wettbewerbsfaktor:**
Für Unternehmen in regulierten Branchen – Finanzdienstleistungen, Gesundheitswesen, öffentliche Verwaltung – sowie für B2B-Anbieter mit strengen Lieferantenanforderungen ist nachweisbarer Datenschutz und Informationssicherheit ein Wettbewerbsvorteil.

1.3 Scope dieses Whitepapers

Dieses Dokument behandelt alle wesentlichen Aspekte des produktiven Betriebs von KI-Systemen aus der Perspektive mittelständischer Unternehmen: Training und Fine-Tuning von KI-Modellen, Inferenz und Model Serving (API-basierte Nutzung, Chatbots, Assistenzsysteme), Retrieval-Augmented Generation (RAG) mit Unternehmensdata sowie MLOps und LLMOps (Versionierung, Deployment-Pipelines, Monitoring, Governance).

Hinweis: Dieses Whitepaper ersetzt keine Rechtsberatung. Es bietet eine technisch-organisatorische Orientierung auf Basis öffentlich verfügbarer Leitlinien und Primärquellen. Für verbindliche rechtliche Einschätzungen empfehlen wir die Hinzuziehung spezialisierter Rechtsberater.

2. KI-Hosting als strategische Entscheidung

2.1 KI als geschäftskritische Infrastruktur

Die Integration von KI in Kerngeschäftsprozesse bedeutet, dass KI-Systeme heute ähnliche Anforderungen an Verfügbarkeit, Datenintegrität, Zugriffsschutz und Ausfallsicherheit stellen wie klassische ERP- oder CRM-Systeme. KI-Anwendungen verarbeiten Kundendaten, Finanz- und Buchhaltungsinformationen, Produktions- und Sensordaten aus industriellen Prozessen, interne Unternehmens-Know-how-Bestände sowie strategisch relevante Entscheidungsgrundlagen.

Fehlentscheidungen beim Hosting können weitreichende Konsequenzen haben: Compliance-Verstöße gegen DSGVO, EU AI Act oder branchenspezifische Regularien, Reputationsschäden durch Datenschutzverletzungen oder Datenlecks, wirtschaftliche Risiken durch Betriebsunterbrechungen oder Cyberangriffe, sowie langfristige strategische Abhängigkeiten von einzelnen Anbietern oder Technologieplattformen. Hosting-Entscheidungen für KI sind damit keine rein technischen Infrastrukturentscheidungen, sondern haben direkten Einfluss auf das Risikoprofil, die regulatorische Compliance und die strategische Unabhängigkeit eines Unternehmens.

2.2 Der regulatorische Druck als Treiber

Mit der Einführung neuer europäischer Regelwerke steigen die Anforderungen an den nachweisbaren, dokumentierten und auditierbaren Betrieb von KI-Systemen erheblich:

DSGVO (Datenschutz-Grundverordnung):

Die DSGVO stellt grundlegende Anforderungen an die Verarbeitung personenbezogener Daten. Für KI-Systeme bedeutet dies: Nachweis der Rechtmäßigkeit der Datenverarbeitung für Training und Inferenz, klare Auftragsverarbeitungsverträge (Art. 28) mit Hosting-Dienstleistern, Implementierung angemessener technischer und organisatorischer Maßnahmen (Art. 32), sowie Durchführung von

Datenschutz-Folgenabschätzungen (DSFA, Art. 35) bei risikoreichen KI-Anwendungen.

EU AI Act (Verordnung EU 2024/1689):

Der AI Act etabliert einen risikobasierten Regulierungsrahmen. Hochrisiko-KI-Systeme (Annex III) unterliegen ab August 2026 umfangreichen Pflichten: Risikomanagementsystem, Datenverwaltung, technische Dokumentation, Transparenz gegenüber Nutzern, menschliche Aufsicht und Genauigkeitsanforderungen. Die Hosting-Entscheidung beeinflusst direkt, welche dieser Anforderungen erfüllbar sind.

NIS2-Richtlinie (umgesetzt im NIS2UmsuCG):

Für Unternehmen, die als wesentliche oder wichtige Einrichtung eingestuft sind, entstehen strenge Anforderungen an Risikomanagement, Incident Reporting und Lieferketten-sicherheit. KI-Systeme als Bestandteil kritischer Geschäftsprozesse fallen in den Scope dieser Anforderungen.

DORA (Digital Operational Resilience Act):

Im Finanzsektor gelten seit Januar 2025 verschärfte Anforderungen an die digitale operationale Resilienz, einschließlich strenger Anforderungen an IKT-Drittanbieterisiken und -verträge.

> Die Konsequenz:

Unternehmen müssen für jeden KI-Workload nachweisen können, wo Daten verarbeitet werden, wer Zugriff auf diese Daten hat, welche technischen Schutzmaßnahmen implementiert sind und wie Risiken systematisch minimiert werden. Hosting-Entscheidungen sind damit Governance-Entscheidungen.

2.3 Die Souveränitätsfrage

Der Begriff "digitale Souveränität" wird im Kontext von KI-Hosting häufig verwendet, aber selten präzise definiert. Für die praktische Entscheidungsfindung empfiehlt sich eine operationale Definition entlang von vier Dimensionen:

Datensouveränität

Verbleiben personenbezogene und geschäftskritische Daten unter der Kontrolle des Unternehmens oder seiner EU-basierten Auftragsverarbeiter? Gibt es nachweislich keine Datenweitergabe an Drittstaaten ohne Einwilligung oder rechtliche Grundlage?

Schlüsselhoheit:

Wer kontrolliert die kryptographischen Schlüssel für die Datenverschlüsselung? Customer Managed Keys (CMK) versus Provider Managed Keys sind ein wesentlicher Unterschied.

Zugriffshoheit:

Wer hat technische und administrative Zugriffsrechte auf Infrastruktur, Daten und Modelle? Sind Remote-Zugriffe aus Drittstaaten durch Supportpersonal ausgeschlossen oder dokumentiert?

Rechtliche Kontrollierbarkeit:

Unterliegt der Hosting-Dienstleister ausschließlich deutschem und europäischem Recht, oder gibt es Drittstaaten-Rechtsbindungen (etwa durch US-amerikanische Muttergesellschaften)?

Wichtiger Hinweis: Die Datenschutzkonferenz (DSK) betont explizit in ihrer Entschliebung zu Confidential Cloud Computing (Juni 2025), dass Marketingbegriffe wie Souveräne Cloud oder Confidential Cloud Computing kritisch gegen realistische Angreifermodelle, technische Grenzen und Nachweise geprüft werden müssen. Standort-Aussagen allein sind kein Nachweis von Souveränität.



3. Hosting-Optionen im Vergleich

3.1 Überblick: Die fünf Hosting-Modelle

Für den Betrieb von KI-Workloads stehen grundsätzlich fünf Hosting-Modelle zur Verfügung, die in der Praxis häufig kombiniert werden. Jedes Modell hat charakteristische Stärken, Einschränkungen und typische Einsatzmuster:

On-Premises:

Betrieb auf eigener Hardware im eigenen Rechenzentrum oder in einem Colocation-Rechenzentrum. Höchste Kontrolle über Hardware und Daten, aber hoher Investitionsaufwand (CAPEX), begrenzte Skalierbarkeit und hohe Anforderungen an interne Expertise.

Private Cloud (DE/EU):

Dedizierte, mandantenspezifische Cloud-Infrastruktur bei einem deutschen oder europäischen Anbieter. Gute Balance zwischen Kontrolle und Managierbarkeit, höhere Flexibilität als On-Premises bei guter Compliance-Positionierung.

Managed GPU-Hosting in Deutschland:

Spezialisierte Managed Services für GPU-basierte KI-Workloads in deutschen Rechenzentren. Schneller Time-to-Value, geringerer interner Betriebsaufwand, bei entsprechendem Vertragswerk hohe Compliance-Tauglichkeit.

Hyperscaler (AWS, Azure, GCP):

Globale Cloud-Anbieter mit umfangreichem Serviceangebot. Sehr hohe Skalierbarkeit und breites KI-Ökosystem, aber regulatorische Komplexität durch Drittstaaten-Rechtsbindungen und eingeschränkte Datensouveränität.

Edge / On-Device:

Ausführung von KI-Modellen direkt auf Endgeräten oder Edge-Infrastruktur in der Produktion. Relevant für Echtzeitanwendungen und Offline-Szenarien, aber hohe Betriebskomplexität durch dezentrales Fleet-Management.

Kriterium	On-Prem	Private Cloud (DE/EU)	Managed GPU (DE)	Hyperscaler	Edge
Datenkontrolle	Sehr hoch	Hoch	Mittel-Hoch	Mittel	Hoch (lokal)
Time-to-Value	Langsam	Mittel	Schnell	Sehr schnell	Mittel
GPU-Skalierung	Begrenzt	Gut	Gut	Sehr gut	Begrenzt
Drittlandrisiko	Niedrig	Niedrig	Niedrig-Mittel	Mittel-Hoch	Niedrig
Betriebsaufwand	Hoch	Mittel-Hoch	Mittel	Niedrig-Mittel	Hoch
Lock-in-Risiko	Niedrig	Mittel	Mittel	Hoch	Mittel
CAPEX	Hoch	Mittel	Niedrig	Niedrig	Mittel
Compliance-Klarheit	Hoch	Hoch	Hoch (vertragsabh.)	Komplex	Hoch (lokal)

3.2 KI-Hosting Made in Germany: Was das wirklich bedeutet

Der Begriff KI-Hosting Made in Germany beschreibt ein Hosting-Modell, das folgende Charakteristika erfüllt: Betrieb in deutschen Rechenzentren mit physisch verifizierbarem Standort, Unterwerfung unter deutsches und europäisches Recht ohne Drittstaaten-Rechtsbindung, DSGVO-konforme Datenverarbeitung mit belastbaren Auftragsverarbeitungsverträgen, transparente Sicherheitsarchitektur mit nachweisbaren technischen und organisatorischen Maßnahmen sowie Zertifizierungen (ISO 27001, BSI C5-Testat), und kein Remote-Administratorzugriff aus Drittstaaten ohne dokumentierte Ausnahme.

Die Vorteile dieses Modells für mittelständische Unternehmen in Deutschland sind substantiell: maximale Rechtssicherheit unter DSGVO und EU AI Act, nachweisbare digitale Souveränität gegenüber Kunden und Aufsichtsbehörden, Vertrauenswürdigkeit als Wettbewerbsargument in regulierten Branchen, Unterstützung durch deutsche Behörden und Zertifizierungsstellen sowie kulturelle und sprachliche Nähe für Support und Eskalationsprozesse. Gleichzeitig sind potenzielle Einschränkungen realistisch einzuplanen: begrenzte globale Skalierungsoptionen im Vergleich zu Hyperscalern, teilweise geringere Service-Vielfalt bei nischenhaften KI-Diensten, sowie höhere Anforderungen an Due Diligence bei der Anbieterauswahl, da der Markt weniger transparent ist als bei etablierten Hyperscalern.

3.3 Hybridstrategien: Der pragmatische Ansatz für den Mittelstand

In der Praxis entscheiden sich viele mittelständische Unternehmen für hybride Ansätze, die die Stärken verschiedener Modelle kombinieren:

- > Daten in DE/EU, Compute elastisch:
Sensible Daten, Vektorindizes und Logs verbleiben in deutschen Rechenzentren. Kurzfristige Compute-Spitzen für nicht-sensible Workloads werden bei Bedarf extern abgerufen – mit klarer Transfer-Governance und Datenklassifizierung.
- > Produktive Inferenz in Deutschland, Experimentierphase extern:
Frühe Proof-of-Concepts und Prototypen werden in flexiblen, kostengünstigen Umgebungen entwickelt. Für die Produktivumgebung gelten harte Compliance- und Sicherheitsanforderungen an einen deutschen Hosting-Anbieter.
- > Fine-Tuning mit öffentlichen Daten extern, Inferenz mit Unternehmensdaten intern:
Das Basis-Training oder Fine-Tuning mit nicht-sensiblen oder öffentlichen Datensätzen kann extern erfolgen. Der produktive Inferenzbetrieb mit Unternehmensdaten findet ausschließlich in kontrollierter Umgebung statt.



4. Technische Referenzarchitekturen

4.1 Workload-Profile: Technische Anforderungen verstehen

Eine fundierte Hosting-Entscheidung erfordert das Verständnis der technischen Profile verschiedener KI-Workloads. Unterschiedliche Workloads stellen grundlegend verschiedene Anforderungen an Compute, Netzwerk, Storage und Betrieb:

Training (From-Scratch):

Ressourcenintensivster Workload, primär relevant für Organisationen mit sehr großen proprietären Datensätzen. Der Engpass liegt typischerweise im GPU-Interconnect-Netzwerk (NVLink/NVSwitch, InfiniBand) und im Storage-I/O für Trainingsdaten. Checkpointing-Strategien und paralleles Filesystem sind kritisch für Effizienz und Ausfallsicherheit.

Fine-Tuning (LoRA/QLoRA):

Für den Mittelstand der praktisch relevantere Workload. Parameter-Efficient Fine-Tuning (PEFT) Methoden wie LoRA oder QLoRA ermöglichen domänenspezifische Modellanpassungen mit deutlich geringerem GPU-Speicher- und Rechenaufwand. Mixed-Precision-Training (BF16/FP16) und effiziente Datenlader-Implementierungen sind für TCO-Optimierung entscheidend.

Inference / Model Serving (Realtime):

Für produktive KI-APIs sind p95/p99-Latenzanforderungen der primäre Optimierungsparameter. Der Serving-Stack (Triton Inference Server, vLLM, KServe), Batching-Strategien, Caching (KV-Cache für Transformer-Modelle) und Autoscaling sind kritische Designentscheidungen. Token-Throughput (Tokens/sec) ist ein wichtiger SLI.

RAG (Retrieval-Augmented Generation):

Kombiniert LLM-Inferenz mit Retrieval über Vektorsuche auf Unternehmensdaten. Neben der LLM-Latenz sind Retrieval-Latenz (p95 über Vektordatenbank-Query), Index-Freshness und die korrekte Implementierung von

Zugriffskontrolle auf Dokumentenebene (Row-Level Security im Index) kritische technische und datenschutzrechtliche Anforderungen.

4.2 GPU-Architekturen und Multi-Tenancy

Für Entscheider ist das Verständnis der GPU-Architektur relevant, weil es direkt Kosten, Leistung und Compliance-Anforderungen beeinflusst:



Multi-Instance GPU (MIG):

NVIDIAs MIG-Technologie (verfügbar bei A100, H100, H200) ermöglicht die physische Partitionierung einer GPU in bis zu sieben isolierte Instanzen mit dedizierten Compute-, Cache- und Memory-Ressourcen. Dies ist für mandantenfähige KI-Plattformen der stärkste Isolationsbaustein, da echte Hardware-Isolation statt Soft-Virtualisierung realisiert wird. MIG ist besonders für Compliance-sensible Workloads relevant, weil "Noisy Neighbor"-Effekte und Cross-Tenant-Datenlecks auf GPU-Ebene technisch ausgeschlossen werden.

GPU Time-Slicing:

Ermöglicht Oversubscription über zeitliche Aufteilung der GPU zwischen mehreren Workloads. Einfacher zu konfigurieren als MIG, aber ohne physische Ressourcenisolation. Geeignet für leichtere Inferenz-Workloads in weniger sensiblen Umgebungen, aber für hochsensible Daten oder strenge Compliance-Anforderungen nicht empfehlenswert.

Thermal Design Power (TDP) als Infrastrukturtreiber:

Moderne KI-Beschleuniger wie die NVIDIA H100 können je nach Konfiguration bis zu 700 Watt TDP erreichen. Ein 8-GPU-Server kann damit eine IT-Last von 5,6 kW allein für GPUs erzeugen. Zusammen mit CPU, RAM, Netzwerk und Storage sowie einem Power Usage Effectiveness (PUE) Faktor von 1,5 bis 1,6 ergibt sich eine Facility-Last von 12 bis 15 kW pro Server. Diese Realität bestimmt Rechenzentrumsauswahl, Kühlanforderungen und Energiekosten maßgeblich.

4.3 Netzwerkarchitektur: Die unterschätzte Dimension

Für skalierbare KI-Workloads ist Netzwerk nicht "nur Bandbreite", sondern oft der limitierende Faktor:

NVLink/NVSwitch für GPU-to-GPU-Kommunikation:

NVIDIAs DGX-Plattformen nutzen NVLink für GPU-to-GPU-Kommunikation innerhalb eines Servers (bis zu 900 GB/s je nach Architekturgeneration). NVSwitch erweitert dies auf Multi-Node-Topologien. Diese Bandbreiten sind für verteiltes Training und Modelle, die nicht in den Speicher einer einzelnen GPU passen, entscheidend.

GPUDirect RDMA für GPU-Storage-Integration:

GPUDirect RDMA ermöglicht den direkten Datentransfer zwischen GPU-Speicher und Netzwerkkarte (InfiniBand oder RoCE) ohne Umweg über Host-CPU und -RAM. Dies reduziert Latenz und CPU-Overhead bei dataintensiven Training-Workloads erheblich. Für produktive KI-Plattformen in Deutschland ist die Verfügbarkeit von RDMA-Netzwerken und dem entsprechenden Betriebswissen beim Anbieter ein wichtiges Auswahlkriterium.

Netzwerksegmentierung für Security:

Kubernetes NetworkPolicies, separate VLANs für GPU-Nodes, Control-Plane und Datenpipelines sowie egress-Kontrollen sind nicht nur Performance-, sondern Sicherheitsanforderungen. East-West-Traffic im Cluster muss durch Micro-Segmentierung eingeschränkt werden, um Lateral Movement bei einem Sicherheitsvorfall zu begrenzen.

4.4 Storage-Architektur: S3 als Standard-Ankerpunkt

KI-Plattformen benötigen eine strukturierte Storage-Strategie über mehrere Ebenen:

> Hot Storage (NVMe/NVMe-oF):

Für aktive Modell-Weights, Intermediate Checkpoints und hochfrequenten Cache-Zugriff. Lokale NVMe-SSDs direkt in GPU-Nodes oder NVMe-over-Fabrics (NVMe-oF) für shared High-Performance Storage.

> Warm Storage (Shared Filesystem):

Für Checkpoints, Artefakte und gemeinsam genutzte Modell-Repositories. Oft als NFS oder paralleles Filesystem realisiert.

> Cold/Tiered Storage (S3 Objektspeicher):

Als "Data Lake"-Schicht für Trainingsdaten, fertig trainierte Modelle, Feature Stores, Experiment-Logs und Backups. S3-API ist de-facto-Standard in KI-Tooling – nahezu alle relevanten Frameworks und Tools (PyTorch, TensorFlow, Spark, MLflow, DVC) unterstützen S3-Protokoll nativ.

Für deutsche KI-Hosting-Umgebungen relevante S3-kompatible Systeme sind Ceph RADOS Gateway (RGW) und MinIO, die S3-Kompatibilität auf Ebene der Core-API-Operationen dokumentieren. Für produktiven MLOps-Betrieb müssen jedoch Edge-Cases – Multipart Upload, Bucket Policies, KMS-Integration, Event Notifications – in einem Pilot getestet und als Kompatibilitätsmatrix dokumentiert werden.

4.5 Referenzarchitektur: Private/Managed GPU-Hosting in Deutschland

Die folgende Referenzarchitektur beschreibt eine produktive KI-Plattform auf Basis von Kubernetes, GPU-Nodes, S3-Objektspeicher und KServe/Triton als Model-Serving-Stack: Architekturschichten (von außen nach innen):

- > User-/App-Layer:
Business-Anwendungen, Chat-UI, externe APIs – kommunizieren ausschließlich über TLS-verschlüsselte Verbindungen mit dem API Gateway.
- > API Gateway / Ingress:
TLS-Terminierung, Rate-Limiting, Anfrage-Routing. Mutual TLS (mTLS) für Service-to-Service-Kommunikation im Cluster.
- > Auth & Policy Layer:
OIDC/SSO-Integration (Unternehmens-IdP), Policy Engine (z.B. OPA/Gatekeeper) für Zugriffskontrolle auf Namespace- und Ressourcenebene.
- > KServe InferenceService:
Kubernetes-CRDs für Model Serving, Autoscaling (HPA/KPA), Canary Deployments, Health Checks und Rollback.
- > Triton Inference Server / vLLM:
Model-Runtime-Schicht. Triton für allgemeines Multi-Framework-Serving (TensorRT, ONNX, PyTorch), vLLM für optimiertes LLM-Serving mit PagedAttention-Optimierung.
- > GPU-Nodes:
Dedizierte Nodes mit NVIDIA GPU Operator (MIG-Konfiguration, Time-Slicing-Policy, Monitoring). RDMA-fähiger High-Speed-Fabric für Training-Workloads.
- > S3 Objektspeicher:
Datasets, Modell-Artefakte, Feature Stores, Logs. Verschlüsselung mit Envelope-Encryption: Daten-DEK pro Bucket/Prefix, KEK im HSM/KMS.
- > KMS/HSM:
Schlüsselverwaltung für alle Datenverschlüsselungsoperationen. Customer Managed Keys (CMK) für maximale Schlüsselhoheit. Auditierbares Key Access Logging.
- > Observability Stack:
Prometheus/Grafana für GPU-Metriken (Utilization, Memory, Temperature), OpenTelemetry für verteiltes Tracing, zentralisiertes Log-Management mit SIEM-Integration.



4.6 RAG-Architektur mit Datenschutz-Governance

Retrieval-Augmented Generation (RAG) ist für den Mittelstand besonders relevant, weil es LLMs mit proprietären Unternehmensdaten kombiniert, ohne aufwändiges Fine-Tuning. Die Architektur eines datenschutzkonformen RAG-Systems umfasst:

Dokumentenvorverarbeitung (ETL-Pipeline):

Ingestion von Unternehmensdokumenten aus verschiedenen Quellen (SharePoint, Confluence, SAP, Filesystem). Chunking-Strategien (semantisches Chunking vs. Fixed-Size-Chunking), Embedding-Generierung und Indexierung in Vektordatenbank. Kritisch: Jedes Dokument muss mit Metadaten versehen werden, die Access Control und Lösbarkeit ermöglichen.

Vektordatenbank mit Row-Level Security:

Die Vektordatenbank (z.B. Weaviate, Qdrant, pgvector) muss dokumentenebene Zugriffskontrolle implementieren. Nutzer dürfen im Retrieval-Schritt nur auf Dokumente zugreifen, die sie auch in der Quellsystem-Berechtigung hätten. Dies ist aus datenschutzrechtlicher und sicherheitstechnischer Perspektive nicht verhandelbar.

Prompt Guard / Policy Engine:

Zwischen API Gateway und LLM-Inferenz muss eine Policy-Engine implementiert werden, die Prompt Injection Attacks erkennt, Output Filtering (PII-Erkennung in Antworten) durchführt und Zugriffskontrolle auf Retrieval-Quellen erzwingt.

Audit Logging:

Prompts, abgerufene Dokument-Chunks und generierte Antworten müssen protokolliert werden (sofern datenschutzrechtlich zulässig), um bei Sicherheitsvorfällen und regulatorischen Anfragen Auskunftsfähigkeit zu gewährleisten. Logs müssen unveränderlich (immutable) gespeichert und selbst durch KMS geschützt werden.

- > Die Datenschutzkonferenz (DSK) hat spezifische Orientierungshilfen zu RAG-Systemen (DSK OH RAG, Oktober 2025) veröffentlicht, die Architekturfragen explizit adressieren: Index-Lifecycle, Lösbarkeit, Auskunftsfähigkeit, Zugriffskontrollen. Diese Orientierungshilfe sollte als Pflichtlektüre in Architektur-Reviews eingehen.



5. Compliance und Rechtssicherheit

5.1 DSGVO: Die Basis für KI-Hosting-Compliance

Die DSGVO stellt für KI-Systeme, die personenbezogene Daten verarbeiten, umfangreiche Anforderungen. Im KI-Hosting-Kontext sind folgende Aspekte besonders relevant:

Auftragsverarbeitung (Art. 28 DSGVO):

Bei jeder Form von Managed Hosting, bei dem ein externer Dienstleister im Auftrag des datenverarbeitenden Unternehmens tätig wird, ist ein Auftragsverarbeitungsvertrag (AVV) zwingend erforderlich. Der AVV muss den Gegenstand, die Dauer und den Zweck der Verarbeitung, die Art der personenbezogenen Daten, die Kategorien der betroffenen Personen sowie die Pflichten und Rechte des Verantwortlichen beschreiben. Besonderes Augenmerk gilt der Subunternehmer-Regelung: Alle Sub-Auftragsverarbeiter müssen benannt und genehmigt werden.

Technische und organisatorische Maßnahmen (Art. 32 DSGVO):

Der Stand der Technik muss bei der Implementierung von Sicherheitsmaßnahmen berücksichtigt werden. Für KI-Systeme bedeutet dies: Verschlüsselung von Daten in Ruhe und in Transit, Mechanismen zur Sicherstellung von Vertraulichkeit, Integrität, Verfügbarkeit und Belastbarkeit der Verarbeitungssysteme, Wiederherstellungsverfahren bei physischen oder technischen Zwischenfällen sowie regelmäßige Überprüfung der Wirksamkeit der Maßnahmen.

Datenschutz-Folgenabschätzung (DSFA, Art. 35 DSGVO):

Viele KI-Anwendungen erfüllen die Tatbestandsmerkmale für eine DSFA-Pflicht: Profiling, umfangreiche Verarbeitung besonderer Datenkategorien, systematische Überwachung öffentlich zugänglicher Bereiche oder neuartige Verarbeitungsmethoden mit möglicherweise hohem Risiko. Eine DSFA muss die Notwendigkeit und Verhältnismäßigkeit der Verarbeitung beurteilen, die Risiken für die Rechte und

Freiheiten der betroffenen Personen bewerten und angemessene Abhilfemaßnahmen festlegen.

Betroffenenrechte im KI-Kontext:

Auskunftsrecht, Recht auf Löschung ("Recht auf Vergessenwerden") und Recht auf Einschränkung der Verarbeitung müssen auch für KI-Systeme technisch implementierbar sein. Dies stellt besondere Anforderungen an RAG-Systeme: Wenn ein Dokument mit personenbezogenen Daten gelöscht werden muss, muss dies auch den Vektorindex betreffen – eine technische Herausforderung, die häufig Neuindexierungsprozesse erfordert.

5.2 Drittlandtransfers und Souveränitätsnachweise

Für das Hosting-Modell "Made in Germany" ist die Frage nach Drittlandtransfers zentral. Personenbezogene Daten dürfen nur dann in Drittländer übermittelt werden, wenn ein Angemessenheitsbeschluss der EU-Kommission vorliegt, geeignete Garantien (z.B. Standard-Datenschutzklauseln, SCC) vereinbart sind oder Ausnahmen nach Art. 49 DSGVO greifen.

Der EU-U.S. Data Privacy Framework (DPF, Implementing Decision EU 2023/1795) erleichtert Datenübermittlungen in die USA für zertifizierte US-Unternehmen. Dieser Angemessenheitsbeschluss ersetzt jedoch nicht die Notwendigkeit einer sorgfältigen Governance für Subunternehmerketten und schützt nicht gegen den US CLOUD Act.

Der EDPB-Empfehlungskatalog 01/2020 zu Supplementary Measures ist maßgeblich für alle Fälle, in denen Transfers auf SCC basieren und ergänzende technische oder organisatorische Maßnahmen erforderlich sind. Für KI-Hosting bedeutet dies: Transfer Impact Assessments (TIA) für alle Datenströme in Drittländer, Verschlüsselungsmaßnahmen, die sicherstellen, dass auch der Empfänger im Drittland keinen Klartextzugriff hat (End-to-End-Verschlüsselung mit im EU-Raum verbleibenden Schlüsseln), sowie vertragliche und

technische Dokumentation aller Subunternehmerketten. Für die Prüfung eines potentiellen Hosting-Anbieters gelten folgende kritische Fragen: Ist der Anbieter oder ein zentraler Subunternehmer eine Gesellschaft, die außerhalb der EU rechtlich kontrolliert wird? Gibt es Support- oder Administrationszugriffe aus Drittstaaten? Sind Telemetrie-Daten oder Metadaten von außereuropäischen Muttergesellschaften zugänglich?

5.3 EU AI Act: Relevanz für Hosting-Entscheidungen

Der EU AI Act (Verordnung EU 2024/1689) hat direkte Implikationen für KI-Hosting-Entscheidungen. Der Act unterscheidet verschiedene Rollen: Anbieter (Provider), die KI-Systeme entwickeln und in Verkehr bringen; Betreiber (Deployer), die KI-Systeme im eigenen beruflichen Kontext nutzen; sowie Importeure und Bevollmächtigte in internationalen Kontexten.

Für mittelständische Unternehmen als Deployer von Hochrisiko-KI-Systemen (Annex III) gelten ab August 2026 unter anderem: Sicherstellung, dass der Anbieter des KI-Systems die Konformitätsbewertungspflichten erfüllt hat, Durchführung einer Folgenabschätzung für Grundrechte (sofern zutreffend), Implementierung angemessener Monitoring- und Oversight-Mechanismen sowie Dokumentations- und Protokollierungspflichten.

Für Unternehmen, die KI-Systeme unter eigenem Namen bereitstellen – sei es an externe Kunden oder intern – entstehen als Provider umfangreichere Pflichten, die direkte Anforderungen an Infrastruktur und Betrieb stellen: Risikomanagementsystem, Qualitätsmanagementsystem, technische Dokumentation und automatische Protokollierung (Logging). Diese Anforderungen sind in einer kontrollierten deutschen Hosting-Umgebung deutlich einfacher zu erfüllen als in einer Hyperscaler-Umgebung mit komplexen Subunternehmerketten.

> Praxishinweis: Die Rollenklärung nach EU AI Act sollte für jeden KI-Use-Case vor der Infrastrukturentscheidung erfolgen. Erstellen Sie eine Rollenlandkarte (Provider / Deployer / Importer) pro Anwendungsfall, um Pflichten frühzeitig zu identifizieren und Hosting-Anforderungen davon abzuleiten.

Compliance-Bereich	On-Prem	Private Cloud(DE)	Managed GPU (DE)	Hyperscaler
AVV Art. 28 (Klarheit)	Hoch	Hoch	Hoch	Hoch (aber komplex)
Subunternehmertransparenz	Hoch	Mittel-Hoch	Mittel	Niedrig-Mittel
Drittlandtransferrisiko	Sehr niedrig	Niedrig	Niedrig-Mittel	Mittel-Hoch
DSFA-Handhabbarkeit	Hoch	Mittel-Hoch	Mittel-Hoch	Mittel
AI Act Nachweisfähigkeit	Hoch	Hoch	Mittel-Hoch	Mittel
RAG-spezifische Governance	Hohe Kontrolle	Hohe Kontrolle	Gute Kontrolle	Eingeschränkt

6. Security und Governance für KI-Hosting

6.1 KI-spezifische Bedrohungslandschaft

KI-Systeme erweitern die klassische IT-Bedrohungslandschaft um spezifische Angriffsvektoren, die in traditionellen Security-Konzepten nicht ausreichend adressiert sind. Das Bundesamt für Sicherheit in der Informationstechnik (BSI) und die European Union Agency for Cybersecurity (ENISA) haben umfangreiche Leitfäden zu KI-Sicherheitsrisiken veröffentlicht:

Data Poisoning / Backdoor-Angriffe:

Durch gezielte Manipulation von Trainingsdaten können Angreifer KI-Modelle mit "Hintertüren" versehen, die unter bestimmten Eingabebedingungen fehlerhafte oder schadhaft beeinflusste Ausgaben produzieren. Besonders relevant für Unternehmen, die Trainingsdaten aus externen Quellen beziehen oder Pre-Trained Models Dritter einsetzen.

Prompt Injection:

Angreifer können durch manipulierte Eingaben versuchen, die Anweisungen und Sicherheitsgrenzen eines LLM zu überwinden ("Jailbreaking"). In Enterprise-Kontexten kann dies zu unerwünschten Datenabflüssen, Umgehung von Zugriffskontrollen oder Manipulation von automatisierten Entscheidungsprozessen führen.

Model Extraction / Model Stealing:

Durch systematische Abfragen einer Inferenz-API können Angreifer proprietäre Modelle und deren Trainingsdatencharakteristika mit relativ geringem Aufwand rekonstruieren. Für Unternehmen mit wertvollem KI-Know-how ist API-Ratenlimiting, Abfrage-Monitoring und Anomalieerkennung ein wichtiger Schutzmechanismus.

Membership Inference / Model Inversion:

Angreifer können mit hoher Wahrscheinlichkeit bestimmen, ob ein bestimmter Datenpunkt im Trainingsdatensatz enthalten war (Membership Inference) oder Teile der Trainingsdaten aus Modell-Ausgaben rekonstruieren (Model Inversion). Dies ist ein erhebliches Datenschutzrisiko, wenn

Modelle mit personenbezogenen Daten trainiert wurden.

Supply Chain Attacks:

KI-Systeme sind komplex gestapelt: CUDA-Bibliotheken, Container Images, Pre-Trained Model Weights aus öffentlichen Repositories, Python-Pakete, Kubernetes-Operator. Jede dieser Ebenen kann kompromittiert sein. Das BSI benennt Supply-Chain-Angriffe als eine der am stärksten wachsenden Bedrohungskategorien.

6.2 Infrastruktur-Härtung nach BSI-Standards

Für eine rechtssichere, "Stand der Technik"-konforme Security-Posture empfiehlt sich die Orientierung an BSI-Standards:

IT-Grundschutz-Bausteine:

Das BSI IT-Grundschutz-Kompodium enthält spezifische Bausteine für Serverinfrastruktur (SYS.1.1 Allgemeiner Server), Containerisierung (SYS.1.6) und Kubernetes-Betrieb (APP.4.4). Diese Bausteine definieren konkrete Anforderungen für Härtung, Zugriffskonzepte, Logging, Patch-Management und sichere Konfiguration. Sie sind zugleich Referenzrahmen für Auditoren und Datenschutzbehörden in Deutschland.

BSI C5

(Cloud Computing Compliance Criteria Catalogue):

C5 ist in Deutschland der etablierte Mindestanforderungskatalog für Cloud-Sicherheit. Ein C5-Testat eines Hosting-Anbieters liefert strukturierten Nachweis über Security-Maßnahmen in 17 Domänen, von Organisationssicherheit über Kryptographie bis hin zu Incident Management. C5 wird von deutschen Aufsichtsbehörden als geeigneter Nachweis anerkannt.

BSI TR-02102

(Technische Richtlinien für Kryptographie):

TR-02102-1 empfiehlt kryptographische Verfahren und Schlüssellängen nach aktuellem Stand (Version 2026-01).

TR-02102-2 enthält TLS-spezifische Empfehlungen zu Protokollversionen und Cipher Suites. Diese Richtlinien sind für die Konfiguration von API-Endpoints, Datenbankverschlüsselung und KMS-Anbindung maßgeblich.

6.3 Schlüsselmanagement und Verschlüsselungsarchitektur

Schlüsselmanagement ist das Herzstück einer souveränen KI-Hosting-Architektur. Ohne kontrollierte Schlüsselhoheit ist Datensouveränität technisch nicht realisierbar:

- > Envelope Encryption als Standardmuster:
Datenverschlüsselungsschlüssel (Data Encryption Key, DEK) werden pro Objekt, Bucket oder Datenshard generiert. Der DEK wird mit einem Schlüsselverschlüsselungsschlüssel (Key Encryption Key, KEK) im Hardware Security Module (HSM) oder Key Management Service (KMS) verschlüsselt und zusammen mit dem verschlüsselten Datum gespeichert.
- > Customer Managed Keys (CMK):
Der entscheidende Unterschied zwischen echter Schlüsselhoheit und "Provider manages keys" ist, ob der Kunde den KEK im KMS/HSM kontrolliert. CMK bedeutet, dass der Kunde den Schlüssel generiert, rotiert und im Notfall widerrufen kann – womit der Anbieter auch bei physischem Zugriff auf Speichersysteme keinen Klartext-Datenzugriff hat.
- > Key Access Logging (unveränderlich):
Jeder Schlüsselzugriff muss unveränderlich protokolliert werden. Audit Logs des KMS müssen durch WORM-Mechanismen (Write Once Read Many) gegen Manipulation geschützt sein und selbst verschlüsselt gespeichert werden.
- > Schlüsselrotation und Revocation:
Automatisierte Schlüsselrotationsprozesse (mindestens jährlich, bei Sicherheitsvorfällen sofort) und getestete Revocation-Prozesse (Schlüsselentzug bei Anbieterwechsel oder Kompromittierung) müssen vertraglich und technisch sichergestellt sein.

6.4 Mandantentrennung in GPU-Clustern

In mandantenfähigen GPU-Hosting-Umgebungen ist die Isolation zwischen Tenants auf mehreren Ebenen zu implementieren:



Compute-Isolation (GPU):

MIG-Partitionierung als stärkster Baustein für GPU-Level-Isolation mit dedizierten Compute Engines, VRAM-Partitionen und L2-Cache-Isolation. Time-Slicing als Alternative für weniger sensitive Workloads.

Control-Plane-Isolation (Kubernetes):

Getrennte Kubernetes-Namespace mit strikten RBAC-Konfigurationen (Least Privilege), Namespace-übergreifende NetworkPolicies (kein East-West-Traffic zwischen Tenants ohne explizite Freigabe), dedizierte Node Pools pro Sensitivitätsklasse und ggf. dedizierte Cluster für hochkritische Workloads.



Daten-Isolation (Storage):

Getrennte S3-Buckets mit isolierten IAM-Policies pro Tenant, getrennte Verschlüsselungs-DEKs mit tenant-spezifischen KEKs, keine gemeinsamen Bucket-Prefixes zwischen Tenants.

Netzwerk-Isolation:

VPC/VLAN-Segmentierung zwischen Tenants, strict Egress-Control, keine unautorisierten DNS-Lookups. Service Mesh (z.B. Istio) für mTLS zwischen allen Services und Policy-Enforcement auf Layer 7.

Wichtig:

Lassen Sie Isolationsbehauptungen durch unabhängige Penetrationstests validieren. Fragen Sie Ihren Anbieter explizit nach dem Ergebnis von Tenant-Isolation-Tests und dem Scope von Audit-Berichten in Bezug auf Multi-Tenancy.

6.5 Secure Model Provenance und Supply Chain Security

KI-Systeme erben alle Sicherheitsrisiken ihrer Komponenten – und fügen spezifische KI-Risiken hinzu. Eine belastbare Supply Chain Security für KI-Plattformen umfasst:

Container Image Signierung:

Alle Container Images (Base Images, CUDA-Layer, Anwendungsimages) werden mit notariell signierten Image-Manifesten versehen. Kubernetes Admission Controller (OPA/Gatekeeper, Kyverno) erzwingen, dass nur signierte Images aus vertrauenswürdigen Registries deployed werden.

Software Bill of Materials (SBOM):

Für alle Plattformkomponenten (Kubernetes Add-ons, GPU Operator, Treiber, Serving-Frameworks) wird eine automatisiert generierte SBOM geführt, die für Vulnerability-Scans und Compliance-Nachweise genutzt wird.

Model Weights Provenance:

Vor dem Deployment von Pre-Trained Models aus öffentlichen Repositories (HuggingFace, Model Zoo) sind Herkunft, Lizenz, bekannte Vulnerabilitäten und Backdoor-Tests (z.B. durch Techniken aus der ML-Security-Forschung) zu dokumentieren und zu bewerten.

Model Registry als Governance-Kern:

Ein zentrales Model Registry (z.B. MLflow Model Registry) mit Versions-, Lineage- und Metadatenverwaltung stellt sicher, dass deployed Modelle eindeutig einem Approval-Prozess zugeordnet werden können und unautorisierte Modelle nicht in Produktion gelangen.

7. Betrieb, SRE und Migration

7.1 SRE-Anforderungen für produktive KI-Plattformen

Site Reliability Engineering (SRE) für KI-Plattformen unterscheidet sich wesentlich vom klassischen Web-Application-SRE. Klassische Metriken wie CPU-Utilization und Memory-Swap-Rate sind unzureichend. Für KI-Plattformen sind folgende Service Level Indicators (SLIs) maßgeblich:

- > Inferenz-Performance:
p95/p99 Latenz pro Endpoint (Chat-API, Embedding-API, Retrieval-API), Tokens per Second als Throughput-Metrik, Time-to-First-Token (TTFT) als Nutzererlebnis-Metrik.
- > GPU-Ressourcenmetriken:
GPU Utilization (%), GPU Memory Used/Total, GPU Temperature, Thermal Throttling Events, GPU-to-GPU-Bandbreite bei Multi-GPU-Workloads.
- > Model-Lifecycle-Metriken:
Model Load Time (relevanter als oft angenommen, besonders bei großen LLMs), Model Deployment Success Rate, Rollback-Trigger-Rate.
- > Queue-Metriken:
Request Queue Depth, Queue Wait Time, Batch Size Distribution – besonders relevant für Batch-Inferenz und Throughput-Optimierung.

- > Fehlermetriken:
HTTP-Error-Rate (4xx, 5xx), gRPC-Error-Rate, Inference-Timeout-Rate, Out-of-Memory-Events.
- > RAG-spezifische Metriken:
Retrieval-Latenz (p95), Vector Database Query Time, Index Freshness (Verzögerung zwischen Dokumenten-Update und Index-Verfügbarkeit), Embedding-Throughput.

7.2 SLA-Anforderungen an Managed GPU-Hosting-Anbieter

Für mittelständische Unternehmen ohne spezialisierte GPU-Betriebsteams ist die Qualität des Managed Services entscheidend. SLAs sollten weit über "Uptime" hinausgehen:

- > Verfügbarkeits-SLAs:
Getrennte SLAs für GPU-Nodes, Kubernetes Control Plane, Objektspeicher und Netzwerk. Geplante Wartungsfenster müssen vertraglich definiert und angekündigt werden. Unterscheidung zwischen "Cluster verfügbar" und "GPU-Kapazität für Workload verfügbar".
- > Recovery-Ziele:
RTO (Recovery Time Objective) und RPO (Recovery Point Objective) für Objektspeicher, Model Registry und zentrale Control Plane. Für produktive KI-APIs sind RTO < 4 Stunden und RPO < 1 Stunde als Orientierungswerte zu nennen, müssen aber workload-spezifisch definiert werden.

- > Incident Response:
Definierte Reaktionszeiten nach Schweregrad (P1-P4), dedizierte Support-Kontakte für Sicherheitsvorfälle, Kommunikationspflichten bei Datenpannen (Art. 33 DSGVO: Meldung an Aufsichtsbehörde innerhalb von 72 Stunden).
- > Patch- und Vulnerability Management:
Verbindliche SLAs für kritische Sicherheitspatches (GPU-Treiber, Kernel, Kubernetes-Komponenten). Für kritische CVEs (CVSS >= 9.0) sollten Patches innerhalb von 24-48 Stunden angewendet werden können.
- > Kapazitätsgarantien:
Reserved Capacity für kritische Zeitfenster (Monatsabschlüsse, Kampagnen) als vertraglich zugesicherter Baustein, nicht nur "Best Effort".
- > Auditrechte:
Recht auf Einsicht in CS-Testierungsberichte, ISO 27001-Zertifikate, Penetrationstestergebnisse und Protokolle von Subunternehmeränderungen.

7.3 Migration zu deutschem KI-Hosting

Die Migration von KI-Workloads in eine deutsche Hosting-Umgebung ist ein strukturierter Prozess, der mehrere Integrationsdimensionen umfasst:

CI/CD für Modelle (MLOps-Pipeline):

Model Packaging, Container Builds, Canary-Deployment-Strategie, A/B-Testing für Modell-Updates, automatisiertes Rollback bei Qualitätsverschlechterung. Integration in bestehende DevOps-Toolchains (GitLab, Jenkins, GitHub Actions).



Identity und Access Management Integration:

OIDC/SSO-Anbindung an den Unternehmens-IdP (Active Directory, Azure AD, Keycloak). Service Account Management für Pipeline-Prozesse. Secrets Management (HashiCorp Vault, Sealed Secrets in Kubernetes) für API-Keys, Datenbank-Credentials und Zertifikate.

Datenpipeline-Integration:

S3-Kompatibilitätstest für alle bestehenden Datenpipelines (Spark, dbt, custom ETL). Validierung von Multipart Upload, Bucket Policy-Semantik und Event Notification-Funktionalität. Für RAG-Systeme: Migration von Dokumenten-Store und Index-Rebuild unter datenschutzrechtlichen Anforderungen (Zugriffskontrolle, Metadaten-Vollständigkeit).

Monitoring-Überführung:

Portierung von Dashboards und Alerting-Regeln. Integration von GPU-spezifischen Metriken (DCGM Exporter für NVIDIA) in bestehende Observability-Stacks (Prometheus, Datadog, Splunk).

8. Wirtschaftlichkeit und Total Cost of Ownership

8.1 TCO-Dimensionen im KI-Hosting

Die häufig anzutreffende Vereinfachung "GPU-Stunden-Kosten × Auslastung = TCO" führt zu erheblichen Fehlkalkulationen.

Eine vollständige TCO-Betrachtung für KI-Hosting umfasst mindestens sechs Kostendimensionen:

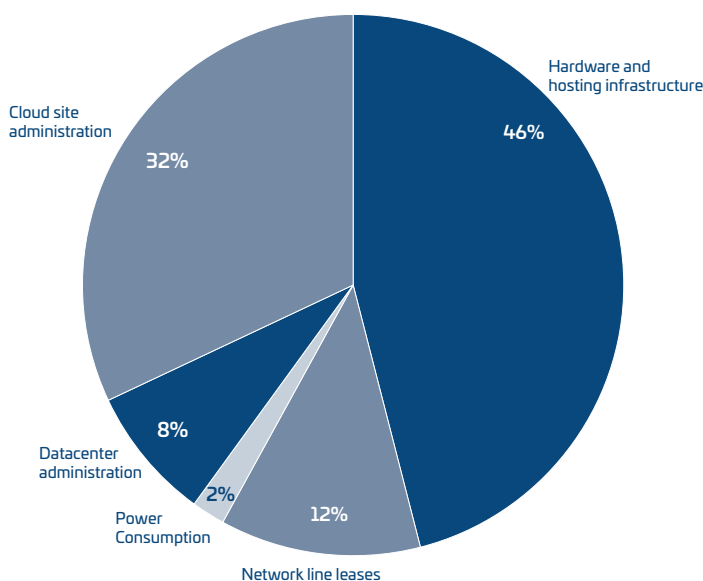
Compute-Kosten:

GPU-Server (CAPEX oder OPEX), CPU, RAM, Netzwerkkarten (insbesondere RDMA-fähige NICs), NVMe-Storage für Hot-Tier.

Facility-Kosten:

Stromkosten (stark abhängig von GPU-TDP und PUE), Kühlung (Luft oder Flüssigkeitskühlung), Colocation-Gebühren (Rack-Miete, Cross-Connects). Bei einem PUE von 1,5 und 0,22 €/kWh Strompreis entstehen für ein 8×H100-System allein an Energiekosten ca. 20.000-25.000 € pro Jahr.

TCO Breakdown:



Netzwerkkosten:

High-Speed-Fabric (InfiniBand, 100G/400G Ethernet), Switches, Verkabelung, Datentransferkosten (Egress bei Cloud-Modellen).

Software- und Lizenzkosten:

Kubernetes Enterprise-Support, Security-Tools (Container Scanning, SIEM), Monitoring (Datadog, Grafana Enterprise), MLOps-Plattform, ggf. Modell-Lizenzen.

Personalkosten:

SRE-Engineering (GPU-Betrieb, Kubernetes, Netzwerk), Security-Engineering (CISO, SecOps), Compliance (DPO, GRC), Data/AI Engineering. Oft der größte und am häufigsten unterschätzte Kostentreiber – für On-Premises GPU-Betrieb sind typischerweise 2-4 Vollzeit-FTEs für Infrastruktur und Security einzuplanen.

Exit- und Migrationskosten:

Datenexportkosten, Re-Architecting für neues Hosting-Modell, Vertragsstrukturierung und -kündigung. Besonders bei Hyperscalern mit hohem Lock-in können Migrationskosten erheblich sein.

8.2 TCO-Vergleich: Hosting-Modelle im Überblick

Kostentreiber	On-Premises	Private Cloud (DE)	Managed GPU (DE)	Hyperscaler
CAPEX Hardware	Sehr hoch	Mittel	Nicht vorhanden	Nicht vorhanden
OPEX Energie/Facility	Mittel-Hoch (direkt)	Im Preis enthalten	Im Preis enthalten	Im Preis (indirekt)
Personalkosten SRE/HPC	Hoch (2-4 FTE)	Mittel-Hoch (1-3 FTE)	Niedrig-Mittel (0,5-1 FTE)	Niedrig-Mittel (0,5-2 FTE)
Lock-in / Exit-Kosten	Niedrig	Mittel	Mittel	Hoch
Kostenvolatilität	Niedrig	Mittel	Mittel	Hoch (Spot-Preise)
Compliance-Zusatzkosten	Niedrig	Niedrig	Niedrig	Hoch (TIA, SCC-Aufwand)

8.3 Benchmark-Methodik für fundierte Entscheidungen

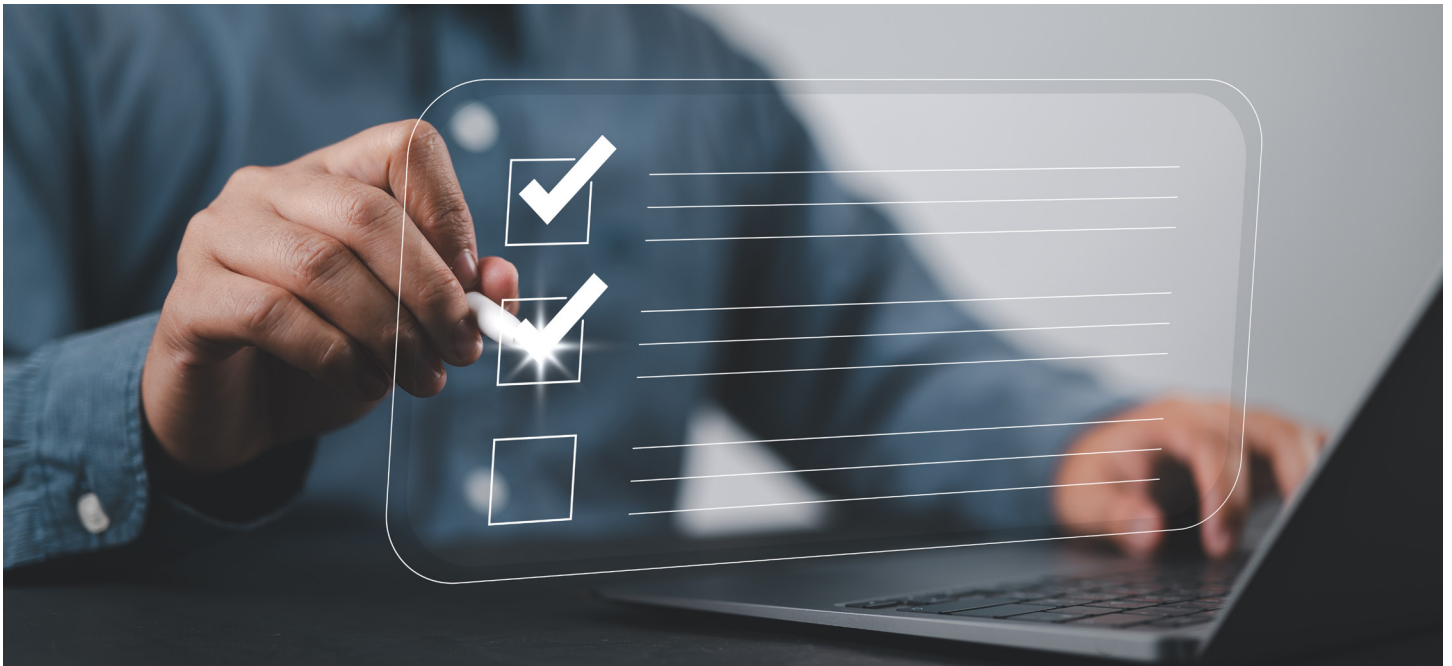
Ein objektiver TCO-Vergleich zwischen Hosting-Modellen erfordert einen workload-basierten Benchmark, kein Preislistenvergleich. Empfohlene Benchmark-Suite:

- > Inferenz-Performance: Tokens/sec bei realistischen Prompt-Längen (kurz/mittel/lang), p95/p99 Latenz bei definierten Concurrency-Stufen (z.B. 10, 50, 200 parallele Anfragen), Time-to-First-Token.
- > RAG-Performance: Retrieval-Latenz über Vektordatenbank (p95), Gesamt-Antwortzeit (End-to-End), Index-Rebuild-Dauer bei Dokumenten-Updates.
- > Betriebs-Benchmark: Model-Deployment-Zeit (Zeit von Git-Commit bis Live-Endpoint), Rollback-Dauer, Recover-Zeit nach simuliertem Node-Failure.
- > Security-Verifikation: IAM-Fehlkonfigurationstest (Tenant-A darf nicht auf Tenant-B-Ressourcen zugreifen), Secret-Scan, Key Access Logging-Verifikation.



Messartefakte müssen reproduzierbar sein: Versionierte Helm-Charts oder Terraform-Konfigurationen, definierte synthetische Test-Datensätze (kein Produktivdaten-Einsatz), protokollierte Ergebnisse mit Zeitstempel und Versionsinformationen aller beteiligten Komponenten (GPU-Treiber, CUDA, Kubernetes-Version, Modell-Revision).

9. Entscheidungsleitfaden



9.1 Entscheidungsdimensionen: Die richtigen Fragen stellen

Bevor ein Hosting-Modell ausgewählt wird, müssen folgende Fragestellungen systematisch beantwortet werden:

Datensensitivität:

Welche Kategorien personenbezogener Daten werden verarbeitet (Art. 9 DSGVO-Kategorien, normale personenbezogene Daten, keine personenbezogenen Daten)? Gibt es Betriebs- oder Geschäftsgeheimnisse, die besonderen Schutz erfordern? Welche Konsequenzen hätte ein Datenleck?

Regulatorisches Umfeld:

Ist das Unternehmen als wesentliche oder wichtige Einrichtung unter NIS2 eingestuft? Gelten branchenspezifische Regularien (DORA für Finanzsektor, ISO 13485 für Medizinprodukte, KRITIS-Anforderungen)? Sind die KI-Systeme als Hochrisiko-KI nach EU AI Act Annex III einzustufen?

Strategische Souveränitätsprioritäten:

Ist Datensouveränität und -unabhängigkeit eine explizite strategische Anforderung der Unternehmensführung? Haben Kunden oder Vertragspartner Anforderungen an

Datenlokalisierung? Gibt es langfristige Ziele bzgl. Unabhängigkeit von US-amerikanischen Technologieanbietern?

Technische Skalierungs- und Latenzanforderungen:

Wie hoch ist das erwartete Request-Volumen? Gibt es harte Latenz-SLAs (z.B. < 500 ms p95)? Sind globale Nutzer oder Multi-Region-Anforderungen relevant?

Interne Betriebsfähigkeit:

Verfügt das Unternehmen über SRE-Expertise für GPU-Infrastruktur und Kubernetes? Kann der Betrieb eines dedizierten GPU-Clusters intern verantwortet werden?

9.2 Entscheidungsmatrix

Entscheidungskriterium	Gewichtung	On-Premises	Private Cloud (DE)	Managed GPU(DE)	Hyperscaler
Compliance/Rechtssicherheit	25%	5/5	5/5	4/5	3/5
Datenkontrolle/Schlüsselhoheit	15%	5/5	4/5	4/5	3/5
Performance/GPU-Skalierung	15%	3/5	4/5	4/5	5/5
Time-to-Value	10%	2/5	3/5	4/5	5/5
Betriebsfähigkeit (Skill-Fit)	15%	2/5	3/5	4/5	4/5
TCO-Vorhersagbarkeit	15%	4/5	4/5	4/5	2/5
Exit- und Portabilität	5%	5/5	4/5	4/5	2/5

Diese Matrix ist als Orientierungsrahmen zu verstehen. Für die konkrete Unternehmensentscheidung empfehlen wir, die Gewichtungen auf die spezifische Unternehmenssituation anzupassen und jeden Score mit einer dokumentierten Begründung zu versehen.

9.3 Empfehlung für den deutschen Mittelstand

Basierend auf der Analyse der Hosting-Optionen und der regulatorischen Anforderungen lässt sich für mittelständische Unternehmen in Deutschland ohne spezialisiertes HPC-Betriebsteam eine klare Empfehlung ableiten:

Für die meisten mittelständischen Unternehmen mit sensiblen Daten, regulatorischen Anforderungen oder strategischen Souveränitätszielen bietet kontrolliertes, in Deutschland betriebenes Managed GPU-Hosting oder eine dedizierte Private Cloud die beste Balance aus Compliance-Sicherheit, Time-to-Value und Betriebsfähigkeit. Voraussetzung ist ein strikter Vertrags- und Sicherheitsrahmen (CS-/ISO-Orientierung, IT-Grundschutz-Bausteine für Container/Kubernetes) und eine klare Exit-Strategie.

Diese Empfehlung gilt insbesondere für folgende Unternehmensprofile: Finanzdienstleister, Versicherungen und Steuerberater mit besonderen Anforderungen aus DSGVO und DORA; Unternehmen der kritischen Infrastruktur (Energie, Wasser, Transport) mit NIS2-Pflichten; Gesundheitsunternehmen, Krankenhäuser und Medizintechnikhersteller mit besonderen Kategorien personenbezogener Daten; sowie alle Unternehmen, die KI-Systeme unter eigenem Namen an externe Kunden als Anbieter nach EU AI Act bereitstellen.

10. Glossar und Referenzen

Glossar

> **AI Act**

Verordnung (EU) 2024/1689 über Künstliche Intelligenz; stufenweise Anwendung ab Februar 2025 bis August 2027.

> **AVV**

Auftragsverarbeitungsvertrag nach Art. 28 DSGVO; regelt die Datenverarbeitung durch Dritte im Auftrag des Verantwortlichen.

> **BSI C5**

Cloud Computing Compliance Criteria Catalogue des BSI; etablierter Mindestanforderungskatalog für Cloud-Sicherheit in Deutschland.

> **CMK (Customer Managed Keys)**

Vom Kunden kontrollierte kryptographische Schlüssel im KMS/HSM; ermöglichen echte Schlüsselhoheit gegenüber dem Hosting-Anbieter.

> **DCGM**

NVIDIA Data Center GPU Manager; liefert GPU-Betriebsmetriken für Monitoring und Observability.

> **Deployer**

Organisation oder Person, die ein KI-System in eigenen Prozessen oder für Dritte einsetzt (Betreiber nach EU AI Act).

> **DSFA/DPIA**

Datenschutz-Folgenabschätzung nach Art. 35 DSGVO; Pflicht bei Verarbeitungen mit voraussichtlich hohem Risiko.

> **DPF**

EU-U.S. Data Privacy Framework; Angemessenheitsbeschluss der EU-Kommission für Datentransfers in zertifizierte US-Unternehmen.

> **Fine-Tuning**

Anpassung eines vortrainierten KI-Modells auf spezifische Daten oder Aufgaben; LoRA/QLoRA sind ressourceneffiziente Methoden.

> **GPUDirect RDMA**

Direkter Datentransfer zwischen GPU-Speicher und Netzwerkkarte ohne CPU-Umweg; relevant für skalierbares Training.

> **HSM**

Hardware Security Module; physisch gesichertes Gerät zur sicheren Speicherung und Nutzung kryptographischer Schlüssel.

> **IAM**

Identity and Access Management; Verwaltung von Identitäten, Zugriffsrechten und Authentifizierung.

> **KMS**

Key Management Service; Dienst oder Software zur Verwaltung kryptographischer Schlüssel.

> **KServe**

Kubernetes-natives Framework für Model Serving mit Autoscaling, Routing und Monitoring.

> **LLM**

Large Language Model; großes Sprachmodell auf Transformer-Basis (z.B. GPT-4, Llama 3, Mistral).

> **LoRA/QLoRA**

Low-Rank Adaptation / Quantized LoRA; Parameter-Efficient Fine-Tuning-Methoden mit reduziertem Speicher- und Rechenaufwand.

> **MIG**

Multi-Instance GPU; NVIDIA-Technologie zur physischen Partitionierung von GPUs in isolierte Instanzen.

> **MLOps**

Machine Learning Operations; Prozesse und Werkzeuge für Entwicklung, Deployment, Monitoring und Governance von ML-Modellen.

> **NVLink/NVSwitch**

NVIDIA-Hochbandbreiten-Interconnect für GPU-zu-GPU-Kommunikation innerhalb und zwischen Servern.

> **OIDC**

OpenID Connect; Authentifizierungsprotokoll auf Basis OAuth 2.0; Standard für SSO-Integration.

> **PUE**

Power Usage Effectiveness; Verhältnis von Gesamt-Facility-Energieverbrauch zu IT-Energieverbrauch; Effizienzmaß für Rechenzentren.

> **RAG**

Retrieval-Augmented Generation; LLM-Architektur, die externe Wissensquellen (Dokumentindex) in die Antwortgenerierung einbezieht.

> **RBAC**

Role-Based Access Control; Zugriffssteuerung auf Basis zugewiesener Rollen.

> **RDMA/RoCE**

Remote Direct Memory Access / RDMA over Converged Ethernet; Netzwerktechnologie für niedrige Latenz und hohen Durchsatz.

> **S3**

Simple Storage Service (AWS); de-facto-Standard-Objektspeicher-API; auch als generischer Begriff für S3-kompatible Systeme.

> **SBOM**

Software Bill of Materials; vollständige Liste aller Software-Komponenten und Abhängigkeiten in einem System.

> **SLI/SLO/SLA**

Service Level Indicator / Objective / Agreement; Metriken, Ziele und vertragliche Zusagen für Service-Qualität.

> **TOM**

Technische und Organisatorische Maßnahmen nach Art. 32 DSGVO; dokumentierte Sicherheitsmaßnahmen für Datenverarbeitung.

> **Triton**

NVIDIA Triton Inference Server; High-Performance Model Serving Framework mit HTTP/gRPC API.

> **vLLM**

Open-Source LLM Inference Engine mit PagedAttention-Optimierung für hohen Durchsatz.

wesentliche Referenzen

Folgende Dokumente sind für die praktische Umsetzung der in diesem Whitepaper beschriebenen Anforderungen maßgeblich:

> **Verordnung (EU) 2024/1689 (EU AI Act)**

Amtsblatt der Europäischen Union, abrufbar über EUR-Lex

> **EU AI Act Service Desk**

Implementierungszeitplan: <https://ai-act-service-desk.ec.europa.eu/en/ai-act/timeline>

> **Datenschutzkonferenz (DSK)**

Orientierungshilfe KI und Datenschutz (06.05.2024): <https://www.datenschutzkonferenz-online.de>

> **Datenschutzkonferenz (DSK)**

Orientierungshilfe zu RAG-Systemen (10/2025): <https://www.datenschutzkonferenz-online.de>

> **Datenschutzkonferenz (DSK)**

Entschiebung Confidential Cloud Computing (06/2025): <https://www.datenschutz-berlin.de>

> **EDPB Recommendations 01/2020**

Supplementary Measures for Data Transfers: <https://www.edpb.europa.eu>

> **EU-Kommission**

Implementing Decision EU 2023/1795 (EU-U.S. Data Privacy Framework): EUR-Lex

> **BSI**

Cloud Computing Compliance Criteria Catalogue C5:2020: <https://www.bsi.bund.de>

> **BSI**

IT-Grundschutz-Kompendium (Bausteine SYS.1.1, SYS.1.6, APP.4.4): <https://www.bsi.bund.de>

> **BSI**

Technische Richtlinien TR-02102-1 und TR-02102-2 (Kryptographie/TLS): <https://www.bsi.bund.de>

> **BSI**

Practical AI Security Guide 2023: <https://www.bsi.bund.de>

> **ENISA**

Securing Machine Learning Algorithms: <https://www.enisa.europa.eu>

> **NVIDIA**

MIG User Guide, GPUDirect RDMA Documentation, H100 Specifications: <https://docs.nvidia.com>

> **Kubernetes**

GPU Scheduling Documentation: <https://kubernetes.io>

> **KServe**

Projektdokumentation: <https://kserve.github.io>

> **MLflow**

Model Registry Documentation: <https://mlflow.org>

> **IEA**

Energy and AI Report 2024: <https://www.iea.org>

> **Uptime Institute**

Global Data Center Survey 2024: <https://datacenter.uptimeinstitute.com>

> **Bitkom**

Digitale Souveränität / Wunsch nach deutscher Cloud: <https://www.bitkom.org>

> **BFDI**

Info 1: Datenschutz und Auftragsverarbeitung: <https://www.bfdi.bund.de>