

DAS KI-KONTINUUM:

WELCHE INFRASTRUKTUR EIGNET SICH AM BESTEN FÜR INFERENZ?

Die Einplanung von KI ist für jede Aktualisierung des Rechenzentrums unerlässlich. GPUs sind wichtig für große KI-Auslastungen, aber die neuesten Generationen von CPUs können neben allgemeinen Auslastungen eine Vielzahl von KI-Aufgaben übernehmen. Beachten Sie diese Aspekte, wenn Sie Ihre wachsenden KI-Inferenzanforderungen ermitteln.

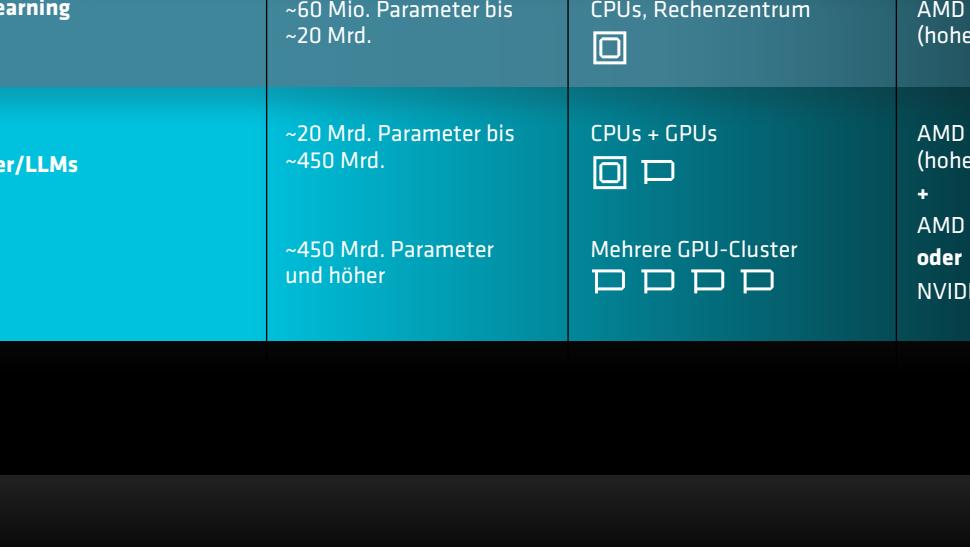
Oft benötigt KI keine Echtzeitergebnisse

Moderne CPUs können kleine bis mittlere KI-Inferenzauslastungen mit einer Latenz von weniger als einer Sekunde ausführen. Wenn KI-Inferenzauslastungen wachsen oder die Reaktionszeiten kürzer werden, müssen Sie möglicherweise diskrete Beschleuniger hinzufügen.

BATCH-VERARBEITUNG	MITTLERE LATENZ	NIEDRIGE LATENZ	NAHEZU ECHTZEIT	ECHTZEIT
Minuten bis Tage	Sekunden bis Minuten	~500 ms bis Sekunden	~100 ms bis ~500 ms	<10 ms bis ~50 ms
ANWENDUNGSFÄLLE				
• Dokumentverarbeitung und -klassifizierung • Data Mining und Analysen • Wissenschaftliche Simulationen	• Übersetzung • Indexierung • Content Moderation • Prädiktive Instandhaltung	• Virtuelle Assistenten • Chatbots • Expert Agents • Videountertitelung	• Betrugerkennung • Entscheidungsfindung • Dynamische Preisfindung • Audio- und Videofilterung	• Finanzhandel • Telekommunikation und Netzwerke • Autonome Systeme
			Mehrere GPU-Cluster	

Wenn KI-Auslastungen ansteigen, werden GPUs immer kostengünstiger

CPUs allein können gemischte Unternehmensauslastungen und KI unterstützen. Wenn Modellgröße, Komplexität und Volumen ansteigen, können GPU-Cluster mehr Performance pro Euro bieten.



Die Grafik dient nur zu Illustrationszwecken. Der Schnittpunkt variiert je nach Auslastungen und Prozessormodellen.

Unterschiedliche Modelle haben einzigartige Verarbeitungsanforderungen

Maschinelles Lernen, Grafikverarbeitung und statistische Methoden werden auf CPUs außergewöhnlich gut ausgeführt. Kleine bis mittlere Large Language Models (LLMs) funktionieren gut auf den neuesten CPUs. Größere Modelle können einen deutlichen Nutzen aus KI-Beschleunigern ziehen.

	MODELLGRÖSSE	PROZESSOR	AMD LÖSUNG
Deep Learning	Klassisches maschinelles Lernen ~1 MB bis ~200 MB	CPUs, integriert in das Rechenzentrum 	AMD Ryzen™ CPUs AMD EPYC™ CPUs
Transformer/LLMs	~60 Mio. Parameter bis ~20 Mrd. ~20 Mrd. Parameter bis ~450 Mrd. ~450 Mrd. Parameter und höher	CPUs, Rechenzentrum CPUs + GPUs Mehrere GPU-Cluster 	AMD EPYC CPUs (hohe Kernzahl) AMD EPYC CPUs (hohe Frequenz) + AMD Instinct™ GPUs oder NVIDIA GPUs

AMD EPYC CPUs brillieren mit KI für Unternehmen

AMD EPYC CPUs der 5. Generation bieten wichtige Performance-Verbesserungen für KI-Auslastungen:

Bis zu 3,8-facher Durchsatz für End-to-End-KI im Vergleich zu CPUs der Konkurrenz¹

Bis zu 90 % schnellerer Durchsatz auf Llama 3.1.8B bei BF16 im Vergleich zu CPUs der Konkurrenz²

Bis zu 86 % schnellere Facebook AI Similarity Search (FAISS) im Vergleich zu EPYC CPUs der vorherigen Generation³

AMDEPYC™ CPUs DER 5. GENERATION: DIE BESTE CPU FÜR UNTERNEHMENS-KI⁴

Erfahren Sie, warum AMD EPYC CPUs der 5. Generation bei KI-Inferenzauslastungen brillieren.

EPYC für KI-Inferenz besuchen

1. KI-Durchsatztest ist vom TPEx-AI-Benchmark abgeleitet und als solcher nicht mit den veröffentlichten TPEx-AI-Ergebnissen vergleichbar, da die Ergebnisse des durchgängigen KI-Durchsatztests nicht der TPEx-AI-Spezifikation entsprechen. 2P AMD EPYC 9656 (384 Kerne gesamt), 12 Instanzen mit 32 Kernen, NPS1, 1,5 TB x 24 x 64 GB DDR5-6400 (bei 6000 MT/s), 1 DPC, 1,0 Gbit/s NetXtreme BCM5720 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled) 2P AMD EPYC 9654 (92 Kerne gesamt) 6 Instanzen mit 32 Kernen, NPS1, 1,5 TB x 24 x 64 GB DDR5-6400 (bei 6000 MT/s), 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.3 LTS, BIOS 1006C (SMT = off, Determinism = Power, Turbo Boost = Enabled) 4 Instanzen mit 32 Kernen, AMX Ein, 1TB x 6 x 64 GB DDR5-5600, 1 DPC, 1,0 Gbit/s NetXtreme BCM5720 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.3 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 4 Instanzen mit 48 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 8 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 12 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 16 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 24 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 32 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 48 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 64 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 96 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 144 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 288 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 576 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 1152 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 2304 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 4608 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 9216 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 18432 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 36864 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 73728 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 147456 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 294912 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 589824 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 1179648 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 2359296 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism = Power, Turbo Boost = Enabled), 4718592 Instanzen mit 16 Kernen, 1,5 TB x 24 x 64 GB DDR5-4800, 1 DPC, 1,0 Gbit/s NetXtreme BCM5719 Gigabit Ethernet PCIe, 3,5 TB Samsung MZWL03T8HCL5-00A07 NVMe®, Ubuntu® 22.04.4 LTS, 6.8.0-40-generic (tuned-adm profile throughput-performance, ulimit -n 198096812, ulimit -n 1024, ulimit -s 8192), BIOS RVT01000C (SMT = off, Determinism