

# Open-Source und Offline-KI

# Beruflicher Hintergrund

## Profil

Johannes Hofmann, selbstständiger Art Director und KI-Manager (IHK), mit jahrelanger nationaler und internationaler Agenturerfahrung und 17 Jahren Selbstständigkeit

## Digital-Expertise und KI

- Einführung des ersten digitalen Schauraums der Audi AG, Prototypen für Content-Module, Technische Spezifikationen
- Prototypen für immersive Gamification-Module im Audi Forum Ingolstadt
- Content für das lokale Digital-Signage-System des Audi Forum Ingolstadt und Einbindung von Echtzeitdaten und Sprachadaptionen von Video-Content
- UX/UI-Design für Augmented Reality-Formate für Audi-Technik
- KI-gestützte Katalogautomatisierung
- Produktivsystem eines Corporate-Voice-LLM für eine deutsche Automotive-Agentur





**“ Wer seine Werkzeuge nicht  
kontrolliert, wird von ihnen  
kontrolliert.”**  
— Open-Source-Prinzip

# Verstehen – Ausprobieren – Reflektieren

1. Das Problem: Warum digitale Souveränität?
2. Die Lösung: Open-Source KI lokal nutzen
3. Live-Demo: OpenWebUI in Aktion
4. Praxis: Was brauche ich wirklich?
5. Next Steps: Heute Abend loslegen



Kann ein Staat souverän  
bleiben, wenn seine KI-Modelle  
aus dem Ausland kommen?

# Was passiert, wenn ChatGPT morgen offline geht?

Rechtliche Gründe? Datenschutz-Urteil?  
Server-Ausfall? Ihre Pressetexte, Förderanträge,  
Content-Planung...  
...alles steht still.

Digitale Souveränität = Handlungsfähig bleiben



# Drei unbequeme Wahrheiten über Cloud-KI

1. Ihre Prompts trainieren fremde Modelle

→ Steht in den AGBs

2. Token-Limits & Kosten

→ 1 Mio. Tokens = 15€

3. Keine Internetverbindung = keine KI



# Cloud-KI vs. Lokale KI: Der ehrliche Vergleich

## **Cloud-KI: Schnell, bequem, abhängig**

- Pro: Beste Modelle, keine Installation
- Contra: Kosten, Datenschutz

## **Lokale KI: Langsamer, souverän**

- Pro: Geringe Kosten, volle Kontrolle
- Contra: Hardware, Setup

# Die Lösung: Open-Source-Modelle lokal nutzen

**Llama, Mistral, Phi**

→ Leistungsstarke Modelle, frei verfügbar

**OpenWebUI, LM Studio**

→ Tools mit Nutzerfreundlichkeit

**Einmal installiert: Unbegrenzt nutzbar**

# Der Automation-Vorteil: Latenz wird irrelevant

**Früher: 4 Tage für 10 Pressetexte**

**Heute: 2 Stunden Arbeitszeit**

- n8n-Workflow läuft nachts
- Morgens: Beste auswählen

**Die Maschine arbeitet 24/7. Sie kuratieren.**





# Live-Demo: Lokale KI in Aktion

## Was Sie gleich sehen:

- OpenWebUI läuft auf diesem Laptop
- Kein Internet, keine Cloud
- Prompt: Presstext
- Side-by-Side: Cloud vs. Lokal

# Tool 1: OpenWebUI – Die All-in-One-Lösung

## Was ist OpenWebUI?

- ChatGPT-ähnliche Oberfläche
- Unterstützt Llama, Mistral, Phi
- Browser: localhost:8080

**Für wen? Fortgeschrittene, Teams**

# Tool 2: LM Studio – Der Einsteiger-Favorit

## Was ist LM Studio?

- Desktop-App (Win/Mac/Linux)
- Ein-Klick-Installation
- Perfekt für erste Schritte

## Für wen? Einsteiger



# Was brauche ich? Hardware-Anforderungen

## Use Case 1: Einfache Texte

→ Llama 3.2 3B | 8 GB RAM

## Use Case 2: Presetexte

→ Llama 3.1 8B | 16 GB RAM

## Use Case 3: Komplexe Analysen

→ Mixtral 8x7B | 32 GB RAM

# 1. Lokale LLM-Nutzung für Text- und Bild-generierung auf Laptop/Desktop-PC/Mac

- Empfohlene Hardware: CPU: Moderne Mehrkern-Prozessoren wie Intel Core i7 oder AMD Ryzen 7
- RAM: Mindestens 32 GB
- GPU: Dedizierte Grafikkarte mit mindestens 12 GB VRAM, z.B. NVIDIA RTX 3060
- Geschätzte Kosten: Desktop-PC: Zwischen 1.500€ und 2.500€, abhängig von den spezifischen Komponenten  
Laptop: Hochleistungs-Laptops mit vergleichbarer Ausstattung liegen zwischen 2.000€ und 3.000€ · Mac mini: 3.500 € · Mac Book Pro: 2000 – 4000 €





## 2. Lokale LLM-Nutzung für kleine Arbeitsgruppen (5–15 Personen) inklusive Automatisierungen

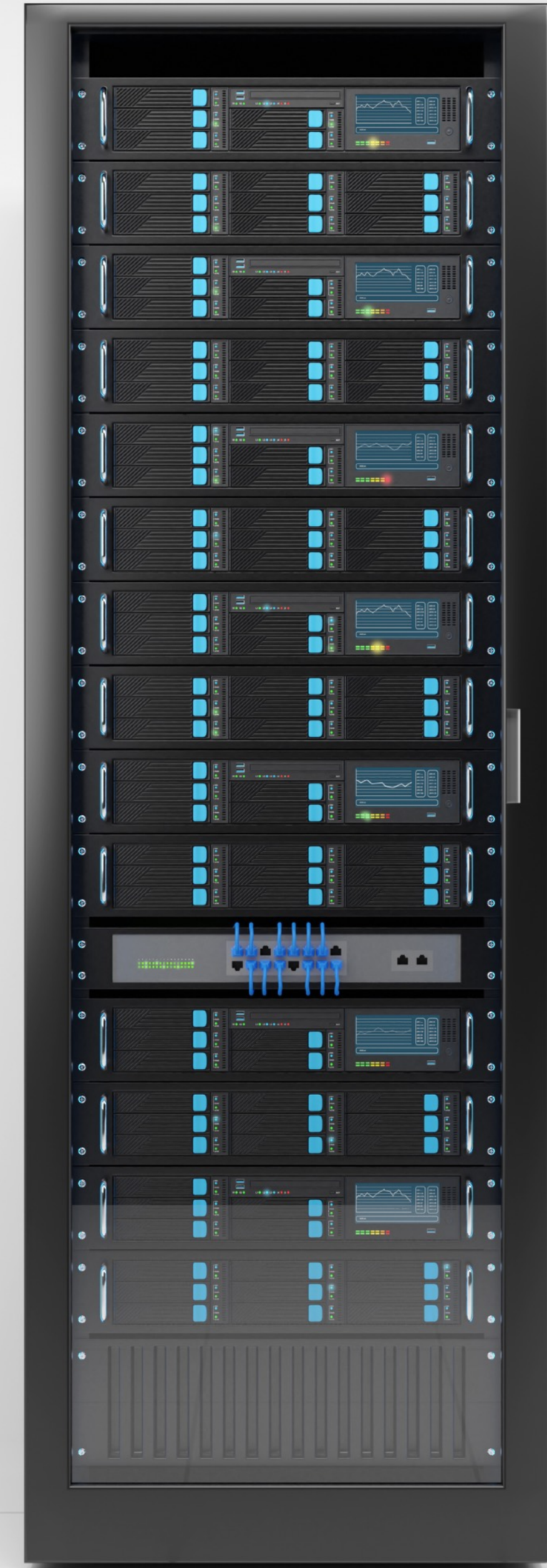
- Empfohlene Hardware: Leistungsstarke Mehrkern-Prozessoren wie AMD Ryzen 9
- RAM: Mindestens 128 GB
- GPU: NVIDIA RTX 4090 mit 24 GB VRAM
- Netzwerk: Gigabit-Ethernet
- Geschätzte Kosten: Desktop-PC: Zwischen 3000 € und 5000 €, abhängig von den spezifischen Komponenten





### 3. Lokale LLM-Nutzung in einem KMU mit ca. 300 Mitarbeitern

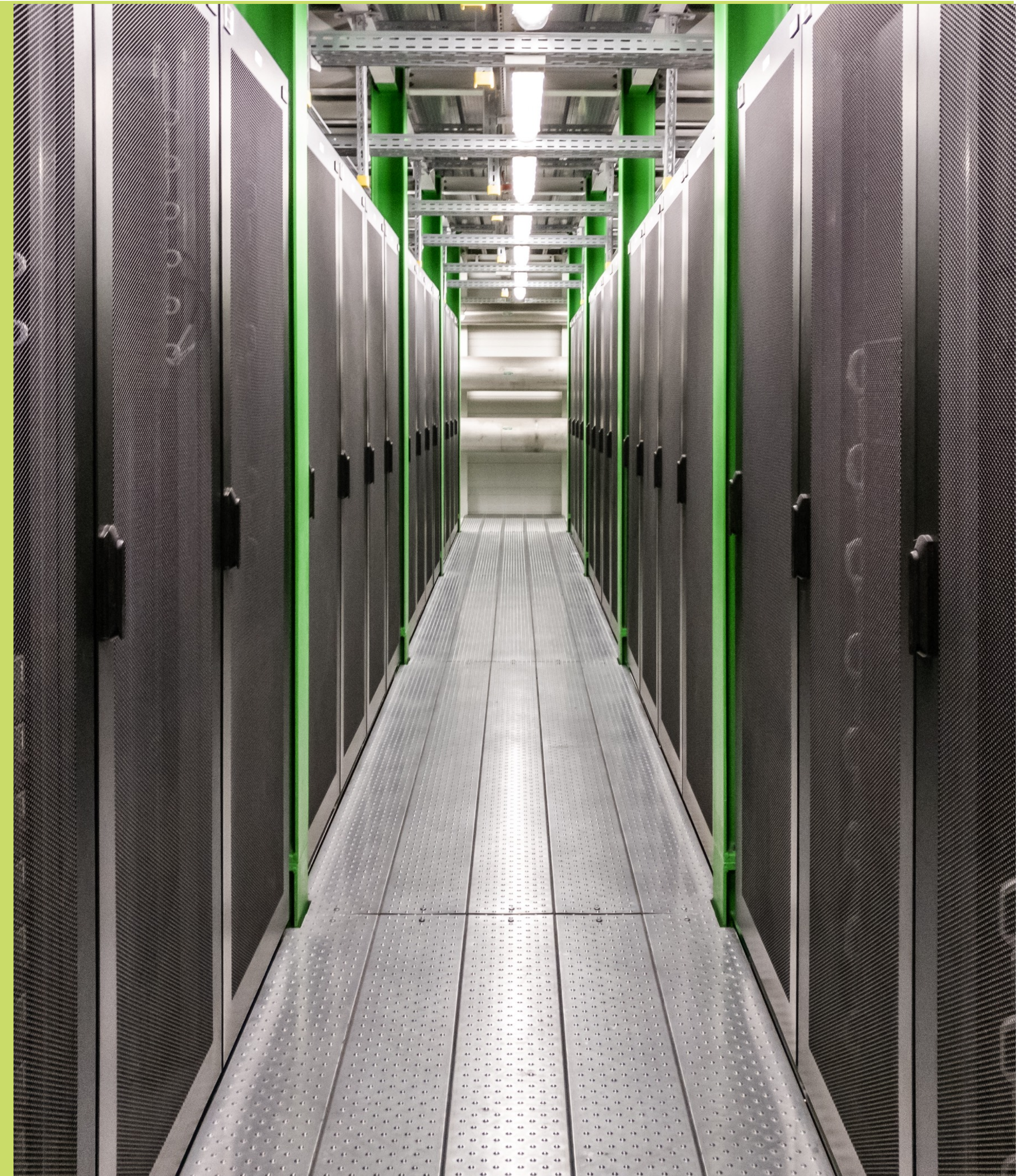
- Empfohlene Hardware: Mehrere Server mit Intel Xeon Gold oder AMD EPYC Prozessoren
- RAM: Pro Server mindestens 256 GB
- GPU: Mehrere NVIDIA A100 oder H100 GPUs mit jeweils mindestens 80 GB VRAM
- Speicher: Schnelle NVMe-SSDs und redundante HDDs/SSDs
- Netzwerk: Gigabit-Ethernet
- Geschätzte Kosten: Pro Server: Zwischen 15.000€ und 30.000€, abhängig von der Konfiguration. Gesamtkosten: Für ein Cluster von 3–5 Servern liegen die Investitionen zwischen 45.000€ und 150.000€





## 4. Enterprise-Lösung für 1.000 Mitarbeiter und mehr

- Empfohlene Hardware: High-End-Server mit Multi-Socket-Motherboards, ausgestattet mit Intel Xeon Platinum oder AMD EPYC Prozessoren
- RAM: Pro Server mindestens 1 TB
- GPU: Mehrere NVIDIA A100 oder H100 GPUs mit jeweils mindestens 80 GB VRAM
- Speicher: Enterprise-Grade NVMe-SSDs und umfangreiche, redundante Speicherlösungen
- Netzwerk: 40-Gigabit-Ethernet oder höher
- Geschätzte Kosten: Pro Server: Zwischen 50.000€ und 100.000€, abhängig von der Konfiguration. Gesamtkosten: Für ein Cluster von 3–5 Servern liegen die Investitionen zwischen 500.000€ und 2.000.000€





# Praxis-Beispiel: Content-Erstellung

**Aufgabe: Presstext für Ausstellung**

**Prompt-Struktur:**

- Wer, Was, Wann, Wo
- Zielgruppe
- Tonalität

**Ergebnis: Entwurf in 30-60 Sek**

# Was funktioniert NICHT gut lokal?

- ✗ Sehr aktuelle Informationen
- ✗ Multimodale Aufgaben (Bild+Text)
- ✗ Team-Kollaboration
- 💡 Hybrid-Ansatz für Spezialfälle



# Installation in 3 Schritten

## Schritt 1: Download

→ lmstudio.ai (~200 MB)

## Schritt 2: Modell laden

→ "Llama 3.2 3B" (~2 GB)

## Schritt 3: Loslegen

→ Chat → Prompt → Enter

# Ihre nächsten Schritte

- Heute Abend: LM Studio installieren
- Diese Woche: Erstes Modell testen
- Nächsten Monat: Workflow automatisieren

 Handout liegt aus!

# Lust auf selbstbestimmtes Handeln?

Heute Abend: Installation starten

## Kontakt:

Johannes Hofmann

info@ki-erlernen.jetzt

Hamburg

<https://ki-erlernen.jetzt>



Link zum Handout



Handout mitnehmen!



# Industrialisierung vs. KI-Revolution: Ähnlichkeiten und Unterschiede

- **Die Industrialisierung** revolutionierte die Produktion von Waren, während die generative KI **Revolution in der kreativen Industrie einsetzt**.
- **Automatisierung kreativer Prozesse:** Industrialisierung ermöglichte die Massenproduktion durch Maschinen. Generative KI-Tools und -Technologien ermöglichen die Automatisierung von kreativen Prozessen, wie z.B. der Erstellung von Inhalten.
- **Demokratisierung der Produktion:** Die Industrialisierung machte Produkte für eine breitere Bevölkerungsschicht zugänglich. Generative KI-Tools und -Technologien machen es auch Nichtexperten ermöglichen, kreative Inhalte zu erstellen, wird quasi demokratisiert.

# Lokale Software zur Nutzung der Modelle (nicht vollständig)

Programm	Mac OSX	PC Windows
LM Studio	Ja	Ja
GPT4All	Ja	Ja
<b>Open WebUI</b> (über Terminal, Docker) Besonders für Arbeitsgruppen geeignet Browser-basiert	Ja	Ja



# Lokale Software zur Nutzung der Modelle (nicht vollständig)

Programm	Mac OSX	PC Windows
<b>Diffusion Bee</b> , leicht installierbar über Installer	Ja	Nein
<b>NMKD</b> , leicht installierbar über Installer	Nein	Ja
<b>Pinokio</b> , leicht installierbar über Installer	Ja	Ja
<b>Automatic1111</b> , im Terminal installierbar	Ja	Ja
<b>Foocus</b> , im Terminal installierbar	Ja	Ja
<b>ComfyUI</b> , im Terminal installierbar	Ja	Ja