

## Nextcloud AI Erweiterungen – Konzepte gegen den GPU-Hunger

# Agenda

## 1. Einleitung

2. Schneller **Überblick** über Featureset

3. **GPU as a Service** in der Public und Private Cloud

4. Was ist der **Plan**?

5. **Fazit**

## Neben Bullshit-Bingo – Wovon reden wir?

- Nextcloud 27 bringt **neue Features** mit, die auf **AI Algorithmen** basieren
- Meine Aussage Nextcloud Advisory Board: „Für mich ist **wichtiger**, **welche Use-Cases** ich mit neuen Features abdecken kann, **als** auf **welcher Algorithmenklasse** sie basieren.“ (sinngemäß)
- Spannend ist jedoch die Frage: **Wie skaliert das?**
  - Auf meinem Mac mini in Ubuntu-Nachnutzung sieht man nichts.
  - Aber mit 4.000+ parallelen Sessions kann das anders aussehen
- Was machen wir also **architekturell**? Und **wie** können wir evtl. **mit** der **Last umgehen**, wenn sie denn signifikant wird?

# Ethical AI Rating

## Drei Bedingungen

1. Ist die Software Open Source (Nutzung und Training)
2. Gibt es ein freies Modell zum Selbst-Hosting?
3. Sind die Trainingsdaten frei und transparent?

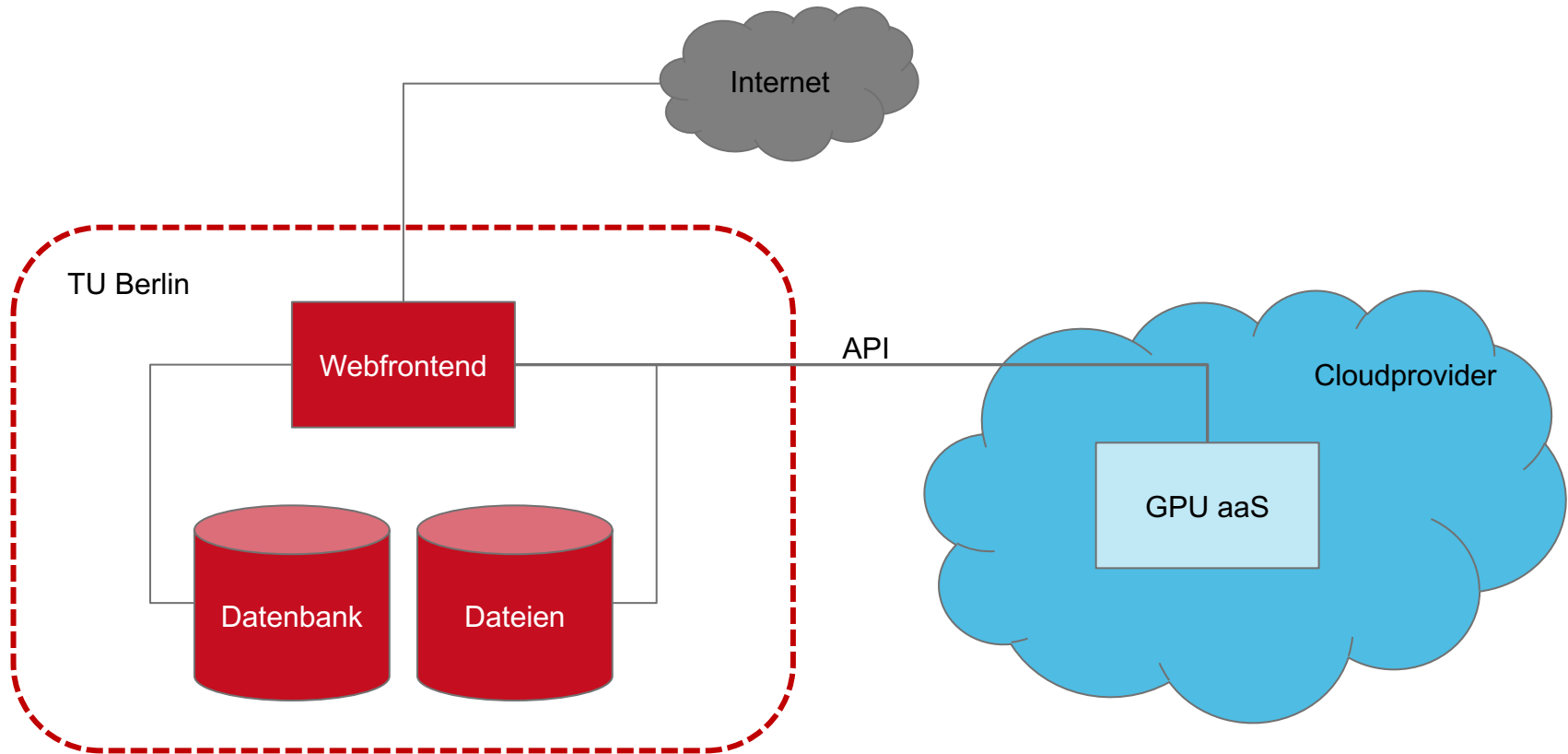
## Vier Wertungen

1. **Grün:** Alle Bedingungen erfüllt.
2. **Gelb:** Zwei Bedingungen erfüllt.
3. **Orange:** Eine Bedingung erfüllt.
4. **Rot:** Keine Bedingung erfüllt.

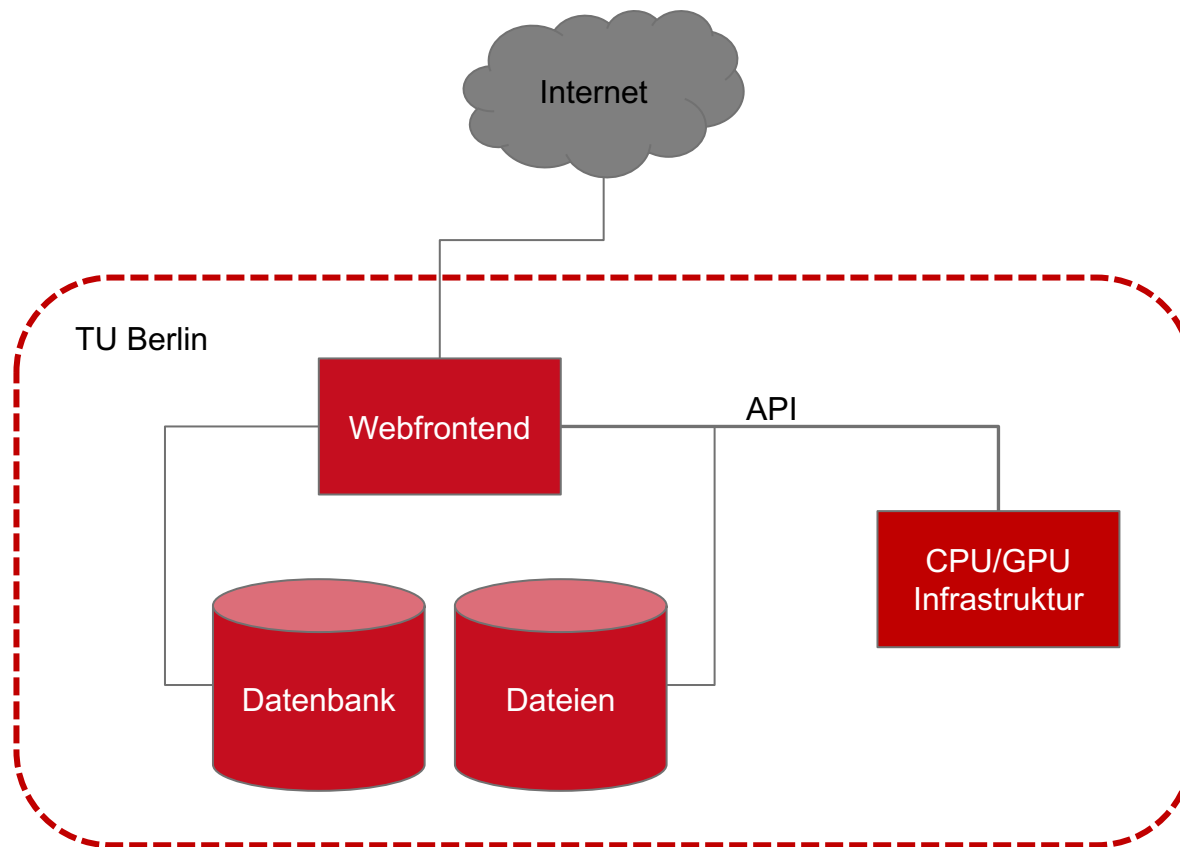
## AI Features beginnend mit Nextcloud 27

Smart inbox	<a href="#">Mail</a>	Green	Image generation	<a href="#">Local Stable Diffusion</a>	Yellow
Image object recognition	<a href="#">Recognize</a>	Green		<a href="#">OpenAI and LocalAI integration (via OpenAI API)</a>	Red
Image face recognition	<a href="#">Recognize</a>	Green		<a href="#">OpenAI and LocalAI integration (via LocalAI)</a>	Yellow
Video action recognition	<a href="#">Recognize</a>	Green		<a href="#">Replicate integration</a>	Yellow
Audio music genre recognition	<a href="#">Recognize</a>	Green	Text generation	<a href="#">Local large language model (via GPT4all Falcon)</a>	Green
Suspicious login detection	<a href="#">Suspicious Login</a>	Green		<a href="#">Local large language model (via Llama 2)</a>	Yellow
Related resources	<a href="#">Related Resources</a>	Green		<a href="#">OpenAI and LocalAI integration (via OpenAI API)</a>	Red
Recommended files	<a href="#">recommended_files</a>	Green		<a href="#">OpenAI and LocalAI integration (via LocalAI)</a>	Green
Machine translation	<a href="#">Translate</a>	Green	Context Chat	<a href="#">Nextcloud Assistant Context Chat</a>	Yellow
	<a href="#">LibreTranslate integration</a>	Green		<a href="#">Nextcloud Assistant Context Chat (Backend)</a>	Yellow
	<a href="#">DeepL integration</a>	Red			
	<a href="#">OpenAI and LocalAI integration (via OpenAI API)</a>	Red			
	<a href="#">OpenAI and LocalAI integration (via LocalAI)</a>	Green			
Speech-To-Text	<a href="#">Whisper Speech-To-Text</a>	Yellow			
	<a href="#">OpenAI and LocalAI integration</a>	Yellow			
	<a href="#">Replicate integration</a>	Yellow			

# GPU as a Service



# GPU as a Service



## Nächste Schritte

1. **Evaluation** im kleinen Kreis:
  - **Features: Use-Cases** an der TU Berlin
2. **PoC** mit Cloud-Provider
3. Tests mit **eigener Infrastruktur**
4. **Auswertung**: Metriken, Nutzungsprofile, Erkenntnisse
5. **Angebot** für DFN-Cloud, evtl. mit eigenem Preisschild



## Fazit

- Neue Nextcloud **AI Features sind** erst einmal nur neue **Features**
- Wir **evaluieren** zunächst, **welche Potentiale** darin stecken.
- Bei den **Tests** an der TU Berlin werden gleich **verschiedene Infrastrukturkonzepte** geprüft.
- **Am Ende berichten wir und** machen ein **DFN-Cloud Angebot** daraus.