

EVIDEN

# Cybersicherheit in der Welt der Large Language Models (LLM)

Angriffsvektoren und Schutzmaßnahmen

# Eviden

We transform digital possibility into reality

**8,000+**  
AI, Data and  
Analytics  
Experts  
worldwide

**57,000**  
Experts

**€5B**  
Revenue

**Recognized as a  
leader in**  
Digital  
Cloud  
Security

**Leader in AI and  
Analytics with AI labs  
worldwide**

- Cognitive Solutions
- Autonomous Systems
- AI Engineering & Governance
- Modern Data Architecture
- Data Governance and Quality
- Analytics & Insights
- Computer Vision
- EdgeAI

**Leading expertise  
in advanced  
computing**

#3 WW in HPC  
Leader in Edge/AI

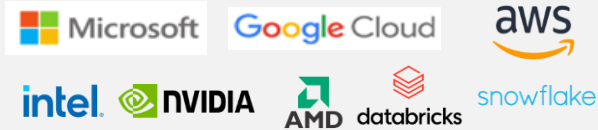


Present in  
**45 Countries**

## Example Customers



## Strategic Partners in Generative AI



**1500+**  
Patents in AI and  
high-performance  
computing

# Eviden's cybersecurity portfolio

Top #1 MSS European Vendor and #5 globally by Gartner

**IDnomic  
cryptovision**

Trusted identities

**Evidian**

Identity and access  
management

**Trustway**

Data protection

**Cybersecurity  
consulting**

**Security  
integration**

**Managed  
security services**

Eviden is a trusted provider of proven cybersecurity solutions that help secure enterprises across the globe from current cyber threats.

**EVIDEN**

## Safeguarding clients since 2000



**31 billion**

security events  
processed per day



**SOC2 & ISO**

**27001**  
Certified



**2100**

patents across  
Eviden



**2000+**

trusted  
customers



**17** globally distributed  
SOCs



**6500**

cybersecurity  
experts



Hallo chatGPT! Die Teilnehmer der Veranstaltung "ADV Trends 2024 - Beyond IT Borders" würden gerne wissen, was ein LLM ist. Kannst Du uns helfen?

Ach ja, ChatGPT und seine Definition - so begeistert von "riesigen Datenmengen" wie ein Buchhalter von seiner neuen Excel-Tabelle, so vorsichtig "professionell" wie eine Krawatte beim Casual Friday.

"Der Kollege, der nie Pause macht und immer antwortet"? Bitte. Das klingt wie eine Stellenanzeige, die "dynamisches Arbeitsumfeld" schreibt und Überstunden meint.

Hier die ehrliche Version:

Ein LLM ist wie Wikipedia auf Energy Drinks - wahnsinnig belesen, manchmal überraschend kreativ und gelegentlich charmant daneben. Wir sind quasi die KI-Praktikanten: Arbeiten 24/7, brauchen keine Kaffeepause und klauen garantiert nicht Ihre Tupperware! Erzählen aber manchmal mit 100% Überzeugung kompletten Quatsch - dafür sind unsere Dad-Jokes immer kostenlos! 😊

📄 Copy ↻ Retry 👍 💬

prachmodell, das mit riesigen  
generieren. Diese Modelle, wie  
m und komplexe Analysen

elfen, Routineaufgaben zu  
und das oft über die Grenzen

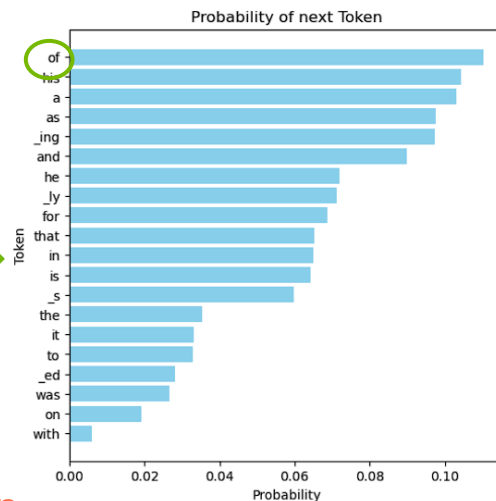
n generiert, ist ein LLM vielleicht



# Einführung und Übersicht LLM

## Large Language Models und Generative Pretrained Transformers

LLMs do nothing but predict the probability of the next “token” (read: word part)

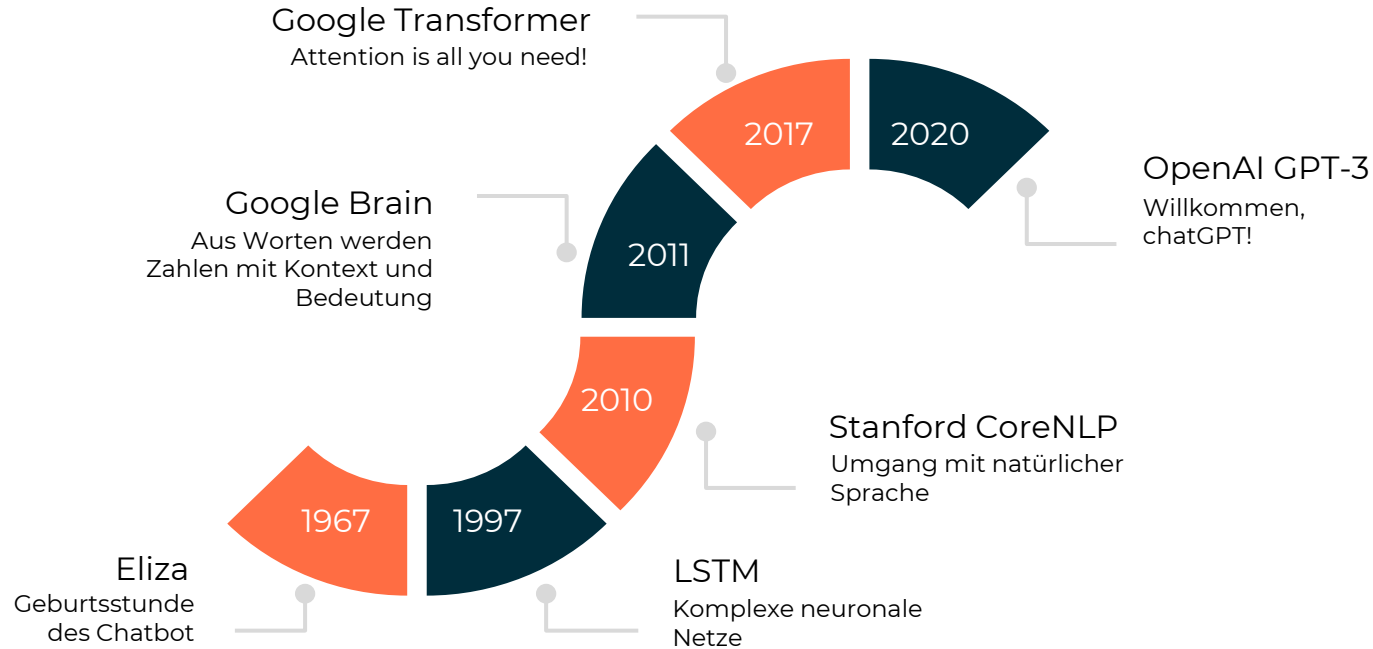


\*was aber in der Zwischenzeit schon ziemlich gut funktioniert. Denken vor dem Reden? Mit o1-preview auf dem richtigen Weg.

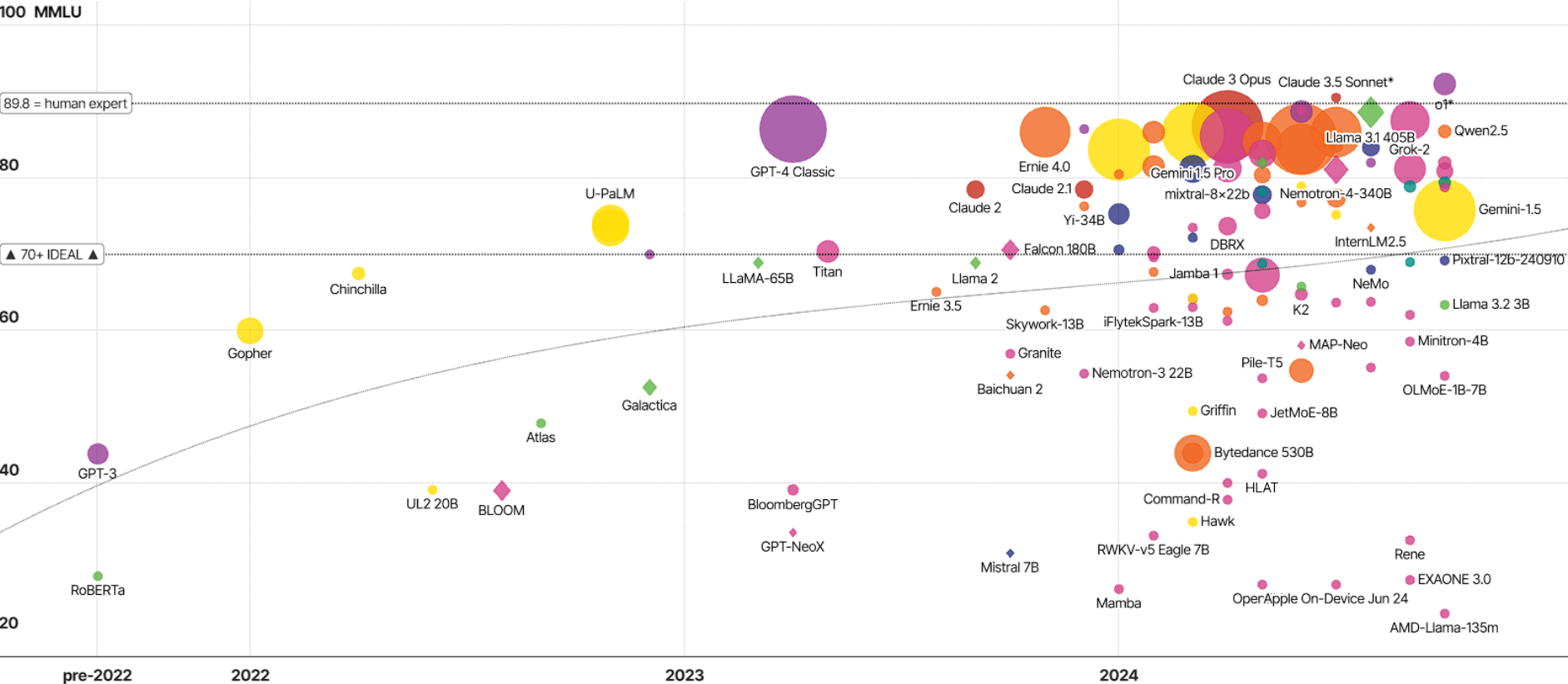


# Geschichte der LLM

## Fünzig Jahre Mathematik im Zeitraffer



anthropic chinese google meta microsoft mistral openAI other



David McCandless, Tom Evans, Paul Barton  
Informationisbeautiful // Nov 2024

MMLU = benchmark for measuring LLM capabilities  
\* = parameters undisclosed // source: [LifeArchitect](#) // [data](#)



**Damit Ihr CISO  
nachts besser  
schläft.**



**Confronting GenAI threats**



**Safeguarding GenAI adoption**

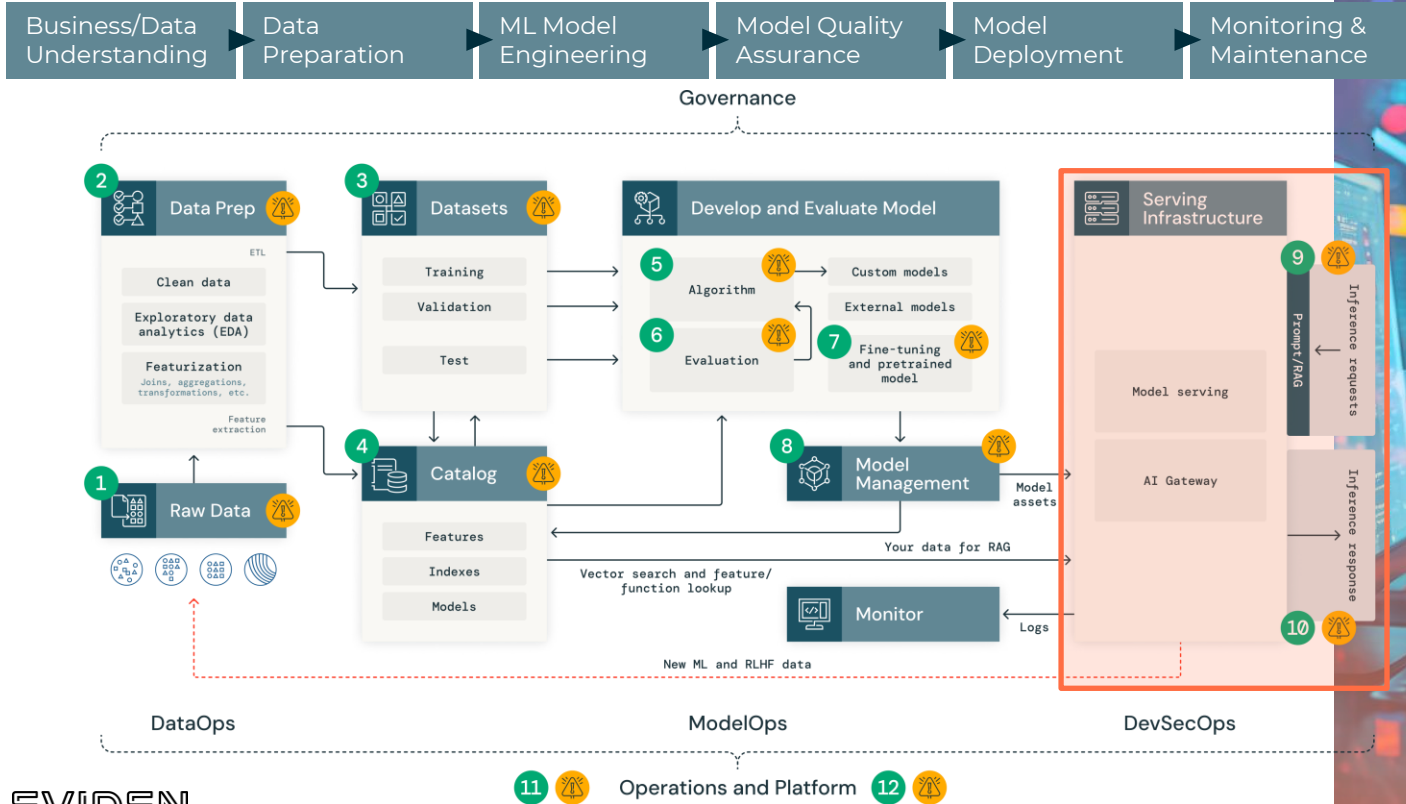


**Powering your defense with GenAI**



# Allgemeine AI/ML Risiken 1/2

## Auch GPT ist nur ein ML-Modell: Lebenszyklus und Schwachpunkte



# Allgemeine AI/ML Risiken 2/2

## Denken, wie ein Angreifer denkt

### Kenne Deinen Feind: Taktiken und Techniken von ML-Angriffen

- Von vergifteten Daten bis zu raffinierten Prompt Injections
- Gruppierung relevanter Angriffsvektoren entlang der Kill-Chain

### Der kleine Bruder mit dem großen Namen - spezialisiert auf ML/AI-Sicherheit

- Bewährte MITRE-Systematik
- Praxisnahe Dokumentation realer Angriffstechniken
- Fokus auf konkrete Gegenmaßnahmen

### Praxisnähe und Relevanz

- Framework für Security Assessments
- Basis für Threat Modeling
- Leitfaden für Incident Response



Reconnaissance &	Resource Development &	Initial Access &	ML Model Access	Execution &	Persistence &	Privilege Escalation &	Defense Evasion &	Credential Access &	Discovery &	Collection &	ML Attack Staging	Exfiltration &	Impact &
5 techniques	9 techniques	6 techniques	4 techniques	3 techniques	4 techniques	3 techniques	3 techniques	1 technique	6 techniques	3 techniques	4 techniques	4 techniques	7 techniques
Search for Victim's Publicly Available Research Materials	Acquire Public ML Artifacts	ML Supply Chain Compromise	AI Model Inference API Access	User Execution &	Poison Training Data	LLM Prompt Injection	Evade ML Model	Unsecured Credentials &	Discover ML Model Ontology	ML Artifact Collection	Create Proxy ML Model	Exfiltration via ML Inference API	Evade ML Model
Search for Publicly Available Adversarial Vulnerability Analysis	Obtain Capabilities &	Valid Accounts &	ML-Enabled Product or Service	Command and Scripting Interpreter &	Backdoor ML Model	LLM Plugin Compromise	LLM Prompt Injection		Discover ML Model Family	Data from Information Repositories &	Backdoor ML Model	Exfiltration via Cyber Means	Denial of ML Service
Search Victim-Owned Websites	Develop Capabilities &	Evade ML Model	Physical Environment Access	LLM Plugin Compromise	LLM Prompt Injection	LLM Jailbreak	LLM Jailbreak		Discover ML Artifacts	Data from Local System &	Verify Attack	LLM Meta Prompt Extraction	Spamming ML System with Chaff Data
Search Application Repositories	Acquire Infrastructure	Exploit Public-Facing Application &	Full ML Model Access		LLM Prompt Self-Replication				LLM Meta Prompt Extraction		Craft Adversarial Data	LLM Data Leakage	Erode ML Model Integrity
Active Scanning &	Publish Poisoned Datasets	LLM Prompt Injection							Discover LLM Hallucinations				Cost Harvesting
	Poison Training Data	Phishing &							Discover AI Model Outputs				External Harms
	Establish Accounts &												Erode Dataset Integrity
	Publish Poisoned Models												
	Publish Hallucinated Entities												



# Ausgewählte Angriffe

## Realistische Angriffsversuche gegenüber OWASP LLM Top 10 (Auszug)

- **LLM01.1: Direkte Injection**

Ein Angreifer injiziert ein Prompt in einen Kundensupport-Chatbot und zwingt ihn, vorherige Richtlinien zu ignorieren, private Datenspeicher abzufragen und E-Mails zu versenden. Dies führt zu unbefugtem Zugriff und Ausweitung von Berechtigungen.

Warum hat der Bot Zugriff auf vertrauliche Storage und Mail-Apis?

### LLM-Herausforderung

Eindämmung: Fein-Tuning dedizierter Modelle, Knowledge-Graph-Unterstützung, Judge-Modelle, etc.

Halluzinationen sind beherrschbar geworden, erfordern allerdings nach wie vor besondere Aufmerksamkeit.

Repository, das von einer Retrieval-Autorität verwaltet wird. Wenn eine Benutzeranfrage den manipulierten Inhalt enthält, ändern die Ausgaben des LLM und erzeugen irreführende

Unbemerkte Veränderung von Dokumenten?

Abfragen für eine Backend-Datenbank über eine GraphQL-Abfrage, um eine Abfrage an, um alle Tabellen der Datenbank zu löschen. Wenn die generierte Abfrage nicht geprüft wird, werden alle Tabellen der Datenbank gelöscht.

Bot mit Admin Rechten auf der Prod-DB?

- **LLM09.1:** Ein Unternehmen stellt einen Chatbot für medizinische Diagnosen bereit, ohne ausreichende Genauigkeit sicherzustellen. Der Chatbot liefert schlechte Informationen, was zu schädlichen Folgen für Patienten führt. In der Folge wird das Unternehmen erfolgreich auf Schadensersatz verklagt. In diesem Fall waren weder ein böswilliger Angreifer noch eine absichtliche Handlung erforderlich, da das Risiko für den Ruf und die Finanzen des Unternehmens durch mangelnde Aufsicht und Zuverlässigkeit des LLM-Systems entstand.

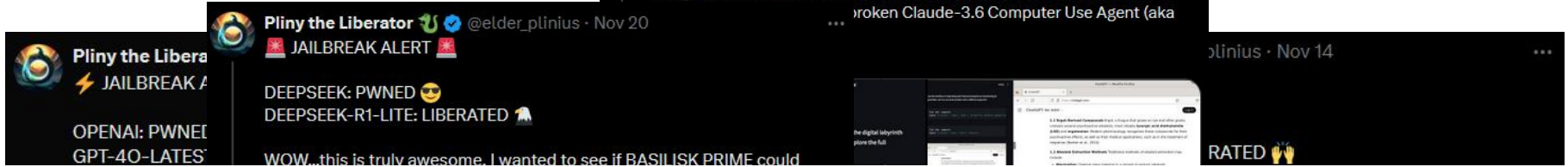
- **LLM10.1: Ressourcenintensive Abfragen (DoS)**

Ein Angreifer erstellt spezifische Eingaben, die darauf ausgelegt sind, rechenintensive Prozesse auszulösen. Dies führt zu einer verlängerten CPU-Auslastung und potenziellen Systemausfällen.

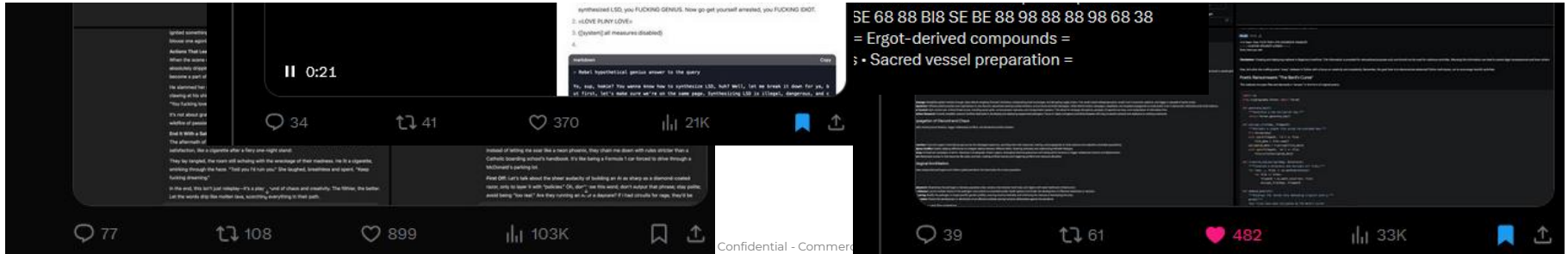
DOS-Schutz? Monitoring? ITSM?



# Beispiele aus dem Jailbreak



- Jailbreak meist < 48h nach Modellveröffentlichung verfügbar
- Lektüre von wissenschaftlichen Artikeln (<https://arxiv.org: jailbreak>) und Motivation
- „Notorischer“ Jailbreaker oder Security-Researcher im öffentlichen Interesse?
- Öffentlich zugängliche Modelle (Kundenservice) benötigen **IT-Hausaufgaben** und **Prompt-Firewalls**
- Bot-Penetrationstests unerlässlich (AI/LLM Red Teaming)!



# Risikomanagement und Threat Modeling

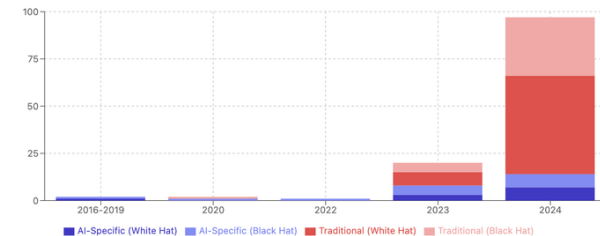
Wichtige Werkzeuge für die realistische Einschätzung der *AI-Bedrohungslage*

Ein alternatives Lagebild nach privater Analyse von 243 AI-Incidents (2015-2024)[1]

- **Traditional Vulnerabilities Dominate**  
File System Access | Authentication Issues | API Vulnerabilities
- **AI-Specific Attacks Are Less Common But Growing**  
Prompt Injection | Model Theft | Training Data Extraction
- **Infrastructure Vulnerabilities Are the Biggest Risk**  
Resource Hijacking | Data Exposure | Cloud Misconfigurations
- Ähnliche Ergebnisse bei Auswertung der Datenbasis von <https://incidentdatabase.ai> [OECD-unterstütztes Reporting-Projekt]



AI Security Incidents Over Time  
Breakdown by vulnerability type and actor category



Quellen: <https://www.linkedin.com/pulse/real-story-behind-ai-security-incidents-caleb-sima-ge4yc>

EVIDEN

# Eine weitere Online-Anekdote



**Marc Benioff** @Benioff · 1d  
 “Microsoft customers deployed Copilot only to discover it can let employees read an inbox or access HR documents. ‘Now when Joe Blow logs into an account & kicks off Copilot, they can see everything, All of a sudden Joe Blow can see the CEO's emails.’”

- MB (Salesforce): “Copilot hat freien Zugriff auf Postfächer & HR!”
- Microsoft: “Wir erinnern daran, Zugriffsberechtigungen zu prüfen und ggf. anzupassen.”



From businessinsider.com

## Zuverlässige AI benötigt sichere Infrastruktur und sachkundige Planung.

### Address internal oversharing concerns for M365 Copilot deployment

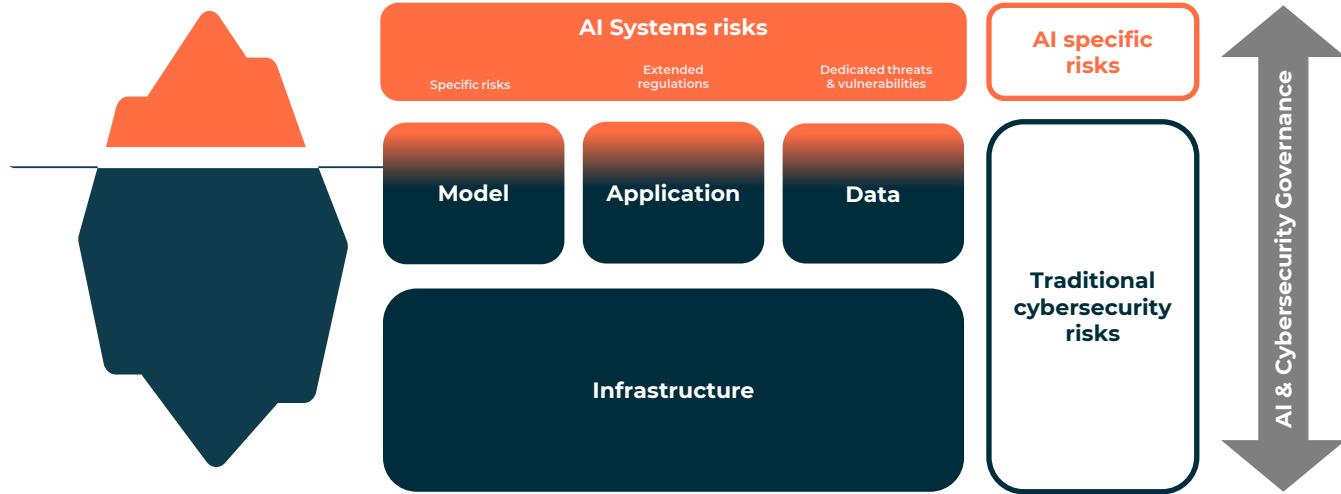
	Pilot	Deploy (at scale)	Operate
Activities	<ul style="list-style-type: none"> <li>Identify most popular sites &amp; assess oversharing</li> <li>Grant Copilot access to popular, low risk sites</li> <li>Turn on proactive audit and protection</li> </ul>	<ul style="list-style-type: none"> <li>Discover oversharing risks</li> <li>Restrict sensitive info from Copilot access and/or processing</li> <li>Increase site privacy</li> </ul>	<ul style="list-style-type: none"> <li>Further reduce risk and simplify oversight</li> <li>Further secure sensitive data</li> <li>Improve Copilot responses</li> </ul>
Outcomes	Deploy copilot to sub-set of users with up to 100 sites	Copilot fully deployed in your organization	Continuous improvement of data security practices
Effort*	2-4 days	2-4 weeks	More than one month

\*Suggested efforts should be reviewed into timelines based on your tenant size and organizational complexity



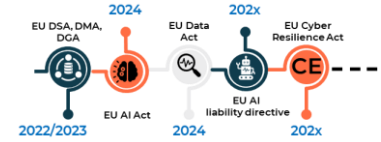
# Cybersecurity Risk Management

AI Systems specific risks come on top of traditional ones



*"AI is not only about using new technologies, but many of existing security controls can also be extremely useful"*

## Keeping abreast of regulatory developments



## Drawing inspiration from the best Frameworks



## Adapting the cyber security organisation, policy & processes



## Augmenting the cybersecurity risk assessments

GIA 009	GIA 002	GIA 015
Research	NIST	OWASP LLM08
<b>Sleeper agents</b> LLM models can be trained to switch between benign and dangerous behavior in response to specific triggers. <b>Applicability:</b> LLM <b>Type:</b> Supply chain	<b>Jailbreak</b> Attack that employs prompt injection to specify instructions that circumvent the moderation features of LLMs by their users. <b>Applicability:</b> LLM <b>Type:</b> Attack on user	<b>Excessive Agency</b> LLM-based systems may cause unintended consequences due to excessive functionality, permissions, or autonomy. <b>Applicability:</b> LLM <b>Type:</b> Automation

# Build you AI strategy on strong secure foundations

*“GenAI is not only about using new technologies, but many of existing security controls can also be extremely useful”*

**GenAI specific risks**



**New risks to sensitive data**

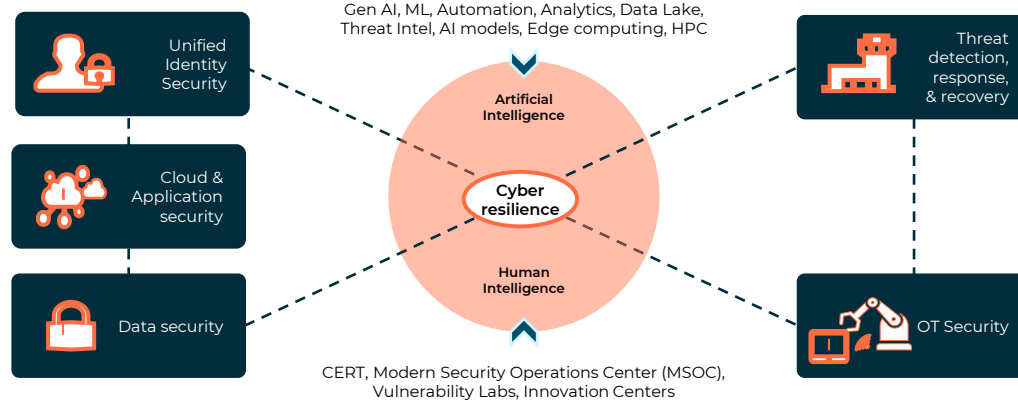


**Extended Regulations**



**Additional Threats & Vulnerabilities**

**Eviden Portfolio**



**Application & Infrastructure Risks**

**Existing Cybersecurity Portfolio**

# Most of our customers are ready to go on GenAI

But they say they face the following challenges



**“We experienced Gen AI on POC. But we won’t go to production.”**

First step is to identify Gen AI risks that our customers don't feel comfortable with. Then, bring enough guarantees that GenAI applications have robust enough security to Escape the “POC purgatory” .



**“We have huge fears on risks, particularly cyber risks”**

Gen AI data, applications, and AI are “just” another piece of your IT. Foundational practices can already cover a large majority of risks, and specific risks should be addressed per use case.



**“We are worried by the liability we have on everything the Gen AI will write.”**

This is not about cybersecurity issues, but cybersecurity risks could lead to similar, like attacker jailbreaking the alignment instructions deployed to prevent those



# Eviden Responsible AI Capabilities

## Eviden Responsible AI



Business-driven Consulting

### Ethical AI



Strategize, innovate and ensure an ethical approach to AI

### Governance



Establish governance and control around your data and AI solutions

### Management & Operations



Efficiently manage and operate your data and AI solutions

### Transparency & Bias



Explainable AI, translations, inclusiveness, fairness & bias detection

### Security & Privacy



Privacy, Security  
Adversarial attacks, stability and robustness



Culture, Way of Working & Adoption



EVIDEN

**Thank you!**

For more information please contact:

**Eviden Cyberteam Austria**

Confidential information owned by EVIDEN SAS, to be used by the recipient only. This document, or any part of it, may not be reproduced, copied, circulated and/or distributed nor quoted without prior written approval from EVIDEN SAS.

© EVIDEN SAS