



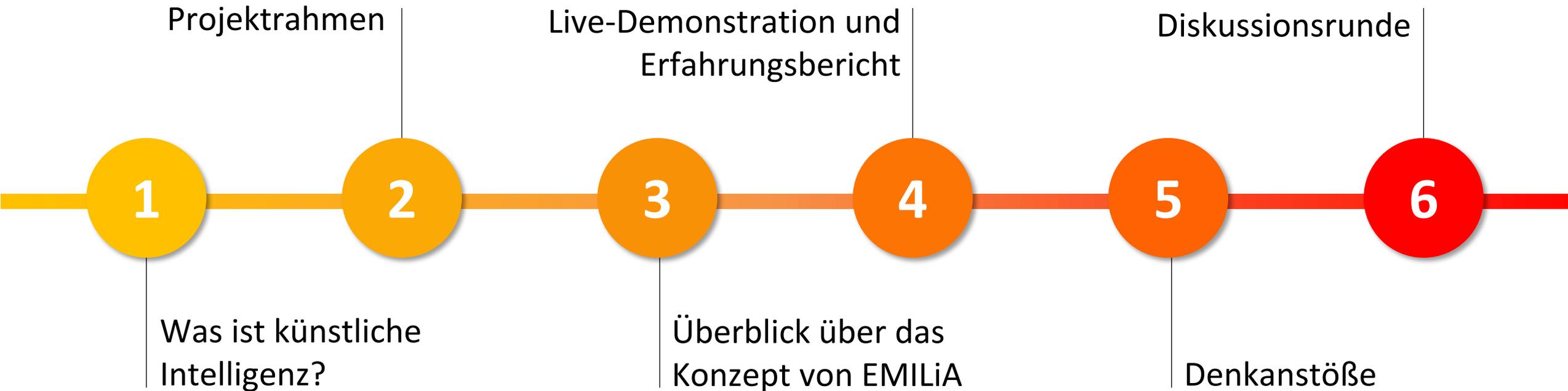
# Wie KI bei der Aufbereitung und Bereitstellung großer Datenmengen helfen kann

Nico Beyer & Felix Gericke

# Ablauf und Organisatorisches

Agenda,

# Was wir heute vorhaben



# Was wir heute nicht vorhaben

1

Computer Vision

2

Training von eigenen KI-Modellen

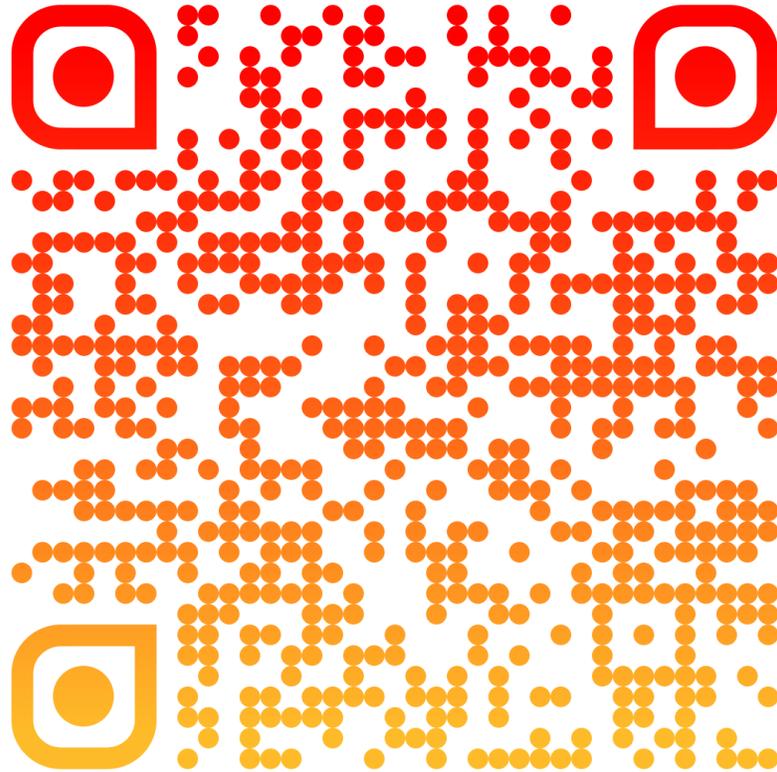
3

Verarbeitung von analogem und digitalisiertem Archivgut

4

Verwendung von Cloud-KI

# Pinnwand für gemeinsames Brainstorming



[url.emilia-archiv.de/pinnwand](http://url.emilia-archiv.de/pinnwand)

# Was ist künstliche Intelligenz?



*„[...] every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it.“*

***McCarthy, J., Minsky, M. L., Rochester, N., & Shannon, C. E. (1955): A Proposal for the Dartmouth Summer Research Project on Artificial Intelligence.***

# Was ist künstliche Intelligenz?



*„Artificial intelligence is the science of making machines do things that would require intelligence if done by men.“*

*Minsky, M.: Computation. Finite and Infinite Machines.*

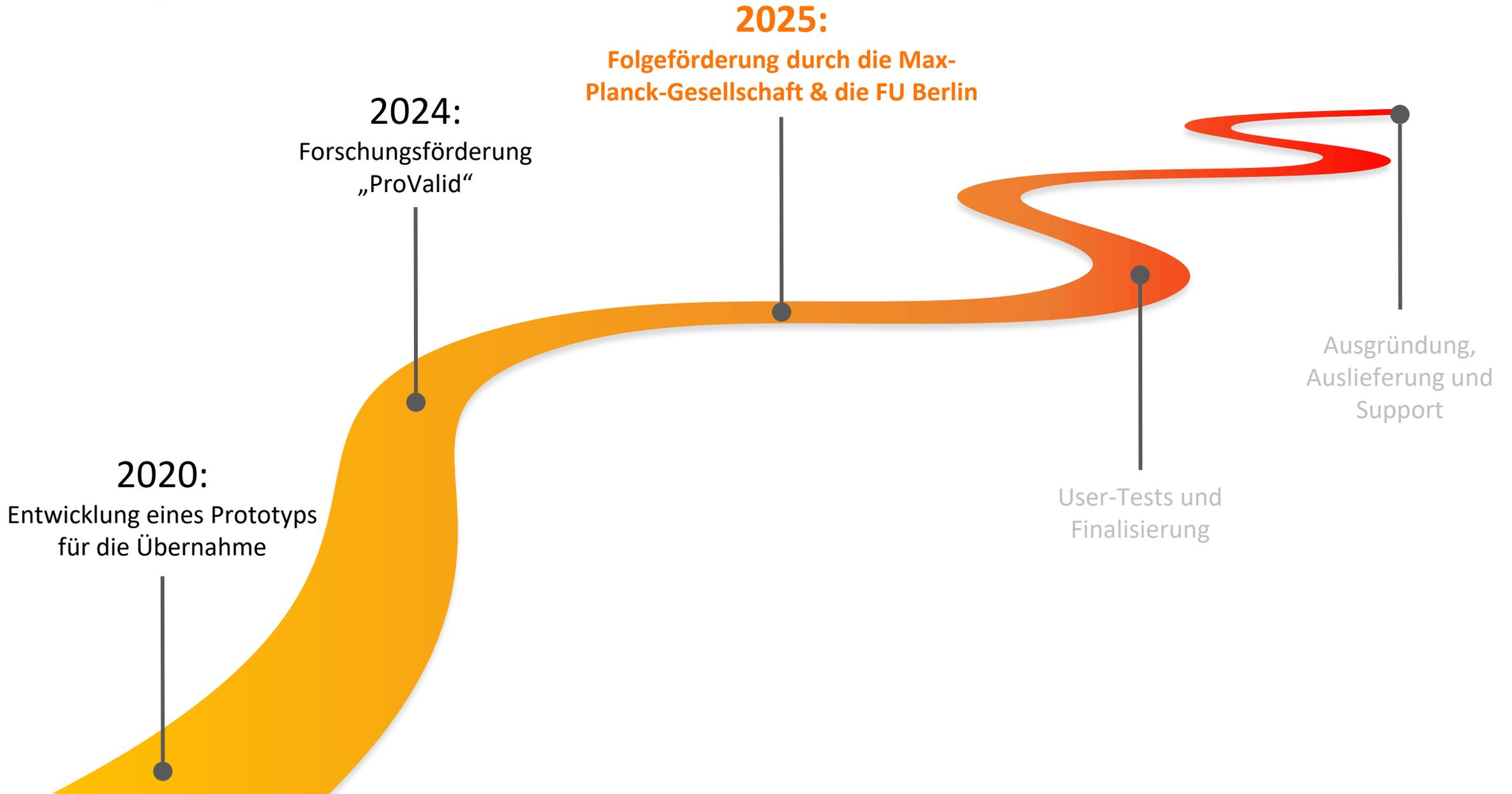
# Was ist künstliche Intelligenz?



*„Künstliche Intelligenz (KI) ist ein Teilgebiet der Informatik. Sie imitiert menschliche kognitive Fähigkeiten, indem sie Informationen aus Eingabedaten erkennt und sortiert.“*

***Fraunhofer-Institut für Kognitive Systeme: Was ist künstliche Intelligenz? Was ist maschinelles Lernen?***

# Projektrahmen



# Ausgangssituation

E-Mails sind für die meisten Menschen zu einem **festen Bestandteil des beruflichen und privaten Alltags** geworden.

E-Mails sind ein zentrales Kommunikationsmedium

Ein großer Teil der meisten Postfächer besteht aus **Spam, Werbung** oder nur **kurzfristig relevanten Informationen**

Viele wertlose Informationen

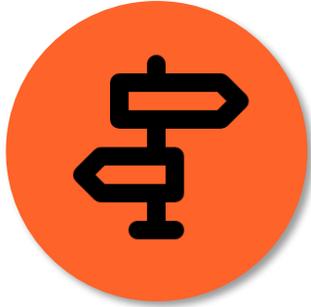
Es gibt jedoch auch **historisch oder rechtlich relevante E-Mails**, die langfristig bewahrt werden sollten

Es muss selektiert werden

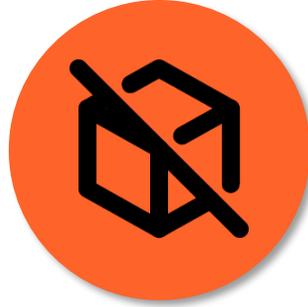
Eine **fachgerechte Auswahl, Archivierung und Auswertung** ist nur mithilfe automatisierter Prozesse möglich

Automatisierung als Chance

# Herausforderungen



E-Mail-Standard macht nur wenige klare Vorgaben



Gängige E-Mail-Container sind nicht gut für die Archivierung geeignet



Anhänge in den unterschiedlichsten Formaten



Signierte und verschlüsselte E-Mails

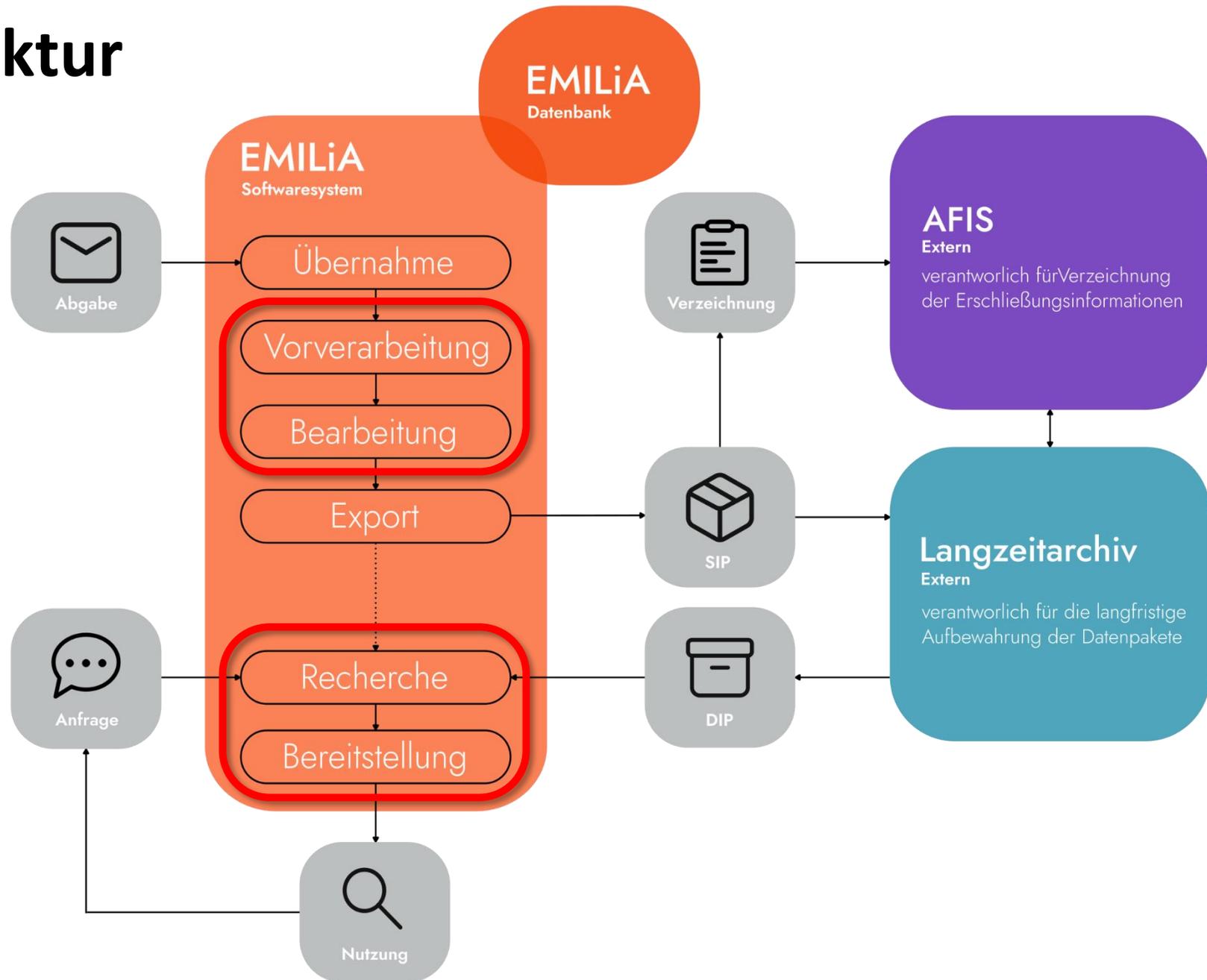


Personenbezogene Daten und urheberrechtlich relevante Dokumente



Archivwürdiger Kern ist nur schwer auffindbar

# Architektur



# Entwicklungsstand

➤ *Mit der aktuellen Version von EMILiA kann ein **vollständiger Archivierungsprozess** durchlaufen werden.*

**Übernahme:** Sicherer Transfer, Konvertierung, Virenprüfung, Formaterkennung, Integritätssicherung

**Bewertung & Erschließung:** Spracherkennung, Spam- und Dublettenerkennung, Themenerkennung, Identifikation personenbezogener Daten, Kassation, Merkliste

**Recherche:** Viewer, Optionale Anonymisierung, Such- und Filterfunktionen, Statistiken und Netzwerke, Zusammenfassung von Threads mittels LLM,

**Export:** MBOX oder strukturierter BagIt-Container mit Metadaten

# Entwicklungsstand

➤ *Mit der aktuellen Version von EMILiA kann ein **vollständiger Archivierungsprozess** durchlaufen werden.*

**Übernahme:** Sicherer Transfer, Konvertierung, Virenprüfung, Formaterkennung, Integritätssicherung

**Bewertung & Erschließung:** Spracherkennung, Spam- und Dublettenerkennung, **Themenerkennung, Identifikation personenbezogener Daten**, Kassation, Merkliste

**Recherche:** Viewer, **Optionale Anonymisierung**, Such- und Filterfunktionen, Statistiken und Netzwerke, **Zusammenfassung von Threads mittels LLM**

**Export:** MBOX oder strukturierter BagIt-Container mit Metadaten

# Live-Demonstration

Vorverarbeitung, Statistik-Dashboard, Viewer, Kassations- und Merkliste

# Experimente

Vergleich verschiedener Modelle, Topic Modelling, Retrieval Augmented Generation (RAG)

# Einfluss der Parameter auf die Performance von LLMs

	Bezeichnung	Parameter	Laufzeit	Laufzeit (E-Mail)
1	Gemma 3	4 Mrd.	00:03:46	~ 1.9 s
2	Gemma 3	12 Mrd.	00:08:17	~ 4.2 s
3	Gemma 3	27 Mrd.	00:13:45	~ 6.9 s

# Einfluss der Parameter auf die Performance von LLMs

## Parameter

## Zusammenfassung

1

4 Mrd.

Eine E-Mail, die vorgibt, von FedEx zu sein, meldet eine fehlgeschlagene Paketzustellung und fordert eine Adressbestätigung an über einen verdächtigen Link. Dies ist höchstwahrscheinlich eine Phishing-Nachricht.

2

12 Mrd.

Diese Nachricht behauptet, von FedEx zu sein und informiert über ein nicht zugestelltes Paket, das eine Adressbestätigung benötigt. Sie enthält Links zu unbekanntem URLs, was auf einen potenziellen Phishing-Versuch hindeutet. Vorsicht ist geboten.

3

27 Mrd.

Diese Nachricht behauptet, eine Benachrichtigung von FedEx bezüglich einer fehlgeschlagenen Paketzustellung zu sein und fordert zur Adressbestätigung auf. Enthält verdächtige Links und eine physische Adresse für die Abmeldung.

# Einfluss der Parameter auf die Performance von LLMs

## Parameter

## Themenextraktion

1

4 Mrd.

['Paketzustellung', 'Betrug', 'E-Commerce']

2

12 Mrd.

['Paketzustellung', 'FedEx', 'Adressbestätigung']

3

27 Mrd.

['Paketzustellung', 'Versand', 'Betrug', 'Phishing']

# Retrieval Augmented Generation



*Retrieval Augmented Generation (RAG) ist eine Technologie, die es ermöglicht, Large Language Models mit externen Wissensquellen zu kombinieren.*

# Retrieval Augmented Generation

Frage:

Antwort:

Ja, einige Namen kommen besonders häufig vor, zum Beispiel █████, █████,  
█████, █████, █████, █████ und █████.

# Retrieval Augmented Generation

Frage:

Antwort:

Ich weiß es nicht.

# Retrieval Augmented Generation

Frage: Bezieht sich der Bestand eher auf wissenschaftliche oder organisatorische Themen? Absenden

Antwort:

Der Bestand bezieht sich sowohl auf wissenschaftliche als auch auf organisatorische Themen. Es gibt Einladungen zu wissenschaftlichen Vorträgen und Diskussionen ( [REDACTED] ) sowie Ankündigungen von Veranstaltungen und organisatorischen Details ( [REDACTED] ).

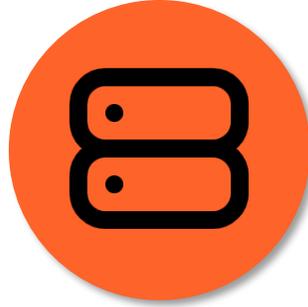
# Denkanstöße

Zeitaufwand, Nachhaltigkeit, Datengrundlage, Hardwareanforderungen und Kosten, Datenschutz, Qualität der Ergebnisse, Vibe-Coding

# Fallstricke



KI-Berechnungen benötigen Zeit



Die Anschaffung geeigneter Hardware  
ist mit Kosten verbunden



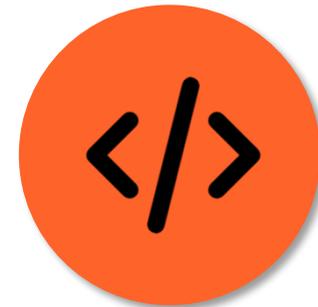
KI-Anwendungen verbrauchen sehr  
viel Energie



Qualität der Ergebnisse kann stark  
variieren

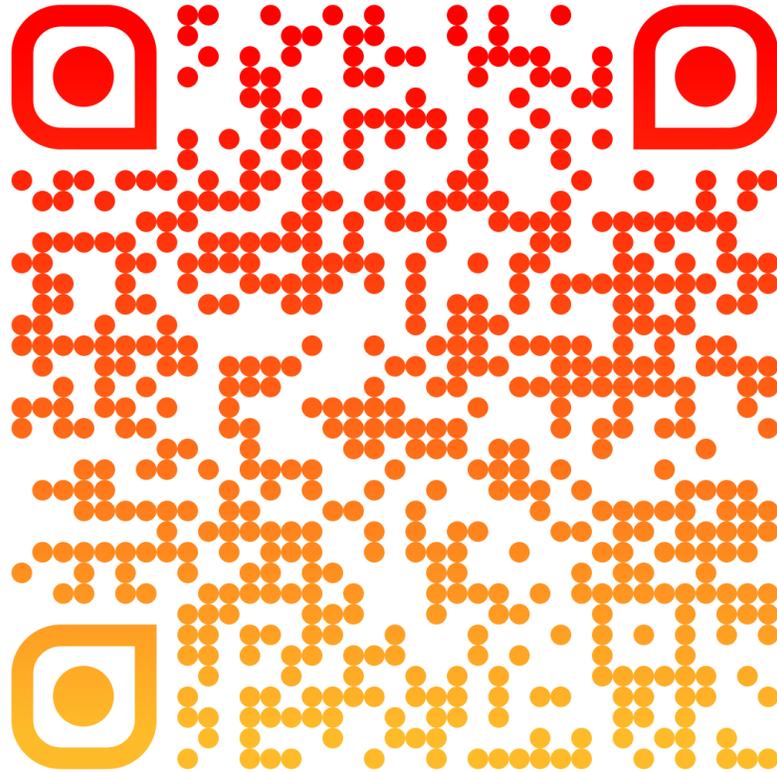


Datenschutzrecht muss beachtet  
werden



„Vibe Coding“  
ist kein Ersatz für geschultes  
Personal

# Pinnwand für gemeinsames Brainstorming



[url.emilia-archiv.de/pinnwand](http://url.emilia-archiv.de/pinnwand)

# Vielen Dank für Ihre Aufmerksamkeit!

E-Mail: [info@emilia-archiv.de](mailto:info@emilia-archiv.de)

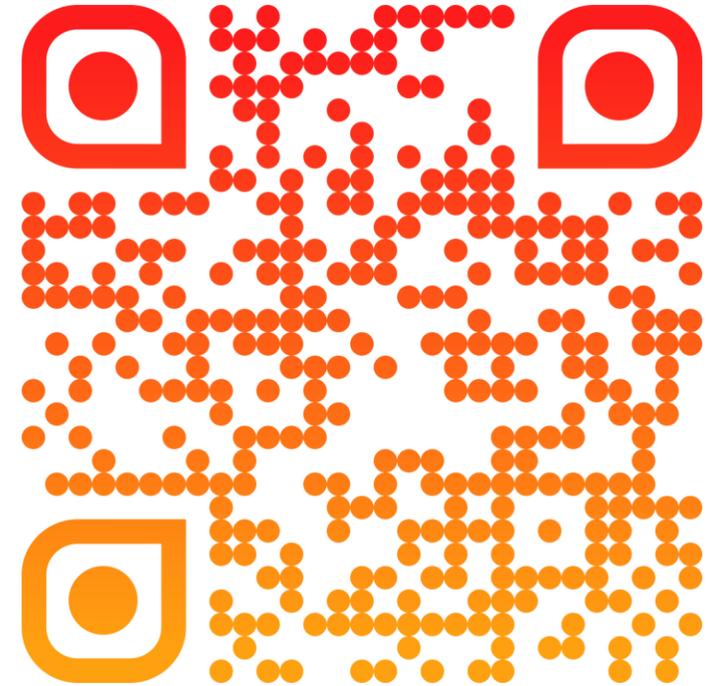
Telefon: +49 30 841 337 15

Archiv der Max-Planck-Gesellschaft

EMILiA-Projekt

Boltzmannstraße 14

14195 Berlin-Dahlem



[www.emilia-archiv.de](http://www.emilia-archiv.de)