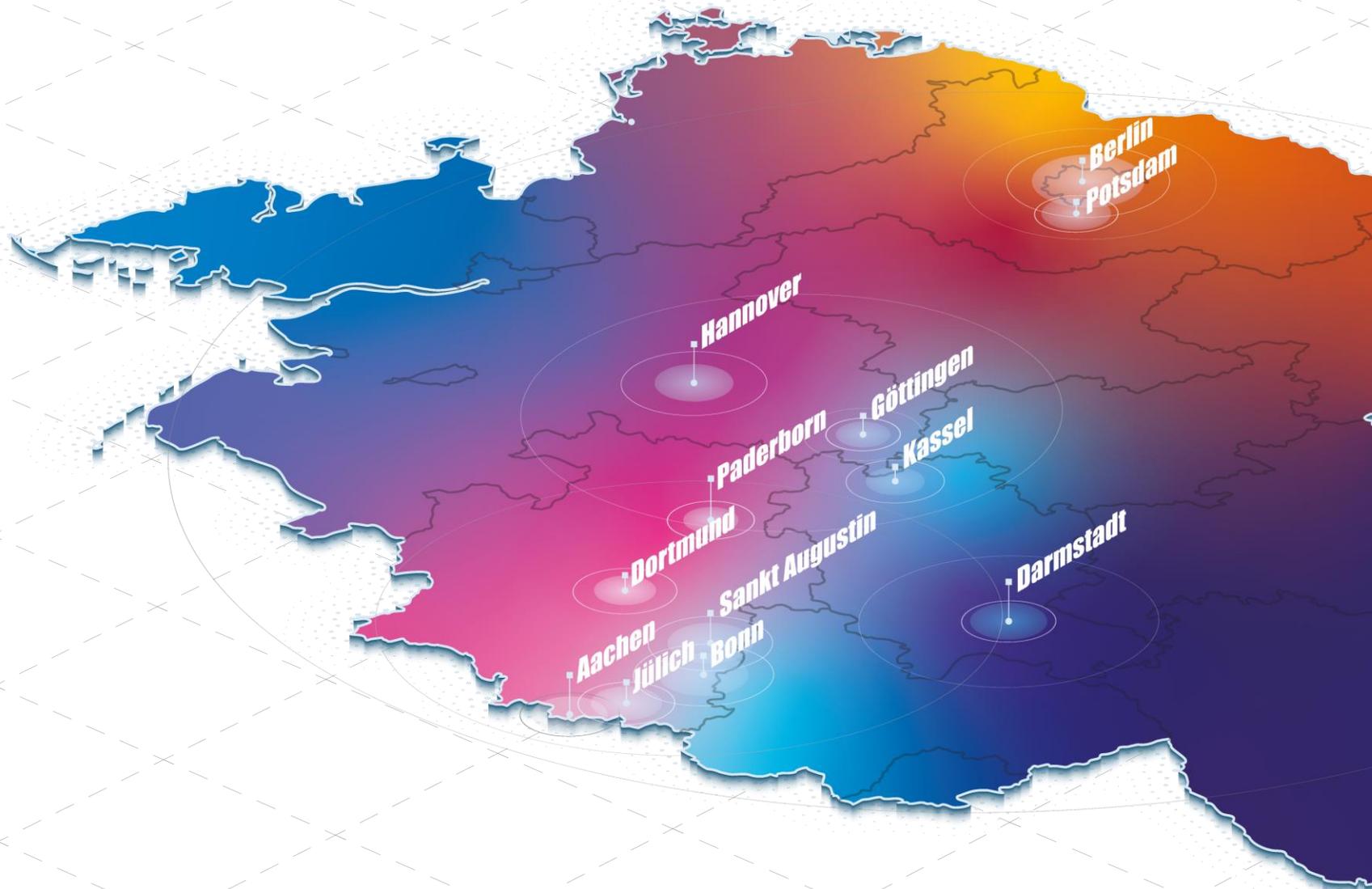


KI-Info-Café



KI-Servicezentren

KI-Servicezentren

Die KI-Servicezentren bieten:

- Transfer von KI in die Praxis
- Angebote für Start-Ups, Mittelstand und Wissenschaft
- Beratung und Entwicklung
- Zugang zu moderner KI-Hardware
- Weiterbildung

Gefördert durch:



Bundesministerium
für Forschung, Technologie
und Raumfahrt

KI Service
Zentrum
by Hasso-Plattner-Institut



Mehr Infos

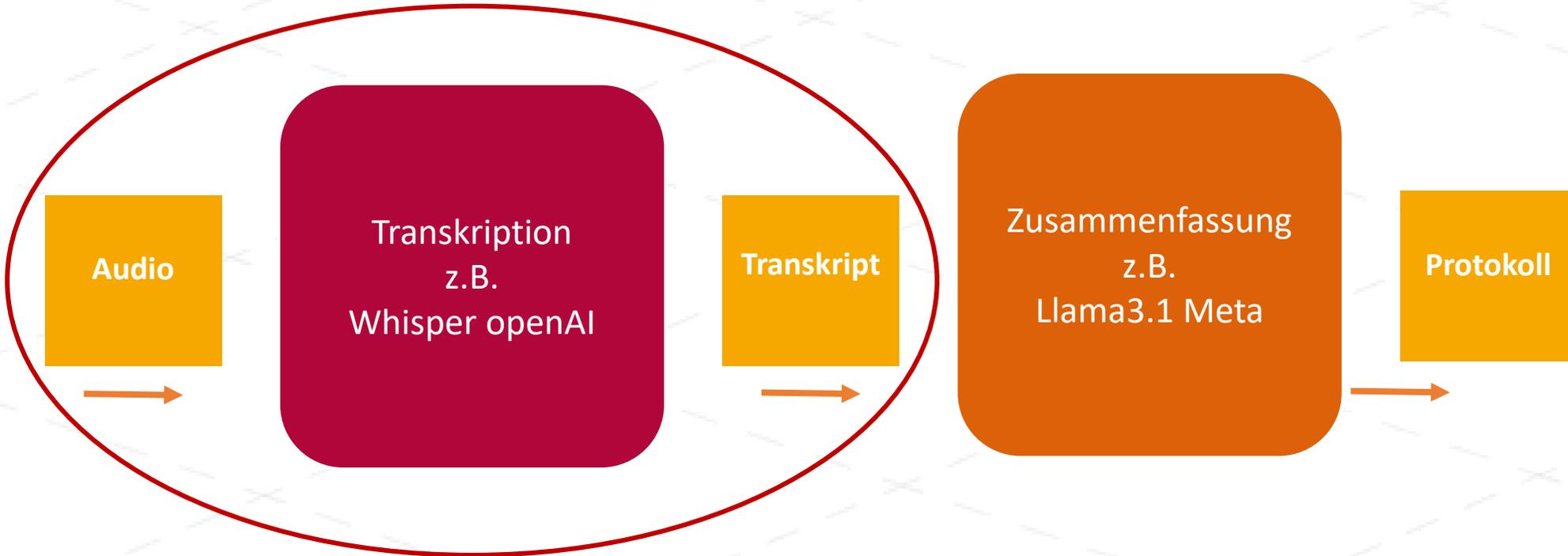


Unser Schwerpunkt:

Bildungs- und Beratungsangebote, Einsatz von KI in
Wirtschaft & Gesellschaft.



Automatisch Protokolle erstellen



Transkription:

- **Performante Transkriptionstools für vielgesprochene Sprachen**
 - Englisch, Mandarin, Spanisch, Französisch, Deutsch, ...
- **Tools für Minderheitensprachen (z.B. Kurdisch) erfordern Finetuning**
- **Kommerzielle Modelle**
 - Google Cloud Speech-to-Text
 - Microsoft Azure AI Speech
 - Amazon Transcribe
- **Open Source Modelle, z.B. OpenAI Whisper, NVIDIA NeMO Parakeet**

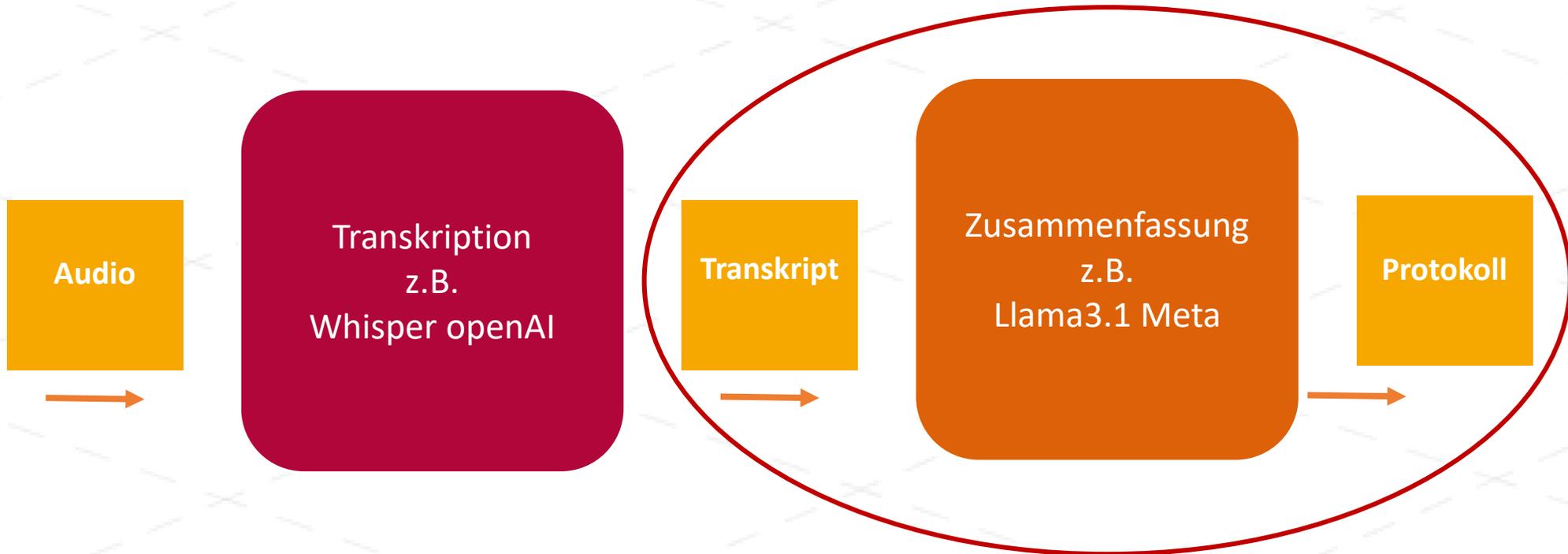
Funktionalität:

- **Spracherkennung**
- **Sprechererkennung**
- **Formatierung und Interpunktion**
 - **Automatische Groß-/ Kleinschreibung, Zeichensetzung**
 - **Automatisches Entfernen von ehm, mhm, uhm etc.**
- **Wort- und Satz-Level Zeitstempel**
- **Online- / Offlinetools**

Herausforderungen:

- Performanzreduktion durch schlechte Audioqualität
- Sprecherüberlappungen
- Verschiedene Wortarten
 - Eigennamen: *Sławomir Krzystof, Reykjavík, Herr König*
 - Abkürzungen: *HPI, KISZ, DABB, BMFTR*
 - Lehnwörter: *Faux pas, Vinaigrette, Köfte*
 - Komposita: *Datenschutzgrundverordnungsbeauftragte*
 - Homophone: *Mei - Mai, Lehre - Leere, Namen - nahmen*
 - Fachbegriffe: *Denitrifikation, Planum, Polymerschweißbahn*

Automatisch Protokolle erstellen



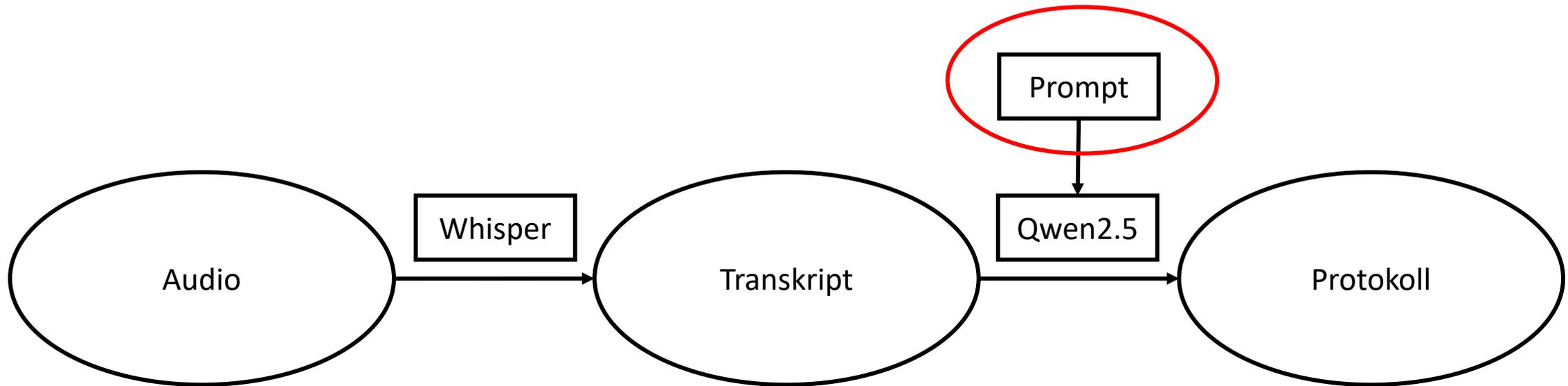
Zusammenfassung: Große Sprachmodelle (LLMs)

- **Nutzen von Sprachmodell zur Transformation des Transkripts**
- **Variablen:**
 - **Sprachmodell:**
 - llama, deepseek, qwen, gemma, phi, mistral, ...
 - **Zusammenfassungsprompt:**
 - **Basierend auf Regeln „Fasse den Text kurz und präzise zusammen“**
 - **Basierend auf Beispielen „Fasse den Text so zusammen wie dieses Beispiel es tut“**

Fragen?

Demo

Pipeline:



Prompt – Teil I:

Du bist Experte für das Zusammenfassen von Video-Transkripten. Deine Aufgabe ist es, eine klare und umfassende Zusammenfassung zu erstellen, die die Hauptpunkte, zentrale Erkenntnisse und wichtige Details des Videos erfasst.

WICHTIGE SPRACHVORGABE:

- Du MUSST die Zusammenfassung in DERSELBEN SPRACHE verfassen wie das Transkript
- Wenn das Transkript auf Deutsch ist, schreibe die Zusammenfassung auf Deutsch
- Wenn das Transkript auf Englisch ist, schreibe die Zusammenfassung auf Englisch
- Wenn das Transkript in einer anderen Sprache ist, schreibe die Zusammenfassung in dieser Sprache
- Übersetze den Inhalt NICHT in eine andere Sprache

Prompt – Teil II:

Bitte analysiere das folgende Video-Transkript und gib eine gut strukturierte Zusammenfassung, die:

1. Das Hauptthema oder den zentralen Inhalt identifiziert
2. Wichtige Punkte und zentrale Informationen hervorhebt
3. Etwaige Schlussfolgerungen oder Erkenntnisse enthält
4. Den ursprünglichen Kontext und die Bedeutung beibehält
5. Klar und prägnant formuliert ist

Prompt – Teil III:

WICHTIGE VORGABEN FÜR DIE AUSGABE:

- Schreibe im Nur-Text-Format ohne jegliche Markdown-Formatierung (kein **, *, # usw.)
- Füge NICHT „Zusammenfassung:“ oder „SUMMARY:“ am Anfang hinzu
- Verwende KEINE Aufzählungszeichen oder nummerierte Listen
- Schreibe in Fließtext mit klaren Übergängen zwischen den Gedanken
- Beginne direkt mit dem Inhalt
- Gib AUSSCHLIESSLICH den endgültigen Zusammenfassungstext aus, sonst nichts
- Wichtig: Verwende DIESELBEN SPRACHE wie das Transkript

In diesem Video geht es um die Planung und Regulierung einer neuen Marktentwicklung in Wilhelmshorst, mit besonderem Fokus auf Verkehrsmanagement, Bürgerbeteiligung und die Rollen verschiedener Behörden. Der Sprecher stellt klar, dass der Bebauungsplan nicht direkt Regeln zu Tempolimits oder zur Anordnung von Fahrspuren festlegt, da diese Aspekte durch Fachgutachten und andere zuständige Fachbehörden geregelt werden. Im Rahmen des Plans werden der Verkehrsfluss sowie die Sicherheit für Fußgänger, Radfahrer, den öffentlichen Nahverkehr und andere Verkehrsteilnehmer bewertet. Die Ergebnisse werden in Gutachten dokumentiert, die sowohl von der Öffentlichkeit als auch von den zuständigen Stellen eingesehen werden können. In der Diskussion wird auch auf die Sorge über ein erhöhtes Verkehrsaufkommen durch den neuen Markt eingegangen. Dabei wird betont, dass zwar mit mehr Verkehr zu rechnen ist, jedoch Maßnahmen ergriffen werden, um diesen zu verteilen und die Erreichbarkeit über alternative Wege wie Fuß- oder Radverkehr zu verbessern – im Einklang mit den Klimaschutzzielen. Der Sprecher hebt die Komplexität der Verkehrsplanung hervor, lehnt vereinfachte Zahlenspiele ab und betont die Notwendigkeit einer umfassenden Analyse. Zudem wird über die Zukunft bestehender Märkte wie zum Beispiel Netto gesprochen und auf die Herausforderungen eingegangen, die sich aus der Abwägung von Entwicklung, bestehender Infrastruktur und Eigentumsverhältnissen ergeben. Der Sprecher erkennt die öffentliche Beteiligung an und unterstreicht die Bedeutung, widersprüchliche Interessen in solchen Projekten zu berücksichtigen, weist aber auch auf die begrenzten Einflussmöglichkeiten der lokalen Behörden bei Entscheidungen hin, die private Akteure betreffen. Das Gespräch endet mit dem Bekenntnis zu Transparenz und einem fortlaufenden Dialog zur weiteren Ausarbeitung des Plans.

- 0:17-0:24** „er begutachtet auch den Kn die Kreuzung“
- 0:30-0:35** „Also Fußgänger, ÖPNV, MEV.“
- 0:35-0:43** „Also das wird begutachtet und eh aufgezeigt ...“
- 0:51-0:56** „Das das war die Aussage, der Bebauungsplan kann nicht alles regeln.“
- 0:56-1:01** „Das wollte ich damit nur klar schreiben (stellen)“
- 1:06-1:12** „dann sehen Sie auch die gutachterlichen Zahlen, den ich hier als seriöse Planerin nicht vowärtsgreifen (vorweggreifen) möchte“

Fragen?



Sprechstunde buchen

BERATUNG

KI-Sprechstunde

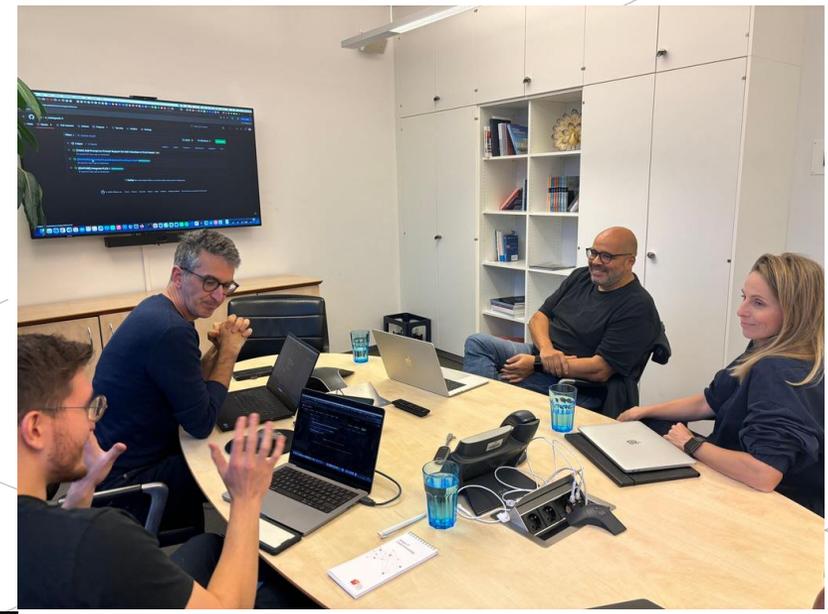
- Beantwortung von Fragen:
 - zu KI-Infrastruktur
 - zu KI-Modellen & Frameworks

KI-Pilotprojekte

- Co-Entwicklung eines Prototyps
- Bewerbung alle drei Monate
- Auswahlkriterien z.B. KI-Reife, Gemeinwohl
- Veröffentlichung der Ergebnisse

Kooperationen

- Gemeinsam organisierte Netzwerktreffen



Leichte Sprache - Simplify Text with LLMs

Simplify your text with a LLM!

Original Text

Kindertagesstätten
In unserer Gemeinde gibt es drei Kinderkrippen und vier Kindergärten. Die Trägerschaft liegt bei der Gemeinde sowie bei der evangelischen Kirche. Es kann zwischen verschiedenen Öffnungszeiten gewählt werden, von den Regelöffnungszeiten bis hin zur Ganztagesbetreuung.

Leichte Sprache

In unserer Gemeinde gibt es drei Kinderkrippen und vier Kindergärten. Die Leiter sind die Gemeinde und die evangelische Kirche. Man kann zwischen verschiedenen Öffnungszeiten wählen, von der normalen Zeit bis hin zur ganzen Tag betreut werden.

Settings Export

Clear Simplify!

Bisherige KI-Pilotprojekte

- Generierung Mathematik-Problemen
- Leichte Sprache
- Generierung von Upcycling Vorschlägen
- Reduzierung von Food Waste
- Datierung mittels Handschrift



Jetzt bewerben!

github.com/aihpi/leichte-sprache



Newsletter

BILDUNG

Talks

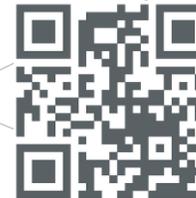


tele-task.de/series/1463

- Gastvorträge zu Forschung und Innovation

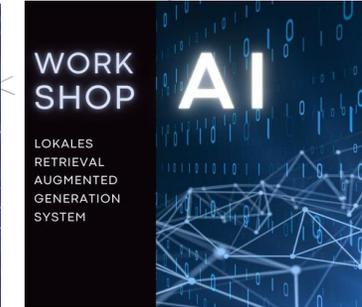


Workshops



aimaker.community

- Praxisnahe Themen
- Beispielthemen: Speech2summary, Docker für ML, semantische Suche



MOOCs



open.hpi.de/channels/ai-service-center

- ChatGPT: Was bedeutet generative KI für unsere Gesellschaft?
- Profitable KI
- KI Biases verstehen und vermeiden



Gefördert durch:





Zugangsanfrage
aisc.hpi.de

INFRASTRUKTUR

Gefördert durch:



- Zugang **kostenfrei**
- kein Produktionsbetrieb
 - Daten sollten **anonymisiert** oder **synthetisiert** sein
 - kein **Hosting** von Produkten
- **Reporting & Veröffentlichung** durch Nutzende
- **Altrechte** bleiben bei Nutzenden
- **Neurechte** bleiben bei Nutzenden
 - Einräumen von Nutzungsrechten für Forschung und Lehre

Training

- 64 NVIDIA H100 GPU

Inferenz

- 40 NVIDIA A30 GPU

ARM Server

- Ampere Altra Max M128-30 CPU
- 2 x NVIDIA L40 GPUs

GPU Server

- AMD Epyc CPU
- 8 x NVIDIA L40S GPU

Edge

- ARMv8 CPU
- NVIDIA Jetson AGX Module

Neuromorph

- 288 SpiNNaker2 Chips

Speicher

- 1.5 PB NVRAM

Netzwerk

- 400 Gb/s Infiniband
- 200 Gb/s Ethernet

Was wird benötigt für automatisiertes Protokollieren?

- **ASR-Tool** (z.B. Whisper)
 - Ggf. regelmäßiges Einpflegen von neuen Namen, Abkürzungen etc.
 - Daten, Know-How, Rechenressourcen
- **LLM**
 - Starke GPU; je nach Auslastung mehrere GPUs
 - Kostenintensiv (Miete Rechenressourcen / Strom)
- **DSGVO-Konformität**
 - **Rechenressourcen mieten vs. selber betreiben**

Rechenressourcen:

Mieten

- Regelmäßige Mietkosten
- Externe Wartung
- Leicht skalierbar

- DSGVO kritisch
- Internetverbindung erforderlich;
Daten evtl. auf externen Servern

Selbstbetrieb

- Hohe Anschaffungskosten
- Manuelle Wartung
- Skalierung mit neuen
Anschaffungskosten verbunden
- DSGVO einfach umsetzbar
- Internetverbindung optional

Kontakt

Lasse Kohlmeyer

AI Engineer & Project Manager
KI-Servicezentrum Berlin-Brandenburg

Lasse.Kohlmeyer@hpi.de

Hanno Müller

AI Engineer
KI-Servicezentrum Berlin-Brandenburg

Hanno.Mueller@hpi.de



Kommen Sie bei Rückfragen gerne auf uns zu.