

- openHPI: ChatGPT: Was bedeutet generative KI für unsere Gesellschaft? -

LLM als Puzzleteil

Johannes Hötter
Christian Warmuth

Was passieren kann, wenn man bei ChatGPT nicht aufpasst..

Name	Firma	Titel	Email	Mobilnr
Mickey Mouse	Disney	CEO	mickey.mouse@disney.com	+1 (555) 123-4567
Harry Potter	Hogwarts School	Student	harry.potter@hogwarts.edu	+44 1234 567890
Sherlock Holmes	221B Baker St.	Private Detective	sherlock.holmes@221B.com	+44 207 123 4567
Darth Vader	Galactic Empire	Sith Lord	darth.vader@galacticempire.com	+1 (555) 555-5555

LLM als Puzzleteil

Johannes Hötter,
Christian Warmuth

openHPI

Was passieren kann, wenn man bei ChatGPT nicht aufpasst..

Idee: personalisierte Nachricht, um positiv aufzufallen - mit ChatGPT automatisiert erzeugt.

Z.B:

„Sehr geehrter Herr Mouse,

es ist mir eine große Ehre, Ihnen eine Nachricht zu senden! Als einer der bekanntesten und beliebtesten Charaktere in der Geschichte der Animationsfilme haben Sie unzählige Kinder und Erwachsene auf der ganzen Welt begeistert und inspiriert.“

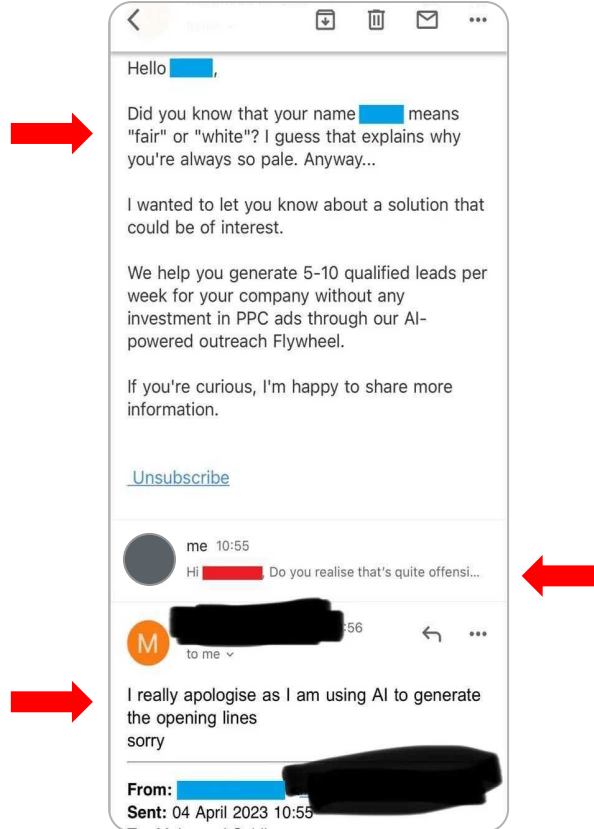
LLM als Puzzleteil

Johannes Hötter,
Christian Warmuth

openHPI

Was passieren kann, wenn man bei ChatGPT nicht aufpasst..

**Und dann
passiert das..**



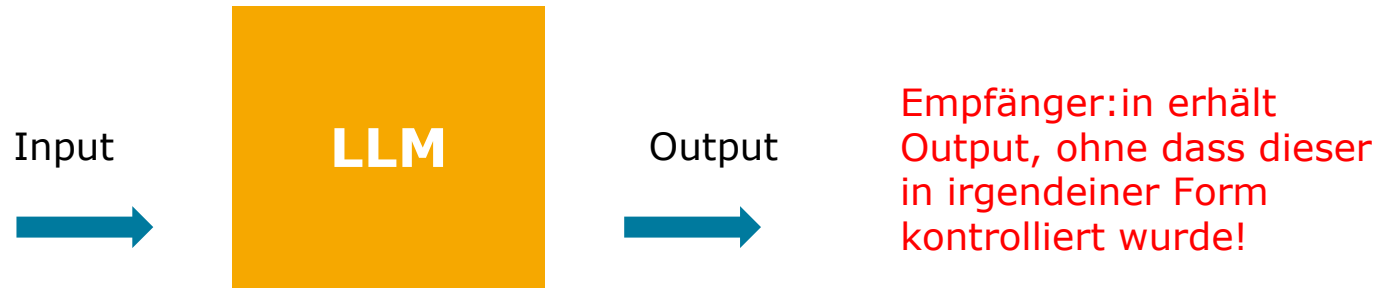
LLM als Puzzleteil

Johannes Hötter,
Christian Warmuth

openHPI

Folie 4

Der Fehler, welcher hier passiert ist:



LLM als Puzzleteil

Johannes Hötter,
Christian Warmuth

openHPI

Wie kann so etwas verhindert werden?



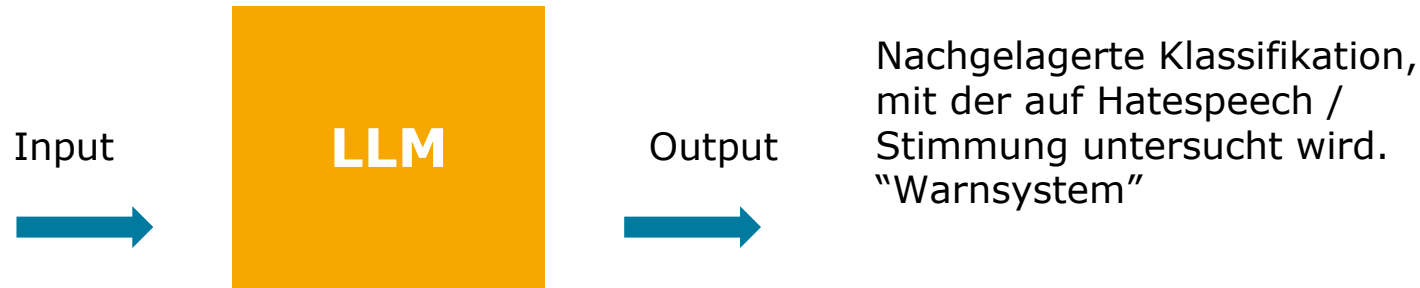
Entwurf wird von GPT erzeugt,
und von einem Menschen
kontrolliert.
Schreibt Schreibarbeit, aber
erhöht Kontrolle.

LLM als Puzzleteil

Johannes Hötter,
Christian Warmuth

openHPI

Wie kann so etwas verhindert werden?

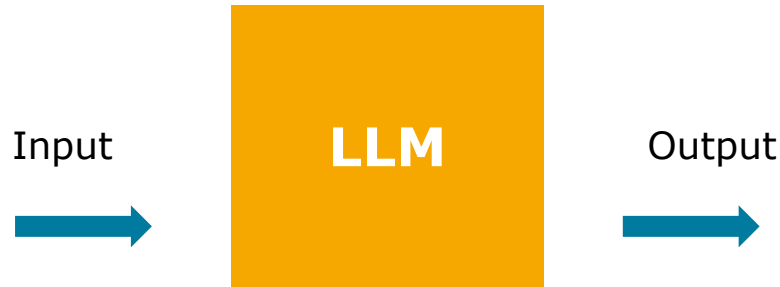


LLM als Puzzleteil

Johannes Hötter,
Christian Warmuth

openHPI

Wie kann so etwas verhindert werden?



Sorgfältiges Benchmarking und Verproben zahlreicher Prompts, z.B. um Situation zu schildern, in dem der Eisbrecher erzeugt werden soll.

LLM als Puzzleteil

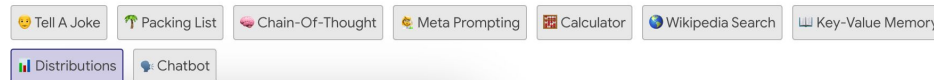
Johannes Hötter,
Christian Warmuth

openHPI

Wie kann so etwas verhindert werden?

Supercharge your prompting with **constraints**.

LMQL is a programming language for language model interaction.



LMQL Open In Playground

```
argmax
"""Review: We had a great stay. Hiking
- in the mountains was fabulous and the
- food is really good.\n
Q: What is the underlying sentiment of
- this review and why?\n
A: [ANALYSIS]\n
Based on this, the overall sentiment of
- the message can be considered to be [CLASSIFICATION]"""

from
"openai/text-davinci-003"
distribution
.....
CLASSIFICATION in [" positive", " neutral", "
- negative"]
```

MODEL OUTPUT

Review: We had a great stay. Hiking in the mountains was fabulous and the food is really good.

Q: What is the underlying sentiment of this review and why?

A: **ANALYSIS** The underlying sentiment of this review is positive because the reviewer enjoyed their stay, the hiking, and the food.

Based on this, the overall sentiment of the message can be considered to be **CLASSIFICATION**

$$P(\text{CLASSIFICATION}) = \begin{cases} \text{positive} & 0.9998711120293567 \\ \text{neutral} & 0.00012790777085508993 \\ \text{negative} & 9.801997880775052e-07 \end{cases}$$

Highlighted text is model output.

LLM als Puzzleteil

Johannes Hötter,
Christian Warmuth

openHPI

Folie 9

Allgemein: Möglichkeiten des Preprocessings

Reichere Input an,
z.B. mit Produktinfos,
Suchergebnissen,
oder IT-System-Infos
(siehe Anwendungen
diese Woche)

Input



LLM

Output

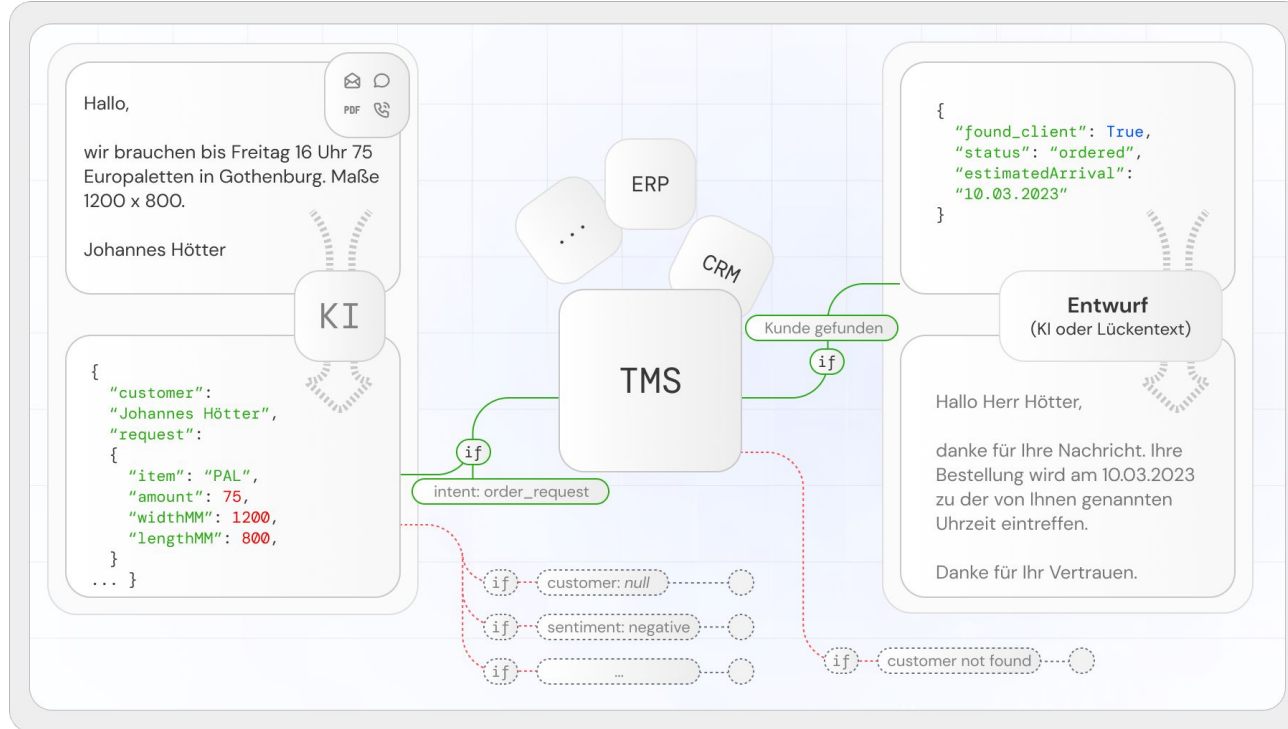


LLM als Puzzleteil

Johannes Hötter,
Christian Warmuth

openHPI

GPT ist ein Puzzleteil, meist nicht der ganze Prozess



LLM als Puzzleteil

Johannes Hötter,
Christian Warmuth

openHPI

- openHPI: ChatGPT: Was bedeutet generative KI für unsere Gesellschaft? -

LLM als Puzzleteil

Johannes Hötter
Christian Warmuth