

# Automatic Term Extraction based on GenAI

Terminology Agent System

Bachelor's Thesis Presentation in Computer Science

Jonas Ludwig Gerg

24.03.2025 - 24.07.2025

Chair of Computer Systems

<https://dse.in.tum.de/>



## Outline

Motivation .....	2
Background .....	4
Goals .....	8
Design .....	10
Implementation .....	14
Testing & Evaluation .....	17
Conclusion .....	21

## Motivation

- Deutsche Bahn is facing a strong demand for qualified workers<sup>1</sup>
- Focus on recruitment from abroad<sup>2</sup>
- Language barriers due to domain-specific terminology
- Provision of explanations in texts and conversations
- Challenge: **How can we efficiently and reliably detect and extract terms and phrases in texts?**

---

<sup>1</sup>[1] “DB Job Portal aktuelle Anzahl an gesuchten Vollzeitstellen.” Accessed: Mar. 12, 2025. [Online]. Available: <https://db.jobs/de-de/Suche>

<sup>2</sup>[2] “Deutsche Bahn rekrutiert Auszubildende im Ausland.” Accessed: Jul. 02, 2025. [Online]. Available: <https://www.sueddeutsche.de/wirtschaft/bahn-azubis-migration-1.5672979>

## Outline

Motivation .....	2
<b>Background .....</b>	<b>4</b>
Goals .....	8
Design .....	10
Implementation .....	14
Testing & Evaluation .....	17
Conclusion .....	21

## Automatic Term Extraction (ATE)

The process of identification and **extraction** of **terminology** specific to a particular **domain** within a **document** using an **automated process**.<sup>1</sup>

---

<sup>1</sup>[3] G. M. Di Nunzio, S. Marchesin, and G. Silvello, “A Systematic Review of Automatic Term Extraction: What Happened in 2022?,” *Digital Scholarship in the Humanities*, vol. 38, no. Supplement\_1, pp. i41–i47, Jun. 2023, doi: 10.1093/lhc/fqad030.

<sup>2</sup>Example using nltk (<https://www.nltk.org>) and the universal tag set

# Automatic Term Extraction (ATE)

The process of identification and **extraction** of **terminology** specific to a particular **domain** within a **document** using an **automated process**.<sup>1</sup>

- Linguistic: “*analyze linguistic properties of language*”<sup>2</sup>

```
[('Auf', 'NOUN'), ('dieser', 'NOUN'), ('Strecke', 'NOUN'), ('musst', 'NOUN'), ('du', 'NOUN'), ('auf', 'NOUN'), ('Sicht', 'NOUN'), ('fahren', 'NOUN'), ('!', '.'), ('.', '.')]
```

- Dictionary based: “*match with known dictionary entries*”
- Statistical: “*extract based on token frequency*”
- Large Language Models: “*extract using pretrained models*”

---

<sup>1</sup>[3] G. M. Di Nunzio, S. Marchesin, and G. Silvello, “A Systematic Review of Automatic Term Extraction: What Happened in 2022?,” *Digital Scholarship in the Humanities*, vol. 38, no. Supplement\_1, pp. i41–i47, Jun. 2023, doi: 10.1093/lc/fqad030.

<sup>2</sup>Example using nltk (<https://www.nltk.org>) and the universal tag set

## Terms and Phrases in Texts

“Auf dieser Strecke musst du **auf Sicht fahren!**”

“**Fahrt auf Sicht** bis Kilometer 103.”

“Ich **fahre** hier seit längerem **auf Sicht.**”

“Ich muss **auf Sicht** für 3km **fahren.**”

## Terms and Phrases in Texts

“Auf dieser Strecke musst du **auf Sicht fahren!**”

“**Fahrt auf Sicht** bis Kilometer 103.”

“Ich **fahre** hier seit längerem **auf Sicht.**”

“Ich muss **auf Sicht** für 3km **fahren.**”

→ **Fahrt auf Sicht**



## Terms and Phrases in Texts

“Auf dieser Strecke musst du **auf Sicht fahren!**”

“**Fahrt auf Sicht** bis Kilometer 103.”

“Ich **fahre** hier seit längerem **auf Sicht.**”

“Ich muss **auf Sicht** für 3km **fahren.**”

→ **Fahrt auf Sicht**

→ **Auf Sicht fahren**

Terms occur in different variations in different languages

## Challenges with ATE<sup>1</sup>

- Variability of Terms and Polysemy
- Domain Dependency
- Cross-Language Compatibility  
→ *different linguistic features*
- Scalability issues  
→ *especially dictionary based methods*

---

<sup>1</sup>[4] R. Mitkov, Ed., *The Oxford Handbook of Computational Linguistics 2nd Edition*, 2nd ed. Oxford University Press, 2014. doi: 10.1093/oxfordhb/9780199573691.001.0001.

## Outline

Motivation .....	2
Background .....	4
<b>Goals .....</b>	<b>8</b>
Design .....	10
Implementation .....	14
Testing & Evaluation .....	17
Conclusion .....	21

## Requirements and Design Goals

- Extensibility for new models, algorithms and data sources
- Reliable and safe definition generation
- Performant extraction and generation

## Requirements and Design Goals

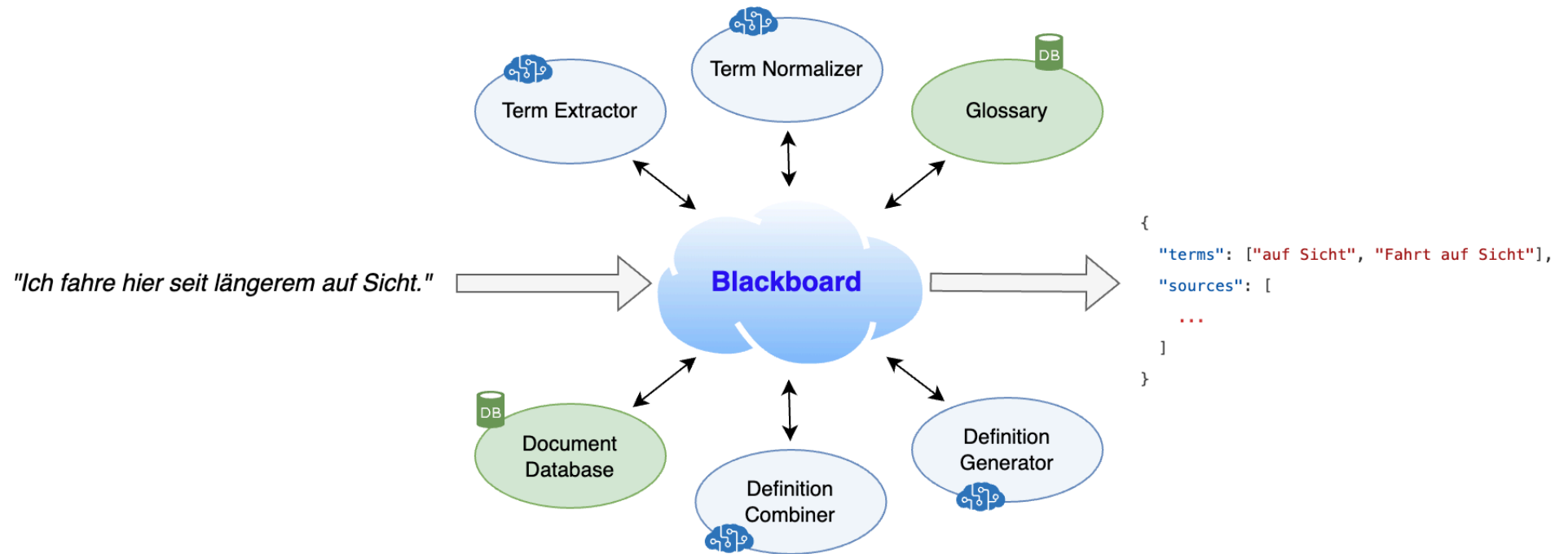
- Extensibility for new models, algorithms and data sources
- Reliable and safe definition generation
- Performant extraction and generation

→ **Main priority on reliability and safety**

## Outline

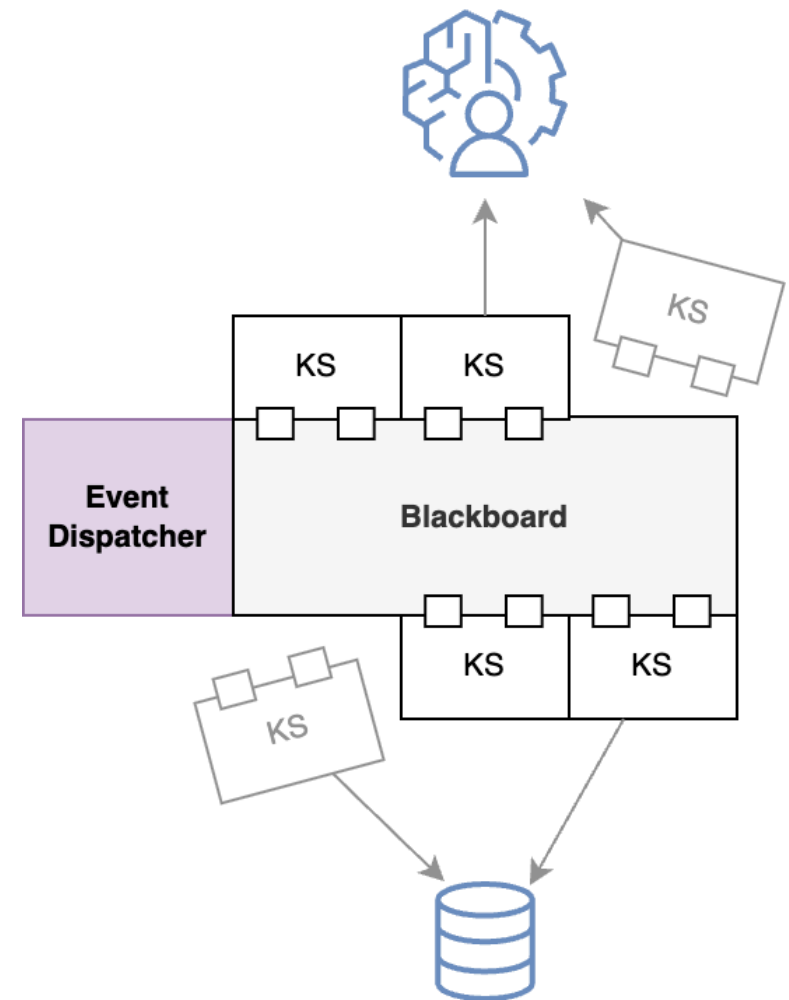
Motivation .....	2
Background .....	4
Goals .....	8
<b>Design .....</b>	<b>10</b>
Implementation .....	14
Testing & Evaluation .....	17
Conclusion .....	21

# Design Overview



## Design Details

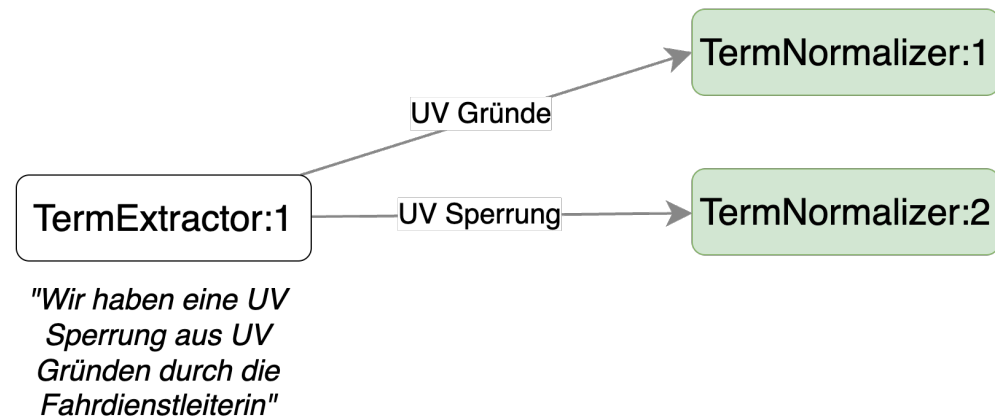
- Event-driven blackboard pattern architecture<sup>1</sup>
- Knowledge Sources (KS) ...
  - ... operate on shared blackboard
  - ... react to and publish new events
  - ... run independently
  - ... employ **LLMs** to solve linguistic tasks
  - ... optionally access external systems



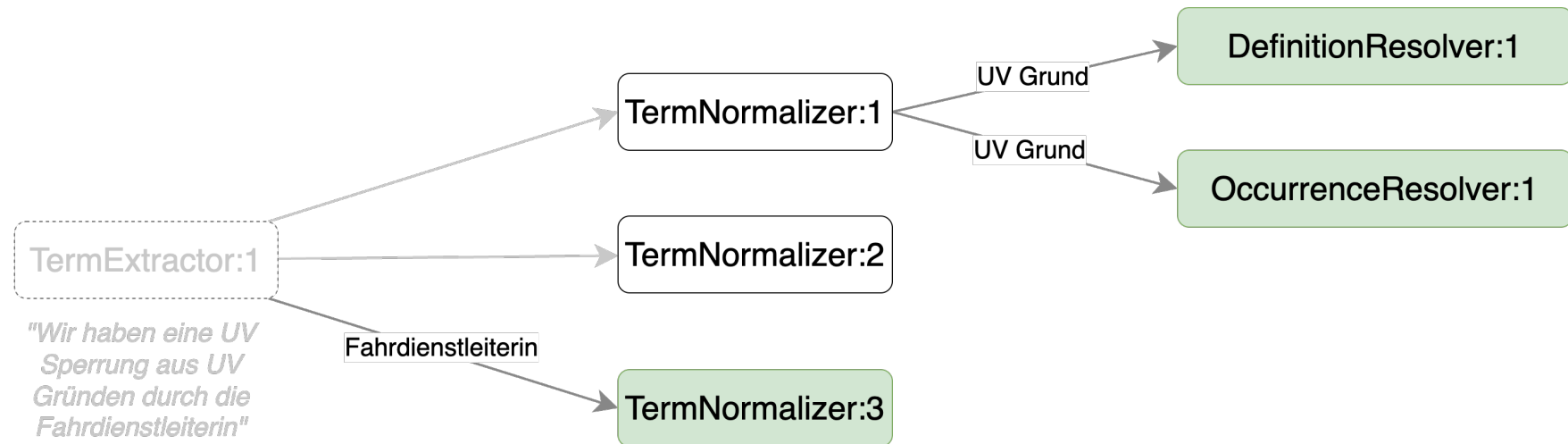
<sup>1</sup>[5] P. Lalanda, "Two Complementary Patterns to Build Multi-Expert Systems."



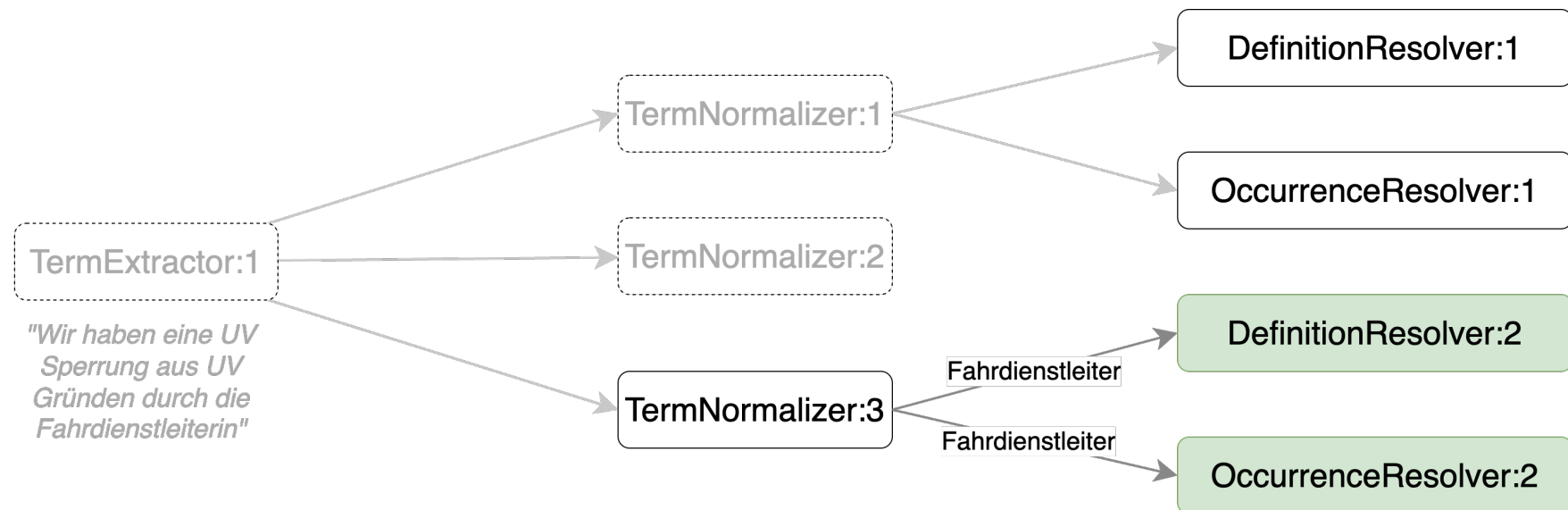
## Example Event Flow



## Example Event Flow



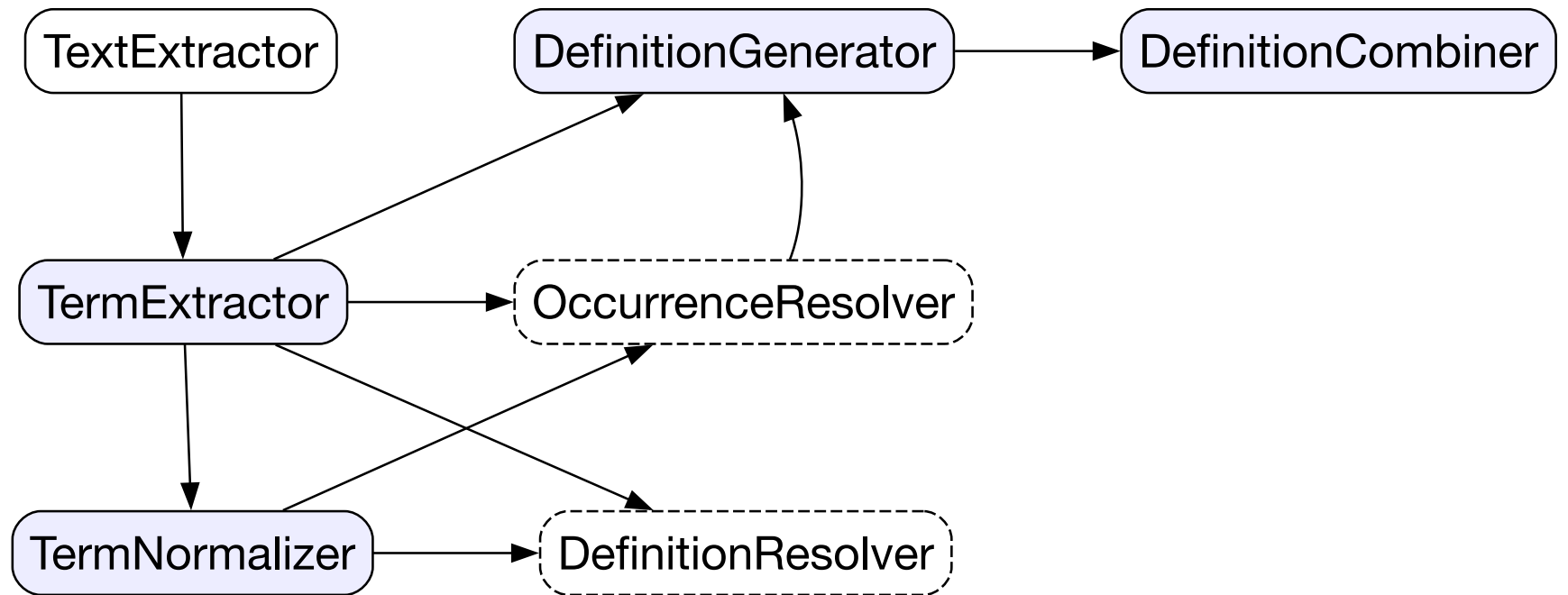
## Example Event Flow



## Outline

Motivation .....	2
Background .....	4
Goals .....	8
Design .....	10
<b>Implementation .....</b>	<b>14</b>
Testing & Evaluation .....	17
Conclusion .....	21

## Knowledge Sources



## LLM based Knowledge Sources

- Use OpenAIs gpt-4o-mini model
- Different Prompt Engineering techniques:<sup>1</sup>
  - Zero-Shot
  - Few-Shot
  - Prompt-Chaining<sup>2</sup>
- temperature=0 → reduce randomness<sup>3</sup>
- Implementation for DB uses internal BahnGPT

TermExtractor

TermNormalizer

DefinitionGenerator

DefinitionCombiner

---

<sup>1</sup>[6] *Prompt Engineering for LLMs*. Accessed: Jun. 11, 2025. [Online]. Available: <https://learning.oreilly.com/library/view/prompt-engineering-for/9781098156145/>

<sup>2</sup>[7] Anthropic, “Chain Complex Prompts for Stronger Performance.” Accessed: Jul. 04, 2025. [Online]. Available: <https://docs.anthropic.com/en/docs/build-with-claude/prompt-engineering/chain-prompts>

<sup>3</sup>[8] M. Renze, “The Effect of Sampling Temperature on Problem Solving in Large Language Models,” in *Findings of the Association for Computational Linguistics: EMNLP 2024*, Miami, Florida, USA: Association for Computational Linguistics, 2024, pp. 7346–7356. doi: 10.18653/v1/2024.findings-emnlp.432.

## Outline

Motivation .....	2
Background .....	4
Goals .....	8
Design .....	10
Implementation .....	14
<b>Testing &amp; Evaluation .....</b>	<b>17</b>
Conclusion .....	21

## Testing LLM Applications

```
1 // probabilities in iteration 1
2 {'FALSE': 0.6224593312018545, 'TRUE': 0.3775406687981454}
3 // probabilities in iteration 2
4 {'FALSE': 0.9149009474519222, 'TRUE': 0.08509905254807776}
5 // probabilities in iteration 3
6 {'TRUE': 0.679178699175393, 'FALSE': 0.32082130082460697}
7 // probabilities in iteration 4
8 {'FALSE': 0.9399133527714579, 'TRUE': 0.060086647228542005}
9 // probabilities in iteration 5
10 {'TRUE': 0.9241418131886822, 'FALSE': 0.0758581868113179}
```

→ Multiple iterations of tests were not consistent, even with deterministic LLM settings



## Challenges with Testing of LLM Applications

1. **Non-Determinism** and Inconsistency of LLM output  
→ *different output for every execution*
2. **Semantic comparison** between oracle and observation  
→ colloquial “**420er**” ↔ official “**420**”
3. **Increased cost** when allowing reasoning  
→ *however: testing was successful afterwards*

## Evaluation

- Limited test set of 21 crafted, short test sentences



## Outline

Motivation .....	2
Background .....	4
Goals .....	8
Design .....	10
Implementation .....	14
Testing & Evaluation .....	17
Conclusion .....	21

## Summary

### Traditional ATE methods not properly applicable

- limited cross-language support
- domain-specificity
- fine-tuning and training required

### TAS:

- Flexible integration of new technologies or data sources
- Out-of-the-box models
- Potentially high recall rate

⚠ Further research required in determinism for testing LLM applications and the performance of TAS

## **Future Research**

- Performance improvements and scalability of TAS
  - Fine-tuning of baseline models
  - Distributed knowledge sources
- Cross-language and cross-domain compatibility
- ChatBot functionality for TAS
- Testing of LLM applications (*“how to deal with non-determinism?”*)

## Sources

- [1] “DB Job Portal aktuelle Anzahl an gesuchten Vollzeitstellen.” Accessed: Mar. 12, 2025. [Online]. Available: <https://db.jobs/de-de/Suche>
- [2] “Deutsche Bahn rekrutiert Auszubildende im Ausland.” Accessed: Jul. 02, 2025. [Online]. Available: <https://www.sueddeutsche.de/wirtschaft/bahn-azubis-migration-1.5672979>
- [3] G. M. Di Nunzio, S. Marchesin, and G. Silvello, “A Systematic Review of Automatic Term Extraction: What Happened in 2022?,” *Digital Scholarship in the Humanities*, vol. 38, no. Supplement\_1, pp. i41–i47, Jun. 2023, doi: 10.1093/llc/fqad030.
- [4] R. Mitkov, Ed., *The Oxford Handbook of Computational Linguistics 2nd Edition*, 2nd ed. Oxford University Press, 2014. doi: 10.1093/oxfordhb/9780199573691.001.0001.
- [5] P. Lalanda, “Two Complementary Patterns to Build Multi-Expert Systems.”
- [6] *Prompt Engineering for LLMs*. Accessed: Jun. 11, 2025. [Online]. Available: <https://learning.oreilly.com/library/view/prompt-engineering-for/9781098156145/>
- [7] Anthropic, “Chain Complex Prompts for Stronger Performance.” Accessed: Jul. 04, 2025. [Online]. Available: <https://docs.anthropic.com/en/docs/build-with-claude/prompt-engineering/chain-prompts>
- [8] M. Renze, “The Effect of Sampling Temperature on Problem Solving in Large Language Models,” in *Findings of the Association for Computational Linguistics: EMNLP 2024*, Miami, Florida, USA: Association for Computational Linguistics, 2024, pp. 7346–7356. doi: 10.18653/v1/2024.findings-emnlp.432.