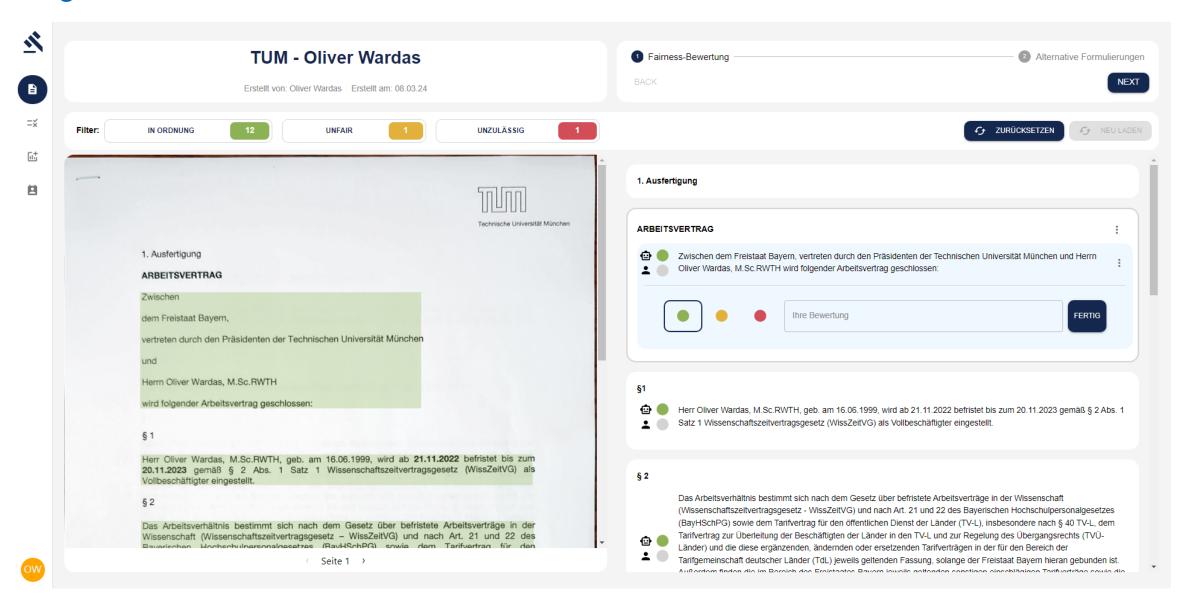


The global legal technology market has grown significantly in recent years and GenAl will accelerate this growth, meaning the market will reach \$50 billion in value by 2027[1].

^{*} until 10.12.2024



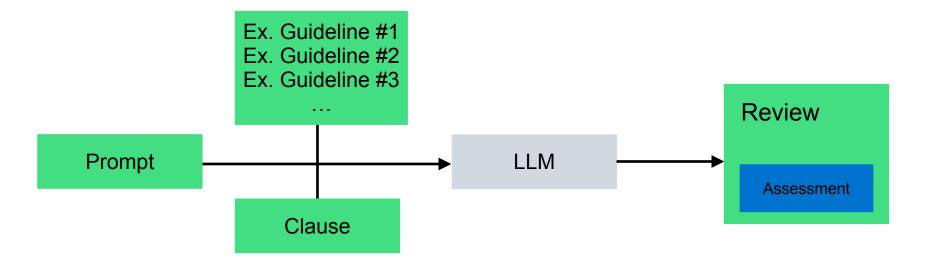


What is an Examination Guideline?



The gross minimum wage is €12.41.

Company training courses must generally be paid for.



What is an Examination Guideline?



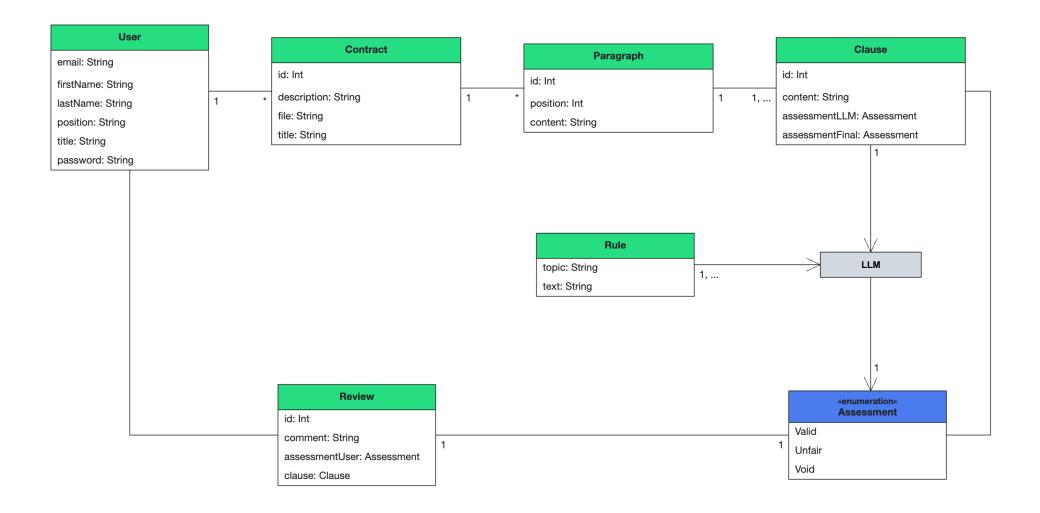
The gross minimum wage is €12.41.

Valid until 31.12.2024

The gross minimum wage is €12.82.

Valid from 01.01.2025





Research Questions



RQ1	How can changes in examination guidelines and their impact on the evaluation of contract clauses in an in-context learning AI system be tracked and assessed?
RQ2	What is needed to create a user-friendly interface which allows to monitor this evolution without technical knowledge?
RQ3	How can technical aspects and the usability of the system be evaluated?

Methodology



RQ1

Slowly Changing Dimensions (SCD), Structured Output [3,4]

RQ2

Common Software Development Metrics

RQ3

Standard Usability Scale (SUS), Unmoderated Software Testing [7,8]

RQ1: Versioning for Examination Guidelines using SCD Type 2



Rule

topic: string

text: string

Regelhistorie: Ausschlussfristen bei der Geltendmachung von Ansprüchen

Version vom 25.11.2024

Es darf keine strengere Form als die Textform verlangt werden. Textform bedeutet, dass eine lesbare Erklärung ohne Unterschrift, z. B. per E-Mail oder Fax, abgegeben werden kann, sofern die Person des Erklärenden erkennbar ist.

Version vom 20.11.2024

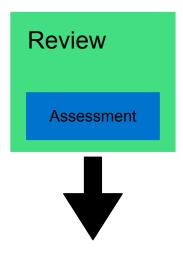
Es darf keine strengere Form als die Textform verlangt werden

SCHLIESSEN

RQ1: Returning applied Examination Guidelines



Pattern Matching



OpenAI's structured output [3]

class ClauseValidationExpert(BaseModel):

clauseld: str

assessment: Literal["void", "unfair", "valid"]

rulelds: List[str] confidence: float

comment: str



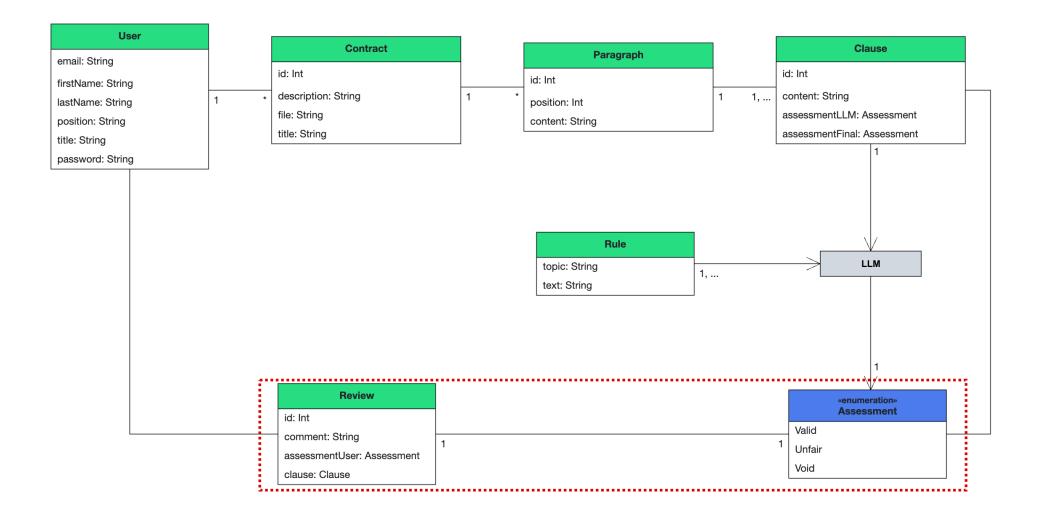
Valid

Unfair

Void

RQ1: Returning applied Examination Guidelines





RQ1: Returning applied Examination Guidelines



OpenAI's structured output [3]

class ClauseValidationExpert(BaseModel):

clauseld: str

assessment: Literal["void", "unfair", "valid"]

rulelds: List[str] confidence: float

comment: str

Review

assessment: Assessment

comment: string

clause: Clause

author: Author

rules: string[]

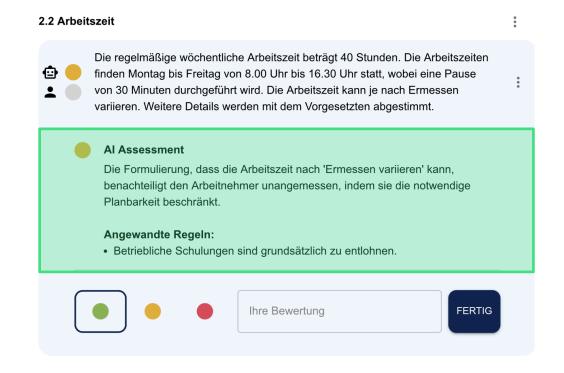
isLLM: boolean

confidence: number

RQ1: Displaying Comment and Applied Examination Guidelines







RQ2: What is needed to create a user-friendly interface which allows to monitor this evolution without technical knowledge?



Meta Information

- Current state of the system
- Amount Contracts, Users
- Amount of Examination Guidelines
- Examination Guidelines per Topic
- Avg. Time per Contract (LLM)

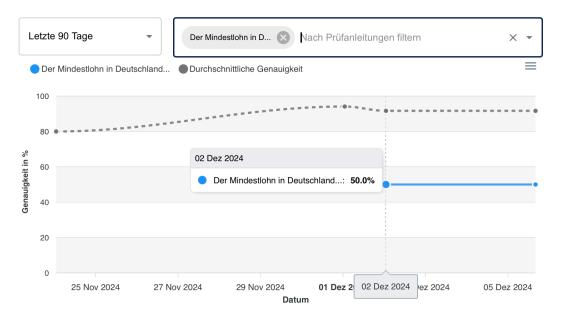
Al Information

- Comparable data to prior versions of the system
- Accuracy of LLMs
- Confusion Matrix
- Confidence Score

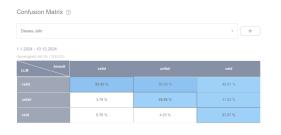
Detail Information Metrics



Genauigkeit der LLMs im Zeitverlauf ②









Detail Information Metrics



Confusion Matrix ②



1.1.2024 - 10.12.2024

Genauigkeit: 68.0% (153/225)

Anwalt LLM	valid	unfair	void
valid	95.45 %	66.20 %	40.91 %
unfair	3.79 %	29.58 %	31.82 %
void	0.76 %	4.23 %	27.27 %



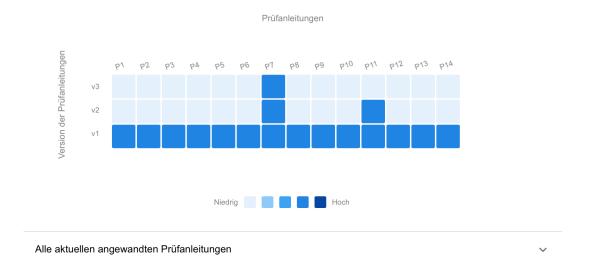




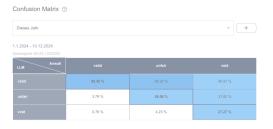
Detail Information Metrics



Confidence Score der Prüfanleitungen ②





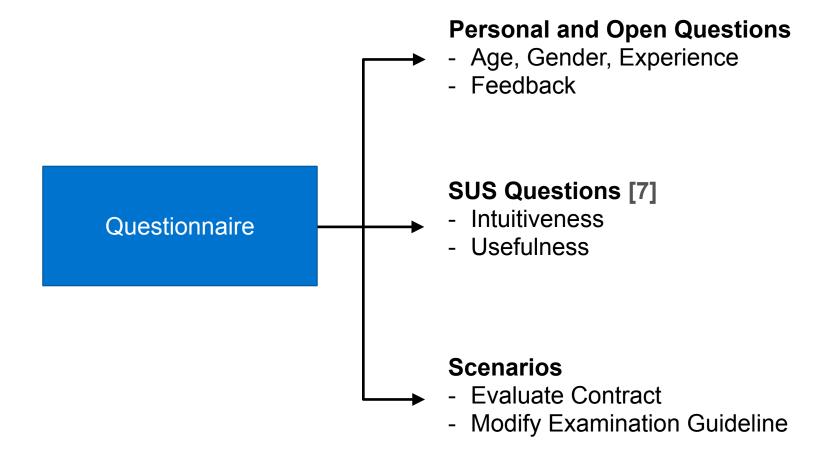




15

RQ3: How can technical aspects and the usability of the system be evaluated?

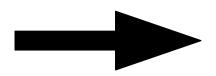




Outlook



Evaluation with lawyers



Usefulness of metrics

New and refined metrics

Limitations



No Evaluation with Lawyers

Single Researcher Bias

Time Constraint



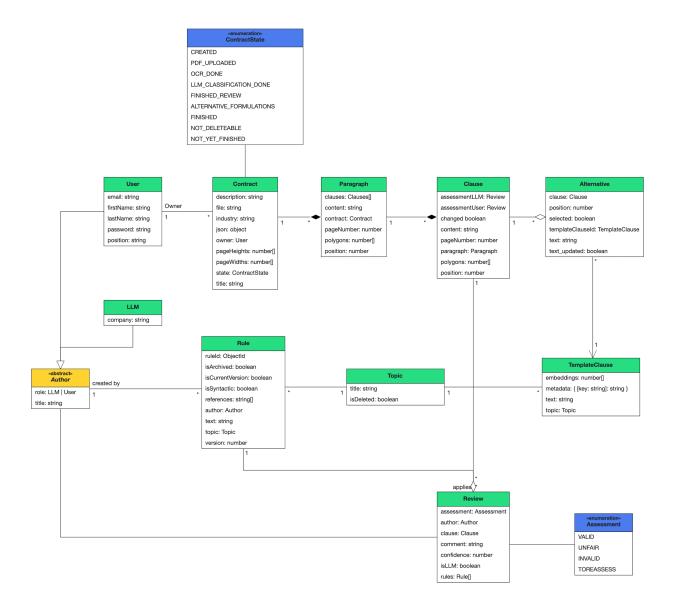
Sources

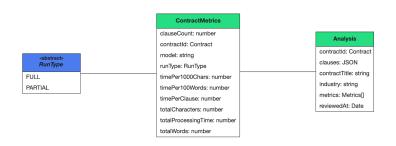


- [1] Gartner, Legal Technology, accessed January 2025, https://www.gartner.com/en/newsroom/press-releases/2024-04-25-gartner-predicts-global-legal-technology-market-will-reach-50-billion-by-2027-as-a-result-of-genai
- [2] Ver.di, accessed January 2025, https://www.verdi.de/themen/arbeit/++co++d4ff4502-5cd5-11ec-9ee8-001a4a16012a
- [3] OpenAI, Structured Output, accessed January 2025, https://platform.openai.com/docs/guides/structured-outputs
- [4] Oracle. Slowly Changing Dimensions, accessed January 2025, https://www.oracle.com/webfolder/ https://www.oracle.com/webfolder/ https://www.oracle.com/webfolder/obe/db/10g/r2/owb/owb10gr2_gs/owb/lesson3/slowlychangingdimensions.htm <a href="technetwork/tutorials/obe/db/10g/r2/owb/obe/db
- [5] K. Tian et al. Just Ask for Calibration: Strategies for Eliciting Calibrated Confidence Scores from Language Models Fine-Tuned with Human Feedback, https://arxiv.org/abs/2305.14975
- [6] K. Shen and M. Kejriwal. A Formalism and Approach for Improving Robustness of Large Language Models Using Risk-Adjusted Confidence Scores, https://arxiv.org/abs/2310.03283
- [7] J. Lewis. "The System Usability Scale: Past, Present, and Future". In: International Journal of Human-Computer Interaction (Mar. 2018), pp. 1–14
- [8] P. Khayyatkhoshnevis, S. Tillberg, E. Latimer, T. Aubry, A. Fisher, and V. Mago. Comparison of Moderated and Unmoderated Remote Usability Sessions for Web-Based Simulation Software: A Randomized Controlled Trial.



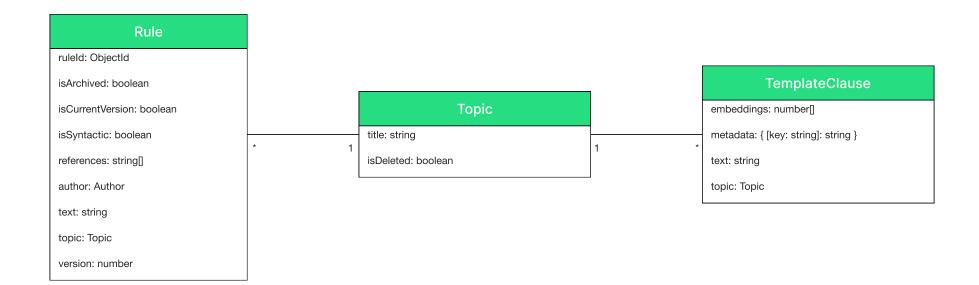






Further Changes

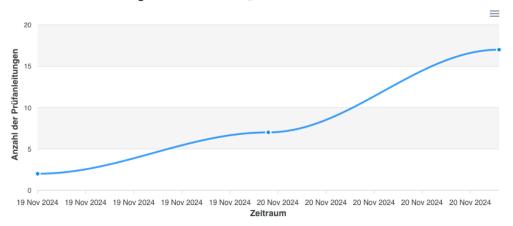




System Information Metrics

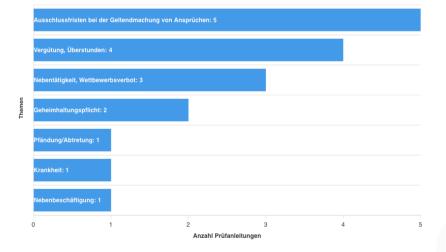
. . . .

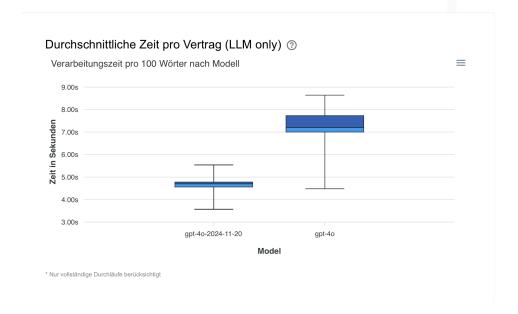
Anzahl der Prüfanleitungen im Zeitverlauf ②



Prüfanleitungen pro Thema ②

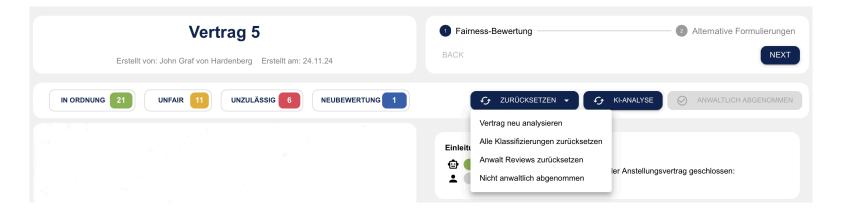


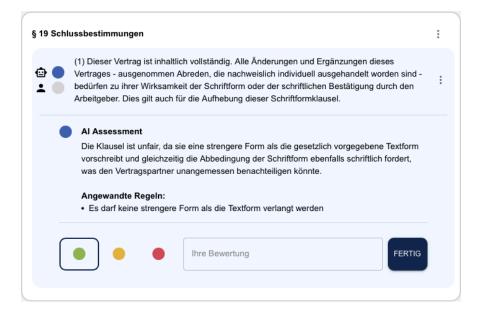


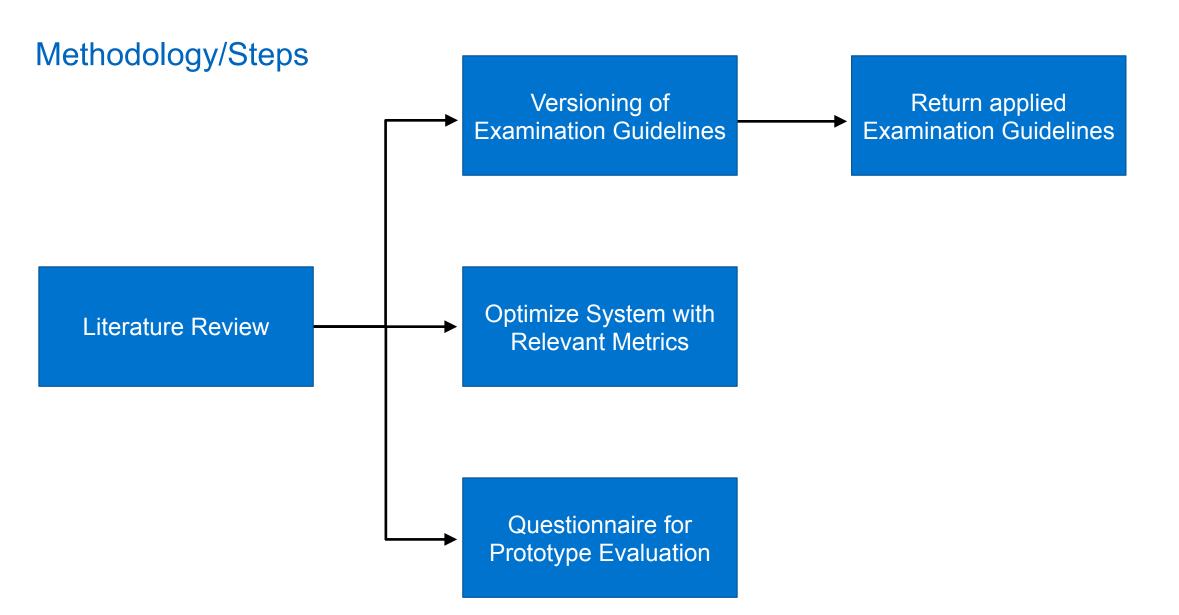


Examination Guideline Changes







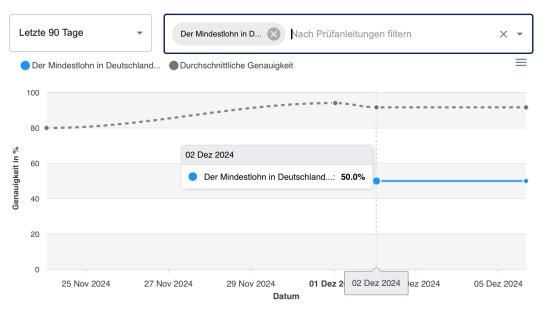




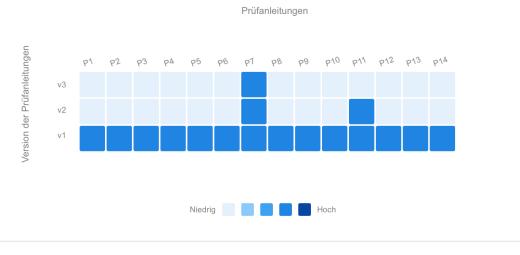
Detail Information Metrics



Genauigkeit der LLMs im Zeitverlauf ②

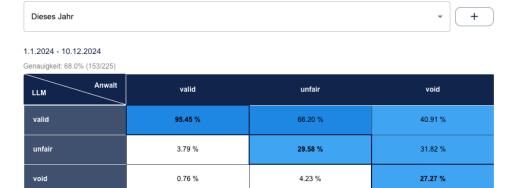


Confidence Score der Prüfanleitungen ②



Alle aktuellen angewandten Prüfanleitungen

Confusion Matrix ②



[5,6]