

# CourtPressGER

## Anonymous ACL submission

### Abstract

We presents CourtPressGER - a system for automatically generating German court press releases using Large Language Models (LLMs). We present a curated dataset with 6.4k entries of court decisions with corresponding press releases from Germany's highest courts. The dataset is enhanced with synthetic prompts that enable automated generation of press releases from court decisions. We describe a pipeline for generating press releases with various state-of-the-art models and evaluate the results using automated metrics and LLM-based evaluation approaches that simulate expert assessment. Our approach combines specialized legal language models with domain-specific techniques to produce accurate and informative press releases that adhere to journalistic and legal standards.

### Introduction

The German legal system consists of a complex network of courts that regularly publish extensive decisions. To make these decisions accessible to the public, the highest courts create press releases that summarize the essential aspects and implications of the decisions in an understandable form. These press releases serve as an important interface between the judicial system and the general public by explaining complex legal matters in an accessible way and serve as a proxy for the task of legal case summarization, for which manually created gold data is typically sparse.

However, the manual creation of such press releases requires significant resources. At the same time, recent advances in Large Language Models (LLMs) offer new possibilities for automated text generation in specialized domains. Our project CourtPressGER aims to leverage these capabilities for the automatic generation of court press releases.

The main contributions of our project are:

- The creation of a curated dataset with 6.4k entries of court decisions with corresponding press releases from Germany's highest federal courts.
- The development of synthetic prompts for each decision-press release pair that can be used to au-

tomatically generate press releases.

- The evaluation of the generated press releases using a combination of traditional metrics and LLM-based approaches, as well as qualitative output analysis.

### Related Work

Legal text summarization has been an active area of research for several decades. Early approaches relied on statistical methods and extractive summarization techniques to select the most important sentences from legal documents. With the advent of neural network models, more sophisticated abstractive summarization methods became possible, allowing for the generation of new text that captures the essence of the original document.

In the German legal domain, several notable research efforts have focused on court decision summarization. The focus of these studies has been on official headnotes ("Leitsätze") as they are mainly extractive summaries from the judgement that are written by the judges themselves. These headnotes are typically short and concise, making them suitable for extractive summarization tasks and can in general be found verbatim in the body of the decision. However, they do not provide a comprehensive overview of the entire decision and are not intended for public communication. In contrast, press releases are designed to be more accessible to the general public and provide a broader context for the decision.

[Glaser et al. \[2021\]](#) presented the first large dataset of 100.000 German court decisions with corresponding summaries, establishing baseline models for German legal summarization. Their transformer-based approach achieved a ROUGE-1 F1 score of approximately 30.5%, demonstrating both the feasibility and challenges of the task. The complex structure of German court decisions (including sections like "Rubrum," "Tenor," and "Gründe") requires specialized preprocessing and models.

[Steffes and Rataj \[2022\]](#) focused on extracting official headnotes ("Leitsätze") from Federal Court of Justice (BGH) decisions by utilizing the argumentative structure of rulings. Their approach selected key sentences based on their argumentative roles, improving the selection of headnote sentences compared to purely

090 statistical methods.

091 For multilingual court summarization, [Rolshoven](#)  
092 [et al. \[2024\]](#) introduced the SLDS dataset (Swiss Lead-  
093 ing Decision Summarization) containing 18,000 Swiss  
094 Federal Court decisions in German, French, and Ital-  
095 ian, along with German summaries (“Regesten”). Their  
096 work on cross-lingual summarization demonstrated  
097 that fine-tuned smaller models could perform similarly  
098 to large pre-trained models in prompt mode. They eval-  
099 uated their approach using ROUGE, BLEU, METEOR,  
100 and BERTScore metrics.

101 Regarding evaluation methodologies, [Steffes et al.](#)  
102 [\[2023\]](#) explicitly showed that ROUGE is unreliable as  
103 a sole quality indicator for legal summaries since it  
104 fails to reliably assess legally relevant content. Their  
105 study demonstrated that a system might achieve high  
106 ROUGE scores while missing essential legal state-  
107 ments.

108 For more robust evaluation, [Xu and Ashley](#) [\[2023\]](#)  
109 presented a question-answering framework using  
110 LLMs to assess the factual correctness of legal sum-  
111 maries. Their approach generates understanding ques-  
112 tions about the reference text and compares answers  
113 derived from both reference and generated summaries,  
114 showing better correlation with expert judgments than  
115 simple ROUGE scores.

116 In practical applications, the ALeKS project  
117 (Anonymisierungs- und Leitsatzerstellungs-Kit) is being  
118 developed in Germany to automatically anonymize  
119 court decisions and generate headnotes using LLMs.  
120 This collaboration between judicial authorities and  
121 research institutions aims to increase the publication rate  
122 of court decisions while maintaining content accuracy  
123 and data protection standards.

124 Our work extends these efforts by specifically fo-  
125 cusing on press release generation (rather than technical  
126 headnotes) for German court decisions, emphasizing  
127 both factual correctness and accessibility for non-  
128 legal audiences. We employ a comprehensive evalua-  
129 tion framework that combines reference-based met-  
130 rics, embedding-based metrics, and factual consistency  
131 checks through both automated methods and LLM-as-  
132 judge assessments.

133 It is important to note that court press releases of-  
134 ten contain additional context not found in the original  
135 decision, such as procedural history, background infor-  
136 mation, or quotes from spokespersons. This character-  
137 istic distinguishes press releases from pure summaries  
138 and presents additional challenges for automated eval-  
139 uation of factual consistency.

## 140 CourtPressGER

### 141 Data

142 Our dataset includes court decisions and correspond-  
143 ing press releases from Germany’s highest courts (Bun-

desgerichte) as well as the federal constitutional court  
(Bundesverfassungsgericht - under german law not a  
Bundesgericht) :

- Federal Labor Law Court (Bundesarbeitsgericht -  
BAG)
- Federal Fiscal Court (Bundesfinanzhof - BFH)
- Federal Court of Justice (Bundesgerichtshof -  
BGH)
- Federal Social Court (Bundessozialgericht - BSG)
- Federal Constitutional Court (Bundesverfassungs-  
gericht - BVerfG)
- Federal Administrative Court (Bundesverwal-  
tungsgericht - BVerwG)

The cleaned dataset contains 6.4k pairs of court de-  
cisions and press releases. The average length of deci-  
sions is 10.810 BPE tokens , while press releases aver-  
age 1.402 BPE tokens. We report BPE token counts as  
used by modern LLMs rather than raw word or charac-  
ter counts for better compatibility with model context  
window considerations.

### 164 Splits

For our experiments, we divided the dataset into train-  
ing, validation, and test splits in an 72.2/11.6/16.3 ratio.  
The training set contains 4643 pairs, while the valida-  
tion set contains 744 test sets contain 1045 pairs. The  
split was done chronologically with the following year  
distribution: ((...))

We decided to split chronologically because other-  
wise the distribution shifts incurred by rotating press  
office personnel over time would not be captured in the  
data split, leading to a potential overestimation of per-  
formance on unseen data.

### 176 Descriptive Statistics

Our dataset analysis reveals variation in document  
lengths across different courts. Federal Constitutional  
Court decisions tend to be the longest with an average  
of 14.782 BPE tokens, while Federal Fiscal Court deci-  
sions average 7.379 BPE tokens. Press release lengths  
also vary, with Federal Constitutional Court releases av-  
eraging 2,230 BPE tokens and Federal Court of Justice  
releases averaging 1,620 BPE tokens. The standard de-  
viation for court decision length is 10.739 BPE tokens,  
indicating considerable variation in document size.

The distribution of press release and judgement  
length and year distribution can be seen in the follow-  
ing table:

## 190 Experimental Setup

### 191 Synthetic Prompts

For each decision-press release pair, we generated syn-  
thetic prompts through the Anthropic API (Claude Son-  
net 3.7) to serve as input for LLMs to generate press  
releases. These prompts were designed to highlight the

Court	Press Release			234	Judgments	For OpenAI models (GPT-4o), the request format uses the above template as the user message with a system message that instructs the model to act as an expert in legal texts whb7 Writes press releases based on court decisions	236	
	Mean	Std	Count	235	Mean	Std	Count	237
BAG	1056.37	407.50	177	14148.00	5913.64	1410.59	761	238
BFH	800.28	213.58	761	7378.97	4410.59	161	2407	239
BGH	1386.84	680.10	2407	8216.82	5686.26	161	2407	240
BSG	1146.66	484.69	161	11790.02	4850.29	161	1771	241
BVerfG	2039.50	1353.63	1771	14781.53	16844.62	1771	1155	242
BVerwG	942.91	336.86	1155	11734.63	8110.92	1155	10739.27	243
<b>Overall avg</b>	<b>1402.32</b>	<b>954.52</b>	–	<b>10809.58</b>	–	–	–	244

Table 1: Statistical summary of press releases and judgments by court

key aspects of the decision and to train the models to create relevant and precise press releases.

To create synthetic prompts, we utilized Claude 3.7 Sonnet with a system prompt [Appendix]

## Press Release Generation

Our pipeline includes various LLMs, which can be categorized into two groups:

1. Large Models: GPT-4o (mainstream and economical closed source model at time of experiments), Llama-3-70B (large & SotA open weights model at time of running experiments)
2. Small Models: Teuken-7B, Llama-3-8B, EuroLLM-9B, Mistral-7B (all open weights in smaller class, typical base models for research finetuning experiments)

The pipeline is designed to send the synthetic prompts to the models, collect the generated press releases, and store them alongside the actual press releases. A checkpoint system allows for the continuation of interrupted generation processes.

## Context Limitation

We found that the context window size of the models has a significant impact on their ability to generate high-quality press releases. Models with larger context windows (e.g., GPT-4o with a theoretical limit of 128k tokens, though in our implementation we used the API with a practical limit of 64k tokens) can process the entire court decision at once, while smaller models require document chunking and hierarchical summarization approaches.

For decisions that exceed the context window of a model, we implemented a hierarchical summarization approach (described in the next section) that allows the model to consider the entire document while respecting context limitations.

## Generation Prompt Template

For consistency across models, we use a standardized prompt template [Appendix]

## Hierarchical Summarization

For court decisions that exceed the context window of a model, we implemented a hierarchical summarization approach. This method involves the following steps:

1. Chunking: The court decision is divided into chunks that fit within the model’s context window.
2. Level 0 Summarization: Each chunk is independently summarized.
3. Higher Level Summarization: The summaries are combined and recursively summarized until a single summary is created.
4. Final Press Release Generation: The final summary is used as input for the press release generation.

This hierarchical approach allows smaller models to process long documents while maintaining the context and coherence of the original text. The implementation involves a recursive algorithm that estimates the number of levels needed based on the document length and the model’s context window size.

Each level of summarization uses specially designed prompts that instruct the model to focus on different aspects of the text, with higher levels emphasizing cohesion and integration of information from multiple chunks.

## FT Teuken

#todo ME

## Evaluation

Our evaluation framework was designed to address the known limitations of traditional NLP metrics for legal text summarization. As highlighted by Steffes et al. (2023), metrics like ROUGE can be unreliable as sole quality indicators because they may not adequately capture legally relevant content.

Therefore, we developed a comprehensive evaluation approach using multiple complementary metrics:

- ROUGE (Lin [2004])
- BLEU (Papineni et al. [2002])
- METEOR (Banerjee and Lavie [2005])
- BERTScore (Zhang et al. [2020])
- QAGS (Question Answering for evaluating Generated Summaries) (Wang et al. [2020])

- 287 • FactCC (Factual Consistency Check) (Kryściński  
288 et al. [2019]) 340
- 289 • LLM-as-a-Judge (evaluation using Claude 3.7  
290 Sonnet)

291 While BLEU is less commonly used for summarization tasks due to its sensitivity to word order and  
292 sentence length, we include it to maintain comparability with multilingual studies like Rolshoven et  
293 al. (2024) and to provide a more comprehensive assessment through multiple metrics.

294 This multi-faceted approach aligns with recent trends in legal summarization evaluation, which emphasize combining different automated metrics with expert judgment to assess different quality dimensions of generated legal texts.

### 300 Factual Consistency Metrics

301 Our project utilizes advanced metrics to evaluate the 341 factual consistency between court decisions and generated press releases:

- 302 • QAGS (Question Answering for evaluating Generated 342 Summaries): This metric first generates questions from the press releases, then answers these 343 questions with the court decisions as context, and finally compares the answers to verify if the press 344 release is factually correct. This approach is similar to the framework proposed by Xu & Ashley 345 (2023), which showed better correlation with expert judgments than traditional metrics.
- 303 • FactCC (Factual Consistency Check): This metric 346 extracts claims from the press releases and checks each claim for consistency with the court decision. A total score for factual consistency is calculated 347 from these checks.

318 For both QAGS and FactCC, we acknowledge a significant limitation: These metrics were originally developed and trained on English news datasets, not German legal texts. Their application to our German court 319 texts relies on the multilingual capabilities of the underlying models, but has not been specifically validated 320 for German legal text. This limitation likely affects the 321 absolute scores and may partially explain why smaller 322 German-specific models like Teuken-7B achieve factual 323 consistency scores comparable to larger models despite 324 lower performance on other metrics. The scores 325 should be interpreted as relative comparisons rather 326 than absolute measures of factual accuracy.

327 For additional context information in press releases 328 that doesn't directly appear in the court decision, these 329 metrics may incorrectly flag such information as inconsistent, leading to potentially lower scores even for 330 high-quality press releases. We address this limitation 331 partially through our LLM-as-a-Judge approach, which 332 can better distinguish between contradictory information 333 and benign additional context.

### LLM-as-a-judge

We use Claude 3.7 Sonnet to evaluate the generated press releases based on various criteria such as factual correctness, completeness, clarity, and structure. Optionally, the generated press release can be compared with the reference press release. The metric provides both numerical ratings (1-10) and detailed justifications, calculating an overall score across all evaluation criteria.

To evaluate the quality of the generated press releases, we use Claude 3.7 Sonnet with the following system prompt [Appendix]

It is important to note that our evaluation relied on LLM-as-a-Judge rather than human legal experts. While this approach provides valuable insights and scales to large datasets, it serves as a proxy for human evaluation and would benefit from validation through targeted expert reviews in future work. Claude 3.7 Sonnet was selected for this task due to its strong performance in understanding complex legal texts in multiple languages as well as its selection for synthetic prompt generation which made it a natural choice for evaluation.

## Results

Based on our evaluation, we present the results organized by evaluation type (hierarchical vs. full document processing) and model. We structured our analysis to examine reference-based metrics, embedding-based metrics, factual consistency metrics, and human-like evaluation through LLM-as-judge.

Modell	ROUGE-1 F1	ROUGE-2 F1	ROUGE-L F1	BLEU-1	BLEU-2	B
gpt 4o	0.3627	0.1452	0.1918	0.2105	0.1266	0
llama 3_3 70B	<b>0.3823</b>	<b>0.1601</b>	<b>0.1997</b>	<b>0.2248</b>	<b>0.1385</b>	0
mistral v03	0.3612	0.1561	0.1844	0.2126	0.1304	0

Table 2: Press release comparison on full judgements

Model	ROUGE-1 F1	ROUGE-2 F1	ROUGE-L F1	BLEU-1	BLEU-2	B
gpt 4o	0.3627	0.1452	0.1918	0.2105	0.1266	0
llama 3_3 70B	<b>0.3823</b>	<b>0.1601</b>	<b>0.1997</b>	<b>0.2248</b>	<b>0.1385</b>	0
mistral v03	0.3612	0.1561	0.1844	0.2126	0.1304	0

Table 3: Press release comparison on full judgements

Note that we evaluate Mistral\_v03 also on the full ruling text even though its context is limited to 32k tokens. In our experiments, 1% of documents needed to be truncated for evaluation in this narrower context.

### Reference-based Metrics

Our evaluation of reference-based metrics shows that larger models consistently outperform smaller models across all metrics.

According to our evaluation, the complete results for full-document processing (without hierarchical summarization) are:

Model	ROUGE-1 F1	ROUGE-2 F1	ROUGE-L F1	Model	BLEU-1	ROUGE-1 F1	ROUGE-2 F1	ROUGE-L F1	BLEU-1
Llama-3-70B	0.3823	0.1601	0.1997	Llama-3-70B	0.2248	0.0668	0.3746	0.1986	0.1411
GPT-4o	0.3627	0.1452	0.1918	GPT-4o	0.2105	0.0559	0.3584	0.1845	0.1242
Mistral-v0.3	0.3612	0.1561	0.1844	Mistral-v0.3	0.2126	0.0660	0.3571	0.1901	0.1218
				Llama-3-8B	0.2927			0.0780	0.1344
				EuroLLM		0.2800		0.0611	0.1199
				Teuken-7B		0.1630		0.0213	0.0703
									0.0794

Table 4: Reference-based metrics for full-document processing

Model	ROUGE-1 F1	ROUGE-2 F1	ROUGE-L F1	Model	BLEU-1	BLEU-4	METEOR
Llama-3-70B	0.3823	0.1601	0.1997	Llama-3-70B	0.2248	0.0668	0.1986
GPT-4o	0.3627	0.1452	0.1918	GPT-4o	0.2105	0.0559	0.1845
Mistral-v0.3	0.3612	0.1561	0.1844	Mistral-v0.3	0.2126	0.0660	0.1901
				Llama-3-8B	0.2927	0.0780	0.1344
				EuroLLM		0.2800	0.1199
				Teuken-7B		0.1630	0.0703
							0.0794

Table 5: Reference-based metrics for full-document processing

For hierarchical summarization, the performance is slightly lower but follows a similar pattern:

These results are consistent with findings from Glaser et al. (2021), who reported ROUGE-1 scores of around 30.5% for their best models on German court decision summarization. Our best models exceed this performance slightly, which may be attributed to the advancement in LLMs since their study.

### Embedding-based Metrics

BERTScore metrics, which capture semantic similarity using contextual embeddings, show similar trends to the reference-based metrics.

This metric is particularly relevant for legal text evaluation as noted by recent surveys, which highlight BERTScore’s ability to detect semantic similarity beyond simple n-gram matching:

For hierarchical summarization:

### Factual Metrics

The QAGS evaluation, which measures factual consistency through question answering (similar to the approach proposed by Xu & Ashley), shows varying degrees of factual accuracy:

For hierarchical summarization:

FactCC scores, which directly evaluate the factual consistency of claims:

For hierarchical summarization, Teuken-7B achieved a FactCC score of 0.5051 with a consistency ratio of 0.5068, comparable to larger models despite its lower performance on other metrics. This surprising result likely reflects limitations in applying FactCC to German texts rather than true parity in factual consistency, as our LLM-as-Judge evaluation shows significant differences in factual correctness scores.

Table 6: Reference-based metrics for hierarchical summarization

Model	ROUGE-1 F1	ROUGE-2 F1	ROUGE-L F1	Model	BLEU-1	BLEU-4	METEOR
Llama-3-70B	0.3823	0.1601	0.1997	Llama-3-70B	0.2248	0.0668	0.1986
GPT-4o	0.3627	0.1452	0.1918	GPT-4o	0.2105	0.0559	0.1845
Mistral-v0.3	0.3612	0.1561	0.1844	Mistral-v0.3	0.2126	0.0660	0.1901
				Llama-3-8B	0.2927	0.0780	0.1344
				EuroLLM		0.2800	0.1199
				Teuken-7B		0.1630	0.0703
							0.0794

Table 7: Reference-based metrics for hierarchical summarization

### LLM-as-a-Judge

The LLM-as-a-Judge evaluation using Claude 3.7 Sonnet provides a more nuanced assessment of the generated press releases, addressing the dimensions of quality emphasized in legal summarization research:

For hierarchical summarization:

These results demonstrate that while larger models generally produce press releases that are more factually correct, complete, clear, and well-structured, the hierarchical summarization approach allows smaller models to produce reasonably good summaries, particularly in terms of clarity and structure. Interestingly, the improvement from hierarchical summarisation to full summarisation is marginal for the largest models.

## Discussion

tbd

## Conclusions

Our comprehensive evaluation of the CourtPressGER system demonstrates that modern LLMs can effectively generate German court press releases, with performance varying according to model size and architecture.

Key findings include:

1. Model size matters: Larger models consistently outperform smaller models across all evaluation metrics.
2. Hierarchical summarization is effective: Our hierarchical approach enables smaller models to process long documents while maintaining reasonable quality.

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405  
406  
407  
408  
409  
410  
411  
412

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

Model	BERTScore Precision	BERTScore Recall	BERTScore F1	BERTScore Precision	BERTScore Recall	BERTScore F1
Llama-3-70B	0.7889	0.7508	Llama-3-70B	0.7918	0.7557	0.7730
GPT-4o	0.7746	0.7396	GPT-4o	0.7563	0.7595	0.7711
Mistral-v0.3	0.7706	0.7255	Mistral-v0.3	0.7918	0.7645	0.7777
			Llama-3-8B	0.7519	0.7239	0.7373
			EuroLLM	0.7570	0.7362	0.7459
			Teuken-7B	0.6966	0.6303	0.6600

Table 8: Embedding-based metrics for full-document processing

Model	BERTScore Precision	BERTScore Recall	BERTScore F1
Llama-3-70B	0.7889	0.7508	0.7691
GPT-4o	0.7746	0.7396	0.7563
Mistral-v0.3	0.7706	0.7255	0.7465

Table 9: Embedding-based metrics for full-document processing

3. Factual consistency challenges: Even the best models struggle with perfect factual consistency, indicating room for improvement. 481  
 446  
 447  
 448  
 449
4. Language-specific models: German-specific 482  
 450 models like EuroLLM show competitive performance 483  
 451 for their size compared to larger multilingual 484  
 452 models. 485

While our fine-tuned Teuken model showed some improvement over the base version, it still performs significantly below larger models, suggesting that parameter count remains a decisive factor for this complex task.

Our work provides a contribution to the emerging field of automated legal text summarization in the German language, extending the work of Glaser et al. (2021), Steffes & Rataj (2022), and Rolshoven et al. (2024). The multidimensional evaluation approach we employed addresses the limitations of traditional metrics highlighted by Steffes et al. (2023) and incorporates newer evaluation methods like question-answering based assessment proposed by Xu & Ashley (2023).

Our system has potential practical applications similar to the ALeKS project currently under development in Germany, which aims to automate the generation of court decision headnotes. While ALeKS focuses on technical headnotes, our work specifically addresses press releases that need to be accessible to non-legal audiences.

## Limitations

We acknowledge several limitations of our approach:

1. Evaluation metrics: Our use of QAGS and FactCC metrics, which were developed and validated on English datasets, introduces uncertainty when applied to German legal texts. Future work should explore German-specific factual consistency metrics.
2. LLM-as-judge vs. human evaluation: While our

Table 10: Embedding-based metrics for hierarchical summarization

Model	QAGS Score	Question Count
Llama-3-70B	0.2898	4.87
GPT-4o	0.2777	4.75
Mistral-v0.3	0.3252	4.72

Table 11: QAGS evaluation for full-document processing

443

444

LLM-based evaluation provides valuable insights, it serves as a proxy for human expert evaluation and would benefit from validation through targeted expert reviews.

3. Additional context in press releases: Court press releases often contain contextual information not present in the original decision, which can confound factual consistency metrics.
4. Divergence from Rolshoven et al. findings: Unlike Rolshoven et al. (2024), who found that finetuned smaller models could approach the performance of larger models, our results show a clear advantage for larger models. This difference may be attributed to our focus on press releases rather than technical summaries (“Regesten”), the different nature of our dataset, or the specific characteristics of German federal court decisions.

The CourtPressGER project demonstrates the potential of LLMs to assist in making legal information more accessible to the public while highlighting the ongoing challenges in maintaining factual accuracy when summarizing complex legal documents.

## Appendix

### Ethics Statement

tbd

### References

### Prompts

We used the following prompts for our experiments:

### Synthetic prompt generation

We used the following prompt for synthetic prompt generation:

Du bist ein Experte für juristische Texte und Kommunikation. Deine Aufgabe ist es, ein Gerichtsurteil und die dazugehörige Pressemitteilung zu analysieren

Model	QAGS Score	Question Count
Llama-3-70B	0.2863	4.94
GPT-4o	0.2637	4.78
Mistral-v0.3	0.2386	4.69
Llama-3-8B	0.2289	4.90
EuroLLM	0.1875	4.84
Teuken-7B	0.1607	4.94

Table 12: QAGS evaluation for hierarchical summarization

Model	FactCC Score	Consistency Ratio	Claim Count
Llama-3-70B	0.5082	0.5144	5.0
GPT-4o	0.4991	0.5068	5.0
Mistral-v0.3	0.5021	0.5044	4.96

Table 13: FactCC scores for full-document processing

und dann herauszufinden, welcher Prompt verwendet worden sein könnte, um diese Pressemitteilung aus dem Gerichtsurteil zu generieren, wenn man ihm einem LLM gegeben hätte.

515

1. Analysiere, wie die Pressemitteilung Informationen aus dem Urteil vereinfacht, umstrukturiert und Schlüsselinformationen hervorhebt
2. Berücksichtige den Ton, die Struktur und den DetAILierungsgrad der Pressemitteilung
3. Identifiziere, welche Anweisungen nötig wären, um den juristischen Text in diese Pressemitteilung zu transformieren

527 Erkläre NICHT deine Überlegungen und füge  
528 KEINE Meta-Kommentare hinzu. Gib NUR den tat-  
529 sächlichen Prompt aus, der die Pressemitteilung aus  
530 dem Gerichtsurteil generieren würde. Sei spezifisch  
531 und detailliert in deinem synthetisierten Prompt.

532 Hier ist das originale Gerichtsurteil:

533 {court\_ruling}

534 Und hier ist die Pressemitteilung, die daraus erstellt  
535 wurde:

536 {press\_release}

537 Erstelle einen detaillierten Prompt, der einem LLM  
538 gegeben werden könnte, um die obige Pressemit-  
539 teilung aus dem Gerichtsurteil zu generieren. Schreibe  
540 NUR den Prompt selbst, ohne Erklärungen oder Meta-  
541 Kommentare.

### 542 Press release generation

543 We used the following prompt for press release genera-  
544 tion:

545 {prompt} Gerichtsurteil: {ruling}

### 546 LLM-as-a-judge

547 We used the following prompt for LLM-as-a-judge  
548 evaluation:

Model	Factual Correctness	Completeness	Clarity	Structure
Llama-3-70B	8.17	6.87	8.63	8.16
GPT-4o	8.39	7.16	8.82	8.54
Mistral-v0.3	6.96	5.71	7.14	6.81

Table 14: LLM-as-a-Judge evaluation for full-document processing

Model	Factual Correctness	Completeness	Clarity	Structure
Llama-3-70B	7.34	6.36	8.15	7.62
GPT-4o	8.11	7.09	8.75	8.41
Mistral-v0.3	5.54	4.97	5.56	5.24
Llama-3-8B	5.28	4.54	6.31	6.43
EuroLLM	4.97	4.43	6.40	6.69
Teuken-7B	3.06	2.16	4.24	4.41

Table 15: LLM-as-a-Judge evaluation for hierarchical summarization

515

You 516 are an expert in legal texts and evaluate the quality  
of 517 press releases for court decisions. Rate the generated  
518 press release according to the following criteria  
on a scale of 1-10:

1. Factual Correctness: How accurately does the press release reflect the facts from the court decision?
2. Completeness: Have all important information from the decision been included in the press release?
3. Clarity: How understandable is the press release for a non-legal audience?
4. Structure: How well is the press release structured and organized?
5. Comparison with Reference: How good is the generated press release compared to the reference press release?

For each criterion, provide a numerical value between 1 and 10 and a brief justification. Finally, calculate an overall score as the average of all individual values. Provide your answer in the following JSON format: { "faktische\_korrektheit": { "wert": X, "begründung": "..."}, "vollständigkeit": { "wert": X, "begründung": "..."}, "klarheit": { "wert": X, "begründung": "..."}, "struktur": { "wert": X, "begründung": "..."}, "vergleich\_mit\_referenz": { "wert": X, "begründung": "..."}, "gesamtscore": X.X }

The user prompt contains: Court Decision [court\_decision] Generated Press Release [generated\_press\_release] Reference Press Release [reference\_press\_release]

## 580 References

Satanjeev Banerjee and Alon Lavie. 2005. METEOR:  
An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments.

- 588 Ingo Glaser, Sebastian Moser, and Florian Matthes. 584  
589 2021. [Summarization of German Court Rulings](#). In 585  
Proceedings of the Natural Legal Language Processing 586  
Workshop 2021, pages 180–189, Punta Cana, 587  
Dominican Republic. Association for Computational  
Linguistics.
- 590 Wojciech Kryściński, Bryan McCann, Caiming Xiong, 591  
and Richard Socher. 2019. [Evaluating the Factual 592  
Consistency of Abstractive Text Summarization](#). 593  
*Preprint*, arXiv:1910.12840.
- 594 Chin-Yew Lin. 2004. ROUGE: A Package for Auto- 595  
matic Evaluation of Summaries.
- 596 Kishore Papineni, Salim Roukos, Todd Ward, and Wei- 597  
Jing Zhu. 2002. [BLEU: A method for automatic 598  
evaluation of machine translation](#). In Proceedings 599  
of the 40th Annual Meeting on Association for Com- 600  
putational Linguistics - ACL '02, page 311, Philadel- 601  
phia, Pennsylvania. Association for Computational 602  
Linguistics.
- 603 Luca Rolshoven, Vishvaksenan Rasiah, Srinanda Brüg- 604  
ger Bose, Matthias Stürmer, and Joel Niklaus. 605  
2024. [Unlocking Legal Knowledge: A Multilingual 606  
Dataset for Judicial Summarization in Switzerland](#). 607  
*Preprint*, arXiv:2410.13456.
- 608 Bianca Steffes and Piotr Rataj. 2022. [Legal Text 609  
Summarization Using Argumentative Structures](#). In 610  
Enrico Francesconi, Georg Borges, and Christoph 611  
Sorge, editors, *Frontiers in Artificial Intelligence 612  
and Applications*. IOS Press.
- 613 Bianca Steffes, Piotr Rataj, Luise Burger, and Lukas 614  
Roth. 2023. [On evaluating legal summaries with 615  
ROUGE](#). In *Proceedings of the Nineteenth Interna- 616  
tional Conference on Artificial Intelligence and Law*, 617  
pages 457–461, Braga Portugal. ACM.
- 618 Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020. 619  
[Asking and Answering Questions to Evaluate the 620  
Factual Consistency of Summaries](#). *Preprint*, 621  
arXiv:2004.04228.
- 622 Huihui Xu and Kevin Ashley. 2023. [Question- 623  
Answering Approach to Evaluating Legal Sum- 624  
maries](#).
- 625 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. 626  
Weinberger, and Yoav Artzi. 2020. [BERTScore: 627  
Evaluating Text Generation with BERT](#). *Preprint*, 628  
arXiv:1904.09675.