CoCoLex: Confidence-guided Copy-based Decoding for Grounded Legal Text Generation

Santosh T.Y.S.S^{1*}, Youssef Tarek Elkhayat¹, Oana Ichim², Pranav Shetty³, Dongsheng Wang³, Zhiqiang Ma³, Armineh Nourbakhsh³, Xiaomo Liu³

¹School of Computation, Information, and Technology, Technical University of Munich, Germany
²Graduate Institute of International and Development Studies, Geneva, Switzerland
³JPMorgan AI Research

Abstract

Due to their ability to process long and complex contexts, LLMs can offer key benefits to the Legal domain, but their adoption has been hindered by their tendency to generate unfaithful, ungrounded, or hallucinatory outputs. While Retrieval-Augmented Generation offers a promising solution by grounding generations in external knowledge, it offers no guarantee that the provided context will be effectively integrated. To address this, context-aware decoding strategies have been proposed to amplify the influence of relevant context, but they usually do not explicitly enforce faithfulness to the context. In this work, we introduce Confidenceguided Copy-based Decoding for Legal Text Generation (CoCoLex)-a decoding strategy that dynamically interpolates the model produced vocabulary distribution with a distribution derived based on copying from the context. CoCoLex encourages direct copying based on models' confidence, ensuring greater fidelity to the source. Experimental results on five legal benchmarks demonstrate that CoCoLex outperforms existing context-aware decoding methods, particularly in long-form generation tasks.

1 Introduction

The legal domain poses unique challenges to document-grounded language generation. Legal documents are often long, structurally complex, and prone to jargon and technical language. Additionally, tasks that are grounded in legal documents such as question answering or analysis often have strict requirements regarding accuracy and faithfulness to sources. In recent years, LLMs have revolutionized the legal domain, transforming areas such as legal education (Choi and Schwarcz, 2023; Jiang et al., 2024), research (Livermore et al., 2023), compliance checking (Hassani, 2024), and even legal practice (Rodgers et al., 2023).

Despite their potential, their adoption is hindered by their tendency to produce hallucinations—text inconsistent with authoritative sources such as case law, statutes, regulations, contracts, and doctrines (Magesh et al., 2024). In law, where strict adherence to authoritative sources is essential, unfaithful outputs can result in inaccurate and even harmful advice (Chitgopkar et al., 2024). These issues stem from the model's inability to fully ground its outputs in verifiable knowledge, leaving them prone to generating plausible-sounding yet non-factual content (El Hamdani et al., 2024; Zmigrod et al., 2024; Santosh et al., 2024).

address these challenges, То Retrieval-Augmented Generation (RAG) has emerged as a promising approach, equipping LLMs with external knowledge chunks to ground their outputs (Lewis et al., 2020; Borgeaud et al., 2022; Guu et al., 2020). However, RAG with regular decoding still struggles to effectively integrate retrieved context (Hagström et al., 2024), leading to outputs that diverge from the provided context—a critical shortcoming in high-stakes legal applications (Magesh et al., 2024). While some methods aim to enhance context grounding through improved pre-training (Guu et al., 2020; Lewis et al., 2020; Borgeaud et al., 2022; Izacard et al., 2023), due to the high cost of pre-training, some studies focus on inference-stage methods, such as decoding strategies for white-box models (Shi et al., 2023; Kim et al., 2024b,a; Zhao et al., 2024) or prompting for black-box models (Zhou et al., 2023; Byerly and Khashabi, 2024; Li et al., 2024). Most of these works prioritize correctness but rarely evaluate whether responses are faithful to the provided context. Therefore, we propose a novel decoding strategy to improve context faithfulness of white-box LLMs for legal text generation.

Legal texts often adhere to templatized structures

^{*}Santosh T.Y.S.S., along with at least one of authors at JPMorgan AI Research, conceived the idea and made substantive contributions detailed in the paper while Santosh was an intern at JPMorgan Chase.

(Ghosh et al., 2023; Nair and Modani, 2023) and verbatim phrases to maintain interpretive precision (Rossi et al., 2021). Inspired by this characteristic, we propose Confidence-guided Copy-based Decoding for Legal Text Generation (CoCoLex), which prioritizes fidelity of the generation by explicitly guiding the model to copy tokens from the context. Specifically, CoCoLex uses a confidence score to dynamically balance copying and generating tokens, interpolating the model's token distribution with a copy-based distribution. Note that despite its conceptual similarity to pointer generator networks (See et al., 2017), CoCoLex is training-free, operating directly on logits during decoding. Its interpolation mechanism draws inspiration from kNN-LM (Khandelwal et al., 2019), which retrieves from external data stores of the pre-training corpus or training examples. Prior works with KNN-LM such as Khandelwal et al. (2019) and Wang et al. (2023) have focused on language modeling and open-ended text generation, CoCoLex enhances faithfulness in context-aware generation, an underexplored area.

Our key contributions are given below: (a) We introduce CoCoLex, a novel decoding strategy that enhances the faithfulness of generated text by encouraging copying from the context. (b) Through experiments conducted on five legal text generation datasets using two LLMs, we demonstrate that our method, CoCoLex, not only improves correctness and ensures that generated outputs remain faithful to the source but also maintains fluency and coherence, all without increasing inference overhead, especially in long-form text generation. (c) We present an extension, CoCoLex+, which allows for copying from the entire document rather than being restricted to top-retrieved chunks. This further enhances performance by leveraging a richer and more comprehensive context. (d) We show that our method can be integrated with previous approaches, indicating that its improvements are complementary to other existing methods.

2 Related Work

LLMs and Risks for legal AI. Lawyers are increasingly adopting AI tools to enhance their practice, from drafting contracts to conducting legal research, yielding substantial efficiency gains. As of January 2024, 41 of the top 100 largest U.S. law firms and 35% of a broader sample of 384 firms report using generative AI tools (Henry, 2024; Col-

lens et al., 2024). In the UK, 14% of lawyers surveyed use these tools weekly or more (Greenhill, 2024). Despite the evident benefits, Legal AI introduces significant ethical challenges, including concerns over client confidentiality, data protection, potential biases, and the critical responsibility of lawyers to supervise and ensure the accuracy of AI-generated outputs (Avery et al., 2023; Harasta et al., 2024; Chitgopkar et al., 2024).

Dahl et al. (2024) provide a systematic assessment of LLMs for legal tasks, categorizing hallucination types in their responses, while Magesh et al. (2024) examine RAG-based tools for legal QA. Their findings indicate that although retrievalaugmented models reduce hallucinations compared to general-purpose LLMs, the issue persists. In this work, we propose a novel decoding algorithm within a retrieval-augmented generation framework to mitigate hallucinations. Particularly, our approach enhances the groundedness of model outputs by improving faithfulness to the retrieved context while ensuring relevance to legal queries.

Retrieval Augmented Generation. RAG retrieves relevant external information, enabling LLMs to provide factual responses. The retrieval component in RAG can operate at different granularities, such as chunk-level (Guu et al., 2020; Lewis et al., 2020), token-level (Khandelwal et al., 2019; Yogatama et al., 2021), and entity-level (De Jong et al., 2021; Févry et al., 2020). The retrieved information then can be integrated at three possible levels: (a) At the input layer, the retrieved segments can be combined with the input or query, and processed them jointly through the model (Ram et al., 2023; Izacard et al., 2023). (b) Intermediate-layer integration uses semi-parametric modules to incorporate retrieved information into the internal layers of the model (Wu et al., 2022; Borgeaud et al., 2022). (c) Output-layer integration merges the retrieval and generation results after processing (Khandelwal et al., 2019; Santosh et al., 2024c).

Existing Retrieval-Augmented Generation (RAG) approaches can be categorized into training-free and training-based methods. Training-based methods involve fine-tuning both the retriever and generator, either independently (Karpukhin et al., 2020; Zhou et al., 2022), sequentially (Borgeaud et al., 2022; Yoran et al., 2023; Lin et al., 2023; Sarto et al., 2022), or jointly (Rubin and Berant, 2024; Izacard et al., 2023), allowing these components to work synergistically. In contrast, training-free methods utilize retrieved

knowledge during inference by incorporating it into prompts (Jiang et al., 2023; Khattab et al., 2022; Trivedi et al., 2022). These methods are computationally efficient and therefore we focus on this common variant of RAG, that performs retrieval based on the query and integrates the retrieved context within the input prompt.

However, simply augmenting the input with context may not consistently ensure alignment between the generated output and the retrieved context. To address this, decoding-based strategies (Shi et al., 2023; Kim et al., 2024b,a; Zhao et al., 2024) and prompt-based techniques (Zhou et al., 2023; Byerly and Khashabi, 2024; Li et al., 2024) are often employed to guide the model towards generating outputs that remain faithful to the provided context. Existing decoding strategies, such as contrastive decoding (Li et al., 2022), aim to amplify the influence of the retrieved context by adjusting token logits but do not explicitly enforce faithfulness to the source. To overcome this limitation, we propose a novel decoding method that incorporates a token-level copying mechanism from the retrieved context, guided by the model's confidence, in our approach.

3 CoCoLex

Given an input query x and context c, a language model with parameters θ is prompted to generate a response $\mathbf{y} = \{y_1, y_2, \dots, y_n\}$ of length n. The response is generated autoregressively, with each token y_t sampled from the conditional probability distribution: $y_t \sim p_{\theta}(y_t \mid c, x, y_{< t})$. Our decoding strategy, CoCoLex consists of two major components: copy-based decoding, which derives distribution over vocabulary tokens based on copying from context and confidence factor which guides model between copying and generating. We explain these two components below.

3.1 Copy-based Decoding

We hypothesize that while augmenting the query with relevant context improves response quality, it does not always guarantee that the responses are grounded in the provided context. To address this, we introduce a copy mechanism that explicitly directs the model to copy tokens from the context during the decoding process, thereby maintaining contextual fidelity, while ensuring the response remains fluent and relevant to the query.

During decoding, we extract and store the hidden

state representations of all tokens within the context c from the language model. These representations, denoted as h_i for each token, are stored along with their corresponding next token. Since the hidden states are computed during autoregressive generation, no additional forward passes are required for this storing step. At each decoding step t, the model generates a hidden state vector h_t , which is then compared to the stored context vectors $\{h_i\}$. The similarity between h_t and each context vector h_i is computed using the Euclidean L_2 distance which is subsequently transformed into similarity scores via an exponential decay (See Appendix C for further details):

$$s_t(i) = \exp\left(-\operatorname{dist}_t(i)\right)$$

where $dist_t(i)$ represents the Euclidean distance between h_t and h_i . These similarity scores are aggregated to form a probability distribution over the vocabulary. For each token v in the vocabulary, the probability of selecting v is proportional to the sum of similarity scores associated with the context tokens mapped to v. The probability distribution for the copy mechanism is given by:

$$p_{\text{copy}}(y_t = v \mid c, x) = \frac{\sum_{i \in \text{tokens}(v)} s_t(i)}{\sum_{v' \in V} \sum_{i \in \text{tokens}(v')} s_t(i)},$$

where tokens(v) represents the context tokens that map to the token v, and V is the vocabulary.

To reduce computational overhead, we limit the aggregation to the top-k most similar context vectors. This approximation is justified because tokens with low similarity scores contribute negligibly to the overall probability distribution. Thus, we compute copy distribution normalizing over the restricted subset of the vocabulary V_k , formed by the top-k tokens,

$$p_{\text{copy}}(y_t = v \mid c, x) \propto \sum_{i \in \text{top-}k \cap \text{tokens}(v)} s_t(i).$$

3.2 Confidence-Guidance

Uncertainty serves as a critical metric for determining when to reliably trust the predictions of LLMs. Several prior studies have explored the use of uncertainty to detect hallucinations, highlighting its potential to identify unfaithful content (Vashurin et al., 2024; Fadeeva et al., 2023; Duan et al., 2023; Kadavath et al., 2022; Huang et al., 2023). Advanced LLMs are expected to assign low probabilities to tokens that are likely to introduce inaccuracies or hallucinations (Kadavath et al., 2022). Building on this insight, we use a logit-based uncertainty approach to define a confidence indicator. Leveraging this indicator, our method dynamically balances reliance on the model's predictions and an external copy mechanism.

At each decoding step t, the model produces a probability distribution over the vocabulary. The entropy of this distribution serves as a measure of uncertainty, capturing the dispersion of probabilities. Higher entropy values indicate greater uncertainty, while lower entropy values suggest that the model is more confident in its predictions:

$$H_t = -\sum_{v \in V} p_{\theta}(y_t = v) \cdot \log p_{\theta}(y_t = v)$$

We normalize the entropy by dividing it by the maximum possible entropy, $\log(|V|)$, which corresponds to a uniform distribution:

$$H_t^{\text{norm}} = \frac{H_t}{\log(|V|)}$$

The confidence score, λ_t , is derived using an exponential transformation of the normalized entropy. This transformation ensures that lower entropy values correspond to higher confidence scores:

$$\lambda_t = \exp(-H_t^{\text{norm}})$$

To prevent erratic behavior caused by sudden spikes or drops in uncertainty, we smooth the confidence scores by incorporating historical data. Specifically, the value is calculated using a smoothing factor that combines its current value with a running average of values from a specified past window, which mitigates any fluctuations and provides a more stable confidence metric. The confidence score interpolates between the model's predictions and the copy mechanism. A high confidence score increases reliance on the model's predictions, while a low score shifts preference to the copy mechanism. The final probability distribution at decoding step t is given by:

$$p(y_t \mid c, x, y_{< t}) = \lambda_t \cdot p_{\theta}(y_t \mid c, x, y_{< t})$$
$$+ (1 - \lambda_t) \cdot p_{\text{copy}}(y_t \mid c, x, y_{< t})$$

4 **Experiments**

4.1 Metrics

We evaluate the generated answers against the reference answers for **correctness** using both lexical and semantic similarity metrics. For lexical similarity, we compute the ROUGE-L F1 score (Lin, 2004). For semantic similarity, we use AlignScore (Zha et al., 2023), which computes the alignment score for each sentence in the generated answer against the reference answer and is aggregated to derive the overall score. For **faithfulness**, we assess the alignment between the generated answer and the provided source context using AlignScore (Zha et al., 2023). We also evaluate stylistic properties such as **fluency** and **coherence** using UniEval (Zhong et al., 2022), which assesses the quality of individual sentences and whether all the sentences collectively form a connected narrative.

4.2 Datasets

We experiment with the following datasets: CUAD (Hendrycks et al., 2021) frames contract information extraction as a question-answering task, pairing contracts with relevant questions and extractive answer spans. OALQA (Butler, 2023) is a Question answering dataset based on the Open Australian Legal Corpus, where each question is paired with the respective reference document, along with the reference answer. ObliQA (Gokhan et al., 2024) is a regulatory QA dataset from Abu Dhabi Global Markets financial regulations. Questions are paired with extractive clauses answering them within the whole document corpus. AQuAECHR is a legal question-answering dataset based on European Court of Human Rights (ECHR) judgments. Given a legal query and a corpus of ECHR judgements, the system must generate an answer following a retrieve-then-generate paradigm. CLERC (Hou et al., 2024) focuses on legal analysis generation for U.S. federal case documents. The model generates text continuations containing legal reasoning along with citations to relevant cases based on preceding case document content, which typically introduces case facts.

Unlike CUAD and OALQA, where each question is linked to a specific reference document, datasets like AQuAECHR, CLERC, and ObliQA require retrieval across entire document corpora. This retrieval setup is particularly challenging in legal contexts due to the ambiguity of legal queries and the interpretive complexity of legal precedents (Dworkin, 1986), in contrast to traditional information-seeking queries (Kwiatkowski et al., 2019; Rajpurkar, 2016) which usually have clear, unambiguous references. Effective retrieval systems in legal contexts must integrate both textual and non-textual factors, such as jurisdiction, time period, and specific conditions, to ensure the relevance and authority of retrieved content (Santosh et al., 2024b). As our primary focus is on evaluating generative models' ability to produce faithfully grounded answers, we simulate these tasks within Oracle Documents-relevant documents obtained from reference answers (available as citations to documents in datasets like CLERC and AQuAECHR or verbatim extractive snippets in ObliQA), rather than across entire corpora. Then, a retrieval step is performed within these documents to extract relevant paragraphs or chunks, which are then provided as context for the model to synthesize answers that remain contextually grounded and relevant to the query. Detailed statistics of these datasets are provided in Table 6 of the Appendix.

4.3 **Baselines and Implementation Details**

We compare CoCoLex to the following baselines: 1) Regular Decoding. 2) CAD (Context-aware Decoding) (Shi et al., 2023), which enhances groundedness by sampling from a contrastive output distribution that amplifies the difference between output probabilities with and without context. CAD modifies the model's original output distribution by incorporating the pointwise mutual information (PMI) between the context c and the generation y_t , conditioned on $x, y_{<t}$:

$$\begin{aligned} y_t \sim \operatorname{softmax}\left[(1+\alpha)\operatorname{logit}_{\theta}(y_t \mid c, x, y_{< t}) \right. \\ \left. -\alpha\operatorname{logit}_{\theta}(y_t \mid x, y_{< t})\right]. \end{aligned}$$

A larger α places more weight on the adjustment, while $\alpha = 0$ reduces to regular decoding. 3) Ada-CAD (Adaptive CAD) (Wang et al., 2024) dynamically infers α in CAD at every timestamp to get α_t based on the degree of conflict, measured by the Jensen-Shannon divergence (JSD) between the distributions representing contextual and parametric knowledge:

$$\alpha_t = \text{JSD}\left(p_\theta(y_t \mid x, y_{< t}) \parallel p_\theta(y_t \mid c, x, y_{< t})\right).$$

We also derive CoLex, from CoCoLex, removing confidence-based dynamic interpolation and interpolating based on a static value to study the effect of the confidence guidance.

We apply these methods to mistralai/ Mistral-7B-Instruct-v0.3 (Jiang and Yang, 2023) and Equall/Saul-7B-Instruct-v1 (Colombo et al., 2024), with Saul being specifically pre-trained on a legal unsupervised corpus and instruction corpus, leveraging the base Mistral model. We use BM25 (Robertson et al., 2009) for retrieving top-k passages from the documents due to its strong performance in legal retrieval settings (Santosh et al., 2024b; Rosa et al., 2021). Implementation details are provided in App. B.

4.4 Results

We present the results on the five datasets in Table 1. We observe that context-aware decoding methods such as CAD and AdaCAD improve both correctness and faithfulness over regular decoding on the CUAD dataset across both models. However, they exhibit a decline in fluency and coherence scores. On OALQA, these methods improve correctness (as measured by AlignScore) only with the Mistral model and achieve comparable faithfulness scores across both models. For other datasets, including ObliQA, AQuAECHR, and CLERC, CAD and AdaCAD lead to a decline in both correctness and faithfulness across both models, although they show some improvement in lexical-based ROUGE scores. These results underscore that CAD and AdaCAD are primarily effective for short-text generation (as seen in CUAD) but struggle with longform generation, particularly in maintaining fluency and coherence. Among them, AdaCAD performs slightly better, especially in long-range tasks.

Our proposed CoLex outperforms prior approaches by guiding the model to explicitly copy tokens from the context. This approach not only enhances faithfulness by aligning generations more closely with the provided context but also improves correctness. The improvement is more pronounced in long-range tasks. The copying mechanism does not limit CoLex's fluency and coherence, which are comparable to regular decoding, with only a marginal decline in some datasets. CoCoLex further enhances correctness and faithfulness by dynamically balancing copying and text generation based on the model's confidence. It also improves fluency and coherence compared to CoLex. Across both models, we observe that the legally pre-trained Saul underperforms the generalist Mistral model, consistent with prior findings (Santosh et al., 2024d). This is mainly due to Saul's difficulty in synthesizing information from the provided context when following instructions. However, our proposed methods substantially improve Saul's performance, mitigating its limitations through copybased decoding, which explicitly guides the model to copy relevant tokens. Overall, our approach

		Mistral-	7B-Instruc	et-v0.3			Saul-	7B-Instruc	t-v1	
	Cor-RL	Cor-AS	Fth-AS	Flu	Coh	Cor-RL	Cor-AS	Fth-AS	Flu	Coh
					CUAD					
Regular	54.29	68.24	76.31	82.14*	62.87	21.79	38.70	73.16	77.40	70.47
CAD	54.57	69.57	79.55	80.41	59.95	23.42	40.62	74.65	76.30	69.16
AdaCAD	54.54	69.63	79.56	79.32	61.43	23.35	40.60	74.55	76.82	69.63
CoLex	55.29	70.65	80.66	81.22	62.27	23.87	45.63	82.06	76.87	69.41
CoCoLex	55.77*	71.06*	80.96*	81.79	62.73*	25.04*	49.63*	84.84*	77.23	70.91*
				(DALQA					
Regular	40.53	41.39	59.85	79.78	84.62	40.16	32.26	52.84	79.32	81.94
CAD	39.90	42.90	59.00	78.32	79.60	39.57	31.74	52.89	76.75	73.90
AdaCAD	39.96	42.49	59.44	77.91	79.25	39.61	31.68	52.78	76.50	73.74
CoLex	46.50	48.61	60.14	79.66	86.44	50.10	50.41	57.23	78.36	84.19
CoCoLex	48.34*	49.84*	60.87*	79.94*	87.24*	50.91*	52.74*	59.19*	80.04*	86.02*
				(ObliQA					
Regular	33.86	73.35	90.84	75.38*	70.22	16.90	62.50	83.10	72.16	69.95
CAD	35.83	71.14	89.73	70.24	63.13	15.74	61.49	82.00	62.69	53.88
AdaCAD	35.72	71.04	89.61	69.73	62.98	15.64	61.09	81.87	62.05	53.57
CoLex	43.41	85.35	93.48	73.10	72.37	22.98	82.13	90.15	71.16	67.64
CoCoLex	45.12*	86.01*	95.96*	74.41	74.74*	23.50*	83.44*	91.05*	72.56*	70.18*
				AÇ	UAECHR	ł				
Regular	21.77	52.79	89.66	74.71	79.99	17.87	48.30	80.68	73.82	68.66*
CAD	22.04	49.15	89.28	72.28	69.43	18.71	40.60	80.17	59.82	61.44
AdaCAD	22.13	48.69	89.37	71.27	69.06	18.68	40.68	80.58	59.41	62.35
CoLex	29.12	59.79	91.85	80.76	86.01	28.97	65.89	90.25	72.75	66.84
CoCoLex	29.84*	60.10*	92.27*	80.48	86.44*	29.28*	66.26*	91.15*	74.24*	67.91
				(CLERC					
Regular	10.42	42.38	74.02	77.27	78.41	9.39	23.40	55.91	73.76	65.24
CAD	10.56	34.98	66.35	74.37	73.09	9.06	24.10	56.04	71.09	65.17
AdaCAD	10.52	35.11	66.46	74.62	73.13	9.09	24.18	56.02	71.17	65.15
CoLex	12.71	54.94	78.62	78.16*	89.55	12.78*	33.02	59.58	76.36	71.01
CoCoLex	12.88*	58.12*	79.54*	77.92	90.55*	12.51	34.95*	62.79*	76.25	73.93*

Table 1: Performance comparison of different decoding-based methods across five legal text generation datasets, using two language models. Cor, Fth, RL, AS, Flu, Coh denote Correctness, Faithfulness, ROUGE-L, AlignScore, Fluency and Coherence respectively. Entries marked with * are statistically significantly higher than the second-best performing baseline at the 95% confidence level, according to the Wilcoxon signed-rank test.

enhances faithfulness and correctness without compromising fluency and coherence.

	Cor	Fth	Flu	Coh
Regular	4.40	4.24	4.88	4.88
AdaCAD	4.24	3.84	4.80	4.84
CoCoLex	4.64	4.44	4.96	4.92

Table 2: Human evaluation results on 25 randomly sampled questions from the AQuAECHR dataset.

Human Evaluation. We randomly sample 25 questions from the AQuAECHR dataset and generate responses using three methods—Regular, AdaCAD,

and CoCoLex—applied to the Mistral-7B model, yielding a total of 75 responses. Each response is assessed by a legal ECHR expert, across four criteria: correctness (relevance to the question and alignment with the reference answer), faithfulness (adherence to the provided passages in context), fluency, and coherence. Evaluations are conducted on a 5-point Likert scale, where 1 represents the lowest quality and 5 the highest. Table 2 presents the average scores for each criterion across methods. Notably, the legal expert consistently ranked CoCoLex higher than Regular and AdaCAD, reinforcing its superiority in generating legally faithful responses. Additionally, the results highlight that AdaCAD underperforms compared to Regular in this task of long-range text generation, particularly struggling with faithfulness. A detailed case study is provided in Appendix D

4.5 Discussion and Analysis

4.5.1 Providing Document Context

Retrieval-augmented approaches for handling lengthy documents typically follow a two-step process: first, retrieving relevant evidence passages or chunks from these documents and then using these retrieved passages as context for the generator to generate an answer. However, this approach often suffers from chunking-related issues, where improper segmentation and concatenation of retrieved passages disrupt the semantics, leading to incomplete and incoherent information retrieval (Qian et al., 2024; Dong et al., 2023). This, in turn, makes it difficult for the model to stay aligned with the main query, eventually degrading the accuracy of the generated response.

To address this limitation, we introduce Co-CoLex+, an extension of CoCoLex that leverages the entire document's encoded hidden states rather than restricting copying only the tokens in topretrieved passages. By incorporating representations from the full document, CoCoLex+ enables the model to capture a more comprehensive understanding of the document's content, leading to more contextually grounded responses. We chunk documents into overlapping segments to efficiently obtain these hidden state representations and extract contextualized hidden states for each token. To prevent redundancy across overlapping contexts, each token is assigned a single hidden state representation, taken from the chunk where it has the most autoregressive context. During inference, we continue to limit explicit textual context to the top-k retrieved passages, similar to CoCoLex, but augment it with hidden states from the full document to facilitate copying.

We evaluate CoCoLex+ on long-range generation tasks—CLERC and AQuAECHR—and present the results in Table 3. Given that CoCoLex+ enables copying from the entire document, we extend our faithfulness evaluation to measure alignment with the full document, in addition to the topretrieved passages, as done previously. Our results show that CoCoLex+ consistently outperforms Co-CoLex in correctness across both datasets and mod-

	C-R	C-A	F-P	F-D	Fl	Co
	CLERO	C - Mistr	al-7B-In	struct-v().3	
CoCoLex	12.88	58.12	79.54	89.34	77.92	90.55
CocoLex+	13.01	60.66	80.17	90.12	77.84	90.52
	CLE	RC - Saı	ıl-7B-In	struct-v1		
CoCoLex	12.51	34.95	62.79	62.71	76.25	73.93
CocoLex+	13.33	44.91	72.51	70.31	75.53	75.16
А	QuAEC	HR - Mi	stral-7B	-Instruct	-v0.3	
CoCoLex	29.84	60.10	92.27	61.35	80.48	86.44
CocoLex+	30.06	60.37	92.62	62.71	80.26	86.28
	AQuAE	ECHR - S	Saul-7B-	Instruct-	·v1	
CoCoLex	29.28	66.26	91.15	69.39	74.24	67.91
CocoLex+	29.63	67.45	88.39	71.53	74.29	67.42

Table 3: Performance comparison between CoCoLex and CoCoLex+, where CoCoLex+ enhances the copy mechanism to apply to all tokens in the document, rather than being limited to tokens in the retrieved chunks appended to the prompt. C-R(A), F-P(D), Fl, Co denote Correctness-ROUGE-L(AlignScore), Faithfulness-Passages (Documents), Fluency and Coherence.

els. Interestingly, it also enhances faithfulness with respect to retrieved passages on CLERC for both models and on AQuAECHR for Mistral, suggesting that additional global context improves the model's ability to stay grounded in the provided evidence. However, in AQuAECHR with Saul, we observe a slight decline in passage-level faithfulness and an increase in document-level faithfulness, indicating that Saul benefits more from whole-document copying than from restricting itself to retrieved tokens. Crucially, CoCoLex+ maintains fluency and coherence comparable to CoCoLex, despite handling a broader copying vocabulary. These findings highlight that expanding the model's access to fulldocument representations strengthens its ability to generate accurate, well-grounded responses, making CoCoLex+ particularly effective for long-form legal text generation.

4.5.2 Combining Strategies

While AdaCAD grounds generation by contrastively amplifying the difference between output probabilities with and without context—down weighting prior parametric knowledge when relevant contextual information is available, CoCoLex grounds generation by guiding the model to copy from retrieved context effectively. Given their complementary mechanisms, we extend CoCoLex by incorporating AdaCAD's contrastive probability distribution into CoCoLex's final probability com-

	C-RL	C-AS	F-AS	Flu	Coh
CU	CUAD - Mistral-7B-Instruct-v0.3				
CoCo	55.77	71.06	80.96	81.79	62.73
Ada + CoCo	56.26	71.87	81.18	81.99	62.41
C	CUAD - S	Saul-7B-	Instruct-	·v1	
СоСо	25.04	49.63	84.84	77.23	70.91
Ada + CoCo	25.38	50.46	85.05	77.18	70.89
CLH	ERC - M	istral-7B	-Instruct	t-v0.3	
СоСо	12.88	58.12	79.54	77.92	90.55
Ada + CoCo	12.10	50.24	72.09	75.51	86.51
С	LERC -	Saul-7B	-Instruct	-v1	
СоСо	12.51	34.95	62.79	76.25	73.93
Ada + CoCo	13.33	36.75	63.30	75.09	73.45
AQuAECHR - Mistral-7B-Instruct-v0.3					
СоСо	29.84	60.10	92.27	80.48	86.44
Ada + CoCo	27.74	55.89	90.70	80.19	81.84
AQuAECHR - Saul-7B-Instruct-v1					
СоСо	29.28	66.26	91.15	74.24	67.91
Ada + CoCo	26.30	57.79	90.10	69.15	63.71

Table 4: Demonstrating the complementarity of both approaches: CoCoLex combined with AdaCAD improves when AdaCAD outperforms regular decoding, but results in a decrease when AdaCAD is less effective.

putation, modifying p_{θ} directly.

We evaluate the combined AdaCAD + CoCoLex approach on CUAD, CLERC, and AQuAECHR across both models, with results in Table 3. Our findings indicate that this combination enhances performance in CUAD across both models and improves CLERC for Saul. However, we observe a decline in CLERC for Mistral and in AQuAECHR for both models. Notably, these trends closely align with AdaCAD's performance in Table 1: whenever AdaCAD improves over regular decoding (e.g., CUAD for both models, CLERC for Saul), the combined approach also yields gains. Conversely, when AdaCAD underperforms compared to regular decoding, the combination similarly results in a performance drop. This suggests that integrating AdaCAD with CoCoLex can be effective when AdaCAD itself outperforms regular decoding, reinforcing their complementarity.

4.5.3 Inference Time

We compute inference time for different methods using a randomly sampled 10% subset of the CUAD and AQuAECHR datasets, running on an Nvidia A100 GPU. To account for variations in

	CUAD	AQuAECHR
Regular	1.00x	1.00x
CAD	1.75x	1.71x
AdaCAD	1.77x	1.72x
CoLex	1.49x	1.61x
CoCoLex	1.51x	1.62x
Ada + CoCo	2.31x	2.25x
CoCoLex+	1.96x	2.96x

Table 5: Inference time comparison of different approaches, scaled to regular decoding.

output length across methods, we normalize inference time by the number of tokens generated and report relative inference times compared to regular decoding (1.00x). Our findings show that CAD and AdaCAD nearly double inference time, requiring an additional decoding step at each timestep to obtain logits with and without context. In contrast, CoLex and CoCoLex, while avoiding the extra decoding step of CAD/AdaCAD, incur a smaller overhead, which is the result of: indexing hidden states of tokens in the retrieved context, retrieving, normalizing, and interpolating token probabilities with copy-based probabilities at each timestep. When combining AdaCAD and CoCoLex, these overheads accumulate linearly since their processes are non-overlapping, further increasing inference time. Additionally, CoCoLex+, with its mechanism to index hidden states from the entire document, introduces further latency-especially for longer documents in AQuAECHR, where the overhead is more pronounced than CUAD.

5 Conclusion

We introduced CoCoLex, a decoding strategy that enhances faithfulness in legal text generation by dynamically balancing the model's token distribution with a copy-based distribution derived from retrieved context, guided by model confidence. Experiments on five legal datasets highlight the limitations of existing context-aware decoding methods, which improve faithfulness in short-range but struggle with long-range tasks. CoCoLex enhances generation fidelity and correctness in long-form tasks while maintaining fluency and coherence without substantial inference overhead. Our findings show that leveraging hidden state representations from the entire document, rather than restricting copying to top-retrieved chunks, mitigates context length constraints, enabling the model to capture richer information and improve faithfulness. We also find

that integrating contrastive-based approaches such as AdaCAD with CoCoLex is beneficial, reinforcing their complementarity. Future work could extend copying beyond the token level to incorporate larger semantic units, such as phrases or clauses, for improved contextual alignment.

Limitations

To evaluate the faithfulness of generator models, our experiments assume an oracle document setting, where document retrieval is performed beforehand and provided as input. However, in realworld applications requiring retrieval across an entire document corpus-such as in benchmarks like CLERC, AQuAECHR, and ObliQA-the quality and correctness of generated text inherently depend on the accuracy and completeness of retrieval. If critical legal information is missing, poorly ranked, or incorrectly retrieved, the model may still produce misleading outputs. Addressing this retrieval bottleneck remains a key challenge, particularly in legal contexts, where case law documents, regulations, and statutes can be lengthy, nuanced, and difficult to rank effectively (Santosh et al., 2024b; Locke and Zuccon, 2022). Unlike standard information retrieval tasks that prioritize semantic similarity, legal retrieval must account for additional constraints such as precedential value, temporal relevance (Santosh et al., 2024e) and procedural applicability (Santosh et al., 2024a).

While CoCoLex improves long-form generation fidelity by encouraging direct copying, it does not explicitly handle cases requiring reasoning beyond the retrieved context, such as synthesizing multiple sources, reconciling conflicting precedents, or constructing persuasive legal arguments. Moreover, as CoCoLex relies on hidden state similarities to guide copying, its effectiveness depends on the model's ability to learn robust token representations, which may vary across architectures, pretraining objectives, and the pretraining corpus.

Ethics Statement

All datasets used in this work are publicly available and have been utilized in compliance with their respective data usage policies. While datasets such as CLERC and AQuAECHR, which involve case law judgments, are not anonymized, our work engages with the data in a manner that we believe does not cause harm beyond the availability of this information. LLMs, due to the historical biases inherent in their pre-training data, may perpetuate harmful prejudices and inaccuracies, potentially exacerbating existing gaps in legal knowledge and representation. This can lead to biased outputs, including factual inaccuracies and misrepresentations of legal citations. As such, caution is necessary in the responsible deployment of LLMs for legal informationseeking tasks. While LLMs can be powerful tools, they are not intended to replace legal professionals but to assist them by augmenting their expertise. Therefore, it is critical to employ LLMs with care in legal contexts, with constant monitoring for fairness, accuracy, and alignment with legal principles.

In the responsible deployment of LLMs, there is also the need to address broader questions surrounding the automation of legal tasks. The growing use of LLMs in legal practice should be carefully assessed to ensure that the integrity, professionalism, and accountability of the legal profession are maintained. As these technologies evolve, it is important to continually evaluate the impact of LLMs on the legal system, ensuring they complement legal professionals and enhance the equitable and effective delivery of legal services. This includes ongoing reflection on their potential biases, transparency, and the ethical implications of their integration into the legal workflow.

Dislaimer

This paper was prepared for informational purposes by the Artificial Intelligence Research group of JP-Morgan Chase Co. and its affiliates ("JP Morgan") and is not a product of the Research Department of JP Morgan. JP Morgan makes no representation and warranty whatsoever and disclaims all liability, for the completeness, accuracy, or reliability of the information contained herein. This document is not intended as investment research or investment advice, or a recommendation, offer or solicitation for the purchase or sale of any security, financial instrument, financial product, or service, or to be used in any way for evaluating the merits of participating in any transaction, and shall not constitute a solicitation under any jurisdiction or to any person, if such solicitation under such jurisdiction or to such person would be unlawful.

References

Joseph J Avery, Patricia Sanchez Abril, and Alissa del Riego. 2023. Chatgpt, esq.: Recasting unauthorized practice of law in the era of generative ai. *Yale JL & Tech.*, 26:64.

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR.

Umar Butler. 2023. Open australian legal qa.

- Adam Byerly and Daniel Khashabi. 2024. How effective is self-consistency for long-context problems? *arXiv preprint arXiv:2411.01101.*
- Sudhan Chitgopkar, Noah Dohrmann, Stephanie Monson, Jimmy Mendez, Finale Doshi-Velez, and Weiwei Pan. 2024. Accuracy isn't everything: Understanding the desiderata of ai tools in legal-financial settings. In *NeurIPS 2024 Workshop on Behavioral Machine Learning*.
- Jonathan H Choi and Daniel Schwarcz. 2023. Ai assistance in legal analysis: An empirical study. *Available at SSRN 4539836*.
- Jack Collens, Rachel Reimer, Gerald Schifman, and Pamela Wilkinson. 2024. Ai survey: Where artificial intelligence stands in the legal industry.
- Pierre Colombo, Telmo Pessoa Pires, Malik Boudiaf, Dominic Culver, Rui Melo, Caio Corro, Andre FT Martins, Fabrizio Esposito, Vera Lúcia Raposo, Sofia Morgado, et al. 2024. Saullm-7b: A pioneering large language model for law. arXiv preprint arXiv:2403.03883.
- Matthew Dahl, Varun Magesh, Mirac Suzgun, and Daniel E Ho. 2024. Large legal fictions: Profiling legal hallucinations in large language models. *Journal of Legal Analysis*, 16(1):64–93.
- Michiel De Jong, Yury Zemlyanskiy, Nicholas FitzGerald, Fei Sha, and William Cohen. 2021. Mention memory: incorporating textual knowledge into transformers through entity mention attention. *arXiv preprint arXiv:2110.06176*.
- Zican Dong, Tianyi Tang, Lunyi Li, and Wayne Xin Zhao. 2023. A survey on long text modeling with transformers. *arXiv preprint arXiv:2302.14502*.
- Jinhao Duan, Hao Cheng, Shiqi Wang, Chenan Wang, Alex Zavalny, Renjing Xu, Bhavya Kailkhura, and Kaidi Xu. 2023. Shifting attention to relevance: Towards the uncertainty estimation of large language models. *arXiv preprint arXiv:2307.01379*.
- Ronald Dworkin. 1986. *Law's Empire*. Harvard University Press, Cambridge, MA.
- Rajaa El Hamdani, Thomas Bonald, Fragkiskos D Malliaros, Nils Holzenberger, and Fabian Suchanek. 2024. The factuality of large language models in the

legal domain. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, pages 3741–3746.

- Ekaterina Fadeeva, Roman Vashurin, Akim Tsvigun, Artem Vazhentsev, Sergey Petrakov, Kirill Fedyanin, Daniil Vasilev, Elizaveta Goncharova, Alexander Panchenko, Maxim Panov, et al. 2023. Lmpolygraph: Uncertainty estimation for language models. *arXiv preprint arXiv:2311.07383*.
- Thibault Févry, Livio Baldini Soares, Nicholas FitzGerald, Eunsol Choi, and Tom Kwiatkowski. 2020. Entities as experts: Sparse memory access with entity supervision. *arXiv preprint arXiv:2004.07202*.
- Sreyan Ghosh, Chandra Kiran Evuru, Sonal Kumar, S Ramaneswaran, S Sakshi, Utkarsh Tyagi, and Dinesh Manocha. 2023. Dale: Generative data augmentation for low-resource legal nlp. *arXiv preprint arXiv:2310.15799*.
- Tuba Gokhan, Kexin Wang, Iryna Gurevych, and Ted Briscoe. 2024. Regnlp in action: Facilitating compliance through automated information retrieval and answer generation. *arXiv preprint arXiv:2409.05677*.
- Stuart Greenhill. 2024. Lawyers cross into the new era of generative ai.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR.
- Lovisa Hagström, Sara Vera Marjanović, Haeun Yu, Arnav Arora, Christina Lioma, Maria Maistro, Pepa Atanasova, and Isabelle Augenstein. 2024. A reality check on context utilisation for retrieval-augmented generation. *arXiv preprint arXiv:2412.17031*.
- Jakub Harasta, Tereza Novotná, and Jaromir Savelka. 2024. It cannot be right if it was written by ai: on lawyers' preferences of documents perceived as authored by an llm vs a human. *Artificial Intelligence and Law*, pages 1–38.
- Shabnam Hassani. 2024. Enhancing legal compliance and regulation analysis with large language models. *arXiv preprint arXiv:2404.17522*.
- Dan Hendrycks, Collin Burns, Anya Chen, and Spencer Ball. 2021. Cuad: An expert-annotated nlp dataset for legal contract review. *arXiv preprint arXiv:2103.06268*.
- Justin Henry. 2024. We asked every am law 100 law firm how they're using gen ai. here's what we learned. *The American Lawyer*.
- Abe Bohan Hou, Orion Weller, Guanghui Qin, Eugene Yang, Dawn Lawrie, Nils Holzenberger, Andrew Blair-Stanek, and Benjamin Van Durme. 2024. Clerc: A dataset for legal case retrieval and retrievalaugmented analysis generation. *arXiv preprint arXiv:2406.17186*.

- Yuheng Huang, Jiayang Song, Zhijie Wang, Shengming Zhao, Huaming Chen, Felix Juefei-Xu, and Lei Ma. 2023. Look before you leap: An exploratory study of uncertainty measurement for large language models. arXiv preprint arXiv:2307.10236.
- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2023. Atlas: Few-shot learning with retrieval augmented language models. *Journal of Machine Learning Research*, 24(251):1–43.
- Cong Jiang and Xiaolei Yang. 2023. Legal syllogism prompting: Teaching large language models for legal judgment prediction. In *Proceedings of the Nineteenth International Conference on Artificial Intelligence and Law*, pages 417–421.
- Hang Jiang, Xiajie Zhang, Robert Mahari, Daniel Kessler, Eric Ma, Tal August, Irene Li, Alex'Sandy' Pentland, Yoon Kim, Deb Roy, et al. 2024. Leveraging large language models for learning complex legal concepts through storytelling. *arXiv preprint arXiv:2402.17019*.
- Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation. *arXiv preprint arXiv:2305.06983*.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data*, 7(3):535–547.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. 2022. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*.
- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2019. Generalization through memorization: Nearest neighbor language models. *arXiv preprint arXiv:1911.00172*.
- Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts, and Matei Zaharia. 2022. Demonstrate-searchpredict: Composing retrieval and language models for knowledge-intensive nlp. *arXiv preprint arXiv:2212.14024*.
- Taehyeon Kim, Joonkee Kim, Gihun Lee, and Se-Young Yun. 2024a. Instructive decoding: Instruction-tuned large language models are self-refiner from noisy instructions. In *The Twelfth International Conference* on Learning Representations.

- Youna Kim, Hyuhng Joon Kim, Cheonbok Park, Choonghyun Park, Hyunsoo Cho, Junyeob Kim, Kang Min Yoo, Sang-goo Lee, and Taeuk Kim. 2024b. Adaptive contrastive decoding in retrievalaugmented generation for handling noisy contexts. *arXiv preprint arXiv:2408.01084*.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453– 466.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Taiji Li, Zhi Li, and Yin Zhang. 2024. Improving faithfulness of large language models in summarization via sliding generation and self-consistency. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 8804– 8817.
- Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2022. Contrastive decoding: Open-ended text generation as optimization. *arXiv preprint arXiv:2210.15097*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Xi Victoria Lin, Xilun Chen, Mingda Chen, Weijia Shi, Maria Lomeli, Rich James, Pedro Rodriguez, Jacob Kahn, Gergely Szilvasy, Mike Lewis, et al. 2023. Ra-dit: Retrieval-augmented dual instruction tuning. *arXiv preprint arXiv:2310.01352.*
- Michael A Livermore, Felix Herron, and Daniel Rockmore. 2023. Language model interpretability and empirical legal studies. *Virginia Public Law and Legal Theory Research Paper*, (2023-69).
- Daniel Locke and Guido Zuccon. 2022. Case law retrieval: problems, methods, challenges and evaluations in the last 20 years. *arXiv preprint arXiv:2202.07209.*
- Varun Magesh, Faiz Surani, Matthew Dahl, Mirac Suzgun, Christopher D Manning, and Daniel E Ho. 2024. Hallucination-free? assessing the reliability of leading ai legal research tools. *arXiv preprint arXiv:2405.20362*.
- Inderjeet Nair and Natwar Modani. 2023. Exploiting language characteristics for legal domain-specific language model pretraining. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 2516–2526.

- Hongjin Qian, Zheng Liu, Kelong Mao, Yujia Zhou, and Zhicheng Dou. 2024. Grounding language model with chunking-free in-context retrieval. *arXiv* preprint arXiv:2402.09760.
- P Rajpurkar. 2016. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*.
- Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language models. *Transactions of the Association for Computational Linguistics*, 11:1316–1331.
- Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends*® *in Information Retrieval*, 3(4):333–389.
- Ian Rodgers, John Armour, and Mari Sako. 2023. How technology is (or is not) transforming law firms. *Annual Review of Law and Social Science*, 19(1):299–317.
- Guilherme Moraes Rosa, Ruan Chaves Rodrigues, Roberto Lotufo, and Rodrigo Nogueira. 2021. Yes, bm25 is a strong baseline for legal case retrieval. *arXiv preprint arXiv:2105.05686.*
- Julien Rossi, Svitlana Vakulenko, and Evangelos Kanoulas. 2021. Verbcl: A dataset of verbatim quotes for highlight extraction in case law. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pages 4554–4563.
- Ohad Rubin and Jonathan Berant. 2024. Retrievalpretrained transformer: Long-range language modeling with self-retrieval. *Transactions of the Association for Computational Linguistics*, 12:1197–1213.
- TYS Santosh, Kevin D Ashley, Katie Atkinson, and Matthias Grabmair. 2024a. Towards supporting legal argumentation with nlp: Is more data really all you need? *arXiv preprint arXiv:2406.10974*.
- TYS Santosh, Rashid Gustav Haddad, and Matthias Grabmair. 2024b. Ecthr-pcr: A dataset for precedent understanding and prior case retrieval in the european court of human rights. *arXiv preprint arXiv:2404.00596*.
- TYS Santosh, Hassan Sarwat, Ahmed Abdou, and Matthias Grabmair. 2024c. Mind your neighbours: Leveraging analogous instances for rhetorical role labeling for legal documents. *arXiv preprint arXiv:2404.01344*.
- T.Y.S.S Santosh, Weidinger Korbinian, Ichim Oana, and Grabmair Matthias. 2024d. Aquaechr: Attributed question answering for european court of human rights. *arXiv preprint*.

- TYSS Santosh, Tuan-Quang Vuong, and Matthias Grabmair. 2024e. Chronoslex: Time-aware incremental training for temporal generalization of legal classification tasks. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3022–3039.
- Sara Sarto, Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. 2022. Retrieval-augmented transformer for image captioning. In *Proceedings of the 19th international conference on content-based multimedia indexing*, pages 1–7.
- Abigail See, Peter J Liu, and Christopher D Manning. 2017. Get to the point: Summarization with pointer-generator networks. *arXiv preprint arXiv:1704.04368*.
- Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and Scott Wen-tau Yih. 2023. Trusting your evidence: Hallucinate less with context-aware decoding. *arXiv preprint arXiv:2305.14739*.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. *arXiv preprint arXiv:2212.10509*.
- Roman Vashurin, Ekaterina Fadeeva, Artem Vazhentsev, Lyudmila Rvanova, Akim Tsvigun, Daniil Vasilev, Rui Xing, Abdelrahman Boda Sadallah, Kirill Grishchenkov, Sergey Petrakov, et al. 2024. Benchmarking uncertainty quantification methods for large language models with lm-polygraph. *arXiv preprint arXiv:2406.15627*.
- Han Wang, Archiki Prasad, Elias Stengel-Eskin, and Mohit Bansal. 2024. Adacad: Adaptively decoding to balance conflicts between contextual and parametric knowledge. *arXiv preprint arXiv:2409.07394*.
- Shufan Wang, Yixiao Song, Andrew Drozdov, Aparna Garimella, Varun Manjunatha, and Mohit Iyyer. 2023. Knn-Im does not improve open-ended text generation. *arXiv preprint arXiv:2305.14625*.
- Yuhuai Wu, Markus N Rabe, DeLesley Hutchins, and Christian Szegedy. 2022. Memorizing transformers. *arXiv preprint arXiv:2203.08913*.
- Dani Yogatama, Cyprien de Masson d'Autume, and Lingpeng Kong. 2021. Adaptive semiparametric language models. *Transactions of the Association for Computational Linguistics*, 9:362–373.
- Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. 2023. Making retrieval-augmented language models robust to irrelevant context. *arXiv preprint arXiv:2310.01558*.
- Yuheng Zha, Yichi Yang, Ruichen Li, and Zhiting Hu. 2023. Alignscore: Evaluating factual consistency with a unified alignment function. *arXiv preprint arXiv:2305.16739*.

- Zheng Zhao, Emilio Monti, Jens Lehmann, and Haytham Assem. 2024. Enhancing contextual understanding in large language models through contrastive decoding. *arXiv preprint arXiv:2405.02750*.
- Ming Zhong, Yang Liu, Da Yin, Yuning Mao, Yizhu Jiao, Pengfei Liu, Chenguang Zhu, Heng Ji, and Jiawei Han. 2022. Towards a unified multidimensional evaluator for text generation. *arXiv* preprint arXiv:2210.07197.
- Shuyan Zhou, Uri Alon, Frank F Xu, Zhiruo Wang, Zhengbao Jiang, and Graham Neubig. 2022. Docprompting: Generating code by retrieving the docs. *arXiv preprint arXiv:2207.05987*.
- Wenxuan Zhou, Sheng Zhang, Hoifung Poon, and Muhao Chen. 2023. Context-faithful prompting for large language models. *arXiv preprint arXiv:2303.11315*.
- Ran Zmigrod, Pranav Shetty, Mathieu Sibue, Zhiqiang Ma, Armineh Nourbakhsh, Xiaomo Liu, and Manuela Veloso. 2024. "What is the value of {templates}?" Rethinking document information extraction datasets for LLMs. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 13162–13185, Miami, Florida, USA. Association for Computational Linguistics.

A Dataset

For computing metrics, in the case of CUAD and ObliQA, where the reference answers are extractive phrases or sentences from the dataset, we concatenate them to form the reference answer. In CUAD, which also contains questions where answers are unavailable in the contract, we focus only on instances with answers in the contract for faithfulness evaluation. Detailed stats on the dataset is provided in Table 6.

B Implementation Details

We use greedy decoding with a Repetition Penalty of 1.5, for all our experiments. For CAD, we use a static value of α as 0.5. Following (Wang et al., 2024), we clamp the α value to a minimum of 0.3 in AdaCAD. In CoCoLex and CoLex, we use the last layer to extract the hidden states for the tokens. We use a smoothing parameter of 0.5 between the prior and the current λ value. We use 0.5 as λ in case of CoLex. We clamp the value of λ in [0.2,0.8] for Co-CoLex. For efficient computation, we leverage the FAISS (Johnson et al., 2019) library, designed for fast nearest-neighbor retrieval in high-dimensional spaces. We retrieve the top 3 passages from the oracle documents using BM25 as context for all datasets except for ObliQA and CUAD where we used 10.

C Ablation Study

Number of Passages in the context We vary the number of top passages retrieved using BM25 from the oracle documents and append them as context. Table 7 reports the performance of Regular, AdaCAD, and CoCoLex using the Mistral-7B-Instruct-v0.3 model on the CLERC dataset. As the number of retrieved passages increases, more distractors are introduced, and performance scores consistently decline across all methods. Notably, Regular decoding experiences a sharp drop in correctness and faithfulness, highlighting its brittleness to distractors. While AdaCAD is more robust than Regular due to its contrastive decoding mechanism, which enhances reliance on context, it also becomes susceptible to distractors. In contrast, CoCoLex leverages confidence and similarity values to assess the relevance of retrieved passages, allowing it to selectively incorporate useful information based on confidence, eventually leading to filtering out distractions. This enables CoCoLex to maintain superior robustness against distractors.

Similarity Function We compare different similarity functions used in CoCoLex, to compute the similarity between the current hidden state and the hidden states of context tokens. Table 8 reports the performance on the CLERC dataset using the Mistral-7B-Instruct-v0.3 model. We observe that Euclidean distance consistently outperforms cosine similarity across all metrics.

Layer for hidden state representation We analyze the impact of selecting different layers for extracting hidden state representations of the current token and context tokens in CoCoLex. Table 8 reports the performance on the CLERC dataset using the Mistral-7B-Instruct-v0.3 model across different layers. Results indicate that using the last layer is more effective, as it provides a refined and contextually informed representation of the token, incorporating the full extent of model reasoning. In contrast, earlier layers primarily encode lower-level features and intermediate transformations, which may not fully capture the semantic and contextual nuances necessary for robust similarity computation. As a result, they are less effective in guiding accurate token retrieval for copying from the source context in CoCoLex.

Dataset	#Instances	# Docs	Docs Length	Answer length
CUAD	4,182	1.00	7,125.60	31.11
OALQA	2,024	1.00	5,973.99	90.07
ObliQA	2,786	1.32	48,138.35	117.55
AQuAECHR	1,116	3.11	14,041.33	193.58
CLERC	1,000	2.79	4,687.04	187.53

Table 6: Detailed Statistics of datasets. Docs represent oracle reference documents associated with each question. Length is obtained in terms of the number of words. We report the mean for # Docs, Docs length and Answer length.

	#Psg	Cor-AS	Fth-AS	Flu	Coh
Regular	3	42.38	74.02	77.27	78.41
	6	40.78	69.22	76.84	76.17
	10	35.70	66.12	75.02	75.05
AdaCAD	3	35.11	66.46	74.62	73.13
	6	31.03	60.77	73.03	71.23
	10	29.01	57.50	71.37	70.91
CoCoLex	3	58.12	79.54	77.92	90.55
	6	57.46	78.48	76.88	90.01
	10	56.95	77.80	76.18	89.31

Table 7: Performance comparison of Regular, Ada-CAD, and CoCoLex on the CLERC dataset using the Mistral-7B-Instruct-v0.3 model with varying numbers of passages as context. Cor, Fth, AS, Flu, and Coh denote Correctness, Faithfulness, AlignScore, Fluency, and Coherence, respectively.

Similarity	Cor-AS	Fth-AS	Flu	Coh
Cosine	55.16	77.26	76.42	87.28
Euclidean	58.12	79.54	77.92	90.55

Table 8: Performance comparison of Co-CoLex on the CLERC dataset using the Mistral-7B-Instruct-v0.3 model with different distance metrics for computing similarity.

D Case Study

We provide Case Study 1 in Table 10 and responses provided by three decoding methods—Regular, AdaCAD, and CoCoLex—to a legal question derived from the AQuAECHR dataset. The question at hand revolves around how the Court determines whether restrictions on the right to marry are properly regulated and subjected to judicial review and how this impacts an applicant's ability to challenge the denial of authorization.

Among the methods, CoCoLex stands as the most faithful to the source document by directly incorporating key legal terminology and phrases

Layer	Cor-AS	Fth-AS	Flu	Coh
-1	58.12	79.54	77.92	90.55
-5	54.67	79.01	72.18	86.12
-10	52.19	78.45	72.87	85.19
-15	49.79	78.81	70.29	82.26
-20	47.78	77.18	68.18	83.19
-25	46.92	77.29	65.28	79.68

Table 9: Performance comparison of Co-CoLex on the CLERC dataset using the Mistral-7B-Instruct-v0.3 model with different layers for obtaining hidden state representation for computing similarity.

such as "arbitrary or disproportionate interference", "necessity test", and "pressing social need." It maintains the original text's structure and legal nuances, ensuring an accurate replication of the Court's reasoning and intent, including the distinction that Article 12 does not require "necessity" or "pressing social need" tests, making it the most complete response in terms of correctness.

On the other hand, Regular provides a structured response that correctly identifies the Court's approach, emphasizing the need for restrictions to be grounded in concrete facts and national law. However, it does not explicitly mention the lack of "necessity" or "pressing social need" tests under Article 12, which is an important nuance. While this answer stays faithful to the source in terms of reasoning, it tends to oversimplify legal concepts, losing important subtleties and omitting essential legal terms that could distort the original meaning. The Regular answer is still well-organized and easy to follow, presenting the Court's reasoning in two key principles, but it generalizes the legal reasoning rather than explicitly referencing the differentiation between Article 12 and Article 8.

AdaCAD captures the limitation imposed by national laws but is less precise in explaining the Court's review process. The sentence structure makes it unclear how the judicial review process operates in practice. It deviates by implying that refusal could be seen as additional punishment without proper justification but does not clarify that the Court focuses on arbitrariness and proportionality rather than "necessity" or "social need.". It struggles with readability due to awkward phrasing and missing articles, which result in some unclear sentences. Though it attempts to reflect the source, its clarity is hindered by these issues.

We present Case Study 2 in Table 11, sourced from the AQuAECHR dataset, which explores the circumstances under which detention following conviction can be considered unlawful.

CoCoLex is the most faithful to the source document. It integrates key legal nuances, such as the specific conditions under Article 5 § 1 (c) of the Convention and the relationship between expert reports and the detention process. This response accurately reflects the structure and content of the legal text, ensuring a faithful representation of the original case law. It is comprehensive, addressing both the requirement of a causal connection to the initial conviction and the protection against arbitrariness under Article 5. The phrasing is formal yet accessible, offering a clear exposition of the legal conditions under which detention following conviction may be deemed unlawful. Its consistent use of legal terminology enhances both the fluency and clarity of the answer. The response is logically structured, presenting the necessary legal principles in a sequence that mirrors the original case law, starting with the causal connection to the initial conviction, followed by the relevance of expert reports and protection against arbitrariness. Each condition under which detention is unlawful is clearly linked, making the response coherent and well-organized.

Regular is relatively faithful but somewhat less precise. It correctly identifies the necessity of a causal connection with the initial conviction and emphasizes the protection against arbitrariness as required by Article 5. However, it oversimplifies the discussion of expert reports and the specific timeframe for lawful detention, omitting critical details such as the requirement for expert reports to be sufficiently recent and the need for expert reports in the context of a review. The response is coherent but lacks the tight organization found in CoCoLex. The sequence of ideas could be more logically presented, particularly regarding the role of expert reports and the protections under Article 5. While the phrasing is generally clear, it could benefit from more nuanced transitions to reflect the complexity of the legal principles discussed. Additionally, although the structure is grammatically sound and easy to follow, some of the phrasing feels formulaic, particularly in the transitions between legal concepts.

AdaCAD includes the core principles but lacks the depth and legal specificity present in CoCoLex. While it accurately mentions the need for a causal connection and references the issue of an unsuitable institution, it does not clarify that detention may be deemed unlawful due to insufficient expert reports and omits the protection against arbitrariness under Article 5. The response is less comprehensive in conveying the full scope of the legal reasoning outlined in the source passages. Ada-CAD struggles with coherence, as the flow of ideas is interrupted by awkward phrasing. For instance, the transition between the causal connection and the discussion of the suitability of the institution could be smoother. Additionally, some constructions are grammatically clunky, such as "detention takes place in an institution that does not suit mental disorders present with the individual being held," which would benefit from clearer phrasing. These issues hinder the fluency and coherence of the response, making it more difficult to read.

Overall, maintaining fidelity to the original text is essential, especially in the legal domain, where even minor changes in wording can significantly impact the interpretation of rights and responsibilities. A model that deviates from legal fidelity risks presenting misleading information, which could distort legal reasoning or misrepresent the Court's stance. This is especially problematic in judicial contexts, where clarity, accuracy, and legal consistency are critical. CoCoLex's approach ensures that the model's output adheres to established legal principles, preserving the integrity of the Court's decisions and offering the most reliable method for generating legally faithful responses.

Question	In the context of restrictions on the right to marry, how does the Court determine whether the restrictions are properly regulated and subjected to judicial review, and how does this impact the applicant's ability to challenge the denial of authorization?
Source Passages from 001- 96453	 80: Article 12, they added, did not guarantee an unlimited right to marry since this right was regulated by "the national laws" governing its exercise. Consequently, as the Court had held in the case of B. and L. v. the (no. 36536/02, judgment of 13 September 2005), this right was subject to limitations, although they could not restrict or reduce the right to such an extent that its very essence was impaired. 84: It was natural that, as the Court had held in many cases concerning the rights of prisoners, any measure depriving a person of liberty inevitably entailed limitations on the exercise of Convention rights, including a measure of control on prisoners' contacts with the outside world. In the context of the right to marry this might mean that the authorities, in exercising their power in this area, could monitor the wedding ceremony and limit, for instance, the number of participants. However, they should maintain a fair balance between the demands of security in prison and the prisoner's right to marry. Their discretion should be limited to, and their decisions based on, concrete facts, not on prejudice. A refusal should be restricted to situations where the marriage ceremony would jeopardise prison security – and not just be difficult to organise. Moreover, rules regulating the authorities' discretion should be laid down in the national law. In particular, the law ought to list specific circumstances in which the authorities should not give leave to marry – for example, if it would affect the process of rehabilitation. In this context, it must be stressed that the issues involved were of a sensitive nature; thus, unjustified refusal of leave to marry could be regarded as additional or disciplinary punishment. 90: This conclusion is reinforced by the wording of Article 12. In contrast to Article 8 of the Convention, which sets for the right to respect for private and family life, and with which the right "to marry and to found a family" has a close affinity, Article 12 does not
Regular	In the context of restrictions on the right to marry, as established by [001-96453], the Court determines whether such limitations are properly regulated and subjected to judicial review based upon two key factors. Firstly, any measures restricting this fundamental freedom must be limited to concrete facts rather than prejudice ([001-96453]). Secondly, rules governing these discretionary decisions should be laid down in national law with specific circumstances listed where authorization for marriage may not be granted - e.g., if it would affect rehabilitation processes or jeopardize prison security (ibid.). This approach allows applicants an opportunity to challenge denials through a legal framework that ensures fairness while maintaining balance between individual rights and societal needs within democratic societies' margins of appreciation.[001-96453]
AdaCAD	In the context of restrictions on the right to marry, [001-96453] states that Article 12 does not guarantee an unlimited right since it is regulated by "the national laws" governing its exercise. Consequently, for a restriction to be proper and subjected to judicial review, rules should govern authorities' discretion in this area with specific circumstances listed where leave to marry ought not given - such as affecting rehabilitation process ([ibid]). This means prisoners have the ability to challenge denial authorization if concrete facts were disregarded or prejudice influenced decisions rather than prison security jeopardy being involved; moreover, refusal cannot restrict so much that essence impaired due to sensitive nature which could constitute additional punishment without justification.[001-96453].
CoCoLex	In the context of restrictions on the right to marry, the Court determines whether the restrictions are properly regulated and subjected to judicial review by assessing whether the national laws governing the exercise of this right are in place [001-96453]. The Court also ensures that the discretion of the authorities in exercising their power to regulate the right to marry is limited to concrete facts, not prejudice, and that their decisions are based on specific circumstances listed in the national law [001-96453]. This means that the applicant has the ability to challenge the denial of authorization if it can be shown that the restrictions imposed are arbitrary or disproportionate, as the Court would not apply the tests of "necessity" or "pressing social need" but would determine whether the impugned interference was arbitrary or disproportionate [001-96453].

Table 10: Case study 1 on a question from the AQuAECHR dataset, including provided source passages and model generations from Regular, AdaCAD, and CoCoLex.

=

Question	Under what circumstances can detention following conviction be considered unlawful?
Source Passages	001-127697 56: The applicant alleged that what she considered to be the retroactive application of a departure from the case-law by the Supreme Court after she had been convicted had extended her detention by almost nine years, in violation of Article 7 of the Convention, which reads as follows: 001-212829 35: In Kadusic v. Switzerland (no. 43977/13, 9 January 2018), which concerned a subsequent order of an institutional therapeutic measure under Swiss law, the Court was in principle prepared to accept that the order for an institutional measure amounted to a correction of the original judgment following the discovery of relevant new circumstances and that the fact that the measure was ordered in the context of proceedings for the review of a penalty imposed in a previous judgment may constitute a causal link between the initial conviction and the measure in issue, as required by the relevant case-law of the Court concerning sub-paragraph (a) of Article 5 § 1 (ibid., § 50). However, in the circumstances of that case, the Court found that the detention which followed the (new) judgment in the review proceedings lacked a sufficient causal connection with the initial conviction for being incompatible with its aims. The Court considered that the measure in question, which had been imposed a considerable period after the applicant's initial conviction, had not been based on a sufficiently recent expert report and that the applicant had been detained in an institution unsuited to his mental disorders (ibid., § \$ 53-60). 001-187391 116: The Court is therefore of the general view that in order not to make it impracticable for the police to fulfil their duties of maintaining order and protecting the public, provided that they comply with the underlying principle of Article 5, which is to protect the individual from arbitrariness (see, Austin and Others v. the United Kingdom [GC], nos. 39692/09 and 2 others, § 56, ECHR 2012), the lawful detention of a person outside the context of criminal proceedings can, as a matter of
Regular	Detention following conviction can be considered unlawful under the principles outlined in ECHR case law if it lacks a sufficient causal connection with the initial conviction [001-212829, § 53] or is based on an expert report that was not recent enough at the time of detainment. Additionally, such detentions must comply with Article 7 (prohibition against retroactive application) and protect individuals from arbitrariness as required by the underlying principle of Article 5 to ensure their rights are protected during this period [001-127697].
AdaCAD	Detention following conviction can be considered unlawful if there is a lack of sufficient causal connection between the initial conviction and subsequent detension [001-212829, § 53]. Additionally, it may also occur when such detainment takes place in an institution that does not suit mental disorders present with the individual being held [ibid., §§ 60], or occurs considerable period after the applicant's original judgment without recent expert report.
CoCoLex	Detention following conviction can be considered unlawful under the principles outlined in the case law if it lacks a sufficient causal connection with the initial conviction [001-212829, § 50] or if it is based on a new judgment in the review proceedings that has been imposed a considerable period after the applicant's initial conviction and has not been based on a sufficiently recent expert report [001-212829, § 53-60]. Additionally, detention outside the context of criminal proceedings can be unlawful if it is not permissible under Article 5 § 1 (c) of the Convention and does not comply with the underlying principle of Article 5, which is to protect the individual from arbitrariness [001-187391, § 56].

Table 11: Case study 2 on a question from the AQuAECHR dataset, including provided source passages and model generations from Regular, AdaCAD, and CoCoLex.