

STARD: A Chinese Statute Retrieval Dataset with Real Queries Issued by Non-professionals

Anonymous ACL submission

Abstract

Statute retrieval aims to find relevant statutory articles for specific queries. This process is the basis of a wide range of legal applications such as legal advice, automated judicial decisions, legal document drafting, etc. Existing statute retrieval benchmarks focus on formal and professional queries from sources like bar exams and legal case documents, thereby neglecting non-professional queries from the general public, which often lack precise legal terminology and references. To address this gap, we introduce the STATute Retrieval Dataset (STARD), a Chinese dataset comprising 1,543 query cases collected from real-world legal consultations and 55,348 candidate statutory articles¹. Unlike existing statute retrieval datasets, which primarily focus on professional legal queries, STARD captures the complexity and diversity of real queries from the general public. Through a comprehensive evaluation of various retrieval baselines, we reveal that existing retrieval approaches all fall short of these real queries issued by non-professional users. The best method only achieves a Recall@100 of 0.907, suggesting the necessity for further exploration and additional research in this area.

1 Introduction

Statutes are written laws formally created and approved by a legislative body, such as a parliament or congress (Livingston, 1990). They set out specific rules and guidelines within a certain area or jurisdiction. Therefore, statutes are the primary source of legal authority in civil law countries and also play a significant role in common law jurisdictions.

Statute retrieval involves finding relevant statutory articles or sections of laws for a specific query. This process is vital in the legal field and supports a wide range of applications, including legal advice

¹All the codes and datasets are available at: <https://anonymous.4open.science/r/STARD/>

Question: Is disclosing the medical case information of a patient considered an invasion of privacy?

Relevant Statute Articles

Personal Information Protection Law, Article 28: Sensitive personal information refers to information that, if leaked or illegally used, could easily harm an individual's dignity or endanger their personal or property safety. This includes biometric data, religious beliefs, specific identities, medical health, financial accounts, tracking information, and personal information of minors under the age of fourteen.

Civil Code, Article 1032: Individuals have the right to privacy. No organization or individual may infringe upon another's privacy rights through snooping, harassment, disclosure, or publicization. Privacy encompasses the tranquility of an individual's private life and the private spaces, activities, and information they wish to keep unknown to others.

Civil Code, Article 1226: Medical institutions and their medical personnel must keep patients' privacy and personal information confidential. Those who disclose patients' private and personal information or publish their medical records without the patient's consent must bear infringement liability.

Table 1: An example of the query and relevant statute articles in the STARD dataset.

services, automated judicial decisions, and logical legal analysis. This task is challenging for the following reasons: **(1)** Statutes use complex legal terminology and linguistic structures rarely found in open-domain corpus. As a result, traditional retrieval models that lack domain-specific knowledge may struggle to accurately capture the meanings of these specialized terms. **(2)** The criteria for assessing information relevance in the legal domain differ greatly from those used in open-domain search tasks. General search tasks focus mainly on textual similarity, while legal tasks involve legal reasoning that requires the understanding of different areas of law, the relations between them, as well as the relevance of specific legal principles and their practical applications.

Due to the challenging nature of statute retrieval and its paramount importance in civil law systems, significant progress has been made in this field. For example, the annual COLIEE competitions introduce a series of statute retrieval tasks using the

061	questions extracted from the Japanese legal bar	developing more accessible and efficient legal systems.	113
062	exams (Goebel et al., 2023; Kim et al., 2022; Ra-		114
063	belo et al., 2022). These tasks aim to retrieve rel-	In conclusion, the contributions of this paper are	115
064	evant statute law from the Japanese Civil Code	as follows:	116
065	Article according to the question from bar exams.		
066	AILA (Bhattacharya et al., 2019) competitions also	• We propose STARD, a statute retrieval dataset de-	117
067	introduce a series of statute retrieval datasets. The	derived from real-world legal consultation posed by	118
068	queries from AILA are case documents that were	non-professionals, with 1,543 queries and their	119
069	judged by the Supreme Court of India. The candi-	corresponding relevant statutes.	120
070	date statutes are part of the set of statutes from		
071	Indian law.	• We propose a comprehensive annotation frame-	121
072	Despite these advancements, a significant gap	work specifically designed for the statute retrieval	122
073	persists in addressing real queries from non-	task based on non-professional queries, which	123
074	professional people, who represent a large pop-	provides references and insights for future anno-	124
075	ulation of legal advice service users. The current	tation in the legal field.	125
076	statute retrieval benchmarks are primarily based on		
077	queries from formal legal documents, such as bar	• We conduct experiments on a wide range of re-	126
078	exam questions or Supreme Court case documents,	trieval baselines and find that statute retrieval	127
079	which differ significantly from the everyday lan-	with queries issued by non-professionals is still a	128
080	guage used by the general public. However, queries	difficult task that requires further investigation.	129
081	from non-professionals often lack precise legal ter-		
082	minology and may include ambiguous references	• We present experiments on LLMs solving legal	130
083	to legal concepts, which significantly complicate	tasks with and without the STARD dataset. Ex-	131
084	the task of statute retrieval.	periments show that STARD can notably enhance	132
085	To address the limitations of existing bench-	the performance of LLMs in legal tasks.	133
086	marks, we propose STAtute Retrieval Dataset		
087	(STARD) i.e. STARD, a Chinese statute retrieval	2 Problem Formulation	134
088	dataset based on real legal consultation questions		
089	from the general public. The STARD dataset com-	2.1 Statute of Civil Law System	135
090	prises 1,543 query cases collected from real-world		
091	legal consultations and 55,348 candidate statutory	Civil law is a legal system primarily based on codi-	136
092	articles extracted from all official Chinese legal reg-	fied laws rather than case precedents, making writ-	137
093	ulations and judicial interpretations. Table 1 shows	ten statutes the main source of legal authority. This	138
094	an example of our dataset. To the best of our knowl-	contrasts with common law systems, where previ-	139
095	edge, STARD is the first statute retrieval dataset	ous judicial decisions also play a central role. In	140
096	where queries are from real-world legal consulting	civil law, statutes are created and enacted by legisla-	141
097	proposed by the general public.	tive bodies, such as parliaments, and are organized	142
098	We conduct experiments on a wide range of in-	into systematic collections known as codes, which	143
099	formation retrieval (IR) baselines on the STARD	cover various areas of law like contracts, torts, and	144
100	dataset, including traditional lexical matching mod-	property. A statute is a formal written law that pro-	145
101	els, open-domain neural retrieval models, legal do-	vides specific rules and guidelines to be followed	146
102	main neural retrieval models, and a dense retriever	within a jurisdiction. Within statutes, there are sec-	147
103	trained with data annotated by GPT-4. The experi-	tions known as statutory articles, which detail indi-	148
104	mental results show that all existing baselines fall	vidual provisions or clauses of the law, addressing	149
105	short of accurately and comprehensively retrieving	particular aspects or requirements. These statutes	150
106	the relevant statutes, leaving significant room for	and their articles are fundamental in civil law sys-	151
107	future work. Additionally, our experimental results	tems to ensure that the legal framework is clear,	152
108	show that employing STARD as an external knowl-	predictable, and accessible, thereby facilitating or-	153
109	edge source for Retrieval-Augmented Generation	der and defining societal rights and responsibilities.	154
110	(RAG) significantly enhances the performance of		
111	large generative language models (LLMs) on le-	2.2 Definition of Statute Retrieval	155
112	gal tasks. This indicates that STARD is useful for		
		The statute retrieval task aims to accurately retrieve	156
		relevant statutory articles in response to a query. To	157
		be specific, given a query q that describes a legal	158
		issue or situation and a corpus of statutory articles	159

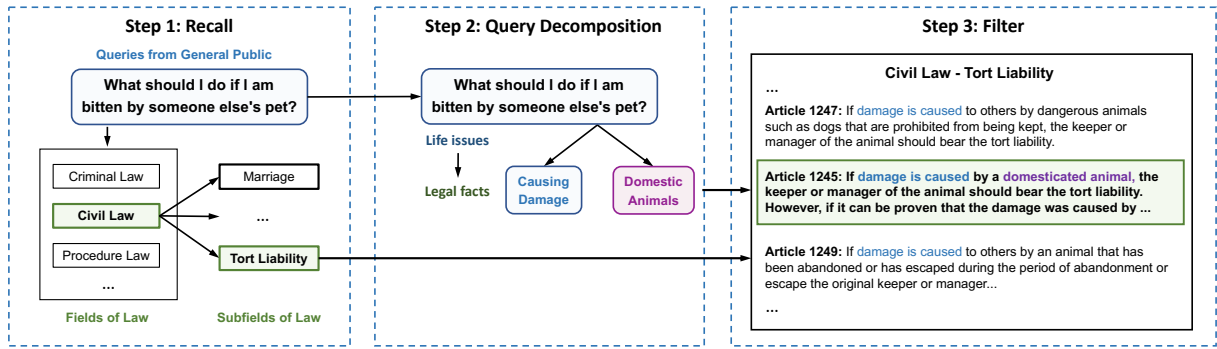


Figure 1: A schematic of our annotation framework with green boxes highlighting query-relevant elements.

160 $S = \{s_1, s_2, \dots, s_n\}$, $n \in \mathbb{N}^+$. For each statute
 161 s_i in the corpus, there is a Bernoulli variable r_i
 162 indicating whether s_i is relevant² to the query q .
 163 The goal of the statute retrieval task is to retrieve
 164 a set of statutes $R = \{s_j | r_j = 1\}$, including all
 165 statutes relevant to the query.

166 3 Annotation Framework

167 This section explains how annotators transform
 168 general questions into professional legal questions
 169 submitted by non-professionals and then identify
 170 the most relevant legal statutes to support these
 171 questions. To be specific, annotators use a three-
 172 step method: recall, query decomposition, and fil-
 173 tering (illustrated in Figure 1). This method mirrors
 174 the structured approach commonly used in legal
 175 reasoning, which involves three logical steps: es-
 176 tablishing a broad legal principle (major premise),
 177 applying it to the specific facts of a case (minor
 178 premise), and then concluding. This section is or-
 179 ganized into three subsections, each detailing a part
 180 of the annotation process that is designed to mirror
 181 these logical steps in legal reasoning.

182 3.1 Step 1: Recall

183 When initiating the annotation of legal statutes per-
 184 tinent to a query, our annotators first narrow down
 185 the scope of the relevant statutes. Specifically, they
 186 start by identifying the most pertinent areas of law
 187 within the entire legal system. The process uses a
 188 top-down refining method. Annotators begin with
 189 broad departmental categories of law, such as civil,
 190 criminal law, and administrative law. Upon encoun-
 191 tering a specific issue, annotators first determine
 192 which category of departmental law it falls under,
 193 then progressively refine the issue to more specific
 194 aspects of the law. For instance, if the issue pertains
 195 to civil law, the annotator assesses whether it relates

²The definition of “relevant” is discussed in detail in Section 3.

196 to contract law or tort law. If it is a matter of con-
 197 tract law, a further determination is made regard-
 198 ing the specific type of contract involved. Similarly, for
 199 tort law, the specific type of tort is identified. This
 200 step effectively narrows the scope of legal statute
 201 retrieval to particular chapters within the relevant
 202 departmental law.

203 3.2 Step 2: Query Decomposition

204 Given the specialized nature of legal knowledge,
 205 individuals without a formal education in law often
 206 frame their queries with informal language rather
 207 than professional legal terminology. These queries
 208 typically consist of straightforward semantic ex-
 209 pressions that do not directly correspond to estab-
 210 lished legal norms. For instance, consider the ques-
 211 tion “What should I do if I am bitten by some-
 212 one’s pet?”. Here, “pet bite” represents a common,
 213 non-technical description of an incident. Search-
 214 ing for legal norms based solely on such descrip-
 215 tions might lead to irrelevant or imprecise results.
 216 Therefore, when annotators perform legal statute
 217 retrieval, they should transform the informal fact
 218 descriptions written by the questioner into legal
 219 facts through interpretation. This is the step to find
 220 the minor premise in the legal logic syllogism. In
 221 this transformation process, the annotator evaluates
 222 the life facts according to the provisions of the law
 223 and selects the legal norms corresponding to these
 224 life facts. For example, for the aforementioned
 225 issue of a pet biting a person, the annotators will
 226 transform “pet bites a person” into the legal fact of
 227 “causing damage to other” and “domestic animals”
 228 according to the provisions of Chapter 9 of the Tort
 229 Liability Compilation of the Civil Code.

230 3.3 Step 3: Filter

231 The filtering process is a critical step where anno-
 232 tators refine and finalize the selection of relevant
 233 legal statutes. This is accomplished by employ-
 234 ing a “subsumption” method, integral to the syllo-

gistic reasoning in law. In this method, the legal facts, which have been interpreted and transformed from real-life scenarios in the previous steps, are matched against the smallest possible subset of legal statutes that adequately address the query.

To be specific, consider a set of legal statutes $S = \{S_1, S_2, S_3\}$ recalled in the first step. Through the transformation process, the query is deconstructed into distinct legal facts F_1, F_2, F_3 . Each fact corresponds to a subset of statutes that it implies, denoted as $S_{F_1} = \{S_1, S_4, S_5\}$, $S_{F_2} = \{S_1, S_5\}$, and $S_{F_3} = \{S_3, S_6\}$. The objective in the filtering stage is to intersect these subsets with the initially recalled set S to determine the most relevant statutes. This is represented as $S_{Golden} = (S_{F_1} \cup S_{F_2} \cup S_{F_3}) \cap S$, yielding $S_{Golden} = \{S_1, S_3\}$.

These statutes in S_{Golden} are considered the “golden” legal statutes for the dataset, as they encompass all the legal implications drawn from the facts of the query. This step ensures that the selected statutes are not only relevant but also comprehensive in covering the legal issues presented in the query, thereby providing a solid legal foundation to support the resolution of the query.

3.4 Generalizability of Our Framework

In this section, we discuss the generalizability of the STARD dataset and our annotation framework, discussing the following two research questions (RQs): **RQ1:** Can the STARD dataset be applied to the legal systems of other countries through direct translation of our dataset? **RQ2:** Can our Three-Step Annotation Framework be applied to other legal systems?

For RQ1, directly translating the STARD dataset into other languages does not guarantee its applicability in foreign legal systems. Each country possesses unique legal statutes; articles selected from one jurisdiction may not exist or may have entirely different implications in another. Thus, the nuances of local laws must be considered, making straightforward translation inadequate for cross-national applications. For RQ2, our proposed Three-Step Annotation Framework is potentially generalizable to other countries under the civil law system. Countries with civil law systems, such as Germany, France, and Japan, typically share a similar process for retrieving law statutes. This process can generally be structured into three steps: Recall, Query Decomposition, and Filtering. Therefore, our framework could be adapted to these environments, supporting the construction of statute re-

trieval datasets and the application of legal statutes across various civil law jurisdictions.

4 Dataset Construction

4.1 Data Sources

All queries in our dataset derive from real legal consultations. Specifically, our legal team creates legal questions from the 12348 China Legal Service Website³, followed by a manual anonymization of each question, which involved removing any potential identifiers associated with entities, corporations, or individuals.

To obtain the 55,348 candidate statutory articles, our legal team conducted extensive research and discussions to compile a comprehensive list of currently valid Chinese statutory laws and regulations⁴. We then manually downloaded the most up-to-date versions of these laws from the government’s official website. These laws were subsequently divided into the smallest searchable units based on articles using automated scripts.

4.2 Recruitment and Payment of Annotators

For recruitment, we sourced annotators from prominent law schools. The annotation team initially consisted of 16 members. Although three members departed during the project, their positions were quickly filled to maintain the team size. Our salary plan remunerates participants based on the number of annotations they complete, with a fixed rate of approximately 10 CNY per annotation. On average, annotators processed four queries per hour, resulting in an average hourly wage of 40 CNY. This pay rate significantly exceeds the minimum hourly wage mandated in Beijing.

4.3 Annotation Process

Annotators are tasked with identifying relevant articles of statutes in response to actual legal queries posed by the general public. The specifics of the annotation framework are detailed in Section 3. Additionally, annotators are instructed not to use generative models, such as ChatGPT, for assistance. The annotation process starts with the manual anonymization of each question within the STARD dataset, involving the removal of any potential identifiers associated with entities, corporations, or individuals. Subsequently, annotators are

³This is the Chinese government’s official website for online legal services: <http://www.12348.gov.cn/>

⁴The entire list of statutes we selected can be found on our official GitHub.

required to locate relevant statutes for each question, following the three-step principle introduced in Section 3.

4.4 Annotation Consistency

For each question, two annotators were assigned. **The final gold standard for each question was established only when both annotators agreed on the same legal provisions⁵.**

To evaluate the reliability of agreement among human annotators, we utilized Cohen’s Kappa (Cohen, 1960) \mathcal{K} coefficient in a binary classification context. Each query-statute article pair corresponds to a binary classification task, where annotators judge whether the query is related or unrelated to the statute. This analysis, conducted on a dataset comprising 1,543 annotated instances, yielded a \mathcal{K} value of 0.5312. This indicates moderate agreement. Achieving such a \mathcal{K} value is considered satisfactory for a complex task involving fifty thousand classifications with multiple possible correct labels.

4.5 Ethics Discussion

We have thoroughly addressed the following ethical considerations: **(1) Privacy and Anonymity:** Given the sensitive nature of legal consultations, we have rigorously anonymized all queries in the STARD dataset. **(2) Transparency:** To promote reproducibility and transparency, we have made the dataset, associated models, and the codes publicly available⁶. This allows other researchers to verify, replicate, and expand upon our work, advancing the field of legal informatics. **(3) Accountability:** Recognizing the dynamic nature of legal statutes, we commit to regularly updating the STARD dataset to reflect the latest changes in law. This ensures the dataset remains accurate and reliable for ongoing research and application. **(4) Accessibility:** The STARD dataset is freely available for download from the official website under the MIT license, facilitating easy access for researchers and practitioners alike. This promotes broader usage and supports innovation across various fields.

5 Dataset Statistics and Analysis

The basic statistics of our proposed dataset are shown in Table 2. STARD comprises a total of 1,543 queries and a large-scale corpus of 55,348

⁵In cases where annotators had differing opinions, the question would not be included in the final dataset.

⁶<https://anonymous.4open.science/r/STARD/>

Table 2: Basic statistics of our proposed STARD dataset.

Statistic	# Number
Total Queries	1,543
Total Candidate Statutory Articles	55,348
Total Num of Relevant Statutory Articles	1,445
Occurrences of Relevant Articles	2,717
Avg. Relevant Articles per Query	1.76
Avg. Query Length	27.30
Avg. Article Length	119.93

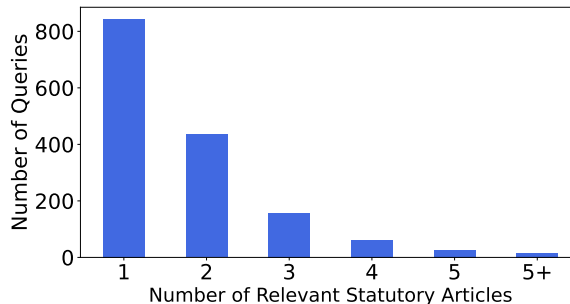


Figure 2: Distribution of relevant statutory article numbers for each query.

candidate statutory articles. Among these candidate statutory articles, 1,445 articles are relevant to at least one of the queries in the dataset. The average length of a query is 27.3 words, and the average length of a statute article is nearly 120 words.

Figure 2 presents the distribution of queries across the number of relevant statutory articles, highlighting the varied complexity within the dataset. A substantial majority of the queries, 843 out of 1,543, correspond to just one relevant statutory article, indicating a significant number of queries can be addressed with a single, specific legal reference. This could suggest that many of the non-professional queries are focused and pertain to specific legal issues that require straightforward statute retrieval. However, 45% of queries require multiple statutory articles which indicates some of the questions are more complex, involving multiple references of law. This diversity in query complexity demonstrates that our dataset is capable of accommodating a wide range of legal questions, from straightforward to highly intricate.

6 Statute Retrieval Experiment

6.1 Selected Retrieval Baselines

We consider four types of baselines for comparison, including traditional IR methods, pre-trained Language models on general domain data, PLMs tailored for IR, and pre-trained language models built with legal documents. **The implementation details of these baselines are provided in Appendix B**

Table 3: The overall experimental results of multiple baselines on STARD. The best results are in bold, and the second-best results are underlined. “R” stands for Recall, and “M” stands for MRR. General PLM and Legal PLM are all in the zero-shot setting. Note that LSI-STARD is a classification model where each statute is treated as a unique label; we report its ranking performance based on the probability for each statute.

		R@5	R@10	R@20	R@30	R@50	R@100	R@200	M@3	M@5	M@10
Lexical Matching	QL	0.3363	0.4020	0.4651	0.4839	0.5537	0.6515	0.7224	0.3052	0.3167	0.3304
	BM25	0.3349	0.3943	0.4504	0.4773	0.5240	0.6493	0.7035	0.3176	0.3251	0.3369
Open-Domain PLM	Roberta	0.3216	0.3908	0.4646	0.5042	0.5715	0.6633	0.7351	0.2766	0.2905	0.3010
	SEED	0.2897	0.3555	0.4264	0.4589	0.4975	0.5626	0.6260	0.2607	0.2708	0.2816
	coCondenser	0.1120	0.1598	0.2223	0.2659	0.3288	0.4292	0.5246	0.0847	0.0922	0.1004
Legal PLM	SAILER	0.2330	0.3050	0.3790	0.4286	0.4885	0.5674	0.6463	0.2006	0.2115	0.2234
	Lawformer	0.2411	0.2989	0.3720	0.4137	0.4733	0.5478	0.6309	0.2205	0.2313	0.2412
Fine-tuned PLM	Dense-STARD	0.5206	0.6061	0.7064	0.7485	0.8107	0.9065	0.9531	0.4372	0.4543	0.4724
	Dense-GPT4	<u>0.4382</u>	<u>0.5174</u>	<u>0.5961</u>	<u>0.6471</u>	<u>0.6810</u>	<u>0.7984</u>	<u>0.8521</u>	<u>0.3842</u>	<u>0.3948</u>	<u>0.4106</u>
	Dense-CAIL	0.0887	0.1272	0.1832	0.2341	0.2712	0.3281	0.3819	0.0660	0.0719	0.0842
	LSI-STARD	0.1861	0.2069	0.2386	0.2564	0.3004	0.3410	0.3956	0.2062	0.2093	0.2156

• Traditional IR Methods

- **QL** (Ponte and Croft, 2017) is a language model based on Dirichlet smoothing and has good performance on retrieval tasks.
- **BM25** (Robertson et al., 2009) is an effective retrieval model based on lexical matching that achieves good performance in retrieval tasks.

• General Domain Pre-trained Models

- **Chinese-RoBERTa-WWM** (Cui et al., 2021) is a language model pre-trained with the Whole Word Masking strategy.
- **SEED** (Lu et al., 2021) is a pre-trained text encoder for dense retrieval that achieves state-of-the-art performance.
- **coCondenser** (Gao and Callan, 2021b) is an enhanced version of Condenser (Gao and Callan, 2021a) that adds an unsupervised corpus-level contrastive loss to warm up the passage embedding space.

• Legal Domain Pre-trained Models

- **Lawformer** (Xiao et al., 2021) apply Longformer (Beltagy et al., 2020) to initialize and train with the MLM task on the legal domain.
- **SAILER** (Li et al., 2023) is a structure-aware pre-trained language model for tailored legal document representation. It utilizes the logical connections between different sections within a legal document.

• Fine-tuned Dense Retrieval Model

- **Dense-CAIL** is a dense retrieval model trained on the CAIL2018 dataset (Xiao et al., 2018). We choose this baseline to verify whether the

existing dataset based on formal professional questions is sufficient for addressing statute retrieval tasks based on non-professional queries.

- **Dense-STARD** employs a five-fold cross-validation technique on the STARD dataset.

We initialize the above two models with Chinese-Roberta-WWM (Cui et al., 2021). For the setting of cross-validation, the dataset is randomly divided into five subsets, where one subset serves as the test set and the remaining four are used as training sets. The details of our fine-tuning process are introduced in Appendix G.

- **Dense-GPT4**: We distill a dense retrieval model from GPT-4. The process involved using GPT-4 to generate legal questions based on statute articles within a given corpus. Specifically, we prompted GPT-4 to create a legal question q that is closely related to a specific statute article a_i^+ , resulting in a query-statute pair (q, a_i^+) . Then, we employ a contrastive learning approach utilizing these query-statute pairs to train the dense retriever. Details are provided in Appendix H.

- **LSI-STARD** is a Transformer based classifier fine-tuned on STARD. In the Legal Statute Identification (LSI) field (Zhong et al., 2018; Paul et al., 2022; Chalkidis et al., 2021), the statute retrieval task is approached as a classification problem, where each statute is treated as a unique label. This method transforms the task into classifying legal documents or queries against a set of labels, each representing a different statute. Following this methodology, we finetune a transformer-based classification model on the STARD dataset, employing the same five-fold cross-validation setting. We initialize the transformer-based model

475	with Chinese-Roberta-WWM (Cui et al., 2021)	525
476	and randomly initialize the outermost MLP Layer.	526
477	Details are provided in Appendix I.	527
478	6.2 Evaluation Metrics	528
479	We use Mean Reciprocal Rank and Recall as eval-	529
480	uation metrics. By using both MRR and Recall,	530
481	we can gain insights into both the accuracy of the	531
482	top-ranked results and the comprehensiveness of	532
483	the relevant statutory articles retrieved by the re-	533
484	trieval model. Detailed definitions of these metrics	534
485	are provided in Appendix C.	535
486	6.3 Experimental Results	536
487	In this subsection, we provide a detailed analysis	537
488	of the performance of various retrieval baselines	
489	evaluated on our proposed STARD dataset. We	
490	have the following insights into the effectiveness	
491	of different retrieval methods:	
492	(1) Under the zero-shot setting, traditional lexi-	
493	cal matching techniques surpass both general and	
494	legal-domain pre-trained language models (PLMs).	
495	This demonstrates that lexical matching methods	
496	are still very strong baselines in retrieval tasks. (2)	
497	Among all the methods that do not use human anno-	
498	tation, the performance of Dense-GPT4 stands out,	
499	exceeding that of all unsupervised methods tested.	
500	This indicates that distilling GPT4 to train task-	
501	specific models is a good choice in scenarios with-	
502	out human annotations. (3) Domain-specific mod-	
503	els like SAILER are optimized for particular tasks,	
504	thus resulting in underperformance compared to	
505	general domain models. Specifically, SAILER is	
506	tailored for legal case retrieval involving long doc-	
507	uments as queries. Consequently, it struggles with	
508	tasks that involve short queries and medium-length	
509	articles, unlike the model STARD. (4) The retrieval	
510	model fine-tuned on the CAIL2018 dataset per-	
511	formed sub-optimally on the STARD dataset. This	
512	suggests significant differences between the non-	
513	professional queries in STARD and the formal leg-	
514	al queries in existing datasets. Consequently, it	
515	underscores the unique nature of STARD, neces-	
516	sitating specialized models for effective statute re-	
517	trieval. (5) While the LSI classifier performs well	
518	in existing studies for tasks involving the classi-	
519	fication of a few dozen statutes, it struggles with	
520	the STARD dataset, which contains over 50,000	
521	labels, resulting in suboptimal performance. As a	
522	result, retrieval methods are more effective than the	
523	LSI approach for large-scale statute retrieval tasks.	
524	(6) The performance of both the lexical matching	
	method and the non-finetuned models is less effec-	525
	tive than that of the Dense-STARD model. This	526
	arises because the former models lack the capaci-	527
	ty to interpret life issues as legal facts, a capabili-	528
	ty that Dense-STARD has acquired through fine-	529
	tuning. It has been trained to associate the life	530
	issues presented in queries with relevant legal arti-	531
	cles. However, Dense-STARD’s training set is con-	532
	fined to just over one thousand query-article pairs.	533
	Consequently, its recall rates remain suboptimal,	534
	with Recall@100 at only 90.65%. These findings	535
	underscore the necessity for further exploration in	536
	this field.	537
	7 Retrieval Augmented Generation	538
	Experiment	539
	7.1 Selected Benchmark	540
	We select two datasets encompassing three tasks	541
	for our RAG experiment:	542
	• JecQA (Zhong et al., 2020) is the most ex-	543
	tensive multiple-choice dataset within the Chi-	544
	nese legal field. This dataset includes two dis-	545
	tinct tasks: Knowledge-Driven Questions (KD-	546
	questions) and Case-Analysis Questions (CA-	547
	questions), encompassing a total of 26,365 ques-	548
	tions. All the questions are multi-select, meaning	549
	that more than one option can be correct.	550
	• CAIL 2018 (Xiao et al., 2018) is a large-scale	551
	Chinese legal dataset designed for judgment pre-	552
	diction with over 2.6 million criminal cases. This	553
	dataset contains detailed annotations of judgment	554
	results, including applicable law articles, specific	555
	charges, and prescribed prison terms. We select	556
	the Charge Prediction task of CAIL 2018 and use	557
	prediction Accuracy as the evaluation metric.	558
	7.2 Selected LLMs and Settings	559
	Our selected LLMs are introduced in Appendix D.	560
	The generation configuration is detailed in Ap-	561
	pendix E. The prompt template for LLMs is de-	562
	tailed in Appendix F.	563
	7.3 Experimental Results	564
	Table 4 presents the results of the LLM’s per-	565
	formance with and without the use of Retrieval-	566
	Augmented Generation (RAG). In the scenario	567
	without RAG, the LLM directly outputs the correct	568
	options based on the question. In the RAG scenario,	569
	the retrieval model (BM25 or Dense-STARD) re-	570
	calls the top 10 relevant statutory articles from the	571

Table 4: The overall experimental results of three LLMs on the JecQA benchmark. We report accuracy as the evaluation metric. The best results are in bold and the second best results are underlined.

	Retriever	JQA-CA	JQA-KD	CAIL
Baichuan	w/o RAG	0.231	0.266	0.850
	BM25	<u>0.233</u>	<u>0.288</u>	0.766
	Dense-STARD	0.238	0.291	<u>0.816</u>
chatGLM	w/o RAG	0.185	0.194	0.636
	BM25	0.189	<u>0.224</u>	<u>0.646</u>
	Dense-STARD	0.200	0.237	0.684
chatGPT	w/o RAG	0.187	0.206	0.496
	BM25	0.233	0.293	0.528
	Dense-STARD	<u>0.193</u>	<u>0.252</u>	<u>0.503</u>

corpus based on the question. The retrieved statutory articles are then integrated into a meticulously designed prompt template (detailed in Appendix F).

The experimental results reveal that using the STARD corpus as the external knowledge base for the RAG significantly enhances the performance of large language models (LLMs) and underscores the value of our proposed dataset in improving the effectiveness of LLMs on legal tasks. The results also reveal that different LLMs have unique preferences for retrievers. For the Baichuan and ChatGLM models, a fine-tuned dense retriever surpasses BM25, indicating that these models benefit from dense retrievers’ high recall rates. However, this advantage is not observed with the ChatGPT model, where BM25 outperforms the fine-tuned dense retriever. This suggests that the performance of RAG is highly dependent on the preferences of the LLM regarding the retriever. The experimental results on the CAIL 2018 dataset align with those observed for JecQA, with one notable exception: the performance of the Baichuan model without RAG. In this setting, Baichuan’s performance is markedly superior to that of chatGLM, chatGPT, and Baichuan with RAG. We hypothesize that this exception arises from the Baichuan model’s utilization of the CAIL 2018 dataset during its pre-training phase, leading to a direct answer accuracy rate that is even 81% higher than that of chatGPT.

8 Related Work

CAIL 2018 (Xiao et al., 2018; Zhong et al., 2018) competitions conduct law statute retrieval work using formal legal judgment documents. The queries in the dataset originate from the “Court’s Findings” part of the judgments, and the candidates are statute articles of Chinese Criminal Law. The annual COLIEE competitions introduce a series

of statute retrieval datasets using the questions extracted from the Japanese legal bar exams (Goebel et al., 2023; Kim et al., 2022; Rabelo et al., 2022). These tasks aim to retrieve relevant statute law from the Japanese Civil Code Article according to questions from bar exams. AILA (Bhattacharya et al., 2019) competitions also introduce a series of statute retrieval datasets. The queries from AILA are legal judgment documents from the Supreme Court of India. The candidate statutes are part of the set of statute articles from Indian law. BSARD (Louis and Spanakis, 2021) is a statutory article retrieval dataset in French with candidate articles from a 22,600+ Belgian law articles corpus.

In the studies of the Legal Statute Identification (LSI) (Zhong et al., 2018; Paul et al., 2022; Chalkidis et al., 2021), finding the relevant statute is approached as a classification problem, where each statute is treated as a unique label. LADAN (Xu et al., 2020) is an LSI method that uses a graph neural network and attention mechanism to distinguish confusing law articles. LeSICiN (Paul et al., 2022) utilizes both textual content and legal citation networks to identify relevant legal statutes.

Legal QA tasks also aim to fulfill the public’s demand for legal information (Do et al., 2017). LLeQA (Louis et al., 2024) is a French long-form legal QA dataset comprising 1,868 expert-annotated legal questions. GerLayQA (Büttner and Habernal, 2024) is a question-answering dataset comprising 21k laymen’s legal questions paired with answers from lawyers and grounded in concrete law book paragraphs.

9 Conclusion

We present STARD, a new benchmark consisting of 1,543 questions from the general public. To the best of our knowledge, STARD is the first Chinese statutes retrieval dataset tailored for the general public. Moreover, we propose an annotation framework to improve the accuracy and relevance of statute retrieval annotation, which offers valuable guidelines for future legal annotations. Our experiments across various retrieval models highlighted the complexities of non-professional statute retrieval, indicating the necessity for further exploration. Additionally, we demonstrated that integrating the STARD dataset significantly boosts the performance of LLMs in legal tasks, showcasing its potential to enhance legal AI applications.

10 Limitations

We acknowledge the limitations of this paper. One of the primary limitations is that our dataset is specifically designed around the Chinese legal system, inherently limiting its direct applicability to legal systems outside of this context. Despite our discussions on potential methodologies for adapting STARD to other civil law systems, such an expansion necessitates creating and annotating new datasets tailored to those systems' distinct legal frameworks and statutes. Thus, our future work will be dedicated to developing additional datasets that encompass a broader range of civil law systems. This endeavor aims to extend the utility of our work and foster further research and development in the domain of legal statute retrieval, ensuring broader applicability and relevance across different legal landscapes.

11 Ethics Statement

In the framework of this research, ethical considerations have been paramount from the initial stages, underscoring our commitment to the responsible advancement and application of artificial intelligence technologies. Our adherence to the principles of open research and the critical importance of reproducibility have compelled us to make all associated models, datasets, and codebases publicly available on GitHub.

Moreover, in the development of our dataset, we have paid scrupulous attention to privacy and respect for individuals' rights. Given the inherently sensitive nature of legal consultations, we have diligently anonymized every query within the STARD dataset. This process involved the removal of any potential identifiers related to entities, corporations, or individuals, thereby safeguarding privacy and preempting the possibility of data misuse.

References

- Qingyao Ai, Keping Bi, Jiafeng Guo, and W Bruce Croft. 2018. Learning a deep listwise context model for ranking refinement. In *The 41st international ACM SIGIR conference on research & development in information retrieval*, pages 135–144.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.
- Paheli Bhattacharya, Kripabandhu Ghosh, Saptarshi Ghosh, Arindam Pal, Parth Mehta, Arnab Bhattacharya, and Prasenjit Majumder. 2019. Fire 2019

- aila track: Artificial intelligence for legal assistance. In *Proceedings of the 11th annual meeting of the forum for information retrieval evaluation*, pages 4–6.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Marius Büttner and Ivan Habernal. 2024. Answering legal questions from laymen in german civil law system. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2015–2027.
- Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. 2007. Learning to rank: from pairwise approach to listwise approach. In *Proceedings of the 24th international conference on Machine learning*, pages 129–136.
- Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael Bommarito, Ion Androutsopoulos, Daniel Martin Katz, and Nikolaos Aletras. 2021. Lexglue: A benchmark dataset for legal language understanding in english. *arXiv preprint arXiv:2110.00976*.
- Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46.
- Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, and Ziqing Yang. 2021. Pre-training with whole word masking for chinese bert. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:3504–3514.
- Phong-Khac Do, Huy-Tien Nguyen, Chien-Xuan Tran, Minh-Tien Nguyen, and Minh-Le Nguyen. 2017. Legal question answering using ranking svm and deep convolutional neural network. *arXiv preprint arXiv:1703.05320*.
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. Glm: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 320–335.
- Luyu Gao and Jamie Callan. 2021a. Condenser: a pre-training architecture for dense retrieval. *arXiv preprint arXiv:2104.08253*.
- Luyu Gao and Jamie Callan. 2021b. Unsupervised corpus aware language model pre-training for dense passage retrieval. *arXiv preprint arXiv:2108.05540*.
- Luyu Gao, Zhuyun Dai, and Jamie Callan. 2021. Rethink training of bert rerankers in multi-stage retrieval pipeline. In *European Conference on Information Retrieval*, pages 280–286. Springer.

A License and Permissions

STARD is available under the MIT License. This permissive license was chosen to encourage the widespread use and adaptation of our resources, allowing for both academic and commercial applications without significant restrictions. For detailed terms and conditions, including how the dataset, code, and models can be used, modified, and shared, please refer to the documentation provided in our GitHub repository⁷.

B Implementation Details of Retrieval Baselines

- For the implementation of traditional IR methods QL and BM25, we use the Pyserini toolkit: <https://github.com/castorini/pyserini>.
- For the implementation of Chinese-RoBERTa-WWM, we directly use their models released on Huggingface⁸. As SEED and Condenser have no available Chinese versions, we reproduce their work on the Chinese Wikipedia based on their open-source training code and follow all settings provided in their paper (Lu et al., 2021; Gao and Callan, 2021a).
- For the implementation of Lawformer (Xiao et al., 2021) and SAILER (Li et al., 2023), we directly use the checkpoints released on the official GitHub⁹.

C Evaluation Metrics

We use Mean Reciprocal Rank and Recall as evaluation metrics.

The Mean Reciprocal Rank is a statistical measure used to evaluate the performance of a query-based system, where the primary goal is to retrieve the highest-ranked item. MRR calculates the average of the reciprocal ranks of results for a sample of queries. The reciprocal rank of a query response is the multiplicative inverse of the rank of the first correct answer:

$$\text{MRR} = \frac{1}{Q} \sum_{i=1}^Q \frac{1}{\text{rank}_i} \quad (1)$$

where Q is the number of queries, and rank_i is the rank position of the first relevant document for the i -th query.

⁷<https://anonymous.4open.science/r/STARD/>

⁸<https://huggingface.co/hfl/chinese-roberta-wwm-ext>

⁹<https://github.com/CSHaitao/SAILER/>,
<https://github.com/thunlp/LegalPLMs>

Recall measures the ability of a model to retrieve all relevant instances in a dataset. It is defined as the ratio of the number of relevant items correctly retrieved to the total number of relevant items in the database, which is critical in scenarios where missing any relevant item could be costly:

$$\text{Recall} = \frac{\text{Number of relevant items retrieved}}{\text{Total number of relevant items}} \quad (2)$$

D Selected LLMs

Our selected LLMs are listed as follows:

- **Baichuan** (Yang et al., 2023) is a series of large-scale multilingual language models, trained from scratch on 2.6 trillion tokens. We choose the **Baichuan-2-Base-13B** model which is widely used in bilingual Chinese-English scenarios.
- **ChatGLM** (Du et al., 2022) is a series of generative language models optimized for Chinese question answering and dialogue. We choose **ChatGLM3-6B** with 6.2 billion parameters.
- **ChatGPT** (Brown et al., 2020) is a series of large language models developed by OpenAI, including several versions. Among these, we choose **GPT-3.5-turbo**, which is identified as the most advanced GPT-3.5 model.

E Generation Configuration

We obtain responses from chatGPT by accessing its official API¹⁰. For Baichuan and chatGLM, we directly download model parameters from each model’s official Hugging Face repositories and use the official Python code provided by Hugging Face to obtain the response. We use the official default configurations provided by each model for the generation configuration.

F Prompt Template for RAG

In our RAG experiments, we employed the following prompt template for the LLM:

¹⁰<https://platform.openai.com/docs/guides/text-generation/chat-completions-api>

Prompt 1

Please answer the question based on the following statute articles:

Article 1: [Content]

.....

Article 10: [Content]

Please answer the following question based on the provided articles and your knowledge, prioritizing the provided knowledge. Note that the provided articles might not include those relevant to the question.

Question: xxx

G Fine-tuning Process

We initialize the model with Chinese-Roberta-WWM (Cui et al., 2021). We use the dual-encoder architecture (Karpukhin et al., 2020) to compute the dot product between two embedding vectors as the relevance score:

$$X(c) = [CLS]_q[SEP], \quad (3)$$

$$X(s) = [CLS]_s[SEP], \quad (4)$$

$$Emb(X) = transformer_{[CLS]}(X), \quad (5)$$

$$S(q, s) = Emb(X(q))^T \cdot Emb(X(s)), \quad (6)$$

where q is the query, s is the statute, $transformer_{[CLS]}(\cdot)$ outputs a contextualized vector for each token and we select the "[CLS]" vector as the embedding vector of the input. In Equation 6, we regard the inner products of embeddings as the relevance score S .

For the loss function, we use the Softmax Cross Entropy Loss (Cao et al., 2007; Ai et al., 2018; Gao et al., 2021) to optimize the retrieval model, which is defined as:

$$\begin{aligned} \mathcal{L}(Q, s^+, N) \\ = -\log \frac{\exp(S(Q, s^+))}{\exp(S(Q, s^+) + \sum_{s^- \in N} \exp(S(Q, s^-))}, \end{aligned} \quad (7)$$

where S is the relevance score function which is defined in Equation 6. Q is the query, s^+ is the relevant statute and N is the set of irrelevant statutes randomly sampled from the corpus.

H Training Process of the Dense Retrieval Model Distilled from GPT-4

We introduce an approach utilizing GPT-4 to generate labels for question-article pairs. Our methodology leverages GPT-4’s capabilities to autonomously generate non-professional legal questions from statutory articles, thus enabling the pairing of these questions with their corresponding articles without the need for human supervision.

The process begins by selecting statutory articles from the corpus of STARD. GPT-4 is then tasked with generating a legal question based on the content of each article. This is achieved by providing GPT-4 with a specific prompt designed to simulate a scenario in which an individual without prior legal knowledge seeks advice. The prompt instructs GPT-4 to formulate a question that such an individual might ask, ensuring that the question is directly related to and explainable by the content of the statutory article provided. The prompt used in this study is structured as follows:

Prompt 1

Given the following known statutory article:

[Content of the statutory article]

Imagine a scenario in which a person without legal knowledge is seeking legal advice. Please generate a question that this party might ask.

Note: The question must be fully explainable using the statutory article mentioned above, and remember that the person who proposes this question has never read the legal articles mentioned before.

Each interaction with GPT-4 results in the creation of a query-statute pair (q, a_i^+) , where q is the generated question and a_i^+ is the positive statute article to which the question is relevant.

Following the generation of query-statute pairs, we employ a contrastive learning framework to train a dense retriever model. We use the same relevance scoring function S , as detailed in Equation 6, which assesses the relevance of articles to the questions.

In the training phase, for each query Q paired with a positive article a_i^+ , we also sample 8 negative articles from the corpus. These negative sam-

990 ples are not relevant to the query and serve as the
 991 negative set. The loss function employed, repre-
 992 sented by Equation 8, is designed to maximize the
 993 score of the positive article relative to the scores of
 994 the negative samples, effectively training the model
 995 to distinguish between relevant and non-relevant
 996 articles accurately.

$$997 \quad \mathcal{L}(Q, a_i^+, N) = -\log \frac{\exp(S(Q, a_i^+))}{\exp(S(Q, a_i^+) + \sum_{a^- \in N} \exp(S(Q, a^-))}, \quad (8)$$

998 where S is the relevance score function, which is
 999 defined in Equation 6, and N is the set of irrelevant
 1000 statutes randomly sampled from the corpus.

1001 I Training Process of the LSI Classifier

1002 We apply a fine-tuned classifier approach to eval-
 1003 uate the performance of Legal Statute Identifi-
 1004 cation (LSI) methods, as defined in previous
 1005 works (Zhong et al., 2018; Paul et al., 2022). LSI
 1006 is framed as a classification task where each legal
 1007 statute is treated as a distinct label. This transforma-
 1008 tion allows for the classification of legal documents
 1009 or queries by associating them with the relevant
 1010 statutory labels.

1011 Our methodology utilizes a transformer-based
 1012 classification model, specifically fine-tuned on the
 1013 STARD dataset within a five-fold cross-validation
 1014 framework. We initiate our model using the
 1015 Chinese-Roberta-WWM (Cui et al., 2021) for the
 1016 transformer’s parameters, while the parameters for
 1017 the outermost Multi-Layer Perceptron (MLP) layer
 1018 are initialized randomly. The process of input trans-
 1019 formation and subsequent classification is defined
 1020 by the following equations:

$$1021 \quad X(q) = [CLS]q[SEP], \quad (9)$$

$$1022 \quad L(q) = MLP(\text{transformer}_{[CLS]}(X(q))), \quad (10)$$

1023 where q represents a query from the STARD
 1024 dataset. The function $\text{transformer}_{[CLS]}(\cdot)$ first
 1025 encodes the input using the transformer architec-
 1026 ture, focusing on the output of the [CLS] token’s
 1027 embedding vector. The $MLP(\cdot)$ then maps this
 1028 embedding onto the space of statutory labels L .