

# LegalDiscourse: Interpreting When Laws Apply and To Whom

Anonymous ACL submission

## Abstract

While legal AI has made strides in recent years, it still struggles with basic legal concepts: *when* does a law apply? *Who* it applies to? *What* does it do? We take a *discourse* approach to addressing these problems and introduce a novel taxonomy for span-and-relation parsing of legal texts. We create a dataset, *LegalDiscourse* of 602 state-level law paragraphs consisting of 3,715 discourse spans and 1,671 relations. Our trained annotators have an agreement-rate  $\kappa > .8$ , yet few-shot GPT3.5 performs poorly at span identification and relation classification. Although fine-tuning improves performance, GPT3.5 still lags far below human level. We demonstrate the usefulness of our schema by creating a web application with journalists. We collect over 100,000 laws for 52 U.S. states and territories using 20 scrapers we built, and apply our trained models to 6,000 laws using U.S. Census population numbers. We describe two journalistic outputs stemming from this application: (1) an investigation into the increase in liquor licenses following population growth and (2) a decrease in applicable laws under different under-count projections.

## 1 Introduction

AI practitioners have long explored how to analyze legal documents – i.e. laws, court opinions and regulations (Mehl, 1958). Automatic legal question answering<sup>1</sup>, document generation<sup>2</sup>, and motion-filing (Gibbs, 2016) are in commercial use (Dale, 2019) and the legal reasoning capabilities of next generation of NLP models are being assessed (Guha et al., 2023), including by testing against the bar exam (Katz et al., 2023).

However, as noted by Dehio et al. (2022), GPT3 models fail when confronted with simple, yet ambiguous conditions present in legal rules (Bommasani et al., 2021), a challenge well-documented in earlier

<sup>1</sup><https://www.chatbotsecommerce.com/nrf-launches-parker-first-australian-privacy-law-chatbot/>

<sup>2</sup><https://legal.thomsonreuters.com.au/products/contract-express/>, <https://turbotax.intuit.com/>

...in counties having a metropolitan form of government and in counties having a population of not less than three hundred thirty-five thousand (335,000) nor more than three hundred thirty-six thousand (336,000), according to the 1990 federal census or any subsequent federal census, the magistrate or magistrates shall be selected and appointed by and serve at the pleasure of the trial court judge...

Figure 1: Paragraph from a sample law, Tennessee § 36-5-402. The colored blocks represent the following legal discourse elements from our schema: PROBE, TEST, SUBJECT, CONSEQUENCE, OBJECT (see Section 2). We train LLMs to identify these spans and build a web application to aggregate these span tags across state-level laws.

Transformer-based models like BERT as well (Zhong et al., 2020; Holzenberger et al., 2020). Additionally, the majority of legal study has been focused a few domains, like contracts (Koreeda and Manning, 2021; Hendrycks et al., 2021), privacy policy (Wilson et al., 2016; Zimmeck et al., 2019), and corporate law (Wang et al., 2023), and the kinds of tasks heretofore studied have been highly domain specific<sup>3</sup>. Other domains, such as state-level administrative law, remain relatively understudied, despite their importance to policy makers, journalists and academics.

We see the need to introduce a unified mode of study that can quickly incorporate new areas of law. In this work, we develop a uniform discourse schema for characterising a legal text. Discourse analyses, or the study of functional role of text and its relations within in a document (Carlson et al., 2003; Prasad et al., 2008), has been successfully applied to areas like argumentation (Eckle-Kohler et al., 2015), dialogue (Chen and Yang, 2021) and journalism (Spangher et al., 2022, 2021). In journalism, for instance, Choubey et al. (2020) use a unified discourse schema to describe textual relations between diffuse domains of journalism.

In this work, we develop a *legal discourse* schema to address this need, which we apply to state-level legal

<sup>3</sup>An example of a domain-specific task: “Classify if the clause limits the ability of a party to transfer the license being granted to a third party” from Hendrycks et al. (2021).

All counties in the state having having duly adopted a consolidated or metropolitan form of government pursuant to title 7, chapter 1, and all counties of the state having a population of six hundred thousand (600,000) or more, according to the 1970 federal census or any subsequent federal census, shall institute an inmate incentive program for workhouse prisoners.

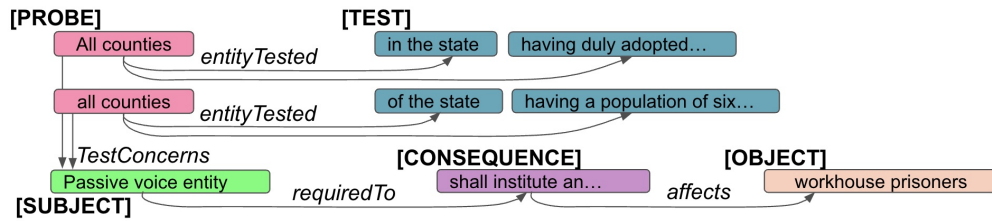


Figure 2: A sample span-and-relation discourse tree generated from a paragraph of legal text. Above, the highlighted text shows the original law text with discourse-spans annotated. Below, relations are drawn between discourse blocks, shown with double-black curved lines and categorically annotated. Note that the **SUBJECT** responsible for carrying out the **CONSEQUENCE** is passively implied.

texts. At the core, our schema seeks to answer the following key questions: (1) When does this law apply? (2) What are the consequences? (3) Who is affected? We show that LLMs struggle to model this schema, yet it is useful for human practitioners. In sum, our contributions are:

In sum, this paper makes three key contributions:

- 1. A Legal Discourse Schema:** We develop a legal discourse schema, consisting of 8 span level and 21 relational classes, some of which are shown in Figure 1 and 2. We annotate 602 state-level laws, with 3,715 spans and 1,671 relations. We show that our schema can be labeled with high inter-annotator agreement, yet fine-tuned LLMs struggle.
- 2. Web Scraping Public Domain U.S. State Law:** We scrape over 100,000 law documents from 52 U.S. states and territories with 20 web-scrappers we build. These scrappers are robustly designed to overcome attempts to block scraping.
- 3. Searching and consuming model output:** We show the practical impact of our work by presenting a web-app we built to help users navigate our dataset. We present two outputs produced by journalists using our interface to study 6,000 laws involving 2020 U.S. Census counts.

We outline our discourse schema and modeling in Section 2. We next discuss our dataset collection process, including the web-scrappers we release for gathering public-domain U.S. state law text (Section 3.1). In Section 3.2 we describe our lightweight and modular span and relation annotation interface which we used to collect data. Next, in Section 6, we describe our web-app, where we surface our model’s output to journalists and engage volunteers to improve our annotations. Finally, we discuss an ongoing use-case to illustrate how one might use our app in Section 6.2.

## 2 A Legal Discourse Schema

A legal rule is a *hypothetical imperative* (Engisch et al., 2018), or a conditional consequence. Reasoning about these rules requires practitioners to understand how and whether conditions of the law are met; what the consequences are (Dehio et al., 2022); and *who* is affected by these consequences.

As shown in Figure 2, modeling the different components of a legal doctrine as *discourse units* and how they interact as *relations* can be an effective way of discern meaning (Carlson et al., 2003; Prasad et al., 2008). Identifying these parts poses a basic test of a model’s legal reasoning and can also lead to practical use-cases (as Spangher et al. (2022) showed in the journalism domain). We introduce the key parts in our schema, starting with span annotations and then relations.

### 2.1 Span-Level Schema

The 8 discourse elements we identify in our schema are **SUBJECT**, **OBJECT**, **PROBE**, **CONSEQUENCE**, **TEST**, **EXCEPTION**, **DEFINITION** and **CLASS**. The first three elements are nearly always be entities (noun phrases), and the rest are predicates (verb phrases) or prepositional phrases.

The first three elements of our schema, **SUBJECT**, **CONSEQUENCE**, and **OBJECT**, capture how law dictates first-degree interactions between entities, inspired by seminal work done by Gardner (1984). We describe each in turn.

- A **SUBJECT** is an entity that gains powers or restrictions under a law. (e.g. “*The trial court judge shall adjudicate property disputes between claimants.*”) Subjects aren’t always explicit, and can be expressed passively<sup>4</sup>.

<sup>4</sup>Example of a Passive **SUBJECT** (and passive **OBJECT**): *Taxes shall be collected at the beginning of every month.* The **SUBJECT** and **OBJECT**, “Tax-collector” and “Tax-payer”, are not actively expressed.

- The **CONSEQUENCE** is the specific power or restriction conferred by the law. Consequences nearly always are attributed to the subject, either passively or explicitly. (e.g. *"The trial court judge shall adjudicate property disputes between claimants."*)
- An **OBJECT** is an entity (noun phrase) affected by the subject, under a law. Typically, when the subject gains powers, the object usually faces more restrictions; if the subject faces restrictions, the object usually faces fewer restrictions. (e.g. *"The trial court judge shall adjudicate property disputes between claimants."*) Like subjects, objects are not always present in the text, or might be expressed passively<sup>5</sup>.

Sometimes, the SUBJECT-CONSEQUENCE-OBJECT involves a longer chain than a 1-hop relationship<sup>6</sup>. In these cases, an entity is both an OBJECT and a SUBJECT. We label this entity as an OBJECT to prioritize the first CONSEQUENCE.

The next three elements in our schema, TEST, PROBE and EXCEPTION, indicate when laws apply.

- A **TEST** is an explicit condition applied to an entity (i.e. an OBJECT, SUBJECT or PROBE) that determines *when* a SUBJECT-CONSEQUENCE-OBJECT relation holds. (e.g. *"In counties with a population above 10,000, the trial court judge shall adjudicate... unless claimants settle."*)
- A **PROBE** is an entity to which a TEST is applied to that is *not* a SUBJECT or an OBJECT. If the TEST is applied to a SUBJECT or an OBJECT, there may not be a need for a PROBE. *"In counties with a population above 10,000, the trial court judge shall adjudicate... unless claimants settle."*
- An **EXCEPTION** is a corollary to a TEST; it specifies when a law does NOT apply. An exception usually modifies a TEST *"In counties with a population above 10,000, the trial court judge shall adjudicate... unless claimants settle."*

Finally, the remaining two classes in our schema, DEFINITION and CLASS, serve to more fully characterize the entities mentioned in legal text. These terms have already been well-described in the literature (Tobia, 2020; Dehio et al., 2022) and incorporated into tasks (Guha et al., 2023). We give definitions in Appendix B. For examples of all span-level discourse types, see Appendix A, Table 6.

<sup>5</sup>Example of a SUBJECT-CONSEQUENCE relation without an OBJECT: *The trial court judge shall begin session at or before 9am.*

<sup>6</sup>Example of a SUBJECT-CONSEQUENCE-OBJECT that is greater than 1-hop: *The magistrate shall designate to the county clerk, who shall adjudicate among taxpayers*

## 2.2 Relational Schema

We define 21 relational categories during our annotation process. There are two categories of relations. (1) The first category occurs between discourse units of *different types*. The type of these relations is usually singular based on the type of the discourse units (e.g. a TEST-PROBE relation means that the TEST is being applied to the PROBE entity), so we do not enumerate them here (we give full definitions in Appendix B). (2) The second category applies between discourse units of the *same type*. These are typically simple grammatical and logical relations. For instance, **sameEntity** indicates that two entities are instances of the same class of entity or the same instance of an entity. **Or**, **And** refers to how two predicate interact (e.g. if test<sub>1</sub> OR test<sub>2</sub> is passed...).

## 3 Dataset Creation

In this section, we describe how we operationalized the schema discussed in Section 2. We scrape a dataset of all state-level laws from 18 states in the U.S., which we discuss in Section 3.1. We then sample a set of paragraphs to annotate. We build an annotation framework, described in Section 3.2, and enlist four annotators, who collectively annotate 602 law paragraphs.

### 3.1 Dataset Construction

Our full legal dataset comprises the more than 100,000 active state-level laws in the United States. We compile this dataset by building a scraper for a public-domain law website called Justia.<sup>7</sup> We then manually audit the output collected by Justia by comparing to state websites and find 19 states where either Justia is incomplete, not updated, or unparsable.<sup>8</sup> We build individual state-level parsers for these states.

State law is public domain,<sup>9</sup> yet it is often inaccessible for bulk downloads and web scraping. For instance, many websites license LexisNexis, a for-profit company, as the official provider for their state codes<sup>10</sup>. Although these websites are publicly accessible, they employ a range of mechanisms (e.g. timeouts, dynamically-generated URLs, cookie-based access) that make them difficult to scrape.<sup>11</sup> To circumvent these, our scrapers are robust and mimic human web-browsing behavior. We develop a generalized scraper for Lexis-

<sup>7</sup><https://www.justia.com/>

<sup>8</sup>Some of the laws provided by Justia, such as those for Colorado, contain data in PDF files (see <https://law.justia.com/codes/colorado/2019/>), which, due to formatting, have a high OCR error rate, so in these cases we extract directly in these cases.

<sup>9</sup><https://fairuse.stanford.edu/overview/public-domain/welcome/>

<sup>10</sup>Ex. Colorado, Georgia and Tennessee: <http://www.lexisnexis.com/hottopics/colorado>, <http://www.lexisnexis.com/hottopics/gacode>, <http://www.lexisnexis.com/hottopics/tncode>

<sup>11</sup>The practical effect of mechanisms to block bulk downloads is the hindrance of law corpora collection for journalistic or academic study.

	% annots	% of docs	# / doc
TEST	28%	91%	2.4
SUBJECT	20%	95%	1.7
CONS.	19%	83%	1.8
OBJECT	15%	69%	1.7
PROBE	9%	46%	1.5
CLASS	6%	34%	1.5
DEF.	2%	11%	1.6
EXC.	1%	6%	1.1

Table 1: The prevalence of different discourse units across our annotated dataset. The left column shows the percentage of units across all annotations. Center shows the percentage of documents in our corpus that have at least one discourse unit. Right shows the average number of units per document, when present.

Nexis Public Access websites using scrapy<sup>12</sup> and selenium-webdriver<sup>13</sup>. In order to scrape Justia, we launch three Google Compute Engine (GCE) instances for a total of 60 compute hours<sup>14</sup>.

### 3.2 Annotation

We recruited 4 annotators, including one former journalist and 2 undergraduate researchers<sup>15</sup>. We trained all of the annotators until they were achieving above an 80% accuracy in both span and relation identification tasks, based on a gold-label set that we constructed. After reaching this agreement level, we begin accepting completed tasks from annotators. Together, the annotators annotated 602 laws, with a 10% overlap, from which we calculated a  $\kappa = .8$

We built a Javascript-based framework to handle span and relation tagging and (1) serve as a standalone web-app for annotators (2) compile to Amazon Mechanical Turk (AMT) tasks<sup>16</sup> (3) integrate into a web-site built for journalists using our work (described in Section 6). Although many NLP-focused annotation tools exist<sup>17</sup> we found that none were flexible enough to be integrated easily into larger websites or automatically generate AMT tasks.<sup>18</sup> We plan to distribute our interface as a

<sup>12</sup><https://scrapy.org/>

<sup>13</sup><https://www.selenium.dev/>.

<sup>14</sup>We will release our code for scraping with Docker images created to perform these scrapes. Given the difficulty in creating this dataset, we believe these routines constitute a considerable resource for academic inquiries into state-level law.

<sup>15</sup>We compensated the undergraduate researchers fairly at a rate of \$20 per hour through AMT, according to University policy

<sup>16</sup>[https://docs.aws.amazon.com/AWSMechTurk/latest/AWSMturkAPI/ApiReference\\_HTMLQuestionArticle.html](https://docs.aws.amazon.com/AWSMechTurk/latest/AWSMturkAPI/ApiReference_HTMLQuestionArticle.html).

<sup>17</sup>There were 87 frameworks as of Neves and Ševa (2021)’s count, including BRAT (Stenetorp et al., 2012), YEDDA (Yang et al., 2017) and WebAnnon (Yimam et al., 2013)

<sup>18</sup>We will release the annotation code as part of this framework

Relation	Percentage
ENTITY $\leftrightarrow$ PREDICATE	61%
ENTITY $\leftrightarrow$ ENTITY	20%
PREDICATE $\leftrightarrow$ PREDICATE	19%

Table 2: Types of relations common in our corpus. ENTITY includes: SUBJECT, OBJECT and PROBE discourse units. PREDICATE includes all others.

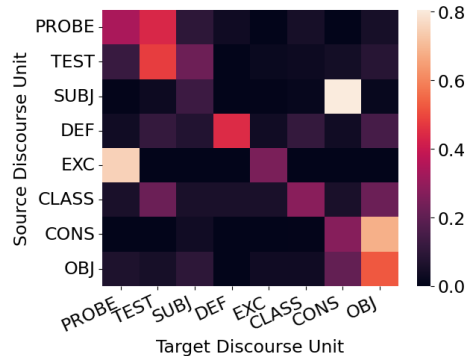


Figure 3: The conditional likelihood of a target discourse class, given a source discourse class. In other words, the color scale is  $p(t|s)$  where  $s$  is the source node and  $t$  is the target node.

stand-alone Javascript package. For more details about the annotation interface, see Appendix C.

### 3.3 Dataset Statistics

**Corpus Description** The length of the legal paragraphs we annotate averages 490 characters. The types of content that we focused on in our sample included topics on Government, education and environment. Certain states in our sample emphasized different topics. For example, California has a higher proportion of laws aimed at Poverty and Development compared with Tennessee, which has a higher proportion of laws focused on Administration (see Appendix A for more information and visualizations).

**Discourse-level Analysis** Discourse unit-level statistics vary widely. As can be seen in Table 1, TEST and SUBJECT are the most common discourse unit, accounting for 48% of all span-level annotations. TEST occurs in 91% documents. Surprisingly, EXCEPTION units were relatively rare, accounting for only 1% of annotations and occurring in only 6% of documents. There are many more TEST units per document, at 2.4 TEST units, than other elements.

**Relation-level Analysis** Next, we analyze the nature of the relations between discourse units. Two discourse spans are much more likely to directly relate if they are closer together in the law text. 62 characters, on average, separate discourse units with relations, while 195 characters, on average, separate all pairs of discourse units without relations. In Section 4.3, we describe how we

275 balance our training datasets to remove this adjacency  
276 bias.

277 Figure 3 shows the likelihood of transitioning to a  
278 target discourse type, given a source discourse type. We  
279 order the  $x$  and  $y$  axes by the most likely starting points  
280 of discourse elements in a document (Note, in Figure  
281 6, that discourse elements that appear first in the  
282 document to be connected with discourse elements later.  
283 See Table 5 in Appendix A for more information). We  
284 see a strong diagonal bias: all discourse elements are  
285 likely to transition to elements of the same type. We also  
286 notice the strong SUBJECT  $\rightarrow$  CONSEQUENCE and  
287 CONSEQUENCE  $\rightarrow$  OBJECT relation, as well as the  
288 PROBE  $\rightarrow$  TEST relation. This reinforces insights by  
289 (Gardner, 1984), (Engisch et al., 2018) and (Dehio et al.,  
290 2022) about the key role of *hypothetical imperative*  
291 language in legal texts (discussed in Section 2).

292 On the other hand, we find that several categories  
293 of relation are simply unlikely to ever occur. For in-  
294 stance, EXCEPTION is almost never applied to CON-  
295 SEQUENCE. We hope in future work to investigate  
296 if these patterns hold up across a wider body of legal  
297 text. See Appendix A for more details. We test the  
298 implication of this in Section 4.3.

## 299 4 Legal Entity and Relational Modeling

300 We frame a new task using the data we collect: Legal En-  
301 tity and Relational Modeling, or *extracting legally signif-*  
302 *icant spans and their relations*. This task is analogous  
303 to end-to-end relation extraction (ERE) (Kameyama,  
304 1997). We will first describe two subtasks that tradi-  
305 tionally compose ERE, and how legal discourse can be  
306 modeled in this framework, then we will discuss meth-  
307 ods, with a particular focus on how we can use this  
308 setup to interrogate the reasoning capabilities of large  
309 language models.

### 310 4.1 Tasks and Datasets

311 **Span-Level Tagging** Given a document  $X$  of  $n$  to-  
312 kens  $x_1, \dots, x_n$ , let  $S = \{s_1, s_2, \dots, s_m\}$  be all possible  
313 spans in  $X$ . Let  $\zeta$  be a set of predefined span-types, in  
314 our case we use a subset of our discourse tags: **SUB-**  
315 **JECT**, **CONSEQUENCE**, **OBJECT**, **TEST**, **PROBE** and  
316 **EXCEPTION**. We focus on these types because they  
317 have within-text consequence, compared with, **DEFINI-**  
318 **TION** and **CLASS**, which are primarily about adding  
319 context and helping to reason across texts (Buey et al.,  
320 2016). Our goal, then, is to predict an entity type  
321  $y_e(s_i) \in \{\zeta, \epsilon\}$ , where  $\epsilon$  is the null class. In legal rea-  
322 soning, this subtask can help test a model’s awareness  
323 of the function of each span of text.

324 We filter our task dataset so that each document has at  
325 minimum two of the primary 6 spans, and we addition-  
326 ally remove spans that are at most one word, as these  
327 were the most ambiguous for our annotators to agree on.  
328 This leaves us with 3,559 spans across 413 documents.  
329 We measure classification accuracy using F1 per class,  
330 and we consider a span to be valid if it contains 80% of

331 more of the same words, after removing stop words and  
332 punctuation, and being no longer than twice in length  
333 as the annotated example.

**Relation Extraction** Let  $R$  be a set of pre-defined  
334 relation types. We seek, for every pair of spans  $s_i \in$   
335  $S, s_j \in S$  to predict a relation-type,  $y_r(s_i, s_j) \in \{R, \epsilon\}$ ,  
336 where  $\epsilon$  is the null class. We consider two versions  
337 of this task: *detection* and *classification*. Detection  
338 involves simply predicting  $y_r(s_i, s_j) \in \{I, \epsilon\}$ , where  $I$   
339 indicates there exists any relation  $\in R$ , and classification  
340 is the classical relation-type detection. This task can  
341 help test a legal model’s ability to identify which spans  
342 are modified by a given span.  
343

344 To construct a challenging legal relation classifica-  
345 tion dataset, we take a subset of relations  $\hat{R} \in R$  that  
346 are observed occurring between span pairs of different  
347 span-types. In other words, we take relations  $r \in \hat{R}$   
348 where  $|\{y_e(s_i), y_e(s_j)\}| > 1 \forall_{i,j} s.t. y_r(s_i, s_j) = r$ .  
349 This allows us to focus less on modeling the semantics  
350 of each span’s type and more on the relation between  
351 them. We additionally sample negative examples, i.e.  
352  $y_e(s_i, s_j) = \epsilon$ .

353 Finally, we notice that discourse units that are more  
354 proximal in the text are more likely to be related, as  
355 noted in Section 3.3. We find in early trials that our  
356 models were overfitting to proximity in text and not  
357 generalizing well to cases where relations are more dis-  
358 tant. So, to make the task more challenging, we sample  
359 negative examples that the same distribution of offsets  
360 our labeled examples. We are left with 1,482 datapoints.  
361 We measure model accuracy using F1, focusing on three  
362 main groupings: relations between entities and entities  
363 (ENT $\leftrightarrow$ ENT), relations between entities and predicates  
364 (ENT $\leftrightarrow$ PRED) and relations between predicates and  
365 predicates (PRED $\leftrightarrow$ PRED).

### 366 4.2 Baselines

367 Relation extraction is a widely studied field, with classi-  
368 cal and current work focusing on modeling each subtask  
369 separately (Sang and De Meulder, 2003; Zelenko et al.,  
370 2003), as well as end-to-end modeling (Li and Ji, 2014).  
371 As such, we build upon two recent methods focused on  
372 each approach:

- 373 • PURE (Zhong and Chen, 2020): separately models  
374 two different embedding spaces, one focused on  
375 span identification and the other focused on rela-  
376 tion extraction, using masked language modeling  
377 (Devlin et al., 2018).
- 378 • ASP (Liu et al., 2022): trains a generative T5 model  
379 (Raffel et al., 2020) to create structured predictions.

### 380 4.3 Generative Modeling

381 Recent work has shown that large language models can  
382 also be effective relation predictors (Wan et al., 2023).  
383 To test this hypothesis, and to add to a growing body of  
384 work focused on benchmarking LLMs for legal tasks  
385 (Guha et al., 2023), we format our tasks as generation

	SUBJECT	CONS	OBJECT	PROBE	TEST	EXC	Macro	Micro
Baselines								
ASP (Liu et al., 2022)	35.7	39.4	26.3	38.9	44.6	33.3	37.7	36.6
PURE (Zhong et al., 2020)	41.5	45.2	25.0	56.1	17.3	36.4	34.3	36.5
GPT3.5								
0-shot	34.4	9.7	14.8	13.4	35.4	54.7	27.1	22.7
3-shot	31.7	23.3	20.4	28.2	43.9	46.2	32.3	30.1
5-shot	30.7	24.1	15.9	30.8	49.8	45.2	32.8	30.8
8-shot	29.7	23.4	15.8	33.5	48.4	53.8	34.1	31.0
GPT fine-tuned	42.1	49.9	35.9	34.9	53.0	56.0	45.3	44.3

Table 3: F1 scores shown for span-identification for our 6 primary discourse elements: SUBJECT, CONSEQUENCE, OBJECT, PROBE, TEST and EXCEPTION. Average Precision, Recall and F1 across all samples are shown. Although fine-tuning improves performance across most categories, leading to +10-point increases in macro and micro f1-scores, although some, like EXCEPTION, are able to be handled relatively well even in zero-shot settings. F1 scores are still below human levels of agreement.

problems and fine-tune GPT3.5 models<sup>19</sup>. For span-prediction, we show that we can recover spans specific to discourse tags, in other words: we model  $p(s|\zeta, X)$ . As each law may contain several discourse elements of each type, we ask the LLM to generate all elements of a certain discourse type in the body of the law (e.g. *You are a legal assistant. I will show you a paragraph of law. Which entities gains powers, restrictions or responsibilities under this law?*) See Appendix F for more examples of prompts. For our evaluation set, we sample such that 30% of data contains no discourse elements of the required type.

For relation detection we ask the model  $p(I|s_1, s_2, X)$ , i.e. whether a relation is present between two spans, and  $p(R|s_1, s_2, X)$ , what the relation is. We still allow  $\epsilon \in R$ , so that our experiments with GPT are comparable to the baseline models. In other words, we construct a prompt where the LLM is given the legal text and two discourse elements, and ask if they are related (e.g. *“Are span A and B related in Law X?”*). We test two different prompt settings. In the first setting, we simply give the two spans of text and the law, and ask the LLM to determine if they are related. In the second setting, we give the LLM the class labels of the discourse units, as well as definitions for what each label means. See Appendix A for examples. We test both tasks in zero-shot, few-shot, and fine-tuned settings and for each test sample, we repeatedly query the LLM for 3 trials, randomizing the few-shot examples it receives.

## 5 Results and Discussion

**Span-Level Tagging** : Table 3 shows F1 scores from our span-tagging experiments. Interestingly, span-tagging appears to be a harder task for GPT: even after fine-tuning, GPT scores below human-level (our anno-

<sup>19</sup>Specifically, we use GPT3.5-turbo as of October 11, 2023.

tators, after conferencing and training). GPT was especially challenged by distinguishing between different entities’ roles: SUBJECT, OBJECT and PROBE (GPT Fine-tuned scores 35-42 F1 on entities, compared with 50-59 F1 for predicates. EXCEPTION stands out as a particular category where even 0-shot GPT performs well.) SUBJECT and OBJECT roles can be particularly ambiguous, as mentioned in Section 2, as there are cases when an entity can be in both a SUBJECT and OBJECT role (we annotated OBJECT, in those cases). Interestingly, too, the gap between GPT and the baseline models is not as large in this task than it is in relational modeling. Perhaps our generative setup for this step,  $p(s|\zeta, X)$ , with 6 different prompts, allowed GPT to generate the same entity for different categories. We might see improvements by disambiguating with another model,  $p(\zeta|s, X)$ , when a single span is generated in multiple categories.

Our broader finding, though, is that this remains a challenging task. Although our task dataset, at 400 documents, is small relative to other language resources, the spans in our schema are syntactically low-level. The spans divide relatively well into different parts of speech, like noun phrases and verb phrases; identifying such chunks in text has long been within the capability of even classical language models (Sang and Buchholz, 2000). Future work either fine-tuning on other resources, or using law-specific models, might show improvements in these areas.

**Relation Identification and Classification** Table 4 show F1 scores from relation detection (Detect) and classification (Class). Relation extraction is a category where fine-tuned GPT performs just as well as our annotators. We notice, too that in some cases GPT does even better on the classification task than it does on the identification task (e.g. ENT $\leftrightarrow$ PRED and ENT $\leftrightarrow$ ENT). It’s possible that the semantics of classification task enforce greater reasoning and justification than the identification

	ENT ↔ PRED		ENT ↔ ENT		PRED ↔ PRED		All (Macro)		All (Micro)	
	Detect	Class.	Detect	Class.	Detect	Class.	Detect	Class.	Detect	Class.
Baselines										
ASP	26.5	14.2	4.5	3.8	4.0	2.2	13.6	6.7	19.5	11.1
PURE	73.9	64.5	15.4	5.3	45.7	38.2	49.5	40.5	63.1	53.9
GPT3.5										
zero-shot	54.9	0.0	42.5	27.1	25.2	23.2	40.8	16.8	48.5	7.2
zero-shot w. def	69.4	0.0	54.2	39.5	60.8	48.2	61.5	29.2	65.1	12.8
few-shot	50.6	55.3	56.8	53.9	40.2	34.2	49.2	47.8	50.5	51.7
few-shot w. def	72.6	60.1	68.5	65.9	65.1	35.2	68.7	53.7	70.8	56.7
GPT finetuned	82.6	85.9	76.5	88.7	81.0	65.9	80.0	80.2	81.1	82.9

Table 4: Relation Detection and Classification F1 score. We examine scores between three categories of relations: ENTITIES ↔ ENTITIES, ENTITIES ↔ PREDICATES, and PREDICATES ↔ PREDICATES. ENTITIES are **SUBJECT**, **OBJECT** and **PROBE**, and PREDICATES are all other discourse types. Classification is only run for discourse-type pairs where more than one relation can exist (see Section 2).

task, like in Wei et al. (2022).

The relation identification task also shows a clear different between the baseline models, which we do not observe in the span-level tagging task. One explanation for the especially poor performance of ASP (Liu et al., 2022) is that the jointly learned model requires the model to make use of more data to fully learn the embedding layers. In fact, tasks that ASP performs well on, like ACE2005 (Sang and De Meulder, 2003), have 10x more documents and annotation than our dataset. We show more details in Appendix E, Figure 7.

## 6 Practical Use Case: Census 2020

To get feedback on our work from a preliminary group of users, we apply our models to a domain of state-level law pertinent to journalists. In 2020, the U.S. Census count faced multiple challenges, and many researchers hypothesized that populations, especially minorities, might be inaccurately counted (Naylor, 2020; Mervis, 2019; Berry-James et al., 2020). Scant insight existed, especially on the state-level, into how population counts were being used in law<sup>20</sup>: the corpus of state-level laws was too large and varied for journalists to parse.

On the other hand, this provided an interesting case for discourse-based reasoning. Population counts typically get used as a relatively unambiguous TEST. Our discourse models help us identify this occurring, and then we can develop ways to parse out the specific ways population is in TEST discourse. We describe the website we built to facilitate different explorations, and then we describe two such explorations that we received permission from the journalists collaborating with us to write about. We will focus on our own contributions in these collaborations.

<sup>20</sup>Besides federal budgeting and Congressional representation, which have already been manually programatized (Reamer, 2018; Berry-James et al., 2020).

### 6.1 Website Design

We design a website to enable exploration of our dataset and modeling output. Users can perform full-text search on all laws in our database, view discourse spans schema by extracting spans across laws and correct or provide new annotations. The website’s overall goal is to facilitate both *deep explorations* and *wide explorations*.

**Going deep** : Going “*deep*” here essentially means *subsetting the laws* first, and then analyzing discourse. The web search functionality<sup>21</sup> helps users do this by exploring a specific term or concept in the law’s plain text or in specific discourse role (e.g. laws affecting OBJECT=“taxpayer”). After the user finds samples of an interesting subset of laws, we do a broader study of them them by using our discourse models to answer: *who* is being affect, under *what conditions*, and *how*?

**Going wide** : Conversely, going “*wide*” means studying discourse units and relations first, then analyzing the laws. The website includes a second functionality: allowing users can view aggregate counts of different discourse units and relations. This helps users notice patterns among the ways in which discourse was being used. After a user notices a specific pattern in discourse roles (e.g. EXCEPTION units modifying TEST units about taxes), then we can analyze the laws that include, or do not include, these elements.

In both flows, visitors can access our annotation framework, described in Section 3.2, which helped us gather more data (to be used in the future). For more on the design of our website, see Appendix D.

### 6.2 Case Study #1: Going Deep (Liquor Store Licenses)

We describe two example articles that are currently being explored by users of our system.

<sup>21</sup>Powered by ElasticSearch (Elasticsearch, 2018)

In the first example, journalists hypothesized that the allocation of new liquor licenses might be population-based. To explore this, they used the search interface; they searched for the term “alcohol OR liquor OR beverage” in the search interface and discovered that interface returned 270 laws. Together, we analyzed the breakdown of liquor-related law by state. We found that the states most likely to base liquor licenses off population counts were Tennessee, New York and Illinois. They then asked us to extract all TESTS from these laws. We found that mid-size cities would be the most likely to be impacted by a 5% or 10% undercount in population. The journalists identified key cities and are seeking sources in these areas.

### 6.3 Case Study #2: Going Wide (Slim Population Thresholds)

In another example, journalists explored the top-level discourse annotations. They noticed that some TESTS are based on explicit population thresholds (ex. Figure 1) and that some of these thresholds were very narrow. Working together, we compiled several keyword filters and regular expressions extract specific population thresholds. We found that, in Tennessee in particular, over 40% of all Census-related laws imposed narrow population tests of fewer than 500 people and 10% imposed tests of fewer than 100 people. This raised questions: what is the purpose of these narrowly targeted laws? Were they trying to target specific counties without mentioning them by name? The journalists are now investigating further by tracking down the authors of these laws.

## 7 Related Work

Although the field of AI-driven legal aids is multifaceted and growing (Kauffman and Soares, 2020), free and open-source frameworks remain few (Morris, 2019; Dale, 2019; Vergottini, 2011). Our discourse-driven web application, designed for legal exploratory analysis is one of the few AI-powered, free applications that exist, and the first to open source tools for legal document collection.

For-profit legal inquiry systems, as mentioned above, are numerous. Bloomberg Law<sup>22</sup>, Westlaw<sup>23</sup>, LexisNexis<sup>24</sup> and Wolters Kluwar<sup>25</sup> are the four main services for legal research (Dale, 2019), which provide subscription-based, Google-style searches. CaseText<sup>26</sup> and Ravel<sup>27</sup> were two upstart case-text search engines (although both have now been aquired); CaseText offered crowdsourced annotations and Ravel linked cases together to create visual maps of important cases (Lee

et al., 2015). We similarly provide a way of collecting user-annotations, and a novel way linking together cases, although ours takes a discourse approach rather than an unsupervised clustering approach.

Various discourse schemas have been developed to understand law texts, including deontological logic-based schemas (Wyner and Peters, 2011; Zeni et al., 2015), and subject matter-specific schemas (Espejo-Garcia et al., 2019). Ours is the first discourse-based approach to take steps towards a big-data approach by setting up a framework for the ingestion of crowdsourced annotations.

Finally, outside of the legal domain, other areas have experienced a growth in academically-oriented systems for human-in-the-loop inquiry. The COVID-19 pandemic has produced a burst in NLP-driven corpora-collection (Wang et al., 2020), demonstrations (Sohrab et al., 2020; Hope et al., 2020; Spangher et al., 2020) and workshops (Verspoor et al., 2020b,a).

Such concerted effort in the NLP domain to expose resources and build open tools for subject matter experts is an inspiring guide for how NLP researchers can contribute to wider inquiries. We hope such efforts expand to other domains as well, forming a common alliance between academics, civil-minded journalists and other researchers and end-users.

## 8 Conclusion

We have sought to take steps towards a semantic understanding of legal texts, a goal long held in computational law (Gardner, 1984). We show that large language models, while achieving impressive results in some parts of our task, show surprisingly weak performance compared to human annotators in others. Language models have an important role to play in interpreting law and lowering the barrier of access to legal systems for citizens, journalists and academics. Our task is an important step towards assessing a sturdy foundation and opening the door to more intensive legal tasks be considered (Guha et al., 2023).

In this work, we have additionally presented three open-source components. (1) A web-app exposing a novel discourse schema and its application to state law referencing U.S. Census counts. (2) A flexible and modular annotation framework that can be seamlessly embedded into web-apps to allow visitors to contribute and update annotations. (3) A set of web-scrapers to help researchers gather public-domain legal text. We demonstrated concrete utility to facilitate journalistic exploration with our discourse schema. Our longer-term goal is to collect feedback and data, and improve our database and machine learning systems. We hope that such efforts can continue to push Legaltech (Hartung et al., 2017) into a more open and accessible domain, and make it easier to understand the laws governing our society.

<sup>22</sup><https://pro.bloomberglaw.com/>  
<sup>23</sup><https://www.westlaw.com/>  
<sup>24</sup><https://www.lexisnexis.com>  
<sup>25</sup><https://www.wolterskluwer.com>  
<sup>26</sup><https://casetext.com/>  
<sup>27</sup><https://home.ravellaw.com/>



## 9 Impact Statement

There were several possible ethical considerations we encountered during this research which we wish to address.

**Dataset Creation:** The creation of our dataset involved scraping numerous websites, including state websites, state-licensed LexisNexis pages and <https://www.justia.com>. In the third case, Justia, we did not violate any terms of service. In fact, Justia’s `robots.txt` file<sup>28</sup> is the most permissive possible, giving unlimited license to any crawler. It is generally accepted that `robots.txt` files are implied licenses of access,<sup>29</sup> and we did not disregard Justia’s file before scraping.

Content derived from the first two categories, state law websites and official, state-licensed websites like LexisNexis are, by law, public domain (Wolfe, 2019; MacWright, 2013). Web-scraping the public domain is neither illegal nor unethical (Mehta, 2021). As we did in the body of the paper, we again emphatically criticize attempts by providers to make web-scraping difficult, and we went to lengths to overcome this.

**Dataset Annotation:** All parties involved in annotating our dataset received valid compensation. We relied entirely on expert researchers to collect our annotations. This included the authors of this paper. All the researchers who provided annotations for us were affiliated with our institution and compensated appropriately by our institution (we leave the determination of “appropriate” for our institution to define.)

Although we describe accommodating AMT tasks in the body of the paper, thus far, we have not used any annotations made by Turkers on AMT or by journalists/researchers using our site. If we do, we will ensure there are no ethical issues by securing university IRB approval or exemption, as deemed fit by the IRB. For the Turkers, we will calculate a payment that equals, on average, \$15 an hour. For the journalists/researchers, we will have exchanged something of value (the use of our web-app) for the annotation.

**Website Usage:** Our website has significant accessibility limitations for the seeing-impaired and for non-English speakers. We have not addressed them in this current version, but are mindful and actively searching for options to expand accessibility.

There are two ways in which seeing-impaired users might suffer. First, blind users will not be able to read any of the site without external tools, as we have not recorded or built in any native audio-scripts, keyboard shortcuts or voice-activated commands. Besides “not containing irrelevant information” (Giraud et al., 2018), we can do more to audit our website (Tosaka, 2005) and organize the flow on our site to increase blind accessibil-

<sup>28</sup>Found here <https://www.justia.com/robots.txt>. Such files govern the site-owners’ standards for scraping and crawling.

<sup>29</sup><https://stackoverflow.com/questions/999056/ethics-of-robots-txt>

ity. Secondly, part of our website introduces users to our discourse schema by introducing them to color-coded segments of text. We are actively investigating color-schemes and other approaches that are more amenable to color-blind individuals, of which there is extensive research (Wakita and Shimamura, 2005; Jambor et al., 2021; Foti and Santucci, 2009). Because of the prototype nature of this website, we have not yet investigated these, but they are crucial next-steps.

Our website focuses on U.S.-based laws and contains only English-language text. We do not attempt, in this version, to perform translations. Our plan in the present iteration of this work was to work with U.S.-based journalists studying U.S.-based law. We have not yet undertaken a study to compare how well our discourse schema would apply to non-U.S. law, be it common or civil (Dainow, 1966). However, if this approach proves useful for journalists and researchers, we will certainly seek to undertake this.

## References

- RaJade M Berry-James, Susan T Gooden, and Richard Gregory Johnson III. 2020. Civil rights, social equity, and census 2020. *Public Administration Review*, 80(6):1100–1108.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.
- María Granados Buey, Angel Luis Garrido, Carlos Bobed, and Sergio Ilarri. 2016. The ais project: Boosting information extraction from legal documents by using ontologies. In *ICAART (2)*, pages 438–445.
- Lynn Carlson, Daniel Marcu, and Mary Ellen Okurowski. 2003. Building a discourse-tagged corpus in the framework of rhetorical structure theory. *Current and new directions in discourse and dialogue*, pages 85–112.
- Jiaao Chen and Diyi Yang. 2021. Structure-aware abstractive conversation summarization via discourse and action graphs. *arXiv preprint arXiv:2104.08400*.
- Prafulla Kumar Choubey, Aaron Lee, Ruihong Huang, and Lu Wang. 2020. Discourse as a function of event: Profiling discourse structure in news articles around the main event. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.
- Joseph Dainow. 1966. The civil law and the common law: some points of comparison. *Am. J. Comp. L.*, 15:419.
- Robert Dale. 2019. Law and word order: NLP in legal tech. *Natural Language Engineering*, 25(1):211–217.

737	Niklas Dehio, Malte Ostendorff, and Georg Rehm.	Markus Hartung, Micha-Manuel Bues, and Gernot Hal-	792
738	2022. Claim extraction and law matching for covid-	bleib. 2017. <i>Legal tech</i> . CH Beck.	793
739	19-related legislation. In <i>Proceedings of the Thir-</i>		
740	<i>teenth Language Resources and Evaluation Confer-</i>	Dan Hendrycks, Collin Burns, Anya Chen, and	794
741	<i>ence</i> , pages 480–490.	Spencer Ball. 2021. Cuad: An expert-annotated	795
		nlp dataset for legal contract review. <i>arXiv preprint</i>	796
		<i>arXiv:2103.06268</i> .	797
742	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	Nils Holzenberger, Andrew Blair-Stanek, and Benjamin	798
743	Kristina Toutanova. 2018. Bert: Pre-training of deep	Van Durme. 2020. A dataset for statutory reasoning	799
744	bidirectional transformers for language understand-	in tax law entailment and question answering. <i>arXiv</i>	800
745	ing. <i>arXiv preprint arXiv:1810.04805</i> .	<i>preprint arXiv:2005.05257</i> .	801
746	Judith Eckle-Kohler, Roland Kluge, and Iryna Gurevych.	Tom Hope, Jason Portenoy, Kishore Vasani, Jonathan	802
747	2015. On the role of discourse markers for discrim-	Borchardt, Eric Horvitz, Daniel Weld, Marti Hearst,	803
748	inating claims and premises in argumentative dis-	and Jevin West. 2020. <i>SciSight: Combining faceted</i>	804
749	course. In <i>Proceedings of the 2015 Conference on</i>	<i>navigation and research group detection for COVID-</i>	805
750	<i>Empirical Methods in Natural Language Processing</i> ,	<i>19 exploratory scientific search</i> . In <i>Proceedings of the</i>	806
751	pages 2236–2242.	<i>2020 Conference on Empirical Methods in Natural</i>	807
		<i>Language Processing: System Demonstrations</i> , pages	808
752	BV Elasticsearch. 2018. Elasticsearch. <i>software</i> ], <i>ver-</i>	135–143, Online. Association for Computational Lin-	809
753	<i>sion</i> , 6(1).	guistics.	810
754	Karl Engisch, Thomas Würtenberger, and Dirk Otto.	Helena Jambor, Alberto Antonietti, Bradly Alicea,	811
755	2018. <i>Einführung in das juristische Denken</i> .	Tracy L Audisio, Susann Auer, Vivek Bhardwaj,	812
756	Kohlhammer Verlag.	Steven J Burgess, Iuliia Ferling, Małgorzata Anna	813
		Gazda, Luke H Hoepfner, et al. 2021. Creating clear	814
757	Borja Espejo-Garcia, Francisco J Lopez-Pellicer, Javier	and informative image-based figures for scientific	815
758	Lacasta, Ramón Piedrafita Moreno, and F Javier	publications. <i>PLoS biology</i> , 19(3):e3001161.	816
759	Zarazaga-Soria. 2019. End-to-end sequence labeling		
760	via deep learning for automatic extraction of agri-	Megumi Kameyama. 1997. Recognizing referential	817
761	cultural regulations. <i>Computers and Electronics in</i>	links: An information extraction perspective. <i>arXiv</i>	818
762	<i>Agriculture</i> , 162:106–111.	<i>preprint cmp-lg/9707009</i> .	819
763	Antonella Foti and Giuseppe Santucci. 2009. Increasing	Daniel Martin Katz, Michael James Bommarito, Shang	820
764	web accessibility through an assisted color specifica-	Gao, and Pablo Arredondo. 2023. Gpt-4 passes the	821
765	tion interface for colorblind people. <i>xD&amp;A</i> , 5:41–48.	bar exam. <i>Available at SSRN 4389233</i> .	822
766	Anne von der Lieth Gardner. 1984. Artificial intelli-	Marcos Eduardo Kauffman and Marcelo Negri Soares.	823
767	gence approach to legal reasoning. Technical report,	2020. <i>AI in legal services: new trends in AI-enabled</i>	824
768	Stanford Univ.	<i>legal services. Service Oriented Computing and Ap-</i>	825
		<i>plications</i> , 14(4):223–226.	826
769	Samuel Gibbs. 2016. <a href="#">Chatbot lawyer overturns 160,000</a>	Yuta Koreeda and Christopher D Manning. 2021.	827
770	<a href="#">parking tickets in london and new york</a> . <i>The</i>	<i>Contractnli: A dataset for document-level natural</i>	828
771	<i>Guardian</i> .	<i>language inference for contracts. arXiv preprint</i>	829
		<i>arXiv:2110.01799</i> .	830
772	Stéphanie Giraud, Pierre Théroutanne, and Dirk D	Katrina June Lee, Susan Azyndar, and Ingrid AB Matt-	831
773	Steiner. 2018. Web accessibility: Filtering redundant	son. 2015. A new era: Integrating today’s next gen	832
774	and irrelevant information improves website usability	research tools ravel and casetext in the law school	833
775	for blind users. <i>International Journal of Human-</i>	classroom. <i>Rutgers Computer &amp; Tech. LJ</i> , 41:31.	834
776	<i>Computer Studies</i> , 111:23–35.		
777	Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher	Qi Li and Heng Ji. 2014. Incremental joint extraction of	835
778	Ré, Adam Chilton, Aditya Narayana, Alex Chohlas-	entity mentions and relations. In <i>Proceedings of the</i>	836
779	Wood, Austin Peters, Brandon Waldon, Daniel N.	<i>52nd Annual Meeting of the Association for Computa-</i>	837
780	Rockmore, Diego Zambrano, Dmitry Talisman, Enam	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages	838
781	Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gre-	402–412.	839
782	gory M. Dickinson, Haggai Porat, Jason Hegland,	Tianyu Liu, Yuchen Jiang, Nicholas Monath, Ryan Cot-	840
783	Jessica Wu, Joe Nudell, Joel Niklaus, John Nay,	terrell, and Mrinmaya Sachan. 2022. Autoregressive	841
784	Jonathan H. Choi, Kevin Tobia, Margaret Hagan,	structured prediction with language models. <i>arXiv</i>	842
785	Megan Ma, Michael Livermore, Nikon Rasumov-	<i>preprint arXiv:2210.14698</i> .	843
786	Rahe, Nils Holzenberger, Noam Kolt, Peter Hender-		
787	son, Sean Rehaag, Sharad Goel, Shang Gao, Spencer	Tom MacWright. 2013. <a href="#">State law is public domain.</a>	844
788	Williams, Sunny Gandhi, Tom Zur, Varun Iyer, and	<a href="#">what’s public domain?</a> <i>macwright.com</i> .	845
789	Zehua Li. 2023. <a href="#">Legalbench: A collaboratively built</a>		
790	<a href="#">benchmark for measuring legal reasoning in large</a>		
791	<a href="#">language models</a> .		

846	Lucien Mehl. 1958. <i>Automation in the legal world</i> .	Alexander Spangher, Nanyun Peng, Jonathan May, and	900
847	National Physical Laboratory.	Emilio Ferrara. 2020. <a href="#">Enabling low-resource transfer</a>	901
848	Ivan Mehta. 2021. <a href="#">Us court says scraping a site without</a>	<a href="#">learning across COVID-19 corpora by combining</a>	902
849	<a href="#">permission isn't illegal</a> . <i>TNW   Security</i> .	<a href="#">event-extraction and co-training</a> . In <i>Proceedings of</i>	903
850	Jeffrey Mervis. 2019. <a href="#">Census citizenship ques-</a>	<a href="#">the 1st Workshop on NLP for COVID-19 at ACL 2020</a> ,	904
851	<a href="#">tion is dropped, but challenges linger</a> . <i>Science</i> ,	Online. Association for Computational Linguistics.	905
852	365(6450):211–211.		
853	Jason Morris. 2019. <a href="#">Making mischief with open-source</a>	Pontus Stenetorp, Sampo Pyysalo, Goran Topić,	906
854	<a href="#">legal tech: Radiant law</a> .	Tomoko Ohta, Sophia Ananiadou, and Jun'ichi Tsujii.	907
855	Lorenda A Naylor. 2020. Counting an invisible class	2012. <a href="#">brat: a web-based tool for NLP-assisted text</a>	908
856	of citizens: The lgbt population and the us census.	<a href="#">annotation</a> . In <i>Proceedings of the Demonstrations</i>	909
857	<i>Public Integrity</i> , 22(1):54–72.	<a href="#">at the 13th Conference of the European Chapter of</a>	910
858	Mariana Neves and Jurica Ševa. 2021. An extensive	<a href="#">the Association for Computational Linguistics</a> , pages	911
859	review of tools for manual annotation of documents.	102–107, Avignon, France. Association for Computa-	912
860	<i>Briefings in bioinformatics</i> , 22(1):146–163.	tional Linguistics.	913
861	Alessandra Potrich and Emanuele Pianta. 2008. L-isa:	Kevin P Tobia. 2020. Testing ordinary meaning. <i>Harv.</i>	914
862	Learning domain specific isa-relations from the web.	<i>L. Rev.</i> , 134:726.	915
863	In <i>LREC</i> .	V Yasuaki Takamoto V Hideki Tosaka. 2005. Web	916
864	Rashmi Prasad, Nikhil Dinesh, Alan Lee, Eleni Milt-	<a href="#">accessibility diagnosis tools</a> . <i>Fujitsu Sci. Tech. J.</i> ,	917
865	sakaki, Livio Robaldo, Aravind K Joshi, and Bon-	41(1):115–122.	918
866	nie L Webber. 2008. The penn discourse treebank	Grant Vergottini. 2011. <a href="#">To go open source or not?</a>	919
867	2.0. In <i>LREC</i> .	Karin Verspoor, Kevin Bretonnel Cohen, Michael Con-	920
868	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	way, Berry de Bruijn, Mark Dredze, Rada Mihalcea,	921
869	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	and Byron Wallace, editors. 2020a. <i>Proceedings of</i>	922
870	Wei Li, and Peter J Liu. 2020. Exploring the limits	<a href="#">the 1st Workshop on NLP for COVID-19 (Part 2) at</a>	923
871	of transfer learning with a unified text-to-text trans-	<a href="#">EMNLP 2020</a> . Association for Computational Lin-	924
872	former. <i>The Journal of Machine Learning Research</i> ,	guistics, Online.	925
873	21(1):5485–5551.	Karin Verspoor, Kevin Bretonnel Cohen, Mark Dredze,	926
874	Andrew Reamer. 2018. Counting for dollars 2020: the	Emilio Ferrara, Jonathan May, Robert Munro, Cecile	927
875	role of the decennial census in the geographic distri-	Paris, and Byron Wallace, editors. 2020b. <i>Proceed-</i>	928
876	bution of federal funds. <i>Initial Analysis</i> , 16.	<a href="#">ings of the 1st Workshop on NLP for COVID-19 at</a>	929
877	Erik F Sang and Sabine Buchholz. 2000. Introduction	<a href="#">ACL 2020</a> . Association for Computational Linguis-	930
878	to the conll-2000 shared task: Chunking. <i>arXiv preprint</i>	tics, Online.	931
879	<i>cs/0009008</i> .	Ken Wakita and Kenta Shimamura. 2005. Smartcolor:	932
880	Erik F Sang and Fien De Meulder. 2003. Introduction	<a href="#">disambiguation framework for the colorblind</a> . In <i>Pro-</i>	933
881	to the conll-2003 shared task: Language-independent	<a href="#">ceedings of the 7th International ACM SIGACCESS</a>	934
882	named entity recognition. <i>arXiv preprint cs/0306050</i> .	<a href="#">Conference on Computers and Accessibility</a> , pages	935
883	Mohammad Golam Sohrab, Khoa Duong, Makoto	158–165.	936
884	Miwa, Goran Topić, Ikeda Masami, and Takamura	Zhen Wan, Fei Cheng, Zhuoyuan Mao, Qianying	937
885	Hiroya. 2020. <a href="#">BANNERD: A neural named entity</a>	Liu, Haiyue Song, Jiwei Li, and Sadao Kurohashi.	938
886	<a href="#">linking system for COVID-19</a> . In <i>Proceedings of the</i>	2023. <a href="#">Gpt-re: In-context learning for relation ex-</a>	939
887	<i>2020 Conference on Empirical Methods in Natural</i>	<a href="#">traction using large language models</a> . <i>arXiv preprint</i>	940
888	<i>Language Processing: System Demonstrations</i> , pages	<i>arXiv:2305.02105</i> .	941
889	182–188, Online. Association for Computational Lin-	Lucy Lu Wang, Kyle Lo, Yoganand Chandrasekhar,	942
890	guistics.	Russell Reas, Jiangjiang Yang, Doug Burdick, Darrin	943
891	Alexander Spangher, Jonathan May, Sz-rung Shiang,	Eide, Kathryn Funk, Yannis Katsis, Rodney Michael	944
892	and Lingjia Deng. 2021. Multitask learning for class-	Kinney, Yunyao Li, Ziyang Liu, William Merrill,	945
893	imbalanced discourse classification. <i>arXiv preprint</i>	Paul Mooney, Dewey A. Murdick, Devvret Rishi,	946
894	<i>arXiv:2101.00389</i> .	Jerry Sheehan, Zhihong Shen, Brandon Stilson,	947
895	Alexander Spangher, Yao Ming, Xinyu Hua, and	Alex D. Wade, Kuansan Wang, Nancy Xin Ru Wang,	948
896	Nanyun Peng. 2022. <a href="#">Sequentially controlled text</a>	Christopher Wilhelm, Boya Xie, Douglas M. Ray-	949
897	<a href="#">generation</a> . In <i>Findings of the Association for Com-</i>	mond, Daniel S. Weld, Oren Etzioni, and Sebastian	950
898	<i>putational Linguistics: EMNLP 2022</i> , pages 6848–	Kohlmeier. 2020. <a href="#">CORD-19: The COVID-19 open</a>	951
899	6866.	<a href="#">research dataset</a> . In <i>Proceedings of the 1st Work-</i>	952
		<a href="#">shop on NLP for COVID-19 at ACL 2020</a> , Online.	953
		Association for Computational Linguistics.	954

955 Steven H Wang, Antoine Scardigli, Leonard Tang,  
956 Wei Chen, Dimitry Levkin, Anya Chen, Spencer  
957 Ball, Thomas Woodside, Oliver Zhang, and Dan  
958 Hendrycks. 2023. Maud: An expert-annotated legal  
959 nlp dataset for merger agreement understanding.  
960 *arXiv preprint arXiv:2301.00876*.

961 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten  
962 Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,  
963 et al. 2022. Chain-of-thought prompting elicits reasoning  
964 in large language models. *Advances in Neural  
965 Information Processing Systems*, 35:24824–24837.

966 Shomir Wilson, Florian Schaub, Aswath Abhilash  
967 Dara, Frederick Liu, Sushain Cherivirala, Pedro Gio-  
968 vanni Leon, Mads Schaarup Andersen, Sebastian  
969 Zimmeck, Kanthashree Mysore Sathyendra,  
970 N Cameron Russell, et al. 2016. The creation and  
971 analysis of a website privacy policy corpus. In *Pro-  
972 ceedings of the 54th Annual Meeting of the Association  
973 for Computational Linguistics (Volume 1: Long  
974 Papers)*, pages 1330–1340.

975 Jan Wolfe. 2019. [U.s. high court to rule on scope of  
976 copyright for legal codes](#).

977 Adam Z Wyner and Wim Peters. 2011. On rule extrac-  
978 tion from regulations. In *JURIX*, volume 11, pages  
979 113–122.

980 Jie Yang, Yue Zhang, Linwei Li, and Xingxuan Li. 2017.  
981 Yedda: A lightweight collaborative text span annota-  
982 tion tool. *arXiv preprint arXiv:1711.03759*.

983 Seid Muhie Yimam, Iryna Gurevych, Richard Eckart  
984 de Castilho, and Chris Biemann. 2013. Webanno: A  
985 flexible, web-based and visually supported system for  
986 distributed annotations. In *Proceedings of the 51st  
987 Annual Meeting of the Association for Computational  
988 Linguistics: System Demonstrations*, pages 1–6.

989 Dmitry Zelenko, Chinatsu Aone, and Anthony  
990 Richardella. 2003. Kernel methods for relation ex-  
991 traction. *Journal of machine learning research*,  
992 3(Feb):1083–1106.

993 Nicola Zeni, Nadzeya Kiyavitskaya, Luisa Mich,  
994 James R Cordy, and John Mylopoulos. 2015. Gaiust:  
995 supporting the extraction of rights and obligations for  
996 regulatory compliance. *Requirements engineering*,  
997 20(1):1–22.

998 Haoxi Zhong, Chaojun Xiao, Cunchao Tu, Tianyang  
999 Zhang, Zhiyuan Liu, and Maosong Sun. 2020. Jec-  
1000 qa: a legal-domain question answering dataset. In  
1001 *Proceedings of the AAAI Conference on Artificial  
1002 Intelligence*, volume 34, pages 9701–9708.

1003 Zexuan Zhong and Danqi Chen. 2020. A frustrat-  
1004 ingly easy approach for entity and relation extraction.  
1005 *arXiv preprint arXiv:2010.12812*.

1006 Sebastian Zimmeck, Peter Story, Daniel Smullen, Ab-  
1007 hilasha Ravichander, Ziqi Wang, Joel R Reidenberg,  
1008 N Cameron Russell, and Norman Sadeh. 2019. Maps:  
1009 Scaling privacy compliance analysis to a million apps.  
1010 *Proc. Priv. Enhancing Tech.*, 2019:66.

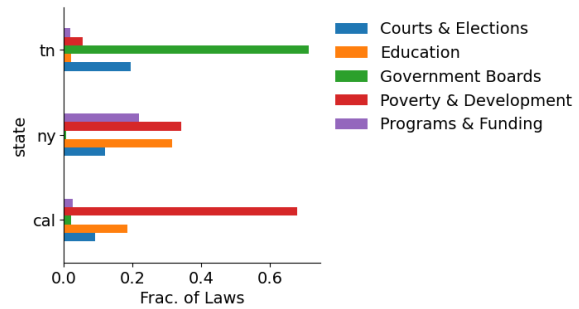


Figure 4: Topic model run over our corpus, showing 3 states. Topics are manually labeled by analyzing top words.

label	start	end
PROBE	28%	32%
TEST	36%	56%
SUBJECT	36%	47%
DEFINITION	41%	55%
EXCEPTION	42%	50%
CLASS	45%	61%
CONSEQUENCE	46%	56%
OBJECT	51%	59%

Table 5: Average start and end character positions of discourse units.

## A Additional Data Analysis

We give analysis of the corpus we collected. In Table 5, we show average character positions of discourse units in the document, as a percentage of the length of the document. PROBE is most likely to occur first in a document, followed by TEST. Discourse units are less likely to occur in the second half of the document.

We examine attributes of relations between discourse elements in Figure 3 and 5. Figure 3 shows the likelihood of transitioning to a target discourse type conditioned on a source type. In Figure 2, we observe there is a strong bias for discourse elements that appear first in the document to be connected with discourse elements later. We order the  $x$  and  $y$  axes by the most likely starting points, as given in Table 5. We see a strong diagonal bias: all discourse elements are likely to transition to elements of the same type. We also notice the strong SUBJECT  $\rightarrow$  CONSEQUENCE  $\rightarrow$  OBJECT relation, as well as the PROBE  $\rightarrow$  TEST relation.

## B Additional Schema Definitions

### B.1 Span-Level Schema: Minor Classes

- A **DEFINITION** is a span of text serving to clarify the ordinary meaning (Tobia, 2020) of a term used in the legal text (e.g. “*Qualified taxpayer*” means a person or entity engaged in a trade or business within...)”)

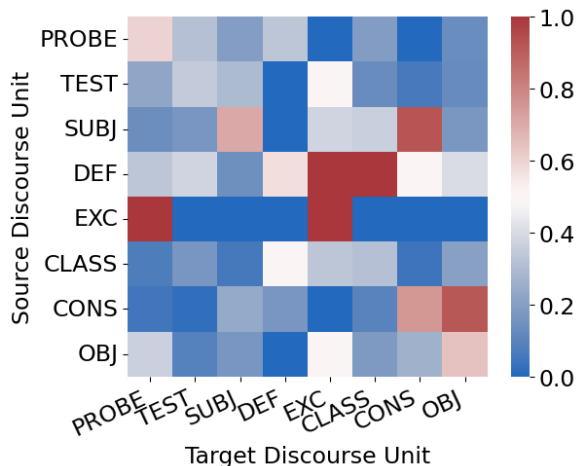


Figure 5: Likelihood of relations, normalized by the random chance of a relation occurring in a document. In other words, if  $n_p$  is the count of all annotated pairs of relations and  $n_r$  is a count of all randomly occurring pairs from a random sample of all  $\binom{n}{2}$  pairs of discourse units in a law text, then the color scale is  $\frac{n_p}{n_p+n_r}$ . Values  $> .5$  are *more* likely to be paired than random chance.

- The **CLASS** of an entity is a modifier that serves to disambiguate the entity from other entities. In knowledge-graph terms, CLASS specifies the source node in an *isA*-type relationship (Potrich and Pianta, 2008); specifically, the entity with a CLASS tag is a subclass of the entity without the class tag (e.g. “The **trial court** judge shall, in...”).

## B.2 Relational Schema

We define 21 relational categories during our annotation process. The first category is of relations that occur between text-spans of different types.

- **EntityEmpoweredTo, EntityRequiredTo:** Indicates which SUBJECT entity (or, in rarer cases, an OBJECT entity) receives powers or responsibilities under a CONSEQUENCE.
- **Affects, AffectedBy:** Indicates which OBJECT entity or entities are affected the CONSEQUENCE of the law, usually mediated through the SUBJECT.
- **TestConcerns, Entitytested:** Indicates which entity a TEST is applied to. This relation typically establishes conditions governing many other entities in the law, not just the entity tested; this is especially the case when the entity is a PROBE.
- **ExceptionAppliesTo, ConditionalTest, ConditionalConsequence:** Indicates which discourse unit (or, in some cases, a second TEST) is applied to. EXCEPTION and multiple levels of TEST are broadly applied to all different kinds of discourse units.

- **Comparison:** A more nebulous relational category, this forms the basis of how the conditional applicability of a law is assessed. Usually found in TEST relations, this relation-type occurs when an attribute of an entity is measured in order to make a determination about whether the conditions are satisfied and the law may be applied. This relation inspired by original attempts to support programmatic legal analysis (Gardner, 1984).

- **EntityHasProperty, PropertyOf:** When any has a particular attribute or CLASS (can be used along with the **Comparison** relation and a TEST).

- **IsDefinedAs, DefinitionOf:** Indicates the entity being defined by a DEFINITION.

The second category of relations typically applies to spans of the same type:

- **SameEntity:** Indicates that two entities are either separate instances of the same class of entity, or they literally refer to the same instance of an entity in legal text.

- **Continuation:** Indicates that two disjoint spans of text refer to the same discourse unit. Can occur when a span is split by another discourse unit.

- **FollowedBy:** When one predicate is conducted or evaluated after the other, in logical order (e.g. in a CONSEQUENCE-CONSEQUENCE relation: “The magistrate must attend the meetings, then they may be seated.”).

- **Or:** Either two predicates or entities are mentioned in the law, but when only one needs to be passed (in the case of a predicate), or only one entity is affected.

- **And:** Either two predicates or entities are mentioned in the law, but both need to be passed (in the case of a predicate), or both entities are affected.

- **SameClass:** Indicates that two discourse units identified as CLASS are the same.

## C Annotation Interface Details

Our annotation tool is we designed a simple and modularized annotation framework in 600 lines of JQuery, Javascript and HTML, with a Datastore backend<sup>30</sup>. Our annotation framework supports span annotation and relation tagging.

The annotation interface itself, shown in Figure 2, is powered by a stateful page object, called PageHandler, that is instantiated with several parameters (page\_height, buttons, relations)

<sup>30</sup>Google Datastore is a NoSQL, scalable JSON store, which is suitable for our usecase. <https://cloud.google.com/datastore>

1113	and handles all of the page interactions. The	<b>E Additional Experimental Results</b>	1164
1114	PageHandler is placed directly in the HTML page	We give more results for the Span-Level tagging task,	1165
1115	containing the text to be annotated, so any service that	reporting on precision and recall as well as F1.	1166
1116	can render text can automatically become an annotation	<b>F Prompt Designs for GPT3.5</b>	1167
1117	service. In our case, we built Jinja templates to render	Here we give sample prompts, along with their true-	1168
1118	our HTML, since our server is coded in Python-Flask.	label completions for each span.	1169
1119	We additionally provide a helper function that, with	<b>F.1 Span Level Tagging Prompts</b>	1170
1120	input data, can compile our Jinja templates as static,	<b>F.1.1 SUBJECT Identification</b>	1171
1121	fully-functional AMT HTMLQuestions.	You are a legal assistant. I will show you a paragraph	1172
1122	We use a Datastore backend to track progress towards	of law. Which entities gains powers, restrictions or	1173
1123	annotation tasks, as shown in Figure ??.	responsibilities under this law? NOT which entities are	1174
1124	We code data entries (the equivalent of MySQL tables) to track helper-	used to test the law, or which entities are affected. In	1175
1125	statistics, helper_summ, how many tasks are left to	other words, what entity is the SUBJECT of the law? It	1176
1126	assign, incomp_tasks, and how many annotations	might not aren't always explicit, and sometimes can be	1177
1127	have been completed, comp_annot. We track these	expressed passively. Restrict your choices to an entity	1178
1128	statistics to ensure that we can obtain multiple annota-	mentioned in the law OR "passive voice entity", if the	1179
1129	tions for each task, and that no helper sees the same	entity is not explicitly mentioned in the text. Enumerate	1180
1130	task more than once. We perform one GET request	all instances of the entity in the text, even if repeated.	1181
1131	at the beginning of each user session to collect user-	If there is no entity that matches this description in the	1182
1132	stats and then use client-side cookies throughout the	text, including "passive voice entity", say "no entity".	1183
1133	session to minimize the number of requests we send	If there are multiple segments of text in the law that apply,	1184
1134	to the back-end. We use a NoSQL database because	join them with a semi-colon. The order of text spans	1185
1135	they are low-latency and designed for streaming, and	does NOT matter. Do NOT say anything else." I will	1186
1136	Datastore because our web-app is hosted on Google	give you 1 examples, and then you will perform the task	1187
1137	App Engine. We include our Datastore management	yourself.	1188
1138	back-end as part of the annotation package. To use our	EXAMPLE: Law: "* 71. Special population census.	1189
1139	tool with other NoSQL providers, <sup>31</sup> a port is necessary.	The expenses incurred by a county, city, town, or village	1190
1140	<b>D Website Design</b>	to conduct a special population census supervised by	1191
1141	In <b>Flow 1</b> , users can use a query box to perform full-text	the United States bureau of the census pursuant to a	1192
1142	and faceted search on laws and then click on and return	contract made pursuant to section twenty of the general	1193
1143	results to read the full text of the law. ElasticSearch	municipal law, three years." Answer: "no entity"	1194
1144	powers both of these endpoints. This flow is useful	NOW IT'S YOUR TURN:	1195
1145	for when journalists want to explore a specific term	Law: "If a vacancy as described in subdivision (d)(1)	1196
1146	or concept irrespective of its discourse role, or simply	occurs after the sixth Thursday before the primary elec-	1197
1147	familiarize themselves with the corpus.	tion in any county having a metropolitan form of gov-	1198
1148	In <b>Flow 2</b> , users can view aggregate counts of dif-	ernment with a population of more than five hundred	1199
1149	ferent discourse elements, by type, across the corpora.	thousand (500,000), according to the 2010 federal cen-	1200
1150	This helps to summarize the corpora from a functional	sus or any subsequent federal census, then the members	1201
1151	standpoint, as described in Section 2. Users navigate	of the county executive committees who represent the	1202
1152	this flow by clicking on one of five buttons to see the	precincts composing such senate district may nominate	1203
1153	counts of each of the five principle discourse spans, then	a candidate to appear on the November election ballot	1204
1154	clicking on any of the returned span results to view all	by any method authorized under the rules of the party."	1205
1155	laws with this span. MySQL serves both of these end-	Answer:	1206
1156	points (and provides additional metrics, such as a map	>> vacancy as described in subdivision (d)(1); mem-	1207
1157	in the about.html page, not shown here.).	bers of the county executive committees	1208
1158	In both flows, visitors can access our annotation	<b>F.1.2 EXCEPTION Identification</b>	1209
1159	framework, described in Section 3.2. From <b>Flow 1</b> ,	You are a legal assistant. I will show you a paragraph	1210
1160	they can click search results to tag a specific paragraph,	of law. What are exception cases when this law does	1211
1161	and from <b>Flow 2</b> they can click to correct an annotated	not apply? Restrict your answer to text in the law. Join	1212
1162	paragraph. Additionally, they can annotate a randomly	non-contiguous segments of text with a semi-colon. If	1213
1163	selected paragraph by clicking "Help Us Tag."	there are no exception cases where this law does not	1214
		apply, say "none". If there are multiple segments of	1215
		text in the law that apply, join them with a semi-colon.	1216
		The order of text spans does NOT matter. Do NOT say	1217

<sup>31</sup>e.g. Amazon DynamoDB – <https://aws.amazon.com/dynamodb/>

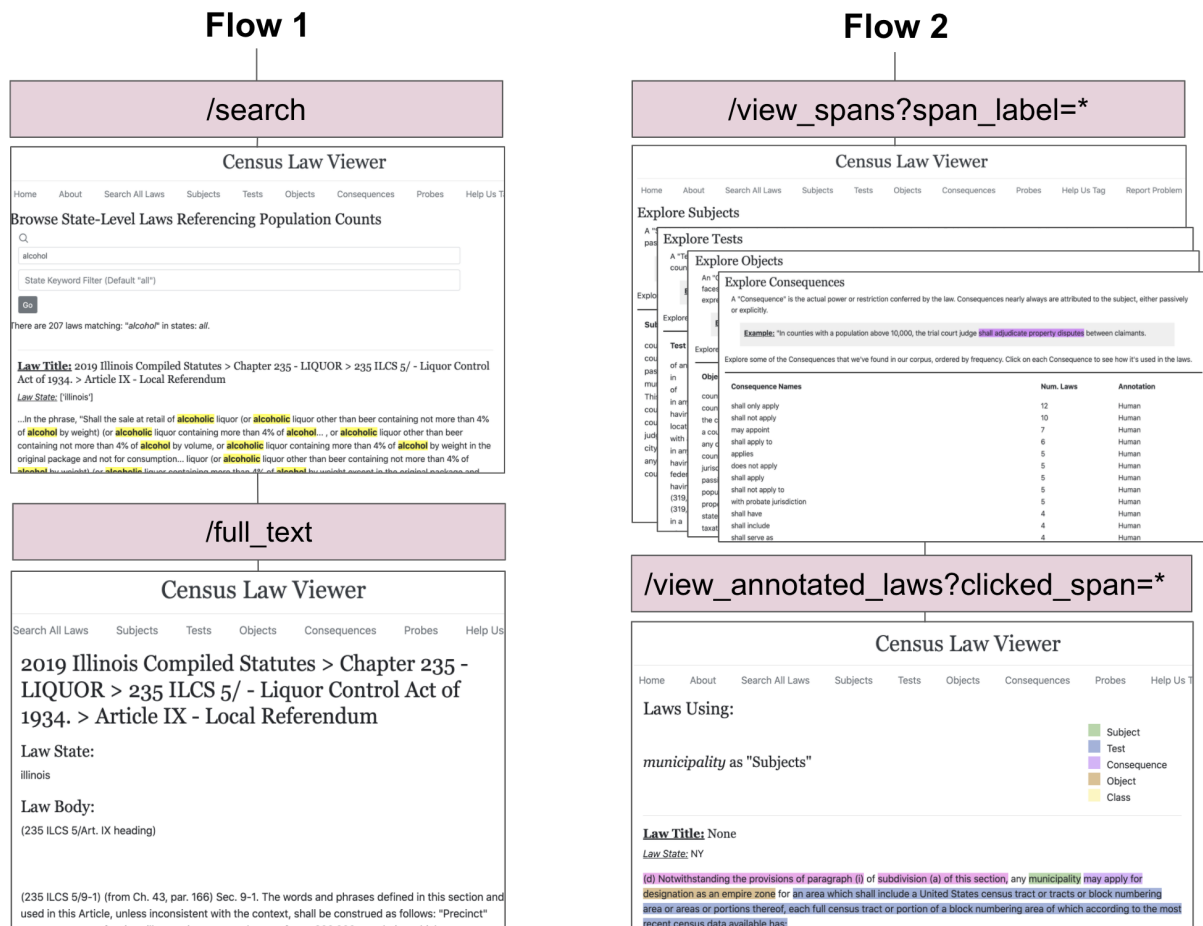


Figure 6: A flow-based sitemap for our website, `statecensulaws.org`, with some details about the back-end and database setup. The left-column shows **Flow 1**, where a user can search and view full-text results. The right-column shows **Flow 2**, where a user can view top law-discourse spans and see all laws these spans are used in. Each flow leads to the annotation framework.

1218 anything else." I will give you 1 examples, and then you  
1219 will perform the task yourself.

1220 EXAMPLE: Law: "Notwithstanding subdivision  
1221 (b)(1), in counties having a population of not less than  
1222 seventeen thousand two hundred fifty (17,250) nor more  
1223 than seventeen thousand five hundred fifty (17,550), ac-  
1224 cording to the 1990 federal census or any subsequent  
1225 federal census, the budget committee shall be composed  
1226 of six (6) members." Answer: "none"

1227 NOW IT'S YOUR TURN:

1228 Law: "In counties having a population of not less than  
1229 three hundred nineteen thousand six hundred twenty-  
1230 five (319,625) nor more than three hundred nineteen  
1231 thousand seven hundred twenty-five (319,725), accord-  
1232 ing to the 1980 federal census or any subsequent federal  
1233 census, a library board of not less than seven (7) mem-  
1234 bers nor more than nine (9) members may be appointed  
1235 by the county legislative body and city governing bod-  
1236 ies which are parties to the agreement, the number ap-  
1237 pointed by each to be determined according to the ratio  
1238 of population in each participating city and in the county  
1239 outside the city or cities, based on the most recent fed-

1240 eral census; provided, that each shall appoint at least  
1241 one (1) member." Answer:

1242 >> provided, that each shall appoint at least one (1)  
1243 member

### 1244 F.1.3 TEST Identification

1245 You are a legal assistant. I will show you a paragraph  
1246 of law. Under what conditions does this law apply? In  
1247 other words, what test is implied by the law? Restrict  
1248 your answer to text in the law. Join non-contiguous  
1249 segments of text with a semi-colon. If there are no  
1250 conditions for this law to apply explicitly stated in the  
1251 text, say "none". If there are multiple segments of text  
1252 in the law that apply, join them with a semi-colon. The  
1253 order of text spans does NOT matter. Do NOT say  
1254 anything else." I will give you 1 examples, and then you  
1255 will perform the task yourself.

1256 EXAMPLE: Law: "(iii) Notwithstanding the forego-  
1257 ing, local governments and voluntary agencies shall be  
1258 granted state aid of one hundred percent of the net oper-  
1259 ating costs expended by such localities and by voluntary  
1260 agencies pursuant to contracts with such local govern-  
1261 ments or with the office of alcoholism and substance

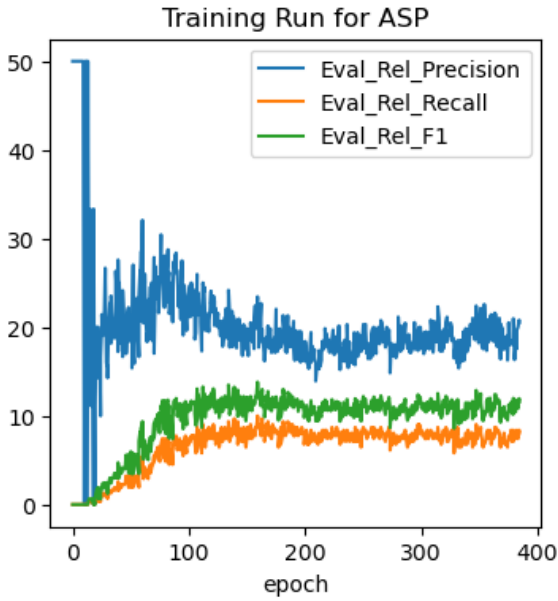


Figure 7: Training run for ASP. The ASP model is learning across epochs, it likely does not have enough data to fully distinguish the embedding space for jointly modeled task.

1262 abuse services for alcohol crisis centers, chemical de-  
 1263 pendency programs for youth, residential services for  
 1264 recovering alcoholics and substance abusers and for al-  
 1265 colism AIDS coordinators. Such state aid may also be  
 1266 granted to programs transferred from the task force on  
 1267 integrated projects for youth and chemical dependency.  
 1268 Such state aid shall also be granted for non-residential  
 1269 services determined to be necessary to serve the public  
 1270 interest by the commissioner of alcoholism and sub-  
 1271 stance abuse services provided by local governments  
 1272 having a population of one hundred twenty-five thou-  
 1273 sand or less as determined by the last preceding federal  
 1274 census, or by voluntary agencies pursuant to contracts  
 1275 with such local governments." Answer: "determined to  
 1276 be necessary to serve the public interest by the com-  
 1277 missioner of alcoholism and substance abuse services;  
 1278 provided by; having a population of one hundred twenty-  
 1279 five thousand or less as determined by the last preceding  
 1280 federal census; pursuant to contracts with; with; trans-  
 1281 ferred from the task force on integrated projects for  
 1282 youth and chemical dependency"

1283 NOW IT'S YOUR TURN:

1284 Law: "(2) If two or more counties included in the  
 1285 measure are required to prepare a translation of ballot  
 1286 materials into the same language other than English,  
 1287 the county that contains the largest population, as de-  
 1288 termined by the most recent federal decennial census,  
 1289 among those counties that are required to prepare a  
 1290 translation of ballot materials into the same language  
 1291 other than English shall prepare the translation, or au-  
 1292 thorize the authority to prepare the translation, and that  
 1293 translation shall be used by the other county or counties,

as applicable." Answer: 1294

»> are required to prepare a translation of ballot ma- 1295  
 terials into the same language other than English; that 1296  
 contains the largest population, as determined by the 1297  
 most recent federal decennial census; among those coun- 1298  
 ties that are required to prepare a translation of ballot 1299  
 materials into the same language other than English 1300

#### F.1.4 OBJECT Identification 1301

1302 You are a legal assistant. I will show you a paragraph of  
 1303 law. Which entities are affected by the powers of this  
 1304 law? NOT which entities gain powers, but who is af-  
 1305 fected by those in power? In other words, what entity is  
 1306 the object of the law? It might not aren't always explicit,  
 1307 and sometimes can be expressed passively. Restrict your  
 1308 choices to an entity mentioned in the law OR "passive  
 1309 voice entity", if the entity is not explicitly mentioned  
 1310 in the text. Enumerate all instances of the entity in the  
 1311 text, even if repeated. If there is no entity that matches  
 1312 this description in the text, including "passive voice en-  
 1313 tity", say "no entity". If there are multiple segments of  
 1314 text in the law that apply, join them with a semi-colon.  
 1315 The order of text spans does NOT matter. Do NOT say  
 1316 anything else." I will give you 1 examples, and then you  
 1317 will perform the task yourself.

1318 EXAMPLE: Law: "This subsection (b) applies only  
 1319 to counties with a metropolitan form of government and  
 1320 to counties having the following populations according  
 1321 to the 1970 federal census or any subsequent federal  
 1322 census:" Answer: "no entity"

1323 NOW IT'S YOUR TURN:

1324 Law: "An authority shall not initiate any redevelop-  
 1325 ment project under this chapter until the governing body,  
 1326 or agency designated by it or empowered by law so to  
 1327 act, of each city or town, herein called "municipalities,"  
 1328 and any county having a population of not less than two  
 1329 hundred seventy-five thousand (275,000) nor more than  
 1330 three hundred twenty-five thousand (325,000), accord-  
 1331 ing to the 1980 federal census or any subsequent federal  
 1332 census, in which any of the area to be covered by the  
 1333 project is situated, has approved a plan, herein called  
 1334 the "redevelopment plan", which provides an outline  
 1335 for the development or redevelopment of the area and is  
 1336 sufficiently complete, to:" Answer:

1337 »> any redevelopment project under this chapter

#### F.1.5 PROBE Identification 1338

1339 You are a legal assistant. I will show you a paragraph  
 1340 of law. Which entities are used to determine when this  
 1341 law applies? NOT which entities gain powers, OR is  
 1342 affected by the law. In other words, which entities are  
 1343 probed by the law? Restrict your answer to text in the  
 1344 law. Join non-contiguous segments of text with a semi-  
 1345 colon. If there is no entity that matches this description  
 1346 in the text, including "passive voice entity", say "no  
 1347 entity". If there are multiple segments of text in the  
 1348 law that apply, join them with a semi-colon. The order  
 1349 of text spans does NOT matter. Do NOT say anything



1350 else." I will give you 1 examples, and then you will 1408  
1351 perform the task yourself. 1409  
1352 EXAMPLE: Law: "2. The commissioner is autho- 1410  
1353 rized to contract to make a state grant, within the limit 1411  
1354 of appropriation therefor, to any planning unit for up 1412  
1355 to ninety percent of the costs to prepare, update or re- 1413  
1356 vise its local solid waste management plan; provided, 1414  
1357 however, that no such grant has been previously made 1415  
1358 to a planning unit which is a part of or is served by the 1416  
1359 planning unit seeking such grant. A planning unit may 1417  
1360 receive a grant pursuant to this subdivision which shall 1418  
1361 not exceed the greater of twenty-five thousand dollars or 1419  
1362 one dollar for each resident of the planning unit, based 1420  
1363 upon the current federal decennial census." Answer: "no  
1364 such grant; a planning unit; the planning unit"  
1365 NOW IT'S YOUR TURN:  
1366 Law: "Such functions may also be delegated by the  
1367 municipality to any not-for-profit corporation acting for  
1368 or on behalf of such municipality; provided, that, ex-  
1369 cept in any county with a metropolitan form of govern-  
1370 ment and having a population of four hundred thousand  
1371 (400,000) or more, according to the 1980 federal census  
1372 or any subsequent federal census, the site selection for  
1373 an energy production facility may be delegated to any  
1374 such not-for-profit corporation, but shall be subject to  
1375 the approval by a two-thirds ( 2/3 ) vote of the legislative  
1376 bodies of the city and the county in which such city is  
1377 located for whom or on whose behalf such not-for-profit  
1378 corporation is acting prior to the purchase of any such  
1379 site." Answer:  
1380 >> any county

1381 **F.1.6 CONSEQUENCE Identification**  
1382 You are a legal assistant. I will show you a paragraph  
1383 of law. What are the powers or obligations granted  
1384 by this law? In other words, what is the law's conse-  
1385 quence? Restrict your answer to text in the law. Join  
1386 non-contiguous segments of text with a semi-colon. If  
1387 there are no powers or obligations explicitly stated in  
1388 the text, say "none". If there are multiple segments of  
1389 text in the law that apply, join them with a semi-colon.  
1390 The order of text spans does NOT matter. Do NOT say  
1391 anything else." I will give you 1 examples, and then you  
1392 will perform the task yourself.  
1393 EXAMPLE: Law: "After January 1, 1980, with re-  
1394 spect to the construction, purchase, or lease of build-  
1395 ings which are located or will be located in a standard  
1396 metropolitan statistical area (SMSA) with a population  
1397 of 250,000 or more according to the most recent decen-  
1398 nial census, which is served by a public transit operator,  
1399 as defined in Section 99210 of the Public Utilities Code,  
1400 the board shall give consideration to the location of  
1401 existing public transit corridors, as defined in Section  
1402 50093.5 of the Health and Safety Code, for the area.  
1403 Construction, purchase, or lease of buildings at loca-  
1404 tions outside of existing public transit corridors may be  
1405 approved after the board has determined: (1) the pur-  
1406 pose of the facility does not require transit access; or  
1407 (2) it is not feasible to locate the facility in an existing  
transit corridor; or (3) the transit operator will provide  
service as needed to effectively serve the facility. The  
board may request the assistance of the transit operator  
in making its determination and shall notify the opera-  
tor of its decision." Answer: "may be approved; shall  
give consideration to; may request the assistance of; in  
making its determination; shall notify; of its decision"  
NOW IT'S YOUR TURN:  
Law: "This part only applies in those counties with a  
metropolitan form of government and in those counties  
with a population according to the 1970 federal census  
or any subsequent federal census of:" Answer:  
>> applies

Discourse Unit	Example
<b>SUBJECT</b>	clerk and master legislative body any person the board of mayor and aldermen “Club”
<b>OBJECT</b>	the library and recreational facilities. to the electors of the county presiding officer the property owners the tenants and their property and the safety and the protection of the premises.
<b>TEST</b>	having a population of not less than eight hundred twenty-five thousand (825,000) nor more than eight hundred thirty thousand (830,000)... upon adoption of a resolution by a two-thirds ( 2/3 ) vote of the county legislative body authorizing the county trustee to collect delinquent property taxes as provided in this subsection who owns real property situated within the corporate limits of such municipality upon entering an order finding it in the best interest of judicial efficiency in areas of historical significance to a locality, the county and the state
<b>CONS.</b>	shall, upon collection of state fines and costs, remit such fines and costs to may be levied be governed by shall make eligible for the waiver be paid from the same fund used for maintaining and operating the county free library.
<b>EXCEPTION</b>	wherever its disapproval of a redevelopment project has been dissolved as prescribed by contracting with other counties and/or cities for joint operation of a free public library except the clerk of the supreme court and chief deputy clerks of the supreme court provided, that each shall appoint at least one (1) member unless the board of supervisors of the county shall, by resolution, provide for fees in excess of that amount
<b>PROBE</b>	county enrolled member and spouse city in Canada an enrolled member of an incorporated volunteer fire company, fire department or incorporated voluntary ambulance service private acts of the state
<b>CLASS</b>	the superior [court] for such county general sessions court [clerk] [the legislative body] of the municipality. the mental health [court] [the commissioner] of mental health,
<b>DEFINITION</b>	shall be determined by the last federal decennial or local special population census... is the same proportion of the total population of the district as each of the other areas. that is the sum of the county public hospital health system’s gross inpatient revenue shall include The Municipality of Metropolitan Toronto and any other similar corporation in Canada means any regular and full-time employee of a county with a metropolitan government

Table 6: Example spans from each discourse type in our annotated dataset.

Source Span	Target Span	Permissible Relations
OBJECT	CLASS OBJECT CONSEQUENCE TEST DEFINITION SUBJECT	hasProperty continuation, And, Or, sameEntity, By, To entityEmpoweredTo, entityRequiredTo entityTested definedAs sameEntity, And, Or, Of
SUBJECT	CLASS OBJECT CONSEQUENCE TEST DEFINITION SUBJECT	hasProperty continuation, And, Or, sameEntity entityEmpoweredTo, entityRequiredTo entityTested isDefinedAs sameEntity, And, Or, Of
TEST	CONSEQUENCE TEST SUBJECT PROBE EXCEPTION	conditionalConsequence continuation, And, Or, followedBy testConcerns testConcerns exceptedBy
CONSEQUENCE	CLASS OBJECT CONSEQUENCE TEST DEFINITION SUBJECT EXCEPTION	hasProperty Affects, comparison continuation, And, Or, followedBy conditionedBy Affects Affects, comparison conditionedBy
CLASS	CLASS OBJECT DEFINITION SUBJECT PROBE	continuation, And, Or, sameClass propertyOf definedAs propertyOf propertyOf
EXCEPTION	OBJECT CONSEQUENCE TEST SUBJECT PROBE EXCEPTION	excepts excepts excepts excepts excepts excepts, continuation, And, Or
PROBE	CLASS TEST PROBE	hasProperty entityTested sameEntity, And, Or, Of
DEFINITION	CLASS OBJECT DEFINITION SUBJECT PROBE	defines defines continuation defines defines

Table 7: All possible relations between discourse units identified in our span-tagging process.

	SUBJECT			CONSEQUENCE			OBJECT		
	P	R	F1	P	R	F1	P	R	F1
0-shot	39.5	30.5	34.4	12.5	7.9	9.7	13.9	15.7	14.8
3-shot	32.1	31.3	31.7	22.4	24.3	23.3	18.3	23.1	20.4
5-shot	27.9	34.1	30.7	23.2	25.1	24.1	13.9	18.4	15.9
8-shot	27.4	32.5	29.7	21.6	25.5	23.4	13.4	19.2	15.8
fine-tuned	41.2	43.1	42.1	51.0	48.8	49.9	38.8	33.3	35.9

Table 8: Precision, Recall and F1 for the first three discourse tags we studied.

	PROBE			TEST			EXCEPTION		
	P	R	F1	P	R	F1	P	R	F1
0-shot	11.3	16.7	13.4	42.8	30.2	35.4	53.3	56.1	54.7
3-shot	23.2	36.0	28.2	45.0	42.8	43.9	45.0	47.4	46.2
5-shot	26.7	36.4	30.8	48.8	50.9	49.8	41.8	49.1	45.2
8-shot	29.1	39.6	33.5	46.7	50.2	48.4	51.6	56.1	53.8
fine-tuned	38.8	31.7	34.9	55.2	51.1	53.0	51.5	61.4	56.0

Table 9: Precision, Recall and F1 for the last three discourse tags we studied.

	Macro			Micro		
	P	R	F1	P	R	F1
0-shot	28.9	26.2	27.1	25.0	21.4	22.7
3-shot	31.0	34.2	32.3	28.8	32.1	30.1
5-shot	30.4	35.7	32.8	28.6	33.6	30.8
8-shot	31.6	37.2	34.1	28.5	34.2	31.0
fine-tuned	46.1	44.9	45.3	45.6	43.3	44.3

Table 10: Macro-average and Micro-averaged Precision, Recall and F1 for all discourse tags we studied.