# Explaining Mixtures of Sources in News Articles

**Anonymous ACL submission**

## Abstract

Human writers plan, *then* write. For large language models (LLMs) to play a role in longer-form article generation, we must understand the planning steps humans make before writing. We explore one kind of planning, source-selection in news, as a case-study for evaluating plans in long-form generation. We ask: why do *specific* stories call for *specific* kinds of sources? We imagine a process where sources are selected to fall into different categories. Learning the article's *plan* means predicting the categorization scheme chosen by the journalist. Inspired by latent-variable modeling, we first develop metrics to select the most likely plan underlying a story. Then, working with professional journalists, we adapt five existing approaches to planning and introduce three new ones. We find that two approaches, or schemas: *stance* (Hardalov et al., 2021) and *social affiliation* best explain source plans in most documents.However, other schemas like *textual entailment* explain source plans in factually rich topics like "Science". Finally, we find we can predict the most suitable schema given just the article's headline with reasonable accuracy. We see this as an important case-study for human planning, and provides a framework and approach for evaluating other kinds of plans, like discourse or plot-oriented plans. We release a corpora, *NewsSources*, with schema annotations for 4M articles, for further study.

## 1 Introduction

Writers use a variety of informational sources to inform storytelling. Consider the following news article, shown in Figure 1. The author shares her planning process[1]:

> *NJ schools are teaching climate change in elementary school. We wanted to understand: how are **teachers** educating*

---
Headline: **NJ Schools Teach Climate Change at all Grade Levels**

**Michelle Liwacz** asked her first graders: what can penguins do to adapt to a warming Earth?  ← *potential labels:* Academic, Neutral

**Gabi**, 7, said a few could live inside her fridge.  ← *potential labels:* Unaffiliated, Neutral

**Tammy Murphy**, wife Governor Murphy, said climate change education was vital to help students.  ← *poten. labels:* Government, Agree

**Critics** said young kids shouldn't learn disputed science.  ← *labels:* Unaffiliated, Refute

A **poll** found that 70 percent of state residents supported climate change being taught at schools.  ← *potential labels:* Media, Agree

---

Table 1: Informational sources synthesized in a single news article. *How would we choose sources to tell this story?* We show two different explanations, given by two competing schema: affiliation and stance. Our central questions: (1) *Which schema best explains the sources used in this story?* (2) *Can we predict, given a topic sentence, which schema to use?*

> *children? How do **parents** and **kids** feel? Is there **pushback**?*

During the planning phase, the journalist chooses different kinds of sources (e.g. teachers, kids) she wishes to use. *Why did she choose these groups?* Was it to capture different sides of an issue (i.e. *stance*-based plan)? Was it to include varied social groups (i.e. *affiliation*-based plan)?

Our motivation for addressing this question is: **(1) As language models (LMs) become more proficient at longer-range generation and incorporate external tools and resources (Schick et al., 2023), evaluating plans is a growing need**. Source selection, as illustrated above, is one of many planning decisions humans make, yet, attempts to instill planning in LMs (Park et al., 2023;

---
[1]Plan: https://nyti.ms/3Tay92f [paraphrased]. Final article: https://nyti.ms/486I11u, see Table 1.
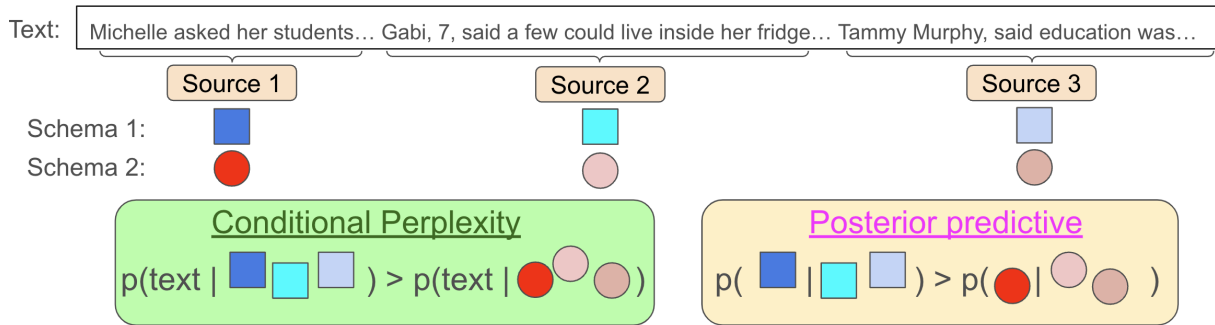
Figure 1: We seek to describe unobserved actions, or *plans*, taken by humans during writing. In this work, *plans* that we study are the choice of sources used during news writing. We adapt 8 schemas to describe different plans and develop two metrics to favor one plan over another: *conditional perplexity* (Airoldi and Bischof, 2016) helps us measure how well a plan corresponds with the observed text and *posterior predictive* (Spangher et al., 2023b) helps us measure how internally coherent a plan is. In the figure above, two *plans* are different sequences of shapes, either squares and circles.

Spangher et al., 2023a) fall short: creative tasks often lack a well-defined objective, so we have no good way of comparing plans, and plans made by humans are mostly not observable in the final text. **(2) Much work exists examining sourcing patterns in journalism (Winter and Krämer, 2014; Hertzum, 2022), however no quantitative measure exists to understand the decisions journalists make on an article-by-article basis.** Sourcing decisions are important for representation and agenda-setting (Manninen, 2017), and provide an important case-study of human planning.

We lay the groundwork for answering both questions by developing novel metrics for comparing document-level plans. Inspired by classical approaches to latent variable modeling in topic modeling (Airoldi and Bischof, 2016), we uncover the optimal *plan* for a document on the following basis: *a plan description is closer the one the original writer followed if it gives more information about the completed document.* We introduce simple metrics for this goal: conditional perplexity and posterior predictive likelihood (Section 2.2).

Then, we work with professional journalists from multiple major news organizations to select different informational schemas which best describe the sourcing plans that journalists make while writing. We adapt an additional five schema from parallel tasks and develop three novel schemas, which we operationalize by annotating over 600 news articles with 4,922 sources. We find that a source's *social affiliation* and *stance* optimally explain plans in most documents. However, for certain kinds of documents, e.g. factually dense topics like "Science", textual entailment (NLI) (Da-

gan et al., 2005) provides a useful structure. The choice of schema, we find, can be predicted with moderate accuracy (ROC=.67) using only the headline of the article, opening the door to different planning approaches for source selection.

Our contributions are threefold:

- We frame *source-type planning* as a lens through which to study planning in writing.

- We collect 8 different plan descriptions, or *schemas* (5 existing and 3 we develop **with professional journalists**). We build a pipeline to extract sources from 4 million news articles and categorize them, building a large public dataset called *NewsSources*.

- We introduce two novel metrics: *conditional perplexity* and *posterior predictive* to compare plans. We find that different plans are optimal for different topics. Further, we show that the right plan can be predicted with .67 ROC given just the headline.

Learning writing plans is parallel to learning *unobserved action sequences* in reinforcement learning setups. By developing metrics to study and compare plans, we open the door for better approaches to plan-based language modeling (Yao et al., 2019; Yang et al., 2022), multi-document retrieval (Pereira et al., 2023; Shapira et al., 2021); and critical media studies (Hernández and Madrid-Morales, 2020). Additionally, taking steps to better predict plans can increase the breath of human-in-the-loop *computational journalism* tools (Wang and Diakopoulos, 2021), e.g. by generating a plan (Yao et al., 2022) for journalists to execute.

2

## 2 Source Categorization

### 2.1 Problem Statement

Our central question is: why did the writer select sources $s_1, s_2, s_3...$ for document $d$? Intuitively, let's say we read an article on a controversial topic. Let's suppose we observe that it contains many opposing viewpoints: some sources in the article "agree" with the main topic and others "disagree". We can conclude that the writer probably chose sources on the basis of their *stance* (Hardalov et al., 2021) (or their opinion-based support) rather than another explanation, like their *discourse* role (which describes their narrative function).

More abstractly, we describe source-selection as a generative process: first, journalists plan *how* they will choose sources (i.e. the *set* of $k$ categories sources will fall into), then they choose sources, each falling into 1-of-$k$ categories. Different plans, or categorizations, are possible (e.g. see Figure 1): *the "right" plan is the one that best predicts the final document.*

Each plan, or categorizations, is specified by a *schema*. For the 8 schemas used in this work, see Figure 2. To apply a schema to a document, we frame an approach consisting of two components: (1) an attribution function, $a$:

$$a(s) = q \in Q_d \text{ for } s \in d \tag{1}$$

introduced in Spangher et al. (2023b), which maps each sentence $s$ in document $d$ to a source $Q_d = \{q_1^{(d)}, ... q_k^{(d)}\}$[2] and (2) a classifier, $c$:

$$c_Z(s_1^{(q)}, ... s_n^{(q)}) = z \in Z \tag{2}$$

which takes as input a sequence of sentences attributed to source $q^{(d)}$ and assigns a type $z \in Z$ for schema $Z$. Taken together, $c_Z$ and $a$ give us a learned estimate of the posterior $p(z|x)$.

This supervised framing is not typical in latent-variable settings; the choice of $z$ and the *meaning* of $Z$ are typically jointly learned without supervision. However, learned latent spaces often do not correspond well to theoretical schemas (Chang et al., 2009), and supervision has been shown to be helpful with planning (Wei et al., 2022). On the other hand, supervised models trained on different schema are challenging to compare, especially when different architectures are optimal for each

schema. A latent-variable framework here is ideal: comparing different graphical models (Bamman et al., 2013; Bamman and Smith, 2014) *necessitates* comparing different schemas, as each run of a latent variable model produces a different schema.

### 2.2 Comparing *Plans*, or Schemas

We can compare plans in two ways: (1) how well they explain each observed document and (2) how structurally consistent they are.

**Explainability**   A primary criterion for a *plan* is for it to explain the observed data well. To measure this, we use *conditional perplexity*[3]

$$p(x|z) \tag{3}$$

which measures the uncertainty of observed data, $x$, given a latent structure, $z$. Measuring $p(x|z)$ for different $z$ (fixing $x$) allows us to compare $z$. Conditional perplexity is a novel metric we introduce, inspired by metrics to evaluate latent unsupervised models, like the "left-to-right" algorithm introduced by (Airoldi and Bischof, 2016). [4]

**Structural Likelihood:**   A second basic criterion for a latent structure to be useful is for it be consistent, which is a predicate for learnability. We assess the consistency of a set of assignments, $z$, by calculating the *posterior predictive*:

$$p(z|z_-, x) \tag{4}$$

Deng et al. (2022) exploring using full joint distribution, $p(z)$, *latent perplexity*, to evaluate the structure text $x$ produced by generative language models ("*model criticism*"). We simplify using the full distribution and instead evaluate the conditional predictive to study document structure. This, we find in early experiments, is easier to learn and thus helps us differentiate different $Z$ better ("*schema criticism*").[5] Now, we describe our schemas.

---

[2]These sources are referenced in $d$. There is no consideration of document-indepedent sources.

[3]We abuse notation here, using $p$ as both probability and perplexity: $p(x) = \exp\{-\mathbb{E} \log p(x_i|x_{<i})\}$.

[4]We note that the term, *conditional perplexity*, was originally introduced by Zhou and Lua (1998) to compare machine-translation pairs. In their case, both $x$ and $z$ are observable; as such, they do not evaluate latent structures, and their usage is not comparable to ours.

[5]Our work is inspired by Spangher et al. (2023b)'s work, where $z$ was the choice of specific action, rather than a general action-type. They had no concept of a "schema" to group actions.

**Affiliation**
*Source's group membership*

| | |
|---|---|
| Academic | Corporate |
| Government | Industry Group |
| Media | NGO |
| Other Group | Political Group |
| Individual | Union |
| Victim | Witness |
| Religious Group | |

**Identity**
*Identifying information*

Named Group
Named Individual
Report/Document
Unnamed Group
Unnamed Individual
Vote/Poll

**Argumentation**
*Type of information*

Anecdote
Assumption
Common-Ground
Other
Statistics
Testimony

**NLI**
*Fact Relation*
Contradiction
Entailment
Neutral

**Stance**
*Opinion Rel.*
Affirm   Discuss
Refute   Neutral

**Role**
*Source's role in group*

Decision Maker
Informational
Participant
Representative

**Retrieval**
*Channel accessed for information*

| | |
|---|---|
| Background | Observation |
| Proposal/Law | Press Report |
| Article | Statement |
| Court Proc. | Email/Social Media |
| Direct/Indirect Quote | |

**Discourse**
*Narrative role of info.*

| | |
|---|---|
| Anecdote | History |
| Consequence | Prev. Event |
| Context | Evaluation |
| Expectations | Main Event |

Figure 2: Label sets of each of the 8 schemas we use to study source categorization. **Extrinsic Source Schemas** Affiliation, role and retrieval-method (Spangher et al., 2023b) capture characteristics of sources *extrinsic* to their usage in the document. **Functional Source Schemas:** Argumentation (), Discourse () and Identity capture functional role of sources for conveying an overall narrative. **Debate-Oriented Source Schemas**: Natural Language Inference (NLI) (Dagan et al., 2005) and Stance (Hardalov et al., 2021) capture the role of sources in broadening the story to encompass multiple sides. Definitions for each label in Appendix D.

## 2.3 Source Schemas

Our schemas, or ways of describing plans taken by journalists, are shown in Figure 2. In this work, we introduce 8 schemas. Each schema provides a set of categories describing the sources used in a news article. Again, the schema that *best predicts the observed text of the article* is the one the journalist most likely adhered to while planning the article. Two of the 8 schemas are **debate-oriented** schemas (i.e. they describe how sources relate to the main topic of the article), three are **functional** (i.e. they describe the role sources serve in the overall flow of the story), and three can be considered **extrinsic** schemas (i.e. they describe how sources fit into societal structures). See Appendix D for more details and definitions for each.

**Debate-Oriented Schemas** Both the *Stance* and *NLI* schemas are debate-orienced schemas. They each capture the relation between two pieces of text. *NLI* (Dagan et al., 2005) captures primarily factual relations between text, while *Stance* (Hardalov et al., 2021) captures opinion-based relations[6]. A text pair may be factually consistent and thus be classified as "Entailment" under a *NLI* schema, but express different opinions and be classified as "Refute" under *Stance*. In our setting, we relate article's headline with the source's attributable information. These schemas assert that a writer uses sources for the purpose of expanding or rebutting information in the narrative.

---

[6]Reddy et al. (2021) views these as the same.

| Schema | Macro-F1 | Schema | Macro-F1 |
|---|---|---|---|
| Argumentation | 68.3 | Retrieval | 61.3 |
| NLI | 55.2 | Identity | 67.2 |
| Stance | 57.1 | Affiliation | 53.3 |
| Discourse | 56.1 | Role | 58.1 |

Table 2: Classification f1 score, macro-averaged, for the 8 schemas. We achieve moderate classification scores for each of schema. In Section 2, when we compare schemas, we account for differences in classification accuracy by introducing noise to higher-performing classifiers.

**Functional Source Schemas** The following schemas: *Argumentation*, *Discourse* and *Identity* all capture the role a source plays in the overall narrative construction of the article. For instance, a source might provide a "Statistic" for a well-formed argument (*Argumentation* (Al Khatib et al., 2016)), or "Background" for a reader to help contextualize (*Discourse* (Choubey et al., 2020)). *Identity*, a novel schema, captures how the reader identifies the source. For example, an "Unnamed Individual" is not identifiable by the reader. b

**Extrinsic Source Schemas** *Affiliation*, *Role* and *Retrieval* schemas serve to characterize attributes of sources external to the news article. Stories often implicate social groups (McLean et al., 2019), such as "academia" or "government." Those group identities are extrinsic to the story's architecture but important for the selection of sources. Sources may be selected because they represent a group

(i.e. *Affiliation*) or because their group position is important within the story's narrative (e.g. "participants" in the events, i.e. *Role*). *Retrieval*, introduced by Spangher et al. (2023b), captures the channel through which the information was found. Although these schema are news-focused, similar ideas can be applied to other fields. For instance, a research article in machine learning might include models from the *open-source*, *academic* and *industry research* communities.

## 3   Building a Silver-Standard Dataset

The schemas described in the previous section give us different theoretical frameworks for identifying the writers' plans. To take further steps towards comparing plans and selecting the plan that *best* describes a document, we must first create a large silver-standard dataset where we have identified and labeled sources in news articles. In this section, we describe how we annotated data and built classifiers for these schema.

### 3.1   Source Extraction

Spangher et al. (2023b) release high-performing source attribution models. With minor modifications[7], we use their models to score a large corpus of $90,000$ news articles from the NewsEdits corpus (Spangher et al., 2022). We find that 47% of sentences in our documents can be attributed to sources, and documents each contain an average of 7.5 +-/5 sources. These statistics are comparable to those reported by Spangher et al. (2023b).

### 3.2   Annotation

We annotate data for our new schemas and evaluate model performance on all schemas. We recruited two annotators, one an undergraduate and the other a former journalist. The former journalist trained the undergraduate for 1 month to identify and label sources, then, they independently labeled 425 sources in 50 articles with each schema to calculate agreement, scoring $\kappa = .63, .76, .84$ on *Affiliation*, *Role* and *Identity* labels. They then labeled 4,922 sources in 600 articles with each schema over 9 months, labeling roughly equal amounts. Finally, they jointly labeled 100 sources in 25 documents with the other schemas for evaluation data over 1 month, with $\kappa \geq .54$.

### 3.3   Training Classifiers for Source Schemas

We train classifiers to assign labels sources under each schema. Unless specified, we use a sequence classifier using RoBERTa-base with self-attention pooling, like in Spangher et al. (2021a); we chose a smaller model that could scale to processing large amounts of articles.

*Affiliation*, *Role*, *Identity*   We use our annotations to train classifiers $p(t|s_1^{(q)} \oplus ... \oplus s_n^{(q)})$, which take as input sentences attributable to source $q$ and output a category in each schema.

*Argumentation*, *Retrieval*, *Discourse*   are labeled on a sentence-level by authors on news and opinion datasets. We use datasets provided by the authors without modification and train classifiers to labels each sentence $s$. For each source $q$, we assign the label $y$ with the most mutual information[8] across attributed sentences $s_1^{(q)}, ... s_n^{(q)}$.

*NLI*   We use an NLI classifier trained by Williams et al. (2022) to label each sentence attributed to source $q$ as a separate hypothesis, and the article's headline as the premise. We use mutual information to assign a single label as above.

*Stance*   We create a news-focused stance dataset by aggregating news-related stance datasets[9]. We filter these training sets to include premises and hypothesis $\geq 10$ words and $\leq 2$ sentences, and distill a T5-based classifier from a fine-tuned GPT3.5-turbo[10] to label news data and label 60,000 news articles. We distill a $T5$ model with this data and achieve comparable performance (Table 2 shows T5's performance).

### 3.4   Classification Results and EDA

As shown in Table 2, we model schemas within a range of f1-scores $(53.3, 67.2)$, showing moderate success in learning each schema. In the next section, we introduce noise (i.e. random label-swapping), to the outputs of these classifiers so that that all have the same accuracy.

---

[7]Described in more detail in Appendix A.

[8]$\arg\max_y p(y|q)/p(y)$

[9]FNC-1 (Pomerleau and Rao, 2017), Perspectrum (Chen et al., 2019), ARC (Habernal et al., 2017), Emergent (Ferreira and Vlachos, 2016) and NewsClaims (Reddy et al., 2021). Data aggregation for stance detection inspired by: (Hardalov et al., 2021; Schiller et al., 2021).

[10]We use OpenAI's GPT3.5-turbo fine-tuning endpoint, as of November 16, 2023.

| Schema | n | Conditional Perplexity $p(x\|z)$ | | | Posterior Predictive $p(\hat{z}\|z_-, x)$ | | |
|---|---|---|---|---|---|---|---|
| | | PPL | Δ base-k (↓) | Δ base-r (↓) | F1 | ÷ base-k (↑) | ÷ base-r (↑) |
| NLI | 3 | 22.8 | 0.62 | -0.08 | 58.0 | 1.02** | 1.01 ** |
| Stance | 4 | 21.5 | -1.71 | -3.21** | 39.1 | 0.88** | 0.83 ** |
| Role | 4 | 22.3 | -0.06 | -0.33** | 38.7 | 1.11** | 1.10 ** |
| Identity | 6 | 21.8 | -0.42 | -0.94 | 25.0 | 1.00 | 1.15 ** |
| Argumentation | 6 | 21.7 | -0.52 | -1.04 | 30.7 | 1.10 ** | 1.12 ** |
| Discourse | 8 | 22.3 | 0.54 | -0.75 | 19.2 | 1.06 ** | 1.08 ** |
| Retrieval | 10 | 23.7 | 1.47 | 0.36 | 15.8 | 1.10 ** | 1.12 ** |
| Affiliation | 14 | 20.5 | -2.11** | -3.04** | 10.5 | 1.26 ** | 1.16 ** |

Table 3: Comparing our schemas against each other. In the first set of experiments, we show *conditional perplexity* results, which tell us how well each schema explains the document text. Shown is PPL (the mean perplexity per schema), $\Delta kmeans$ (PPL - avg. perplexity of kmeans) and $\Delta random$ (PPL - avg. perplexity of the random trial). Statistical significance ($p < .05$) via a $t$-test calculated over perplexity values is shown via **.In the second set of experiments, we show *posterior predictive* results, measured via micro F1-score. We show F1 (f1-score per schema), ÷ kmeans (F1 / f1-score of kmeans), ÷ random (F1 / f1-score of random trial). Statistical significance ($p < .05$) via a $t$-test calculated over 500-sample bootstrapped f1-scores is shown via **.

## 4   Comparing Schemas

We are now ready to explore how well these schemas explain source selection in documents. We start by describing our experiments, then baselines, and finally results. All experiments in this section are based on $90,000$ news articles from NewsEdits (Spangher et al., 2022), labeled as described in the previous section. We split $80,000/10,000$ train/eval.

### 4.1   Metrics

We describe here how we implement the metrics introduced in Section 2.2: (1) *conditional perplexity* and (2) *posterior predictive*.

For an illustration of each metric, please refer to Figure 1. The overall goal of the metrics is to determine *which schema, or labeling of sources, best explains the observed news article*. As the figure shows, if schema A (e.g. in Figure 1: squares) describes an article better than schema B (e.g. in Figure 1: circles), then labels assigned to each source under schema A will outperform labels assigned to each source under Schema B.

**Conditional Perplexity**   To measure *conditional perplexity*, $p(x|z)$, we fine-tune GPT2-base models (Radford et al., 2019) to take *as a prompt a sequence of latent variables, each for a different source and then assess likelihood assigned to observed article text*.[11] This is similar to measuring

*vanilla perplexity* on observed text, except: (1) we provide latent variables as conditioning (2) by fixing the model used and varying the labels, *we are measuring the signal given by each set of different labels*.

Our template for GPT2 is:

$$\langle h \rangle \; h \; \langle l \rangle \; (1) \; l_1 \; (2) \; l_2 ... \langle t \rangle$$
$$(1) \; s_1^{(q_1)} ... s_n^{(q_1)} \; (2) \; ...$$

Where `<tokens>` (e.g. "(1)", "⟨text⟩") are structural markers while variables $l, h, s$ are article-specific. Variables mean the following: $h$ is the headline, $l_i$ is the label for source $i$ and $s_1^{(q_1)} ... s_n^{(q_1)}$ are the sentences attributable to source $i$. Red text is the prompt, or conditioning, and green text to calculate perplexity. *We do not use GPT2 for generation, but to compare the likelihood of observed article text under each schema.* [12]

**Posterior Predictive**   To learn the *posterior predictive* (Equation 4), we train a BERT-based classification model (Devlin et al., 2018) to take the article's headline and a sequence of source-types *with a one randomly held out*. We then seek to predict

---

[11]We note that this formulation has overlaps with recent work seeking to learn latent plans (Deng et al., 2022; Wang et al., 2023; Wei et al., 2022).

[12]Initial experiments show that text markers are essential for the model to learn structural cues. However, they also provide their own signal (e.g. on the number of sources). To reduce the effects of these artifacts, we use a technique called *negative prompting* (Sanchez et al., 2023). Specifically, we calculate perplexity on the *altered* logits, $P_\gamma = \gamma \log p(x|z) - (1 - \gamma) \log p(x|\hat{z})$, where $\hat{z}$ is a shuffled version of the latent variables. Since textual markers remain the same in the prompt for $z$ and $\hat{z}$, this removes markers' predictive power.

*that one*, and evaluate using f1-score. Additionally, we follow Spangher et al. (2023b)'s observation that some sources are *more important* (i.e. have more information attributed). We model the posterior predictive among the 4 sources per article with the most sentences attributed to them.

### 4.2 Baselines

Vanilla perplexity has been criticized for it's use in model comparison (Meister and Cotterell, 2021; Oh et al., 2022) because it can be affected by factors outside goodness-of-fit (e.g. tokenization scheme) can affect the perplexity measurements. We hypothesized that the dimensionality of each schema's latent space might also have an effect (Lu et al., 2017); larger latent spaces tend to assign lower probabilities to each point. Thus, we benchmark each schema against baselines with similar latent dimensions.

**Base-r** Random baseline. We generate $k$ unique identifiers[13], and randomly sample one to each source in each document. $k$ is set to match the number of labels in the schema being compared to.

**Base-k** Kmeans baseline. We first embed sources as paragraph-embeddings using Sentence BERT (Reimers and Gurevych, 2019) [14] Then, we cluster all sources across all documents into $k$ clusters using kmeans (Likas et al., 2003), where $k$ is set to match the number of labels in the schema being compared to. We assign each source the cluster number it was assigned to.

### 4.3 Results and Discussion

As shown in Table 3, the supervised schemas mostly have have lower conditional perplexity than their random and unsupervised kmeans baselines. However, only the *Stance*, *Affiliation* and *Role* schemas improve significantly (at $p < .001$), and the *Role* schema's performance increase is minor. *Retrieval* has a statistically significant *decrease* in explainability. There are two reasons for this: (1) a small number of examples are very high perplexity, and this shifts the distribution significantly (when considering median statistics, as shown in Appendix B, the difference disappears.) (2) We examine examples and find that *Retrieval* does not impact wording as expected: writers make efforts

---

[13]Using MD5 hashes, from python's `uuid` library.
[14]Specifically, `microsoft/mpnet-base`'s model `https://www.sbert.net/docs/pretrained_mo dels.html` given all sentences associated with the source.

to convey information similarly whether it was obtained via a quote, document or a statement.

Interestingly, we *do* observe statistically significant improvements of kmeans over random baselines in all cases (except $k = 3$). In general, our baselines have lower variance in perplexity values than experimental schemas. This is not unexpected: as we will explore in the next section, we expect that schemas will be optimal for certain articles and suboptimal for others, resulting in a greater range in performance. For more detailed comparisons, see Appendix B.

Posterior predictive results generally show improvement across trials, with the *Affiliation* trial showing the highest improvement over both baselines. This indicates that most tagsets are, to some degree, internally consistent and predictable. *Stance* is the only exception, showing significantly lower f1 than even random baselines. This indicates that, although Stance is able to explain observed documents well (as observed by it's impact on conditional perplexity), it's not always predictable how it will applied. Perhaps this is indicative that writers do not know a-priori what sources will agree or disagree on any given topic before talking to them, and writers do not always actively seek out opposing sides.

For another baseline, we implemented latent variable model. In initial experiments, it does not perform well. We show in Appendex G that the latent space learned by the model is sensible. Bayesian models are attractive for their ability to encode prior belief, and ideally they would make good baselines for a task like this, which interrogates latent structure. However, more work is needed to better align them to modern deep-learning baselines.

## 5 Predicting Schemas

Taken together, our observations from (1) Section 3.4) indicate that schemas are largely unrelated and (2) Section 4.3 indicate that *Stance* and *Affiliation* both have similar explanatory power (although *Stance* is less predictable). We next ask: which kinds of articles are better explained by one schema, and which are better explained by the other?

In Table 4, we show topics that have low perplexity under the *Stance* schema, compared with the *Affiliation* schema (we calculate these by aggregating document-level perplexity across keywords assigned to each document in our dataset). As we can

| Stance | Affiliation |
|---|---|
| Bush, George W | Freedom of Speech |
| Swift, Taylor | 2020 Pres. Election |
| Data-Mining | Jazz |
| Artificial Intelligence | Ships and Shipping |
| Rumors/Misinfo. | United States Military |
| Illegal Immigration | Culture (Arts) |
| Social Media | Mississippi |

Table 4: Top keywords associated with articles favored by stance or affiliation. Keywords are manually assigned by news editors

| | | | |
|---|---|---|---|
| Affiliation | 41.7% | Argument. | 1.2% |
| Identity | 22.7% | Discourse | 1.1% |
| Stance | 17.7% | NLI | 1.1% |
| Role | 13.4% | Retrieval | 1.1% |

Table 5: Proportion of our validation dataset favored by one schema, i.e. $\hat{Z} = \arg\max_Z p(x|z)$

see, topics requiring greater degrees of debate, like "Artificial Intelligence", and "Taylor Swift" are favored under the *Stance* Topic, while broader topics requiring many different social perspectives, like "Culture" and "Freedom of Speech" are favored under *Affiliation*. We set up an experiment where we try to predict $\hat{Z} = \arg\min_Z p(x|z)$, the schema for each datapoint with the lowest perplexity. Using perplexity scores calculated in the prior section[15], we calculate the lowest-perplexity schema. Table 5 shows the distribution of such articles. We down-sample the articles until the classes are balanced, and train a simple linear classifier[16] to predict $\hat{Z}$. We get .67 ROC-AUC (or .23 f1-score). These results are tantalizing and offer the prospect of being able to *better plan source retrieval*, in RAG, and computational journalism settings, by helping decide an axis on which to seek different sources. More work is needed to validate these results.

## 6 Related Work

**Latent Variable Persona Modeling** Our work is inspired by earlier work in persona-type latent variable modeling (Bamman et al., 2013; Card et al., 2016; Spangher et al., 2021b). Authors model characters in text as mixtures of topics. We both seek to learn and reason about about latent character-types, but their line of work takes an unsupervised approach. We show that supervised schemas outperform unsupervised.

**Multi-Document Retrieval** In multiple settings – e.g. multi-document QA (Pereira et al., 2023), multi-document summarization (Shapira et al., 2021), retrieval-augmented generation (Lewis et al., 2020) – information *from a single source* is as-

sumed to be insufficient to meet a user's needs. In typical information retrieval settings, the goal is to retrieve a single document closest to the query (Page et al., 1998). In settings where *multiple sources are needed*, on the other hand, retrieval goals are not clearly understood[17]. Our work attempts to clarify this, and can be seen as a step towards better retrieval planning.

**Planning in Language Models** Along the line of the previous point, chain-of-thought reasoning (Wei et al., 2022) and few-shot prompting, summarized in (Sanchez et al., 2023), can be seen as latent-variable processes. Indeed, work in this vein is exploring latent-variable modeling for shot selection (). Our work, in particular the *conditional perplexity* formulation and it's implementation, can be seen as a way of comparing different chain-of-thought plans as they relate to document planning.

**Computational Journalism** seeks to apply computational techniques to assist journalists in reporting. Researchers have sought to improve detection of incongruent information (Chesney et al., 2017), detect misinformation (Pisarevskaya, 2017) and false claims made in news articles (Adair et al., 2017).

## 7 Conclusions

In conclusion, we explore ways of thinking about sourcing in human writing. We compare 8 schemas of source categorization, and adapt novel ways of comparing them. We find, overall, that *affiliation* and *stance* schemas help explain sourcing the best, and we can predict which is most useful with moderate accuracy. Our work lays the ground work for a larger discussion of retrieval aims in multi-document retrieval settings, it also takes us steps towards tools that might be useful to journalists. Naturally, our work is a simplification of the real human processes guiding source selection; these categories are non-exclusive and inexhaustive. We hope by framing these problems we can spur further research in this area.

---

[15]across the dataset used for validation, or 5,000 articles
[16]Bag-of-words with logistic regression

[17]As Pereira et al. (2023) states, *"retrievers are the main bottleneck"* for well-performing multi-document systems.

## 8 Limitations

A central limitation to our work is that the datasets we used to train our models are all in English. As mentioned previously, we used English language sources from Spangher et al. (2022)'s *NewsEdits* dataset, which consists of sources such as nytimes.com, bbc.com, washingtonpost.com, etc.

Thus, we must view our work with the important caveat that non-Western news outlets may not follow the same source-usage patterns and discourse structures in writing their news articles as outlets from other regions. We might face extraction and labeling biases if we were to attempt to do such work in other languages.

## 9 Ethics Statement

### 9.1 Risks

Since we constructed our datasets on well-trusted news outlets, we assumed that every informational sentence was factual, to the best of the journalist's ability, and honestly constructed. We have no guarantees that our classification systems would work in a setting where a journalist was acting adversarially.

There is a risk that, if planning works and natural language generation works advance, it could fuel actors that wish to use it to plan misinformation and propaganda. Any step towards making generated news article more human-like risks us being less able to detect and stop them. Misinformation is not new to our media ecosystem, (Boyd et al., 2018; Spangher et al., 2020). We have not experimented how our classifiers would function in such a domain. There is work using discourse-structure to identify misinformation (Abbas, 2022; ?), and this could be useful in a source-attribution pipeline to mitigate such risks.

We used OpenAI Finetuning to train the GPT3 variants. We recognize that OpenAI is not transparent about its training process, and this might reduce the reproducibility of our process. We also recognize that OpenAI owns the models we fine-tuned, and thus we cannot release them publicly. Both of these thrusts are anti-science and anti-openness and we disagree with them on principle. We tried where possible to train open-sourced versions, as mentioned in the text.

### 9.2 Licensing

The dataset we used, *NewsEdits* (Spangher et al., 2022), is released academically. Authors claim that they received permission from the publishers to release their dataset, and it was published as a dataset resource in NAACL 2023. We have had lawyers at a major media company ascertain that this dataset was low risk for copyright infringement.

### 9.3 Computational Resources

The experiments in our paper required computational resources. We used 64 12GB NVIDIA 2080 GPUs. We designed all our models to run on 1 GPU, so they did not need to utilize model or data-parallelism. However, we still need to recognize that not all researchers have access to this type of equipment.

We used Huggingface models for our predictive tasks, and will release the code of all the custom architectures that we constructed. Our models do not exceed 300 million parameters.

### 9.4 Annotators

We recruited annotators from our educational institutions. They consented to the experiment in exchange for mentoring and acknowledgement in the final paper. One is an undergraduate student, and the other is a former journalist. Both annotators are male. Both identify as cis-gender. The annotation conducted for this work was deemed exempt from review by our Institutional Review Board.

## References

Ali Haif Abbas. 2022. Politicizing the pandemic: A schemata analysis of covid-19 news in two selected newspapers. *International Journal for the Semiotics of Law-Revue internationale de Sémiotique juridique*, 35(3):883–902.

Bill Adair, Chengkai Li, Jun Yang, and Cong Yu. 2017. Progress toward "the holy grail": The continued quest to automate fact-checking. In *Computation+ Journalism Symposium, Evanston*.

Edoardo M Airoldi and Jonathan M Bischof. 2016. Improving and evaluating topic models and other models of text. *Journal of the American Statistical Association*, 111(516):1381–1403.

Khalid Al Khatib, Henning Wachsmuth, Johannes Kiesel, Matthias Hagen, and Benno Stein. 2016. A news editorial corpus for mining argumentation strategies. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 3433–3443.

David Bamman, Brendan O'Connor, and Noah A Smith. 2013. Learning latent personas of film characters.

In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 352–361.

David Bamman and Noah A Smith. 2014. Unsupervised discovery of biographical structure from text. *Transactions of the Association for Computational Linguistics*, 2:363–376.

Ryan L Boyd, Alexander Spangher, Adam Fourney, Besmira Nushi, Gireeja Ranade, James Pennebaker, and Eric Horvitz. 2018. Characterizing the internet research agency's social media operations during the 2016 us presidential election using linguistic analyses.

Dallas Card, Justin Gross, Amber Boydstun, and Noah A. Smith. 2016. Analyzing framing through the casts of characters in the news. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1410–1420, Austin, Texas. Association for Computational Linguistics.

Jonathan Chang, Sean Gerrish, Chong Wang, Jordan Boyd-Graber, and David Blei. 2009. Reading tea leaves: How humans interpret topic models. *Advances in neural information processing systems*, 22.

Sihao Chen, Daniel Khashabi, Wenpeng Yin, Chris Callison-Burch, and Dan Roth. 2019. Seeing things from a different angle: Discovering diverse perspectives about claims. In *Proceedings of NAACL-HLT*, pages 542–557.

Sophie Chesney, Maria Liakata, Massimo Poesio, and Matthew Purver. 2017. Incongruent headlines: Yet another way to mislead your readers. In *Proceedings of the 2017 EMNLP Workshop: Natural Language Processing meets Journalism*, pages 56–61, Copenhagen, Denmark. Association for Computational Linguistics.

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*.

Harald Cramér. 1999. *Mathematical methods of statistics*, volume 43. Princeton university press.

Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The pascal recognising textual entailment challenge. In *Machine learning challenges workshop*, pages 177–190. Springer.

Yuntian Deng, Volodymyr Kuleshov, and Alexander M Rush. 2022. Model criticism for long-form text generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11887–11912.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

William Ferreira and Andreas Vlachos. 2016. Emergent: a novel data-set for stance classification. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: Human language technologies*. ACL.

Ivan Habernal, Henning Wachsmuth, Iryna Gurevych, and Benno Stein. 2017. The argument reasoning comprehension task: Identification and reconstruction of implicit warrants. *arXiv preprint arXiv:1708.01425*.

Momchil Hardalov, Arnav Arora, Preslav Nakov, and Isabelle Augenstein. 2021. Cross-domain label-adaptive stance detection. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9011–9028.

Miriam Hernández and Dani Madrid-Morales. 2020. Diversifying voice, democratizing the news? a content analysis of citizen news sources in spanish-language international broadcasting. *Journalism Studies*, 21(8):1076–1092.

Morten Hertzum. 2022. How do journalists seek information from sources? a systematic review. *Information Processing & Management*, 59(6):103087.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.

Aristidis Likas, Nikos Vlassis, and Jakob J Verbeek. 2003. The global k-means clustering algorithm. *Pattern recognition*, 36(2):451–461.

Kun Lu, Xin Cai, Isola Ajiferuke, and Dietmar Wolfram. 2017. Vocabulary size and its effect on topic representation. *Information Processing & Management*, 53(3):653–665.

Ville JE Manninen. 2017. Sourcing practices in online journalism: An ethnographic study of the formation of trust in and the use of journalistic sources. *Journal of Media Practice*, 18(2-3):212–228.

Kate C McLean, Moin Syed, Kristin Gudbjorg Haraldsson, and Alexandra Lowe. 2019. Narrative identity in the social world: The press for stability. *Handbook of Personality Psychology*.

Clara Meister and Ryan Cotterell. 2021. Language model evaluation beyond perplexity. *arXiv preprint arXiv:2106.00085*.

Byung-Doh Oh, Christian Clark, and William Schuler. 2022. Comparison of structural parsers and neural language models as surprisal estimators. *Frontiers in Artificial Intelligence*, 5:777963.

10

Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. 1998. The pagerank citation ranking: Bring order to the web. Technical report, Technical report, stanford University.

Kyeongman Park, Nakyeong Yang, and Kyomin Jung. 2023. Longstory: Coherent, complete and length controlled long story generation. *arXiv preprint arXiv:2311.15208*.

Jayr Pereira, Robson Fidalgo, Roberto Lotufo, and Rodrigo Nogueira. 2023. Visconde: Multi-document qa with gpt-3 and neural reranking. In *European Conference on Information Retrieval*, pages 534–543. Springer.

Dina Pisarevskaya. 2017. Deception detection in news reports in the Russian language: Lexics and discourse. In *Proceedings of the 2017 EMNLP Workshop: Natural Language Processing meets Journalism*, pages 74–79, Copenhagen, Denmark. Association for Computational Linguistics.

Dean Pomerleau and Delip Rao. 2017. Fake news challenge stage 1 (fnc-i): Stance detection. *Retrieved March*, 15:2023.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Revanth Gangi Reddy, Sai Chinthakindi, Zhenhailong Wang, Yi R Fung, Kathryn S Conger, Ahmed S Elsayed, Martha Palmer, and Heng Ji. 2021. Newsclaims: A new benchmark for claim detection from news with background knowledge. *arXiv preprint arXiv:2112.08544*.

Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*.

Guillaume Sanchez, Honglu Fan, Alexander Spangher, Elad Levi, Pawan Sasanka Ammanamanchi, and Stella Biderman. 2023. Stay on topic with classifier-free guidance. *arXiv preprint arXiv:2306.17806*.

Evan Sandhaus. 2008. The new york times annotated corpus. *Linguistic Data Consortium, Philadelphia*, 6(12):e26752.

Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*.

Benjamin Schiller, Johannes Daxenberger, and Iryna Gurevych. 2021. Stance detection benchmark: How robust is your stance detection? *KI - Künstliche Intelligenz*.

Ori Shapira, Ramakanth Pasunuru, Hadar Ronen, Mohit Bansal, Yael Amsterdamer, and Ido Dagan. 2021. Extending multi-document summarization evaluation to the interactive setting. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 657–677.

Alexander Spangher, Xinyu Hua, Yao Ming, and Nanyun Peng. 2023a. Sequentially controlled text generation. *arXiv preprint arXiv:2301.02299*.

Alexander Spangher, Jonathan May, Sz-Rung Shiang, and Lingjia Deng. 2021a. Multitask semi-supervised learning for class-imbalanced discourse classification. In *Proceedings of the 2021 conference on empirical methods in natural language processing*, pages 498–517.

Alexander Spangher, Nanyun Peng, Jonathan May, and Emilio Ferrara. 2021b. "don't quote me on that": Finding mixtures of sources in news articles. *arXiv preprint arXiv:2104.09656*.

Alexander Spangher, Nanyun Peng, Jonathan May, and Emilio Ferrara. 2023b. Identifying informational sources in news articles. *arXiv preprint arXiv:2305.14904*.

Alexander Spangher, Gireeja Ranade, Besmira Nushi, Adam Fourney, and Eric Horvitz. 2020. Characterizing search-engine traffic to internet research agency web properties. In *Proceedings of The Web Conference 2020*, pages 2253–2263.

Alexander Spangher, Xiang Ren, Jonathan May, and Nanyun Peng. 2022. Newsedits: A news article revision dataset and a novel document-level reasoning challenge. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 127–157.

Timoté Vaucher, Andreas Spitz, Michele Catasta, and Robert West. 2021. Quotebank: a corpus of quotations from a decade of news. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*, pages 328–336.

Xinyi Wang, Wanrong Zhu, Michael Saxon, Mark Steyvers, and William Yang Wang. 2023. Large language models are latent variable models: Explaining and finding good demonstrations for in-context learning. In *Thirty-seventh Conference on Neural Information Processing Systems*.

Yixue Wang and Nicholas Diakopoulos. 2021. Journalistic source discovery: Supporting the identification of news sources in user generated content. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–18.

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

11

Adina Williams, Tristan Thrush, and Douwe Kiela. 2022. Anlizing the adversarial natural language inference dataset.

Stephan Winter and Nicole C Krämer. 2014. A question of credibility–effects of source cues and recommendations on information selection on news sites and blogs. *Communications*, 39(4):435–456.

Kevin Yang, Yuandong Tian, Nanyun Peng, and Dan Klein. 2022. Re3: Generating longer stories with recursive reprompting and revision. *arXiv preprint arXiv:2210.06774*.

Lili Yao, Nanyun Peng, Ralph Weischedel, Kevin Knight, Dongyan Zhao, and Rui Yan. 2019. Plan-and-write: Towards better automatic storytelling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 7378–7385.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*.

GuoDong Zhou and KimTeng Lua. 1998. Word association and MI-Trigger-based language modeling. In *36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics, Volume 2*, pages 1465–1471, Montreal, Quebec, Canada. Association for Computational Linguistics.

# Appendix

In Appendix A, we include more, precise detail about our experimental methods. Then, Appendix B, we present more exploratory analysis to support our experiments, including comparisons between schemas. In Appendix D, we give a more complete set of definitions for the labels in each schema. In Appendix G, we define the unsupervised latent variable models we use as baselines, including providing details on their implementation.

## A  Additional Methodological Details

### A.1  Source Extraction

Before classifying sources, we first need to learn an attribution function (Equation 1) to identify the set of sources in news articles. Spangher et al. (2023b) introduced a large source attribution dataset, but their models are either closed (i.e. GPT-based) or underperforming. So, we train a high-performing open-source model using their dataset. We fine-tune GPT3.5-turbo [18], achieving a prediction accuracy of 74.5% on their test data[19]. Then, we label a large silver-standard dataset of 30,000 news articles and distill a BERT-base span-labeling model, described in (Vaucher et al., 2021), with an accuracy of 74.0%.[20] We use this model to score a large corpus of 90,000 news articles from the NewsEdits corpus (Spangher et al., 2022). We find that 47% of sentences in our documents can be attributed to sources, and documents each contain an average of 7.5 +-/5 sources. These statistics are comparable to those reported by Spangher et al. (2023b).

## B  Exploratory Data Analysis

We explore more nuances of our schemas, including comparative analyses. We start by showing a view of $\hat{Z}$, the conditions under which a schema best explains the observed results. In Tables 6 and 7, we show an extension of Table 4 in the main body: we show favored keywords across all schemas. (Note that in contrast to Table 4, we restrict the keywords we consider to a tighter range). When topics require a mixture of different informa-

---

[18]As of November 30th, 2023.

[19]Lower than the reported 83.0% accuracy of their Curie model. We formulate a different, batched prompt aimed at retrieving more data, see Appendix **??**
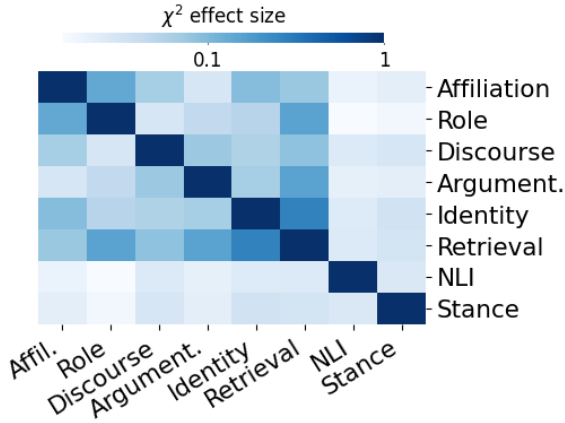
[20]All models will be released.

Figure 3: Correlation between 8 schemas, measured as Cramer's V (Cramér, 1999), or the effect-size measurement of the $\chi^2$ test of independence.
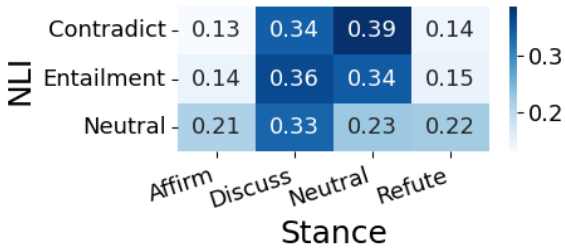


Figure 4: Stance and NLI schema definitions are not very aligned. We show conditional probability of labels in each category, $p(x|y)$ where $x =$ Stance and $y =$ NLI.
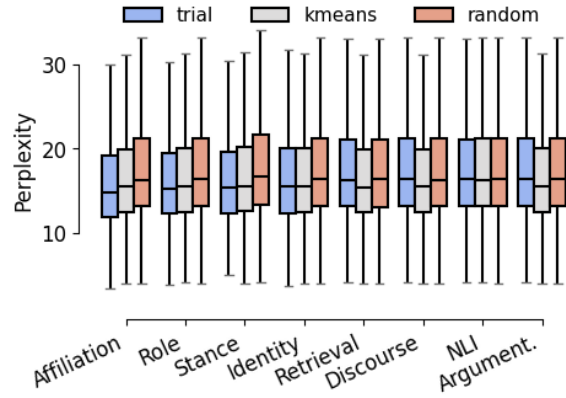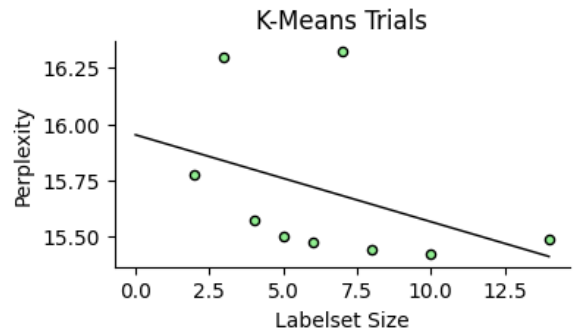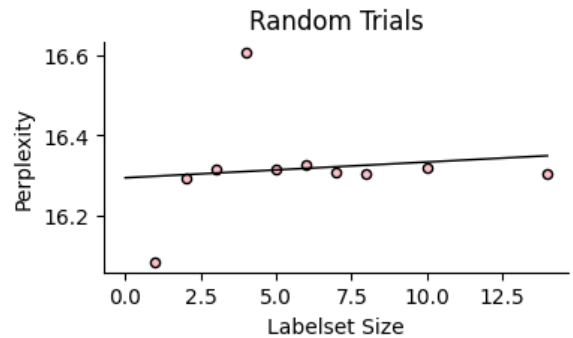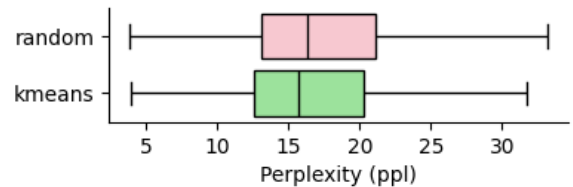


Figure 5: Distribution of conditional perplexity measurements across different experimental groups.



(a) Relationship between the size of the labelset and perplexity for kmeans trials



(b) Relational between the size of the labelset and perplexity for random trials.



(c) Distribution over perplexity scores for all random trials and kmeans trials, compared.

Figure 6: To explore the effects of labelset size, and confirm that conditional perplexity does align with basic intuitions, we compare Random trials and Kmeans trials across all of our labelset sizes.

tion types, like statistics, testimony, etc. *Argumentation* is favored. When story-telling is on topics like "Travel", "Education", "Quarantine (Life and Culture)", where it incorporates background, history, analysis, expectation, *Discourse* is favored. In Table 9, we show the top *Affiliations* per section of the newspaper, based on the NYT LDC corpus (Sandhaus, 2008).

Next, we further explore the relation between different labelsets. In Figure 5, we show the same story as in Table 3 in the Main Body, except with a broader view of the distributional shifts. As can be seen, by comparing differents between the means in Table 3 and the medians in 5, we see that the effect of outliers is quite large, which reduces the significance we observe. In 7, we show the correlation between perplexities across labelsets. We observe clusters in our schemas of particularly high correlation. Interestingly, this stands in contrast to Figure 3, which showed almost no relation between the tagsets. We suspect that outlier effects on perplexity (e.g. misspelled words, strange punctuation)

| Affiliation | Argumentation | Discourse | NLI |
|---|---|---|---|
| Inflation (Economics) | Race and Ethnicity | Travel and Vacations | Deaths (Fatalities) |
| Writing and Writers | Books and Literature | Quarantine (Life and Culture) | Murders, Homicides |
| United States Economy | Demonstrations, Protests and Riots | Education (K-12) | Law and Legislation |
| Race and Ethnicity | Travel and Vacations | Fashion and Apparel | States (US) |
| Disease Rates | Suits and Litigation | Murders, Homicides | Science |
| Real Estate and Housing (Residential) | Senate | Great Britain | Politics and Government |
| China | United States International Relations | Deaths (Fatalities) | Personal Profile |
| Supreme Court (US) | Deaths (Fatalities) | Pop and Rock Music | Children/ Childhood |
| Ukraine | Labor and Jobs | Demonstrations, Protests and Riots | China |

Table 6: Keyword topics that are best explained (i.e. have the lowest conditional perplexity) by the following schemas: Affiliation, Discourse, NLI. Broader topics, like "Inflation" which require sources from different backgrounds, favor Affiliation-based source selection, while topics integrating many different, possibly conflicting, facts, favor NLI-based selection.
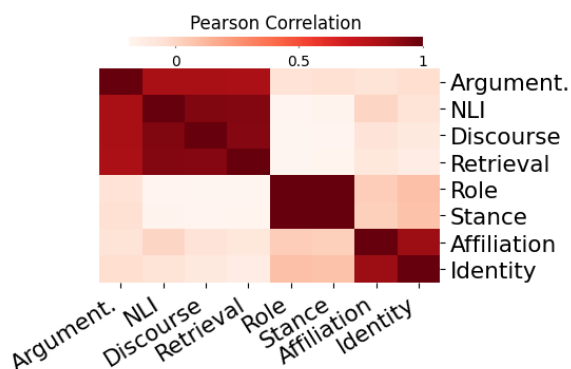


Figure 7: Pearson Correlation between conditional perplexity per document under different schemas.

has a high effect on relating different conditional perplexities, swamping the effects of the schema. This points to the caution in using perplexity as a metric; it must be well explored and appropriately baselined.

In Figure 4, we explore more why NLI and Stance are not very related. It turns out that many of the factual categories can fall in any one of the opinion-based categories. A lot of "Entailing" facts under NLI, for example, might be the the basis of "Discussion" under Stance. This points to the need to be cautious when using NLI as a stand-in for

Stance, as in (Reddy et al., 2021).

In Figures 6, we compare random and kmeans perplexities across the latent dimension size. Our experiments show that indeed, we are learning important cues about perplexity. As expected, "Random" assignments have almost no affect on the perplexity of the document, while "kmeans" assignments do. Increasing the dimensionality space of Kmeans, interestingly, *decreases* the median perplexity, perhaps because the Kmeans algorithm is allowed to capture more and more meaningful semantic differences between sources.

Finally, we discuss label imbalances in our classification sets. We do not observe a strong correlation between the number of labels in the schema and the classification accuracy ($\rho = -.16$). As seen in Table 8, many schema are highly skewed, with, for example, the minority class in Argumentation ("common ground") being present in less than .22% of sources. Using our classifiers to label the news articles compiled in Section A.1, we find that the schemas all offer different information. Figure 3 shows the effect size of the $\chi^2$ independence test, a test ranging from $(0, 1)$ which measures the relatedness of two sets of categorical variables (Cramér, 1999). The schemas are largely uncorrelated, with the highest correspondence be-

14

| Retrieval | Role | Identity | Stance |
|---|---|---|---|
| Actors and Actresses | Inflation (Economics) | United States Economy | Midterm Elections (2022) |
| Fashion and Apparel | House of Representatives | Disease Rates | Presidential Election of 2020 |
| Pop and Rock Music | Presidential Election of 2020 | Real Estate and Housing (Residential) | California |
| Elections | United States Economy | Movies | Storming of the US Capitol (Jan, 2021) |
| Personal Profile | Trump, Donald J | Education (K-12) | Vaccination and Immunization |
| Deaths (Fatalities) | Education (K-12) | Race and Ethnicity | News and News Media |
| Primaries and Caucuses | Elections, House of Representatives | Ukraine | United States Economy |
| Politics and Government | Supreme Court (US) | Trump, Donald J | Defense and Military Forces |
| Regulation and Deregulation of Industry | Computers and the Internet | Presidential Election of 2020 | Television |

Table 7: Keyword topics that are best explained (i.e. have the lowest conditional perplexity) by the following schemas: Retrieval, Role, Identity, Stance. Political topics, like "House of Representatives" which often have a mixture of different roles, favor Role-based source selection, while polarizing topics like "Storming of the US Capitol" favor Stance.

| Schema | n | H | % Maj. | % Min. |
|---|---|---|---|---|
| Affiliation | 14 | 2.2 | 32.9 | 0.46 |
| Role | 4 | 1.0 | 53.3 | 4.61 |
| Identity | 6 | 1.3 | 52.2 | 0.69 |
| Argument. | 6 | 1.1 | 62.9 | 0.22 |
| NLI | 3 | 1.1 | 40.4 | 22.6 |
| Stance | 4 | 1.3 | 34.8 | 15.5 |
| Discourse | 8 | 1.9 | 30.0 | 1.09 |
| Retrieval | 10 | 2.0 | 21.4 | 0.05 |

Table 8: Description of the size of each schema (n) and the class imbalance inherent in it, shown by: Entropy (H), % Representation of the Majority class (% Maj.) and % Representation of the Minority class (% Min.).

ing $\nu = .34$ between "Identity" and "Retrieval". We were surprised that NLI and Stance were not very related, as they have similar labelsets and have been used interchangeably (Reddy et al., 2021). This indicates that significant semantic differences exist between fact-relations and opinion-relations, resulting in different application of tags. We explore this in Appendix B.

## C Article Example

Here is an article example, annotated with different schema definitions, along with a description by the journalist of why they pursued the sources they did.

> *We mined state and federal court paperwork. We went looking for [previous] stories. We called police and fire communications people to determine [events]. We found families for interviews about [the subjects'] lives.*[21]

## D Further Schema Definitions

Here we provide a deeper overview of each of the schemas that we used in our work, as well as definitions that we presented to the annotators during annotation.

- **Affiliation:** Which group the source belongs to.

---

[21]https://www.nytimes.com/2017/01/23/insider/on-the-murder-beat-times-reporters-in-new-yorks-40th-precinct.html

| Newspaper Sections | Proportion of Sources in each Category | | |
| --- | --- | --- | --- |
| Arts | Individual: 0.29 | Media: 0.19 | Witness: 0.17 |
| Automobiles | Corporate: 0.41 | Witness: 0.17 | Media: 0.11 |
| Books | Individual: 0.26 | Media: 0.19 | Witness: 0.18 |
| Business | Corporate: 0.51 | Government: 0.2 | Industry Group: 0.06 |
| Dining and Wine | Witness: 0.28 | Individual: 0.18 | Media: 0.17 |
| Education | Government: 0.36 | Academic: 0.19 | Witness: 0.1 |
| Front Page | Government: 0.5 | Political Group: 0.09 | Corporate: 0.08 |
| Health | Government: 0.33 | Academic: 0.19 | Corporate: 0.12 |
| Home and Garden | Individual: 0.21 | Witness: 0.19 | Corporate: 0.17 |
| Job Market | Corporate: 0.26 | Individual: 0.15 | Witness: 0.14 |
| Magazine | Witness: 0.23 | Media: 0.2 | Individual: 0.18 |
| Movies | Individual: 0.28 | Media: 0.18 | Witness: 0.18 |
| New York and Region | Government: 0.36 | Witness: 0.13 | Individual: 0.12 |
| Obituaries | Government: 0.18 | Individual: 0.18 | Media: 0.16 |
| Opinion | Government: 0.43 | Media: 0.14 | Witness: 0.12 |
| Real Estate | Corporate: 0.33 | Government: 0.21 | Individual: 0.12 |
| Science | Academic: 0.4 | Government: 0.19 | Corporate: 0.1 |
| Sports | Other Group: 0.38 | Individual: 0.15 | Witness: 0.14 |
| Style | Individual: 0.23 | Witness: 0.2 | Corporate: 0.17 |
| Technology | Corporate: 0.41 | Government: 0.17 | Academic: 0.09 |
| The Public Editor | Media: 0.44 | Individual: 0.16 | Government: 0.16 |
| Theater | Individual: 0.34 | Witness: 0.18 | Media: 0.14 |
| Travel | Witness: 0.25 | Corporate: 0.21 | Government: 0.15 |
| U.S. | Government: 0.44 | Political Group: 0.12 | Academic: 0.08 |
| Washington | Government: 0.6 | Political Group: 0.1 | Media: 0.08 |
| Week in Review | Government: 0.37 | Academic: 0.11 | Media: 0.1 |
| World | Government: 0.54 | Media: 0.09 | Witness: 0.09 |

Table 9: Distribution over source-types with different *Affiliation* tags, by newspaper section.

– **Institutional:** The source belongs to a larger institution.

  1. **Government:** Any source who executes the functions of or represents a government entity. *(E.g. a politician, regulator, judge, political spokesman etc.)*

  2. **Corporate:** Any source who belongs to an organization in the private sector. *(E.g. a corporate executive, worker, etc.)*

  3. **Non-Governmental Organization (NGO):** If the source belongs to a nonprofit organization that operates independently of a government. *(E.g. a charity, think tank, non-academic research group.)*

  4. **Academic:** If the source belongs to an academic institution. Typically, these are professors or students and they serve an informational role, but they can be university administrators, provosts etc. if the story is specifically about academia.

  5. **Other Group:** If the source belongs or is acting on behalf of some group not captured by the above categories (please specify the group).

– **Individual:** The source does **NOT** belong to a larger institution.

  1. **Actor:** If the source is an individual acting on their own. *(E.g. an entrepreneur, main character, solo-acting terrorist.)*

  2. **Witness:** A source that is ancillary to events, but bears witness in either an active *(e.g. protester, voter)* or inactive *(i.e. bystander)* way.

16

**Headline: Services failed to prevent crime**

⎵'s voice became a preoccupation of ⎵, who told the police that he heard her calling his name at night. ← Government, Neutral

"Psychotic Disorder," detectives wrote in their report. ← *labels:* Government, Refute

"She had a strong voice," said Carmen Martinez, 85, a neighbor. ← Witness, Neutral

Records show a string of government encounters failed to help ⎵ as his mental health deteriorated. ← *labels:* Government, Agree

"This could have been able to be avoided," said ⎵'s lawyer. ← *labels:* Actor, Agree

Table 10: Informational sources synthesized in a single news article[22]. Source categorizations under two different schema: affiliation and stance. Our central question: *which schema best characterizes the kinds of sources needed to tell this story?*

3. **Victim:** A source that is affected by events in the story, typically negatively.

4. **Other:** Some other individual (please specify).

- **Role:**

  1. **Participant:** A source who is either directly making decisions on behalf of the entity they are affiliated with, or taking an active role somehow in the decision-making process.

  2. **Representative:** A source who is speaking on behalf of a *Participant*.

  3. **Informational:** A source who is giving information on ongoing decisions or events in the world, but is not directly involved in them.

  4. **Other:** Some other role that we have not captured (please specify).

- **Role Status:**

  1. **Current:** A source who is currently occupying the role and affiliation.

  2. **Former:** A source who *used* to occupy the role and affiliation.

  3. **Other:** Some other status that we have not captured (please specify).

We note that **Rote Status** was a schema that we collected, but ultimately did not end up modeling.

# E Example GPT Prompts

We give more examples for prompts.

## E.1 Source Attribution Prompts

In Section A.1, we discuss training a GPT3.5-Turbo model with Spangher et al. (2023b)'s source attribution dataset to create more labeled datapoints, which we then distil into a BERT model. We train a batched model to save on costs. The prompt takes the following form:

```
Input:
1.  <sent 1>
2.  <sent 2>
3.  <sent 3>...
Response:
1.  <attribution 1>
2.  <attribution 2>
...
```

Here is an example:

```
System Prompt:
You are a journalist's
fact-checker who identifies
sources providing information
for each sentence.  The user
will show you a sentences in
an article and you'll respond
with the source of the sentences.
Consider the whole article and be
sure to answer every question.
Answer either by directly
copying text in the article
OR with "passive-voice" when
a canonical source is clearly
consulted OR "journalist" when
a direct observation is made OR
"No source" when no source is
referenced, the information is
vague, or the source is unclear.
Do not make up names, or say
anything that is not in the
article besides those phrases
above.
User Input:  1:  BANGKOK | A
plane carrying key senior Laotian
government officials crashed
Saturday morning, leaving at
least four people dead, Laotian
diplomats said Saturday.
2:  Killed in the crash were
two top figures in the security
apparatus of the authoritarian
```

17

```
Lao government:  the deputy prime
minister, Douangchay Phichit,
and Thongbane Sengaphone, the
minister of public security,
according to two Lao diplomats.
3:  For a Communist party that
relies on force and intimidation
to stay in power, the loss of
what were arguably the two most
powerful people in the security
apparatus was a significant blow.
4:  The governor of Vientiane
province was also killed in the
crash.
5:  In addition to his post
as deputy prime minister Mr.
Douangchay was defense minister
and a member of the Politburo
, the highest decision – making
body of the Communist party.
6:  Mr.  Thongbane, the public
security head, was feared in the
country and was said to be one of
the officials leading a crackdown
against dissent over the past
year and half.
```
**Response:**
```
1:  Laotian diplomats
2:  No Source
3:  No Source
4:  No Source
5:  No Source
6:  No Source
```

### E.2  Stance-Based Prompts

In Section 3.3 we discuss the prompts we formu-
lated to do appropriate transfer learning from the
stance datasets others have annotated to our news
setting. Because in Stance detection, there are usu-
ally many claims made for each hypothesis, we
used batched prompts to save costs, in the follow-
ing form:

```
   Premise:  <premise>
Claim:
1.  <claim 1>
2.  <claim 2>
Response:  1.  <label 1>
2.  <label 2>
...
   Here is an example:
```

```
   System Prompt:  You are a
journalist's assistant who spots
opposing claims.  The user will
give you a premise and 5 claims.
Respond to each one, in numbered
order from 1 to 5, with a choice
from:  ['Neutral', 'Affirm',
'Discuss', 'Refute'].
Don't say anything else, and be
sure to answer each one.
```
**User Prompt**
```
Premise:  3-D printing will
change the world.
Claims:
1:  I can see 3D printing for
prototypes, and some custom work.
However manufacturing industries
use thousands of plastics and
thousands of metal alloys...
2:  Flash backwards to 1972,
Colorado, where the newly
enfranchised...
3:  This is precisely the way I
feel about 3D printers...another
way to fill the world with
plastic junk that will end up
in landfills, beaches, and yes,
mountains and oceans.  ...
4:  I am totally terrified with
the thought of 3-D printed,
non-traceable, guns and bullets
in every thugs hands.  May that
never happen.  But then Hiroshima
did (bad thing)...
5:  Hate to point out an obvious
solution is to tie the tax rate
to unemployment....
```
**Response:**
```
1:  Refute
2:  Neutral
3:  Refute
4:  Affirm
5:  Neutral
```

### E.3  GPT-2 Conditional Perplexity Prompts

In Section 4.1, we discuss crafting prompts for
GPT2-base models in order to calculate conditional
perplexity. We give the outline of our prompt. Here
is an example:

```
   Revelations from the artist's
autobiography threaten to cloud
```

```
her new show at the San Francisco
Museum of Modern Art.
<labels>
(1):  NGO,
(2):  Media,
(3):  Media,
(4):  Media,
(5):  Corporate
<text>
(1):  In a telephone interview
on Tuesday, the museums current
director, Christopher Bedford ,
said he welcomed the opportunity
to "be very outspoken about
the museums relationship to
antiracism" and ...
(2):  Last week a Chronicle
critic denounced the museums
decision to proceed with the
show.
(3):  Its longest-serving
curator, Gary Garrels, resigned
in 2020 soon after a post quoted
him saying, "Dont worry, we will
definitely continue to collect
white artists."
(4):  The website Hyperallergic
surfaced those comments in June .
(5):  And its previous director,
Neal Benezra, apologized to
employees after removing critical
comments from an Instagram post
following the murder of George
Floyd.
(6):  And the San Francisco
Museum of Modern Art has been
forced to reckon with what
employees have called structural
inequities around race.
(7):  The popular Japanese artist
Yayoi Kusama, whose " Infinity
Mirror Rooms " have brought
lines around the block for one
blockbuster exhibition after
another, has...'
```

## F  Combining Different Schema

We show how two schema, *Role* and *Affiliation* may be naturally combined. One function of journalism is to interrogate the organizations powering our society. Thus, many sources are from
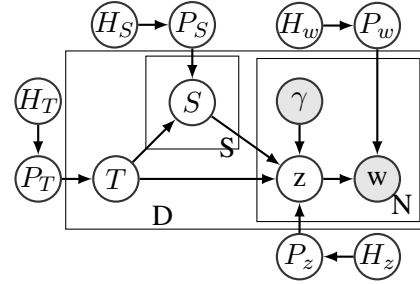


Figure 8: Plate diagram for Source Topic Model

**Affiliations**: *Government*, *Corporations*, *Universities*, *Non-Governmental Organizations* (NGOs). And, they have different *Roles* in these places. Journalists first seek to quote *decision-makers* or *participants*: presidents, CEOs, or senators. Sometimes decision-makers only comment though *Representatives*: advisors, lawyers or spokespeople. These sources all typically provide knowledge of the inner-workings of an organization. Broader views are often sought from *Informational* sources: experts in government or analysts in corporations; scholars in academia or researchers in NGOs. These sources usually provide broader perspectives on topics. Table 11 shows the intersection of these two schema.

## G  Latent Variable Models

As shown in Figure 8, our model observes a switching variable, $\gamma$ and the words, $w$, in each document. The switching variable, $\gamma$ is inferred and takes one of two values: "source word" for words that are associated with a source "background", for words that are not.

The model then infers source-type, $S$, document type $T$, and word-topic $z$. These variables are all categorical. All of the variables labeled $P_.$ in the diagram represent Dirichlet *P*riors, while all of the variables labeled $H_.$ in the diagram represent Dirichlet *H*yperpriors.

Our generative story is as follows:
For each document $d = 1, ..., D$:

1. Sample a document type $T_d \sim Cat(P_T)$
2. For each source $s = 1, ..., S_{(d,n)}$ in document:

    (a) Sample source-type $S_s \sim Cat(P_S^{(T_d)})$

3. For each word $w = 1, ... N_w$ in document:

    (a) If $\gamma_{d,w}$ = "source word", sample word-topic $z_{d,w} \sim Cat(P_z^{(S_s)})$

| | | | **Role** | | |
|---|---|---|---|---|---|
| | | | *Decision Maker* | *Representative* | *Informational* |
| **Affiliation** | *Institutional* | *Government* | President, Senator... | Appointee, Advisor... | Expert, Whistle-Blower... |
| | | *Corporate* | CEO, President... | Spokesman, Lawyer... | Analyst, Researcher... |
| | | *NGO* | Director, Actor... | Spokesman, Lawyer... | Expert, Researcher... |
| | | *Academic* | President, Actor... | Trustee, Lawyer... | Expert, Scientist... |
| | | *Group* | Leader, Founder... | Member, Militia... | Casual, Bystander... |
| | *Individ.* | *Actor* | Individual... | Doctor, Lawyer... | Family, Friends... |
| | | *Witness* | Voter, Protestor... | Spokesman, Poll... | Bystander... |
| | | *Victim* | Individual... | Lawyer, Advocate... | Family, Friends... |

Table 11: Our source ontology: describes the affiliation and roles that each source can take. A *source-type* is the concatenation of *affiliation* and *role*.

(b) If $\gamma_{d,w}$ = "background", sample word-topic $z_{d,w} \sim Cat(P_z^{(T_d)})$

(c) Sample word $w \sim Cat(z_{d,n})$

The key variables in our model, which we wish to infer, are the document type ($T_d$) for each document, and the source-type ($S_{(d,n)}$) for each source. It is worth noting a key difference in our model architecture: Bamman et al. (2013) assume that there is an unbounded set of mixtures over person-types. In other words, in step 2, $S_s$ is drawn from a document-specific Dirichlet distribution, $P_S^{(d)}$. While followup work by Card et al. (2016) extends Bamman et al. (2013)'s model to ameliorate this, Card et al. (2016) do not place prior knowledge on the number of document types, and rather draw from a Chinese Restaurant Process.[23] We constraint the number of *document-types*, anticipating in later work that we will bound news-article types into a set of common archetypes, much like we did for *source-types*.

Additionally, both previous models represent documents solely as mixtures of characters. Ours, on the other hand, allows the type of a news article, $T$, to be determined both by the mixture of sources present in that article, and the other words in that article. For example, a *crime* article might have sources like a government official, a witness, and a victim's family member, but it might also include words like "gun", "night" and "arrest" that are not included in any of the source words.

### G.1 Inference

We construct the joint probability and collapse out the Dirichlet variables: $P_w$, $P_z$, $P_S$, $P_T$ to solve

[23]Card et al. (2016) do not make their code available for comparison.

a Gibbs sampler. Next, we discuss the document-type, source-type, and word-topic inferences.

#### G.1.1 Document-Type inference

First, we sample a document-type $T_d \in 1, ..., T$ for each document:

$$
\begin{aligned}
p(T_d|T_{-d}, s, z, \gamma, H_T, H_S, H_Z) &\propto \\
(H_{TT_d} + c_{T_d,*}^{(-d)}) \times \prod_{s=1}^{S_d} \frac{(H_{Ss} + c_{T_d,s,*,*})}{(c_{T_d,*,*,*} + SH_S)} & \\
\times \prod_{j=1}^{N_T} \frac{(H_{zk} + c_{k,*,T_d,*})}{(c_{*,*,T_d,*} + KH_z)} &
\end{aligned} \tag{5}
$$

where the first term in the product is the probability attributed to document-type: $c_{T_d,*}^{(-d)}$ is the count of all documents with type $T_d$, not considering the current document $d$'s assignment. The second term is the probability attributed to source-type in a document: the product is over all sources in document $d$. Whereas $c_{T_d,s,*,*}$ is the count of all sources of type $s$ in documents of type $T_d$, and $c_{T_d,*,*,*}$ is the count of all sources of any time in documents of type $T_d$. The third term is the probability attributed to word-topics associated with the background word: the product is over all background words in document $d$. Here, $c_{k,*,T_d,*}$ is the count of all words with topic $k$ in document type $T_d$, and $c_{*,*,T_d,*}$ is the count of all words in documents of type $T_d$.

#### G.1.2 Source-Type Inference

Next, having assigned each document a type, $T_d$, we sample a source-type $S_{(d,n)} \in 1, ..., S$ for each source.

$$
\begin{aligned}
p(S_{(d,n)}|S_{-(d,n)}, T, z, H_T, H_s, H_z) &\propto \\
(H_{SS_d} + c_{T_d,S_{(d,n)},*,*}^{-(d,n)}) & \\
\times \prod_{j=1}^{N_{S_{d,n}}} \frac{(H_z + c_{z_j,*,S_{(d,n)},*,*})}{(c_{*,*,S_{(d,n)},*,*} + KH_z)} &
\end{aligned} \tag{6}
$$

20

The first term in the product is the probability attributed to the source-type: $c_{T_d,S_{(d,n)},*,*}^{-(d,n)}$ is the count of all sources of type $S_{(d,n)}$ in documents of type $T_d$, not considering the current source's source-type assignment. The second term in the product is the probability attributed to word-topics of words assigned to the source: the product is over all words associated with source $n$ in document $d$. Here, $c_{z_j,*,S_{(d,n)},*,*}$ is the count of all words with topic $z_j$ and source-type $S_{(d,n)}$, and $c_{*,*,S_{(d,n)},*,*}$ is the count of all words associated with source-type $S_{(d,n)}$.

### G.1.3 Word-topic Inference

Finally, having assigned each document a document-type and source a source-type, we sample word-topics. For word $i,j$, if it is associated with sources ($\gamma_{i,j}$ = Source Word), we sample:

$$
\begin{aligned}
&p(z_{(i,j)}|z^{-(i,j)}, S, T, w, \gamma, H_w, H_S, H_T, H_z) \propto \\
&(c_{z_{i,j},*,S_d,*,*}^{-(i,j)} + H_z z_{i,j}) \times \frac{c_{z_{i,j},*,w_{i,j},*}^{-(i,j)} + H_w}{c_{z_{i,j},*,*,*}^{-(i,j)} + V H_w}
\end{aligned}
\tag{7}
$$

The first term in the product is the word-topic probability: $c_{z_{i,j},*,S_d,*,*}^{-(i,j)}$ is the count of word-topics associated with source-type $S_d$, not considering the current word. The second term is the word probability: $c_{z_{i,j},*,w_{i,j},*}^{-(i,j)}$ is the count of words of type $w_{i,j}$ associated with word-topic $z_{i,j}$, and $c_{z_{i,j},*,*,*}^{-(i,j)}$ is the count of all words associated with word-topic $z_{i,j}$.

For word $i,j$, if it is associated with background word-topic ($\gamma_{i,j}$ = Background), we sample:

$$
\begin{aligned}
&p(z_{(i,j)}|z^{-(i,j)}, S, T, w, \gamma, H_w, H_S, H_T, H_z) \propto \\
&(c_{z_{i,j},*,T_d,*}^{-(i,j)} + H_z z_{i,j}) \times \frac{c_{z_{i,j},*,w_{i,j},*}^{-(i,j)} + H_w}{c_{z_{i,j},*,*,*}^{-(i,j)} + V H_w}
\end{aligned}
\tag{8}
$$

Equation 8 is nearly identical to 7, with the exception of the first term, the word-topic probability term, where $c_{z_{i,j},*,T_d,*}^{-(i,j)}$ refers to the count of words associated with word-topic $z_{i,j}$ in document-type $T_d$, not considering the current word. The second term, the word probability term, is identical.