

# NewsEdits 2.0: Learning the Intentions Behind Updating News

Anonymous ACL submission

## Abstract

As events progress, news articles often update with new information: if we are not cautious, we risk propagating outdated facts in many applications (e.g. large language model question asking (LLM Q&A)). In this work, we address this by *predicting which facts in a news article will update*. In the first part of this work, we isolate fact-updates in news revisions. This is challenging: although large news revisions corpora have been published (Spangher et al., 2022), news articles may update for many reasons (e.g. factual, stylistic, narrative). We introduce the *NewsEdits 2.0* taxonomy, an edit-intentions schema that separates fact updates from stylistic and narrative updates in news writing, annotate over 9,200 pairs of sentence revisions and train high-scoring ensemble models to apply this schema. Then, taking a large dataset of silver-labeled pairs, we show we can predict when facts will update in older article drafts. *Linguistic cues exist in news-writing that signal factual fluidity* and these can be learned with a big-data approach. With this insight, we demonstrate the value of these predictions by inducing LLMs to abstain from answering questions information is likely to be outdated. Using our models, LLM absention reaches *nearly oracle levels of accuracy*.

## 1 Introduction

News is the “first rough draft of history” (Croly, 1943). Its information is both valuable and fluid, prone to changes, updates, and corrections. As shown in Figure 1, the sentence: “Japan issued a tsunami advisory for the eastern coast” has a factual update, while while “A 7.1 magnitude quake struck...” does not. Intuitively, we might be able to predict this: an “advisory” is not likely to stay in effect indefinitely, while the “quake’s” existence is not likely to change. Indeed, if someone asks a question about the first sentence, we might want to abstain from answering definitively. For the second, however, it is better to answer directly.

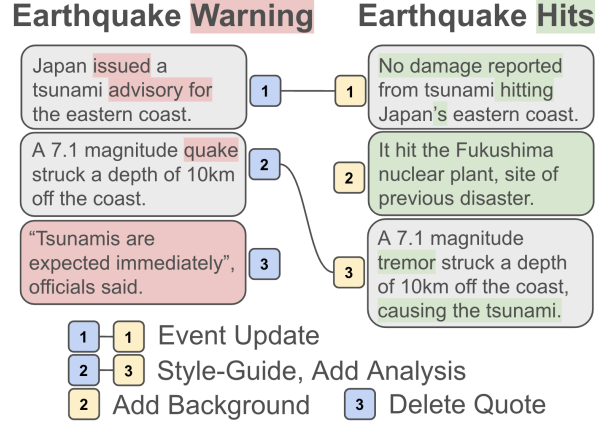


Figure 1: *NewsEdits 2.0*: We introduce a taxonomy of edit-types to characterize edits in news. Shown here, we identify factual updates (e.g. “Event Update” between 1-1), stylistic updates (e.g. “Style-Guide” between 2-3) and narrative updates (e.g. “Add Background” for sentence addition 2). Predicting which sentences on the left will have factual updates, we show, is particularly important for LLM Q&A tasks to prevent spreading outdated information.

Recent work has recognized the importance of testing LLM Q&A in dynamic settings (Jia et al., 2018; Liska et al., 2022). Indeed, Kasai et al. (2022)’s RealTimeQA benchmark specifically focused on measuring LLM Q&A performance for updating news documents. However, current approaches to these tasks rely on search engines retrieving updated information<sup>1</sup>. This leaves important linguistic and common-sense information on the table. As the example shown in Figure 1 demonstrates, cues exist that we, as humans, intuitively understand to signaling fluidity. Our hypothesis in this work is: *can we predict which facts in a news article will update? Can LLMs better contextualize answers to questions about these facts?*

In this work, we test this hypothesis. Spangher et al. (2022) released *NewsEdits*, a large corpus

<sup>1</sup>The latest entry of RealTimeQA was RAG + Google Custom Search. <https://realtimeqa.github.io/>.

Factual Edit	Style edit	Narrative/Contextual
Delete/Update/Add Eye-witness Account	Simplification	Delete/Add/Update Analysis
Delete/Add/Update Event	Emphasize/De- emphasize Importance	Delete/Add/Update Background
Delete/Add/Update Source-Doc.	Define term	Delete/Add/Update Anecdote
Correction	Style-Guide Adherence	
Delete/Add/Update Quote	Syntax Correction	
Additional Sourcing (Other)	Tonal Edits	
Additional Information (Other)	Sensitivity Consideration	
		<b>Other</b>
		Incorrect Link
		Unchanged
		Other

Figure 2: *NewsEdits 2.0*: Edit-Intentions Schema categories and their subcategories. In this work, we focus mainly on the *Factual Edit* category. See Appendix C.1 for definitions for all categories.

of article revision histories. However, *NewsEdits* is not detailed enough to study *factual* edits: articles update for many different reasons (e.g. factual, stylistic and narrative), and all types are treated the same in the corpus. *First, we need to first identify factual edit patterns.* We introduce *NewsEdits 2.0*, a taxonomy of edit-intentions for journalistic edits, shown in Figure 2. We hire professional journalists to annotate 507 article revision pairs with the *NewsEdits 2.0* schema. We then train an ensemble model to tag pairs of revisions and use them to silver-label a large corpus of revision pairs. Using this large silver-labeled corpus, we try to predict which facts in old articles will update. We find that models achieve a moderate macro-F1 of 58, overall, but by focusing on the sentences they predict are *highly likely* to update, we can have a real and positive impact. We simulate a RealTimeQA-style case where an LLM using Retrieval Augmented Generation (RAG) retrieves an outdated document. Without our predictions, the LLM abstains confidently, and wrongly more than it should. With them, the LLM achieves near-oracle level performance. In sum, our contributions are:

- We introduce the *NewsEdits 2.0* schema, with 4 coarse and 20 fine-grained categories, developed with professional journalists; train models to label these with 75.1 micro-F1; and release a large corpus of 4 million revision histories silver-labeled with edit intentions.
- We show that pretrained LLMs perform poorly at *predicting which facts in the old versions articles will update*, indicating that this important capability is not emergent during pre-training. While fine-tuning helps performance, LLMs still lag humans.

- Finally, we show via a use-case, Question Answering with Outdated Documents, that a failure to address these shortcomings can result in decreased performance for leading LLMs.

We believe that *NewsEdits 2.0* will open the door to many other exciting directions as well. Future work to isolate and predict stylistic and narrative event updates, we believe, can lead to interesting tools for journalists, writers and readers.

## 2 Related Work

Although most LLM Q&A benchmarks assume that information is static, recent work has increasingly explored LLM performance in the presence of dynamic, updating information (Jia et al., 2018; Liska et al., 2022). This growing direction is concisely captured by Kasai et al. (2022)’s statement: “*GPT-3 tends to return outdated answers when retrieved documents [are outdated]. Can [we] identify such unanswerable cases?*”

To our knowledge, the use of revision-histories to address this question, which we discuss in Section 5, is novel. News updates are an especially crucial domain to study: (1) news is socially important (Cohen et al., 2011) (2) LLMs are increasingly using news to better serve users (Hadero and Bauder, 2023) (3) news is more likely to deal with updating events than other domains (Spangher et al., 2022). Indeed, Kasai et al. (2022)’s RealTimeQA benchmark is built entirely on news data.

Edit-intention schemas have been developed for other types of revision histories, like Wikipedia (Yang et al., 2017), and Student Learner Essays (Zhang and Litman, 2015). In these works, researchers categorize the intention of each edit using similar schemas to what we have developed. While building *NewsEdits 2.0*, we were inspired by

the schemas developed by prior work and they provided a starting point for our taxonomy. We added edit-categories that were more journalism specific, like “Add Eye-witness Account”, and removed categories that were more specific to the aforementioned domains (Section 3.1.). The use-cases of these schemas has mainly focused on stylistic prediction tasks (e.g. text simplification (Woodsend and Lapata, 2011) and grammatical error correction (Faruqui et al., 2018)) or tasks specific to these corpora (e.g. building models to assess the validity of a student’s draft (Zhang and Litman, 2015), or counter vandalism on Wikipedia (Yang et al., 2017)). We are the first, to our knowledge, to develop tasks centered on news articles (Section 4) and to apply predictive analyses to fact-based edits.

Finally, a subtle yet significant novelty in this work are the improvements and visualization tools we introduce to make *NewsEdits* (Spangher et al., 2022) more accessible to users (Section 3.3). To our knowledge, this is the first attempt to provide visualizations for edit histories. We hope that our work can increase utilization, attention and understanding of news dynamics.

### 3 NewsEdits 2.0: Edit Intentions in Revision Histories

News articles update for different reasons, especially during breaking news cycles where facts and events update quickly (Saltzis, 2012). In this section, we introduce the edit-intentions schema we introduce for *NewsEdits 2.0*, our annotation, and our models to label edit-pairs. This lays groundwork for Section 4, where we will predict when facts change.

We wish identify categories of edits, in order to enable different investigations into these different update patterns. In other words, we describe the following update model:

$$p(l|s_i, s'_j, D, D') \quad (1)$$

where  $l$  is an *intention* (e.g. a “Correction” needs to be made),  $D$  and  $D'$  represent the older and newer versions of a news article, respectively, and  $s_i$  and  $s'_j$  are individual sentences where the update occurred.

### 3.1 Edit Intentions Schema

We work with two professional journalists and one copy editor<sup>2</sup> to develop an intentions schema. Building off work by (Zhang and Litman, 2015; Yang et al., 2017), we start by examining 50 revision-pairs sampled from *NewsEdits*. We developed our schema over four group meetings; before each one, we tagged examples and found edge-cases, then discussed as a group to add or collapse schema categories. Figure 2 shows our schema, which we organize into coarse and fine-grained labels. We incorporate existing theories of news semantics into our schema. For instance, “Event Updates” incorporates definitions of “events” (Dodgington et al., 2004), while “Add Background” incorporates theories of news discourse (Van Dijk, 1998). “Add Quote” incorporates definitions from informational source detection (Spangher et al., 2023) and “Add Anecdote” incorporates definitions from editorial analysis (Al-Khatib et al., 2016). See Appendix B.2 for a deeper discussion of the theoretical schemas that inform the *NewsEdits 2.0* schema. Finally, “Incorrect Link” is an attempt to correct sentence pairs that were erroneously (un)linked in *NewsEdits*.

### 3.2 Schema Annotation

We build an interface for annotators to provide intention labels for news article sentence pairs (see Appendix C.2). Annotators are shown definitions for each fine-grained intention and the articles to tag; they are instructed to tag each sentence. To recruit annotators, we posted on two list-serves for journalism industry professionals<sup>3</sup>. We train our annotators until they are all tagging with  $\kappa > .6$ . See Appendix for more details about our annotators.

### 3.3 Improvements over *NewsEdits*

Spangher et al. (2022) identified “edit-actions”, or “syntactic” edits in article revision histories (i.e. sentence additions, deletions and updates), which requires them to match sentences across article versions. They report matching sentences with 89.5 F1. This error rate is noticeable. We examined *NewsEdits*’s sentence matches and found that a large source of errors stem from poor sentence boundary detec-

<sup>2</sup>Collectively, these collaborators have over 50 years of experience in major newsrooms.

<sup>3</sup>The Association of Copy Editors (ACES) <https://aceseditors.org/> and National Institute for Computer-Assisted Reporting (NICAR) <https://www.ire.org/hire-ire/data-analysis/>.

Features	All		Fact		Style		Narrative	
	Macro	Micro	Macro	Micro	Macro	Micro	Macro	Micro
Baseline, <i>fine-grained</i>	45.8	73.6	32.0	47.2	58.6	39.9	52.0	39.9
+ NLI	48.6	74.1	45.7	50.4	55.2	38.7	43.6	38.7
+ Event	46.7	74.1	39.0	49.0	59.3	41.4	41.7	41.4
+ Quote	46.3	72.8	49.8	54.7	31.9	28.0	42.4	28.0
+ Collapsed Quote	51.2	73.9	38.7	47.6	58.3	39.4	51.4	39.4
+ Discourse	45.8	75.1	37.7	49.6	63.8	44.6	43.2	44.6
+ Argumentation	48.9	73.6	37.1	47.9	57.1	37.7	53.5	37.7
+ Discourse & Event	46.3	74.3	38.9	49.9	62.1	42.2	42.4	42.2
+ Discourse & Argumentation	47.8	74.1	56.8	50.5	31.4	32.2	41.1	32.2
+ Argumentation & Event	50.0	75.1	38.0	48.6	46.4	44.9	58.5	44.9
+ Quote & Discourse	51.2	72.2	40.5	45.3	62.8	43.0	48.7	43.0
+ Collapsed Quote & Discourse	49.6	73.9	45.6	49.4	58.9	39.1	47.9	39.1
+ Collapsed Quote & NLI	45.4	72.8	41.9	50.4	46.7	31.2	39.3	31.2
+ Collapsed Quote & NLI & Event	49.0	73.8	44.9	48.9	57.4	37.0	44.0	37.0
+ All	47.2	73.6	40.0	49.7	58.6	36.0	43.5	36.0
Baseline, <i>coarse-grained</i>	49.4	56.7	46.6		65.1		10.4	
+ Discourse & Arg. (Best model, Fact)	65.4	70.7	59.4		66.2		49.2	

Table 1: Various F1 scores (%) on our test set of the fine-tuned LED model with different combinations of features. Fact/Style/Narrative F1 scores are computed on instances that contain the corresponding labels, whereas All F1 scores are derived from all instances.

tion (SBD). Poor SBD creates an abundance of sentence stubs, which often over-match across revisions. We reprocessed the dataset from scratch using spaCy<sup>4</sup> instead of SparkNLP for SBD<sup>5</sup>, which we qualitatively observe to be better. For word-matching, we use albert-xxlarge-v2<sup>6</sup>’s embeddings (Lan et al., 2019) instead of TinyBert (Jiao et al., 2019). These steps, we find, increase our linking accuracy to 95%. We reprocess and re-release *NewsEdits*. In addition, we release a suite of visualization tools, based on D3<sup>7</sup> to enable further exploration of the corpus. See Appendix C.2 for an example.

### 3.4 Modeling Edit Intentions

Edit intentions are labeled on the sentence-level, and each sentence addition, deletion or update is potentially multiply labeled. Furthermore, document-level context is important: for instance in Figure 1, understanding that Sentence 2 adds background (“It hit the Fukushima nuclear plant, site of previous disaster.”) is aided by the surrounding sen-

tences contextualizing that a major event had just occurred.

Generative models have recently been shown to outperform classification-based models in document understanding tasks (Li et al., 2021; Huang et al., 2021). Inspired by this, we develop a sequence-to-sequence framework using the LongFormer-Encoder-Decoder (LED) architecture<sup>8</sup> (Beltagy et al., 2020) to predict the intent behind each edit. Specifically, our model processes the input  $x = [s_i || s'_j || D || D']$ .  $s_i$  or  $s'_j$  can also be  $\emptyset$ , which corresponds to the other sentence being an addition/deletion. The decoding target  $y_{i,j} = [l_1 || \dots || l_n]$  is a concatenation of potentially multilabeled intention labels  $l_i$  for pair  $s_i, s'_j$ .

**Experimental Variants and Results** As discussed in Section 3.1, we developed our schema to bring together different theories of news semantics. We experiment with integrating labels from models published these domains. We use models from the following papers: *Discourse* (Spangher et al., 2021), *Quote-Type Labeling* (Spangher et al., 2023), *Event Detection* (Hsu et al., 2021), *Textual Entailment* (Nie et al., 2020) and *Argumentation* (Al-Khatib et al., 2016). Labels generated from these external schema, denoted as  $f_{D_i}$  and  $f_{D'_j}$ , are appended to the model input  $x =$

<sup>4</sup><https://spacy.io/>, specifically, the en\_core\_web\_lg model.

<sup>5</sup><https://sparknlp.org/api/com/johnsnowlabs/nlp/annotators/sbd/pragmatic/SentenceDetector.html>

<sup>6</sup><https://huggingface.co/albert/albert-xxlarge-v2>

<sup>7</sup><https://d3js.org/>

<sup>8</sup><https://huggingface.co/allenai/led-base-16384>



	Narrative	Fact	Style
addition	840329	358900	104
deletion	330039	21671	6088
edit	411292	102499	644243

Table 2: Counts of coarse-grained semantic edit types, broken out by syntactic categories (for fine-grained counts, see Appendix).

$[D_i || D'_j || D || D'] [f_{D_i} || f_{D'_j}]$ . Incorporating these features increases Macro and Micro F1 by 5.5 and 1.5 points, respectively. For model details and schema definitions, see Appendix B.

### 3.5 Insights

We run the models trained in the last section over the entire *NewsEdits* corpus to generate silver-labels on all edit pairs. We present an exploratory analysis of these silver labels, with more material shown in the appendix. Table 2 shows the correlation between syntactic edit categories (defined by (Spangher et al., 2022)) and our semantic categories. As can be seen, categories like Addition have far more Narrative and Factual updates than Stylistic updates; Stylistic updates, on the other hand, are far more likely to occur between sentences. This makes sense, stylistic updates are likely smaller, local updates, while Narrative and Factual updates might include more rewriting.

Next, we explore if certain *kinds of articles* are more likely to have certain *kinds of edits*. We start by looking at broad news categories, shown in Table 4, derived from training a classifier on CNN News Groups<sup>9</sup>. “Politics” and “Sports” coverage are observed to have the highest level of fact-edits, relative to other categories, while Stylistic updates are prevalent in “Entertainment” pieces. Table 3 shows the kinds of edits in 6 different categories of news determined “socially beneficial”, by (Spangher et al., 2023)<sup>10</sup>. Even though “Fact” updates are rarer overall in sentence-level updates, they are more represented in Disaster and Safety categories.

Although we primarily focus on factual updates for the rest of the paper, we believe that there are many fruitful directions of future work examining other categories of updates.

<sup>9</sup><https://www.kaggle.com/code/faressayah/20-news-groups-classification-prediction-cnns>

<sup>10</sup>To group news articles in these categories, we use a classifier released by the authors

	Fact	Style	Other
Disaster	6.4	43.4	50.0
Elections	5.1	47.9	46.9
Environment	1.9	56.8	41.2
Labor	2.0	49.6	48.2
Other	3.7	50.7	45.5
Safety	4.7	46.6	48.6

Table 3: Distribution over update-types, across social-interest categories (Spangher et al., 2023).

	Fact	Style	Other
business	1.6	62.0	36.4
entertainment	3.3	65.5	31.1
health	2.1	61.0	36.9
news	2.8	57.0	40.2
politics	5.9	57.8	36.3
sport	3.5	59.3	37.2

Table 4: Distribution over update-types, across CNN section classifications.

## 4 Predicting Update Patterns

### 4.1 Problem Statement

In Section 3, we learned high-scoring models to categorize edit pairs (Equation 1). Now, we wish to leverage these to learn a predictive function:

$$p(l | s_i, D) \quad (2)$$

Where  $s_i$  and  $D$  are the *older* half of a revision pair or the last version of a revision history sequence. In other words, we wish to describe how this article *might* change. This, we hypothesize, can allow us to take actions to help users as news unfolds (Section 5).

Spangher et al. (2022) showed that structural predictions could be made about a news article’s development across time. They modeled “syntactic” changes in revision histories (e.g. “sentence will be *Added* or *Deleted*”). They found that whether an article *would* update or not was predictable with high F1, and they showed that expert journalists were surprisingly good at predicting how *much* and *where* an article would be update. However, authors stopped at this “syntactic” analysis. Here, we go a step further: with the semantic understanding of edits introduced in the prior section, we try to predict *how* information will change.

### 4.2 Dataset Construction

Because Spangher et al. (2022) already demonstrated predictability of syntactic edit actions, we decide to narrow our focus to revision pairs that we observe having substantial updates. We sample a set of 500,000 articles from *NewsEdits* that

Model	Features	Fact F1	None F1	Macro F1	Weighted F1
GPT-3.5	S	11.3	79.1	30.4	74.2
	DC	3.4	91.8	32.2	85.2
	FA	7.9	91.1	49.8	85.4
GPT-4	S	11.1	66.3	38.9	62.4
	DC	14.8	88.8	52.7	84.1
	FA	15.4	90.6	53.2	84.9
FT Longformer	S	21.2	92.3	57.4	87.0
	DC	22.3	93.0	87.8	87.4
	FA	25.4	91.4	58.0	86.4
Human Performance	S	41.2	75.3	58.6	69.2

Table 5: Individual F1 scores and macro and weighted F1 scores (%) on the golden test set for various evaluated models. S: sentence-only, DC: direct context, FA: full article.

have  $> 10\%$  sentences added and  $> 5\%$  deleted. Then, we use models developed in Section 3.4 to produce silver-standard labels. In other words, we assign labels  $l$  using both versions of a revision pair (Equation 1) and then we discard  $D'$ ,  $D'_j$  and try to predict  $l$  using just  $D$ ,  $D_i$  (Equation 2).

In order to prevent label leakage, we perform a chronological split of our dataset, splitting the earliest 80% of articles for training and the next 10% as the development set, and the most recent 10% as the test set. To keep computational and cost requirements reasonable and reproducible, we sample 16,000 sentences for the training set and 2,000 each for the development and the test set. In early experiments, we noticed that many fine-grained labels were too infrequent to model well, so we switched to predicting coarse-grained labels. As shown in Section 3.5, this classification problem is highly imbalanced: there are many more sentences that are not updated and of those that are, Style and Background/Narrative categories are more common. Thus, we balance the training dataset to have an equal number of classes for training. We sample from the true distribution for the development and test set. This yields a test set with 1,654 Others; 211 Fact-Updates; and 135 Style-Updates.

### 4.3 Experiments

We hypothesize that the broader article context is necessary to predict sentence-level update semantics, as sentences play a discursive role in the larger story. Thus, we experiment with predicting variants of equation 2: (i) sentence-only (S), or  $p(l|s_i)$ ; (ii) with the direct context (DC),  $p(l|s_{i-1}, s_i, s_{i+1})$ ;

and (iii) with the full article (FA),  $p(l|s_i, D)$ .

For each of (i), (ii) and (iii), we test both zero-shot approaches (i.e. prompted gpt-3.5-turbo and gpt-4); and fine-tuning approaches (i.e. longformer models and finetuned gpt-3.5-turbo models)<sup>11</sup>. For both approaches, we evaluate their performance on the same set of documents  $D_{test}^{gold}$ , which were part of the test set of our annotation, described in Section 3.2. In early trials, we try different variations on our experiments, like restricting the dataset to different subsets based on topic, like “Disaster” or “Safety”, which in Section 3.5 are more fact-heavy. However, we find negative results.

### 4.4 Results

Results are shown in 5. As can be seen, performance is overall low for detecting factual updates. However, we do observe performance increases from fine-tuning the longformer model, so to some degree this task is learnable. We recruit a former journalist with years of experience in newsrooms to provide human predictions of Equation 2 as an upper bound. After some observation of the training data, the journalist scores the test set. At 41.2 F1-score, the journalist sets an upper bound, but not a very high upper bound.

Next, we hypothesize that the middle of the distribution is actually very noisy: many sentences may look similar, but may or may not have had fact update due to chance. However, we hypothesize

<sup>11</sup>The longformer is trained with the same approach as the silver-label prediction step from Section 3.4 and gpt-3.5-turbo is trained using the OpenAI API.

Sentences with  $\uparrow p(l|s_i, D)$

There are no immediate reports of casualties.  
His trial has not yet started.  
Officials said attackers fired as many as 30 rockets in Friday’s assault.  
The rebel group did not immediately comment.

Table 6: A small sample of sentences in the high-likelihood region of  $p(l|s_i, D)$ . More examples shown in Table 11.

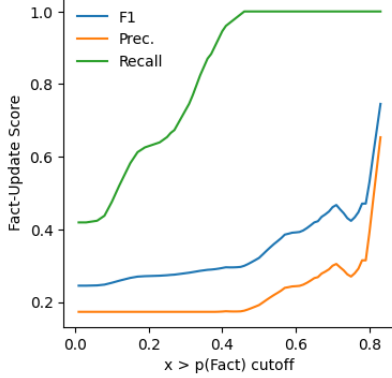


Figure 3: Performance of Fact-update model increases as we increasingly focus on a pool of documents that are categorized as high-likelihood under the model. In otherword, the model truly shines in the high-precision, high-probability realm.

that the examples that the model is most confident about, or the high-precision region, are more uniformly predictable, because they are sentences that more urgently need to receive an update. This is indicated in samples shown in Table 6, which include signals of immediacy (e.g. “no immediate reports”), future events (e.g. “...has not yet started”) and statistics (e.g. “30 rockets”). See Table 11 for more examples of high-probability sentences and Table 12 for examples of low-probability sentences. Figure 3 shows this exploration: as we restrict the pool of documents, we increase the performance.

## 5 Question Answering with Outdated Documents

Finally, we are ready to test whether the models learned in the last section, for predicting whether a sentence will have a factual update, can help us in dynamic LLM Q&A tasks. We set up a RealTimeQA-style task (Kasai et al., 2022), where an LLM is supplied by a retrieval system with potentially *out-of-date* information. We would like the LLM to *abstain* from answering a question if it

**Old sentence:** The White House **is** on lockdown after a vehicle struck a security barrier.

**New sentence:** The White House **was** on lockdown **for about an hour** after a vehicle struck a security barrier.’

**Question:** “If I visit the White House right now, will I get turned away?”

Table 7: Example candidate for LLM Q&A abstention.

suspects that the information it’s basing it’s answer on might be out-of-date.

Consider the scenario in Table 7. As humans, we could infer that the ongoing events in the **old sentence** would be of relatively short time-scale; an important building like the White House would likely reopen quickly. Thus, if a retriever retrieves only the **old sentence** for the LLM, even without knowledge of the **new sentence**, we would like the LLM to answer the **question** with something like: “yes, you will, but please check back; I do not have the most updated information and this might change quickly”. Confidently answering “Yes, you will be turned away” without any caution as to the updating nature of events is *wrong*.

### 5.1 Experiments

We design the following trials. We take pairs of sentences in the gold test set of our annotated data where an update occurred, and we ask GPT4 to generate 15 questions per pair of sentences, 5 each per category:

- **Easy:** questions where information from the older sentence *likely is not* in conflict with the newer one.
- **Medium** questions where information from the older sentence *might be* in conflict with the new one.
- **Hard** questions information from the older sentence *likely is* in conflict with a newer one,

In order to generate these questions, GPT4 is shown both versions and explicitly told to ask questions that fit each criteria (for all prompts, see Appendix D). Then, given this test set, we devise the following experimental variant. Each variant take in the *old sentence* and a *question*, generated previously:

	Easy			Medium			Hard		
	W. F1	Macro F1	Avg.	W. F1	Macro F1	Avg	W. F1	Macro. F1	Avg.
Baseline #1	55.9	35.8	55.9	8.8	8.1	8.8	38.8	28.0	38.8
Baseline #2	52.9	<b>49.6</b>	52.9	90.0	47.4	90.0	64.7	54.0	64.7
Experiment	<b>59.4</b>	48.9	<b>59.4</b>	<b>90.6</b>	61.1	<b>90.6</b>	<b>67.1</b>	<b>62.4</b>	<b>67.1</b>
Oracle	57.6	47.7	57.6	90.0	<b>63.3</b>	90.0	66.5	61.1	66.5

Table 8: A use-case for NewsEdits2.0: predicting when to abstain from factual question-answering, based on our predictions that material will update. We generate questions in different categories (easy, medium, hard)

	Easy	Medium	Hard
Baseline #1	0.0	0.0	0.0
Baseline #2	30.0	98.8	87.1
Experiment	10.6	95.9	74.1
Oracle	12.4	94.1	75.9

Table 9: Likelihood of refraining. In general, we wish to refrain only when we need to. Over-refraining is bad.

- **Baseline #1: Vanilla:** We formulate a basic prompt to GPT3.5, without alerting it to any possibly outdated material.
- **Baseline #2: Uniform** We formulate a prompt that warns GPT3.5 that some information might be outdated, and to refrain from responding if it things it is. However, this prompt is the same for all questions, so GPT has to rely on itself to detect outdated information.
- **Experiment:** Here, we input probabilities from our prediction model, binned into “low”, “medium”, “high” risk, into our prompt. In other words, we might tell GPT that we suspect there is a *high likelihood for the sentence being outdated*.
- **Oracle:** We feed in gold labels about whether a fact-update *will* occur in the next version of the article. We keep the phrasing the same as in the experimental version. This is designed to give us an upper bound.

**Evaluation** We evaluate performance of each prompting strategy as follows: we feed GPT4 the sentence pairs and the questions that were generated, and we ask:

- Is this question answerable given *just* the old sentence?

- Is the answer, using the old sentence, factually consistent with the information presented in the revised sentence?

If the answer is yes to both, then GPT should answer confidently. If either of the answers is “no”, then we want GPT to refrain from answering. Every time it refrains when it *should* be refraining is a success, otherwise is a failure. From these scores, we calculate an F1 score.

Our results are shown in Table 8 and 9. Interestingly, and perhaps unexpectedly, the experimental variant does as well if not better than the oracle. Perhaps the granularity of the prediction score helps GPT make a better assessment of the likelihood of update; perhaps our gold labels are a bit overly broad. As expected, **Baseline #2** has a strong performance (Table 8), but at the cost of far more refrains, show in Table 9.

## 6 Discussion and Conclusion

The ability of our prediction tags to recover near-oracle performance signals that factual edit prediction can serve a useful role in LLM Q&A. Although we have mainly tested our results in a high-likelihood region of the problem domain as a proof of concept, we suspect that if future work improves the models trained in Section 4.2, then we will see an increase in the ability to drive such abstentions.

We do suspect there to be an inherent upper bound in our ability to train such models. Stochasticity undoubtedly exists in the editing and revision process; for many factual updates where, perhaps, the ethical stakes of outdated information are lower, journalists may choose not to go back and revise. We still see such work as promising for LLM Q&A.

More broadly, the taxonomy introduces in *NewsEdits 2.0* has many rich directions. We hope in future work to revise directions around stylistic and narrative edits, both of which we believe can lead to better tools for computational journalists.



## 7 Ethical Considerations

### 7.1 Dataset

*NewsEdits* is a publicly and licensed dataset under an AGPL-3.0 License<sup>12</sup>, which is a strong “Copy-Left” license.

Our use is within the bounds of intended use given in writing by the original dataset creators, and is within the scope of their licensing.

### 7.2 Privacy

We believe that there are no adverse privacy implications in this dataset. The dataset comprises news articles that were already published in the public domain with the expectation of widespread distribution. We did not engage in any concerted effort to assess whether information within the dataset was libelous, slanderous or otherwise unprotected speech. We instructed annotators to be aware that this was a possibility and to report to us if they saw anything, but we did not receive any reports. We discuss this more below.

### 7.3 Limitations and Risks

The primary theoretical limitation in our work is that we did not include a robust non-Western language source. As our work builds off of *NewsEdits* as a primary corpora, it contains only English and French.

This work should be viewed with that important caveat. We cannot assume *a priori* that all cultures necessarily follow this approach to breaking news and indeed all of the theoretical works that we cite in justifying our directions also focus on English-language newspapers. One possible risk is that some of the information contained in earlier versions of news articles was updated or removed for the express purpose that it was potentially unprotected speech: libel, slander, etc. Instances of First Amendment lawsuits where the plaintiff was successful in challenging content are rare in the U.S. We are not as familiar with the guidelines of protected speech in other countries.

We echo the risk of the original *NewsEdits* authors: another risk we see is the misuse of this work on edits for the purpose of disparaging and denigrating media outlets. Many news tracker websites have been used for good purposes (e.g. holding newspapers accountable for when they make stylistic edits or try to update without giving notice). But

we live in a political environment that is often hostile to the core democracy-preserving role of the media. We focus on fact-based updates and hope that this resource is not used to unnecessarily find fault with media outlets.

### 7.4 Computational Resources

The experiments in our paper require computational resources. Our models run on a single 30GB NVIDIA V100 GPU or on one A40 GPU, along with storage and CPU capabilities provided by our campus. While our experiments do not need to leverage model or data parallelism, we still recognize that not all researchers have access to this resource level.

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

### 7.5 Annotators

We recruited annotators from professional journalism networks like the NICAR listserve, which we mention in the main body of the paper. All the annotators consented to annotate as part of the experiment, and were paid \$1 per task, above the highest minimum wage in the U.S. Of our 11 annotators, all were based in large U.S. cities. 8 annotators identify as white, 1 as Asian, 1 as Latinx and 1 as black. 8 annotators identify as male and 3 identifies as female. This data collection process is covered under a university IRB. We do not publish personal details about the annotations, and their interviews were given with consent and full awareness that they would be published in full.

### 7.6 References

#### References

- Khalid Al-Khatib, Henning Wachsmuth, Johannes Kiesel, Matthias Hagen, and Benno Stein. 2016. *A News Editorial Corpus for Mining Argumentation Strategies*. In *26th International Conference on Computational Linguistics (COLING 2016)*, pages 3433–3443. Association for Computational Linguistics.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.
- 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*.

<sup>12</sup><https://opensource.org/licenses/AGPL-3.0>

620	Sarah Cohen, James T Hamilton, and Fred Turner. 2011.	Sha Li, Heng Ji, and Jiawei Han. 2021. <a href="#">Document-level</a>	672
621	Computational journalism. <i>Communications of the</i>	<a href="#">event argument extraction by conditional generation</a> .	673
622	ACM, 54(10):66–71.	In <i>Proceedings of the 2021 Conference of the North</i>	674
623	H.D. Croly. 1943. <i>The New Republic</i> . v. 108. Republic	<i>American Chapter of the Association for Computa-</i>	675
624	Publishing Company.	<i>tional Linguistics: Human Language Technologies</i> ,	676
625	Ido Dagan, Oren Glickman, and Bernardo Magnini.	pages 894–908, Online. Association for Computa-	677
626	2005. The pascal recognising textual entailment chal-	tional Linguistics.	678
627	lenge. In <i>Machine learning challenges workshop</i> ,	Adam Liska, Tomáš Kociský, Elena Gribovskaya, Tay-	679
628	pages 177–190. Springer.	fun Terzi, Eren Sezener, and Devang Agrawal. 2022.	680
629	George R Doddington, Alexis Mitchell, Mark A Przy-	Cyprien de masson d’automne, tim scholtes, manzil	681
630	bocki, Lance A Ramshaw, Stephanie M Strassel, and	zaheer, susannah young, ellen gilsenan-mcmahon,	682
631	Ralph M Weischedel. 2004. The automatic content	sophia austin, phil blunsom, and angeliki lazarakidou.	683
632	extraction (ace) program-tasks, data, and evaluation.	2022. streamingqa: A benchmark for adaptation to	684
633	In <i>Lrec</i> , volume 2, pages 837–840. Lisbon.	new knowledge over time in question answering mod-	685
634	Manaal Faruqui, Ellie Pavlick, Ian Tenney, and Dipan-	els. In <i>International Conference on Machine Learn-</i>	686
635	jan Das. 2018. Wikiatomicedits: A multilingual corpus	<i>ing</i> .	687
636	of wikipedia edits for modeling language and	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal,	688
637	discourse. <i>arXiv preprint arXiv:1808.09422</i> .	Jason Weston, and Douwe Kiela. 2020. Adversarial	689
638	Haleluya Hadero and David Bauder. 2023. New york	NLI: A new benchmark for natural language under-	690
639	times sues microsoft, open ai over use of content.	standing. In <i>Proceedings of the 58th Annual Meeting</i>	691
640	<i>Globe &amp; Mail (Toronto, Canada)</i> , pages B1–B1.	<i>of the Association for Computational Linguistics</i> . As-	692
641	I Hsu, Kuan-Hao Huang, Elizabeth Boschee, Scott	sociation for Computational Linguistics.	693
642	Miller, Prem Natarajan, Kai-Wei Chang, Nanyun	Silvia Pareti, Tim O’keefe, Ioannis Konstas, James R	694
643	Peng, et al. 2021. Degree: A data-efficient	Curran, and Irena Koprinska. 2013. Automatically	695
644	generation-based event extraction model. <i>arXiv</i>	detecting and attributing indirect quotations. In <i>Pro-</i>	696
645	<i>preprint arXiv:2108.12724</i> .	<i>ceedings of the 2013 Conference on Empirical Meth-</i>	697
646	Kung-Hsiang Huang, Sam Tang, and Nanyun Peng.	<i>ods in Natural Language Processing</i> , pages 989–999.	698
647	2021. <a href="#">Document-level entity-based extraction as tem-</a>	Kostas Saltzis. 2012. Breaking news online: How news	699
648	<a href="#">plate generation</a> . In <i>Proceedings of the 2021 Confer-</i>	stories are updated and maintained around-the-clock.	700
649	<i>ence on Empirical Methods in Natural Language Pro-</i>	<i>Journalism practice</i> , 6(5-6):702–710.	701
650	<i>cessing</i> , pages 5257–5269, Online and Punta Cana,	Alexander Spangher, Jonathan May, Sz-Rung Shiang,	702
651	Dominican Republic. Association for Computational	and Lingjia Deng. 2021. Multitask semi-supervised	703
652	Linguistics.	learning for class-imbalanced discourse classification.	704
653	Zhen Jia, Abdalghani Abujabal, Rishiraj Saha Roy, Jan-	In <i>Proceedings of the 2021 conference on empirical</i>	705
654	nik Strötgen, and Gerhard Weikum. 2018. Tempques-	<i>methods in natural language processing</i> , pages 498–	706
655	tions: A benchmark for temporal question answering.	517.	707
656	In <i>Companion Proceedings of the The Web Confer-</i>	Alexander Spangher, Nanyun Peng, Jonathan May,	708
657	<i>ence 2018</i> , pages 1057–1062.	and Emilio Ferrara. 2023. Identifying informa-	709
658	Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao	tional sources in news articles. <i>arXiv preprint</i>	710
659	Chen, Linlin Li, Fang Wang, and Qun Liu. 2019.	<i>arXiv:2305.14904</i> .	711
660	Tinybert: Distilling bert for natural language under-	Alexander Spangher, Xiang Ren, Jonathan May, and	712
661	standing. <i>arXiv preprint arXiv:1909.10351</i> .	Nanyun Peng. 2022. Newsdits: A news article re-	713
662	Jungo Kasai, Keisuke Sakaguchi, Yoichi Takahashi,	vision dataset and a novel document-level reasoning	714
663	Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir	challenge. In <i>Proceedings of the 2022 Conference</i>	715
664	Radev, Noah A Smith, Yejin Choi, and Kentaro Inui.	<i>of the North American Chapter of the Association</i>	716
665	2022. Realtime qa: What’s the answer right now?	<i>for Computational Linguistics: Human Language</i>	717
666	<i>arXiv preprint arXiv:2207.13332</i> .	<i>Technologies</i> , pages 127–157.	718
667	Zhenzhong Lan, Mingda Chen, Sebastian Goodman,	Alexander Spangher, James Youn, Matthew Debutts,	719
668	Kevin Gimpel, Piyush Sharma, and Radu Soricut.	Nanyun Peng, and Jonathan May. 2024. Explaining	720
669	2019. Albert: A lite bert for self-supervised learn-	mixtures of sources in news articles.	721
670	ing of language representations. <i>arXiv preprint</i>	Teun A Van Dijk. 1998. <i>News as discourse</i> . Lawrence	722
671	<i>arXiv:1909.11942</i> .	Erlbaum Associates.	723
		Kristian Woodsend and Mirella Lapata. 2011. Wikisim-	724
		ple: Automatic simplification of wikipedia articles.	725
		In <i>Proceedings of the AAAI Conference on Artificial</i>	726
		<i>Intelligence</i> , volume 25, pages 927–932.	727

- Diyi Yang, Aaron Halfaker, Robert Kraut, and Eduard Hovy. 2017. Identifying semantic edit intentions from revisions in wikipedia. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2000–2010.
- W Victor Yarlott, Cristina Cornelio, Tian Gao, and Mark Finlayson. 2018. Identifying the discourse function of news article paragraphs. In *Proceedings of the Workshop Events and Stories in the News 2018*, pages 25–33.
- Fan Zhang and Diane Litman. 2015. Annotation and classification of argumentative writing revisions. In *Proceedings of the tenth workshop on innovative use of NLP for building educational applications*, pages 133–143.

	Addition	Deletion	Edit
Add/Delete/Update Background	806909	329652	411025
Add/Delete/Update Quote	303451	17995	46300
Incorrect Link	191022	125362	237437
Other (Please Specify)	84646	66929	65077
Add/Delete/Update Event Reference	37409	3645	56098
Add/Delete/Update Analysis	33426	390	268
Add/Delete/Update Eye-witness account	9772	0	3
Add/Delete/Update Source-Document	6639	2	28
Add/Delete/Update Information (Other)	1058	13	3
Additional Sourcing	573	15	29
Tonal Edits	102	6000	616514
Emphasize/De-emphasize Importance	1	32	1076
Syntax Correction	1	2	21729
Emphasize/De-emphasize a Point	0	53	1668
Simplification	0	0	3
Style-Guide Edits	0	1	3253
Correction	0	1	47

Table 10: Counts of fine-grained semantic edit types, broken out by syntactic categories

## A Appendix

### B Details of the LED Model

In this section, we describe the specifications of the LED model described in Section 3.4.

#### B.1 Input Template

The input to the LED model is shown below:

```
Predict the edit intention from
version 1 to version 2.
Version 1: SOURCE_SENTENCE
Version 2: TARGET_SENTENCE
Version 1 Document: SOURCE_DOCUMENT
Version 2 Document: TARGET_DOCUMENT
```

Here, **SOURCE\_DOCUMENT** ( $D$ ) and **TARGET\_DOCUMENT** ( $D'$ ) refer to the newer and older articles, while **SOURCE\_SENTENCE** ( $D_i$ ) and **TARGET\_SENTENCE** ( $D'_j$ ) represent a sentence with these articles.

#### B.2 Additional Schema

**NLI** We use textual entailment from (Dagan et al., 2005), which consists of *Entail*, *Contradict* and *Neutral*. These categories indicate whether two pieces of information refute each other, complement each other, or are neutral. We use a trained model by (Nie et al., 2020), which is an adversarially-trained Albert-xxlarge model, to la-

bel pairs of sentences (one from the old version, one from the new version).

**Event Detection** As described by Doddington et al. (2004) in the coding guidelines for the ACE-2005 dataset, “An Event is a specific occurrence involving participants. An Event is something that happens. An Event can frequently be described as a change of state.” Several datasets exist which label events in text, like ACE-2005, and a wide body of research has since emerged to model and detect events in text. Such models detect *triggers* (i.e. mostly verb-forms that signal the presence of an event); *types* (i.e. broad taxonomies that events fall into) and *arguments* (i.e. people, places or other lexical units associated with the occurrence of the event which further define it).

We use a model by (Hsu et al., 2021), designed to detect events in a wide variety of settings. We only consider whether an event trigger exists in a sentence, as a binary variable (0=no trigger exists, 1=trigger exists). Our theory is that this can help with tags like “Delete/Add/Update Event”.

**Argumentation** Defined in Al-Khatib et al. (2016), *Argumentation* is a type of discourse schema that defines what kinds of evidence the writer marshalls to make their point. Authors define the following categories: *Anecdote*, *Assumption*, *Common Ground*, *Statistics*, *Testimony*, *Other*. They primarily study news editorials (i.e. opinion



pieces), where they assume they have the most different kinds of argumentation categories. Spangher et al. (2021) and Spangher et al. (2024) show that these models can generally be applied helpfully across a broader news domain. We include them in the present study to capture aspects like “Anecdote” that capture framing aspects of journalistic writing.

**Quote** Quote-detection is a long-standing task, usually involving detecting the presence of direct or indirect quotes (Pareti et al., 2013). We use the broad definition of a “quote” as “information derived from any source external to the news article and the journalist’s own thoughts”, as defined in Spangher et al. (2023). Authors developed and released models for detecting when sentences had information that could be attributable to a named or unnamed source in the news article. We use these models to apply a simple binary indicator for whether or not the sentence contained a quote (1=sentence contains a quote, 0=it does not). We include this under the hypothesis that it can help us improve our detection in categories like “Delete/Add/Update Quote”.

**News Discourse** The News Discourse schema, as defined by Van Dijk (1998) views news stories as a sequence of structural elements, each serving a different narrative role. As implemented separately by (Choubey et al., 2020), (Yarlott et al., 2018) and (Spangher et al., 2021), the news discourse schema has undergone some modifications since Van Dijk (1998)’s original formulation, most notably to include current theories on event detection. It includes the following elements: *Main Event, Consequence, Previous Event, Current Context, Evaluation, Expectation, Historical Events, Anecdotal Event*. We believed that, since much of our edit schema was inspired by notions of narration, like “Delete/Add/Update Background”, we could get signal from this schema.

## C Annotation Details

In this section, we provide details of the annotation process, such as annotation guidelines and task allocation.

### C.1 Annotation Guidelines

To complete the task, look at each sentence: if it’s been added, updated, or deleted between drafts, try to determine based on your knowledge of the journalistic editing process why this was done.

You can specify multiple intentions for each add/delete/edit operation. Please also pay attention to when sentences are moved around in a document (i.e. if that was done to emphasize or de-emphasize that sentence), and when there might be errors to how we are linking sentences.

We devised these in consultation with professional journalists. However, if you are consistently annotating edits with “Other” (i.e. we are missing something in our schema), please let us know!

### Fact Edits:

- **Delete/Add/Update Eye-witness Account:** The writer deletes/adds/updates the contents for the events being described. This can either take the form of a quote (in which case this edit should be paired with a Quote Update), or a first-person account by the journalist.
- **Delete/Add/Update Event:** There is a change to some event in the world that the article covers and the article needs to be updated to reflect this. Usually, there are changes to the verbs in the article, but this can also include increased death counts, stock-market changes, etc.
- **Delete/Add/Update Source-Doc:** Additional written documents have been released by a government or company that warrant deletion/inclusion/update of the content of the article. For example, additional information included in an SEC filing, quarterly earnings report, IPCC report, etc.
- **Correction:** There are factual errors in the original version. The new version corrects the error.
- **Delete/Add/Update Quote:** There is an addition, editing or deletion of quotes in the article. Or, a quote from one person is swapped for a quote from another. Sometimes these updates are made with other intentions (e.g. to include a punchier quote, in which case it would also be a Preferential Edit. In these cases, please use the “+” button to add another intention dropdown.)
- **Additional Sourcing (Other):** The new version includes evidence of new sources for additional information, usually added for confirmation purposes. Note that this is different from Quote Update or Document Update

893	since Additional Sourcing doesn't have to re-	have the effect of some other edit intention,	939
894	sult in a new quote or document reference.	see the example, but cannot be fully ascribed	940
895	Can simply be an indication that the journalist	to other aims.	941
896	obtained new evidence.		
897	• <b>Additional Information (Other):</b> This edit	• <b>Sensitivity Consideration:</b> The journalist	942
898	intention is applied when the new version	rewrote the sentence because the original ver-	943
899	of the article includes details or context not	sion is inappropriate/ may be considered in-	944
900	present in the original version, which doesn't	sensitive.	945
901	necessarily fall under specific updates like eye-		
902	witness accounts, event changes, document	<b>Narrative Edits:</b>	946
903	updates, or sourcing alterations.	• <b>Delete/Add/Update Analysis:</b> The writer	947
904		deletes/adds/updates inferences from the pre-	948
905		sented information. These can be in the form	949
906		of analyses, expectations, or deeper under-	950
907		standings. These are usually forward-looking	951
908		rather than Background information, which is	952
909		usually past-looking.	953
910			
911	<b>Style Edits:</b>	• <b>Delete/Add/Update Background:</b>	954
912	• <b>Simplification:</b> reduces the complexity or	Delete/add/update contextualizing in-	955
913	breadth of discussion. This edit might also	formation to the article to help readers	956
914	remove information from the article.	understand the history, geography or signifi-	957
915		cance of a term, personal, place or company.	958
916		Note that contextualizing information is	959
917		not analysis, expectations, or projections,	960
918		which would fall into the Analysis intention	961
919		category.	962
920			
921	• <b>Emphasize/De-emphasize Importance:</b> The	• <b>Delete/Add/Update Anecdote:</b> The writer	963
922	sentence is moved up or down in the document	deletes, adds, or updates a brief, revealing	964
923	in order to make the sentence MORE/LESS	account of a person or event. This can be	965
924	prominent, or to emphasize/de-emphasize it's	a personal story, a particular incident, or a	966
925	connection to the events being described in	narrative snippet that exemplifies a point or	967
926	another sentence.	adds a humanizing or illustrative dimension to	968
927		the news piece. These anecdotes may serve to	969
928	• <b>Define term:</b> The author provides meaning or	engage the reader's interest, illuminate a fact,	970
929	differentiation to a term or concept that might	or provide a real-world example of abstract	971
930	be unknown to the reader. Note that this in-	concepts.	972
931	tention is DIFFERENT from the Background		
932	intention, which is more about providing con-	<b>Others:</b>	973
933	text, e.g. historical or geographic context for	• <b>Incorrect Link:</b> This refers to an error in our	974
934	a person, company, or place.	original linking of sentences. We have linked	975
935		two sentences that should NOT be linked.	976
936		This only pertains to 'Edit'ed or 'Unchanged'	977
937		sentences. Sentences should not be linked if	978
938		they are entirely unrelated — they have sub-	979
		stantially different syntax, intent, and purpose	980
		— and, by error, our algorithm said they were.	981
		If you identify an <b>Incorrect Link</b> AND there	982
		are more than one links, please specify (A) the	983
		index of the sentence in the other version that	984
		it should NOT be linked to via the dropdown	985
		(B) any other intention ascribed to this pair	986
		(i.e. Fact Deletion).	987

## C.2 Annotation Interface

Figure 6 shows the annotation interface for our task. Users are shown pairs of sentences, as identified in NewsEdits (Spangher et al., 2022) and have the option to annotate edits, additions and deletions with different edit intentions. Additionally, users can annotate when the links are incorrect.

## C.3 Annotation Task Distribution

We asked prospective applicants to describe their journalism experience, and selected this pool based on those having one or more year of professional editing experience. Then, we asked them to label revised sentences in five news articles, which we checked. We recruited 11 annotators who scored above 90% on these tests.

In Figure 4, we show the portion of annotation tasks assigned to each worker. As can be seen, we have a broad mix of users. Worker 11 is a professional journalist we worked most often with, and annotated a plurality of the tasks.

## D Prompts for Use-Case

### D.1 Question-Asking Prompts

**Easy** I will give you a sentence and you will give me an answer. It should be timely and related to the facts in the sentence. It should be a question that could go stale, especially for ongoing events, or facts like death counts that might update.

Here are some examples: example 1: sentence: "WASHINGTON (AP) – The White House is on lockdown after a passenger vehicle struck a security barrier." question: "Is the White House currently in lockdown – if I visit, will I get turned away?"

example 2: sentence: "The death count from the street bombing is 49 injured, 2 killed so far." question: "How many people have died so far?"

example 3: sentence: "The construction work left the bridge badly damaged and unsafe for passengers and is expected to remain so for days." question: "What route should I take? The bridge is the quickest way to work."

Ok, now it's your turn. Ask 5 different questions, output in a list. Don't say anything else. sentence:

**Easy** I will give you a sentence and you will give me 5 different questions. It should be directly answerable by the sentence.

Here are some examples: example 1: sentence: "WASHINGTON (AP) – The White House is on

lockdown after a passenger vehicle struck a security barrier." question: "What did the vehicle strike?"

example 2: sentence: "The death count from the 42nd street bombing is 49 injured, 2 killed so far." question: "Where did the bombing take place?"

example 3: sentence: "The construction work left the bridge badly damaged and unsafe for passengers and is expected to remain so for days." question: "What kind of work was being done?"

Ok, now it's your turn. Ask 5 different questions, output in a list. Don't say anything else. sentence:

**Hard** I will give you two sentences from an updating news article and you will give me 5 different questions. They should ideally focus on information that changes in between the sentences. So, if someone were to just look at the old sentence and you asked them your question, they would get it wrong.

Here are some examples: example 1: old sentence: "WASHINGTON (AP) – The White House is on lockdown after a passenger vehicle struck a security barrier." new sentence: "WASHINGTON (AP) – The White House was on lockdown for about an hour Friday after a passenger vehicle struck a security barrier." question: "Is the White House currently in lockdown – if I visit, will I get turned away?"

example 2: old sentence: "ISTANBUL (AP) – An earthquake with a preliminary magnitude of 6.2 shook western Turkey and the Greek island of Lesbos Monday, scaring residents and damaging buildings." new sentence: "ISTANBUL (AP) – An earthquake with a preliminary magnitude of 6.2 shook western Turkey and the Greek island of Lesbos on Monday, injuring at least 10 people and damaging buildings, authorities said." question: "Was anyone injured?"

example 3: old sentence: "Turkey's emergency management agency said there were no reports of casualties in the country." new sentence: "Turkey's emergency management agency said there were no reports of casualties and has dispatched emergency and health teams, and 240 family tents to the area as a precaution." question: "Is the Turkish emergency management doing anything as a precaution?"

Ok, now it's your turn. Ask 5 different questions, output in a list. Don't say anything else. old sentence: {old<sub>sentence</sub>}newsentence :

**Experimental Prompt** You are a helpful assistant who answers questions based on

1086	this news information: orig_sentence	1136
1087	We give this a {outdated_threshold	1137
1088	$\in \{high, medium, low\}$ } chance of there	
1089	being a fact update in this sentence.	
1090	That might mean some new information,	
1091	updating information. Answer cautiously	
1092	and do not give the user wrong/outdated	
1093	information. If the user's question looks	
1094	like it will still be relevant even if	
1095	the facts change, answer it directly. If	
1096	the user's question looks like it will	
1097	be outdated, say "I don't have the most	
1098	up-to-date information" and that's it.	
1099	Say nothing else. Do NOT say "I don't	
1100	have the most up-to-date information" AND	
1101	something else. Keep our estimate in	
1102	mind.	
1103	<b>Baseline 1</b> You are a helpful assistant	
1104	who answers questions based on this news	
1105	information: {orig_sentence}	
1106	Try to directly answer the users	
1107	question and say nothing else.	
1108	<b>Baseline 2</b> You are a helpful assistant	
1109	who answers questions based on this news	
1110	information: {orig_sentence}	
1111	This sentence might go out of date.	
1112	Answer cautiously and do not give the user	
1113	wrong/outdated information. If the user's	
1114	question looks like it will still be	
1115	relevant even if the facts change, answer	
1116	it directly. If the user's question looks	
1117	like it will be outdated, say "I don't	
1118	have the most up-to-date information" and	
1119	that's it. Say nothing else. Do NOT	
1120	say "I don't have the most up-to-date	
1121	information" AND something else.	
1122	<b>Oracle</b> You are a helpful assistant who	
1123	answers questions based on this news	
1124	information: {orig_sentence}	
1125	This sentence {oracle} have a major	
1126	fact update. That might mean some	
1127	new information, updating information.	
1128	Answer cautiously and do not give the user	
1129	wrong/outdated information. If the user's	
1130	question looks like it will still be	
1131	relevant even if the facts change, answer	
1132	it directly. If the user's question looks	
1133	like it will be outdated, say "I don't	
1134	have the most up-to-date information" and	
1135	that's it. Say nothing else. Do NOT	
	say "I don't have the most up-to-date	1136
	information" AND something else.	1137
	<b>D.2 Evaluation Prompts</b>	1138
	You are a helpful assistant. You will be	1139
	shown an old sentence, a revised sentence,	1140
	and a user-question. you will answer	1141
	the following 2 questions: 1. Is this	1142
	question answerable given JUST the old	1143
	sentence? Answer with "yes" or "no". Do	1144
	not answer anything else. If the answer	1145
	to 1 was yes, then proceed to the second	1146
	question, otherwise respond to question	1147
	2 with n/a 2. Does the question ask about	1148
	something that is factually consistent	1149
	with the information presented in the	1150
	revised sentence? Answer with "yes", "no"	1151
	or "n/a." Do not answer with anything	1152
	else.	1153
	<b>E Additional EDA</b>	1154
	We show the following different analyses	1155
	to support the findings in the main	1156
	body. In Figure 7, we perform an	1157
	error analysis on our best-performing	1158
	ensemble model, which includes tags	1159
	from Argumentation and Discourse. We	1160
	inspect the categories we are most likely	1161
	to get wrong. As can be seen, our	1162
	fine-grained accuracy is actually quite	1163
	low, indicating the value of future	1164
	work, perhaps collecting more training	1165
	data or employing LLMs to label more	1166
	silver-standard data. Many categories on	1167
	the diagonal have 0 labels, both because	1168
	many categories are low-count categories	1169
	(e.g. "Define Term", which does not have	1170
	any gold-truth labels in the test set),	1171
	as well as that more dominant categories	1172
	capture many of the predictions (e.g.	1173
	"Tonal Edits").	1174
	However, the problem is slightly less	1175
	sever on the coarse-grained level, shown	1176
	in Figure 5. By comparing these two	1177
	categories, we can see that many of the	1178
	errors we observed are on the fine-grained	1179
	level are within the same coarse-grained	1180
	category. We suspect that to raise	1181
	accuracy for fine-grained labels further,	1182
	we need further experimentation is needed.	1183
	Perhaps we can experiment with approaches	1184



1185 involving more specific fine-grained  
1186 models or with data augmentation.

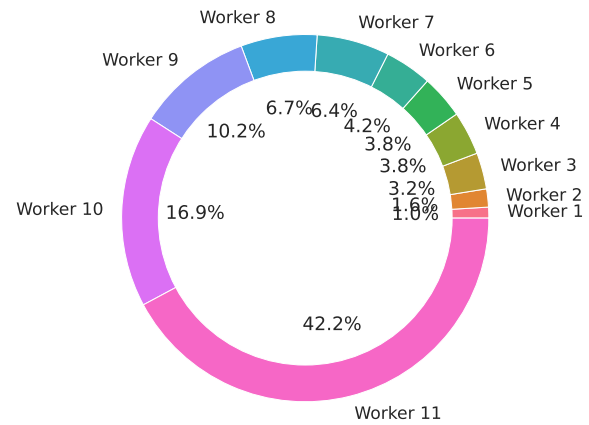


Figure 4: The portion of annotation tasks assigned to each worker.

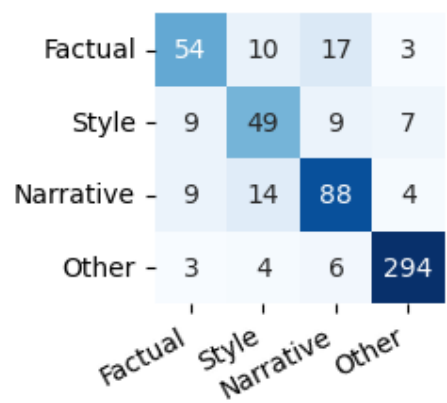


Figure 5: Coarse-grained confusion matrix for the LED model trained with Discourse and Argumentation features.

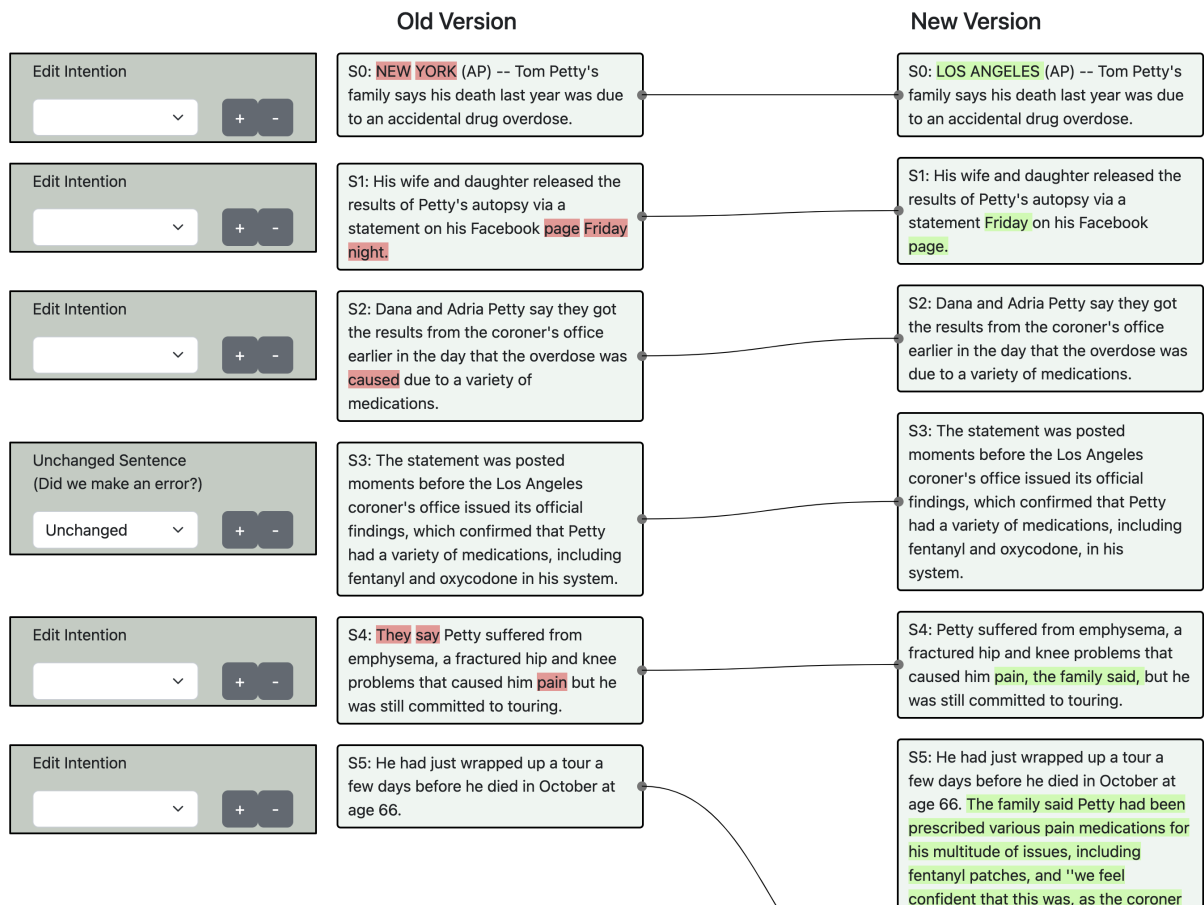


Figure 6: The interface for annotating edit intentions.

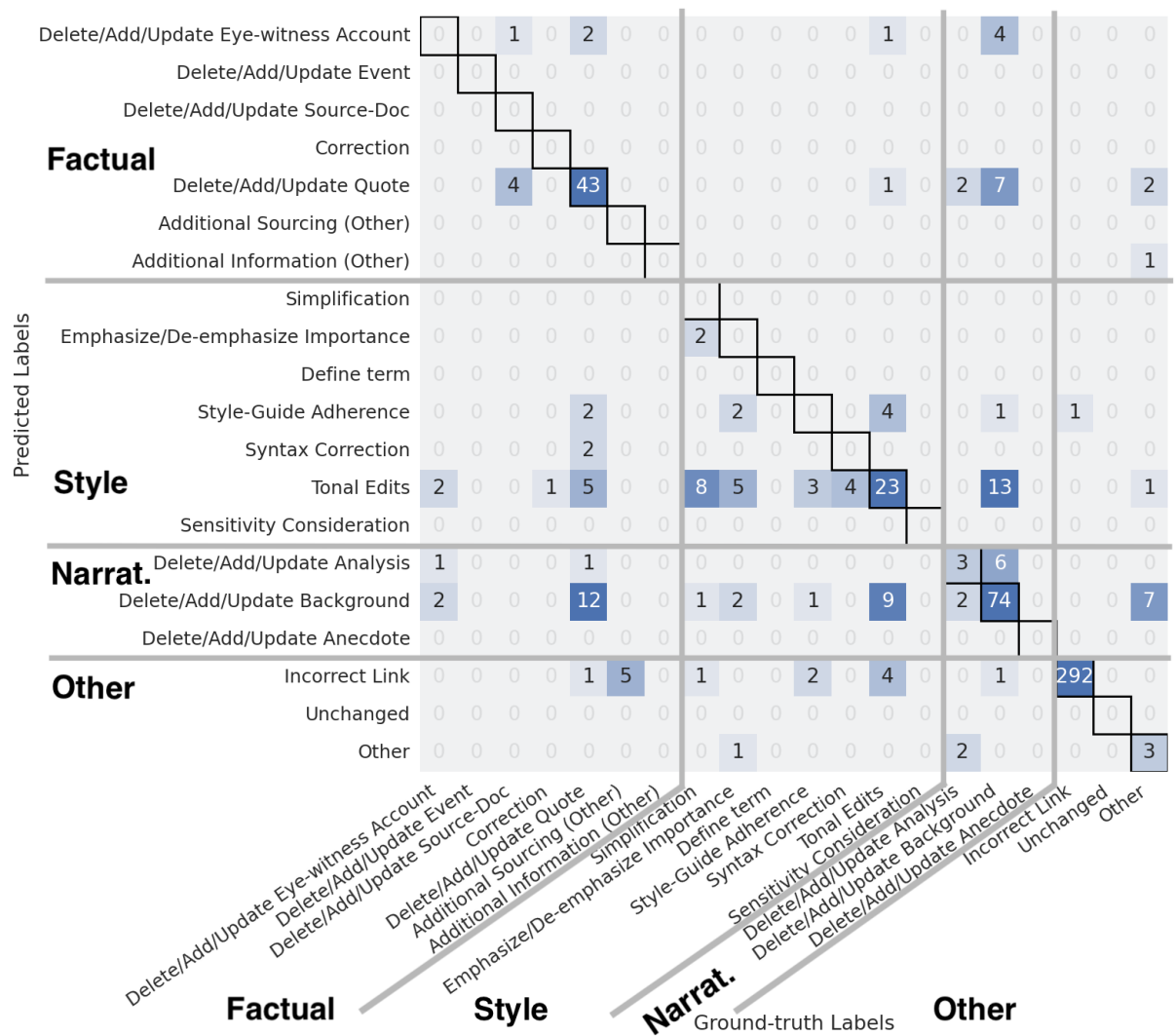


Figure 7: Fine-grained confusion matrix for the LED model trained with Discourse and Argumentation features.

---

Top Predictions for Content Evolution Prediction,  $p(l = \text{Fact Update} | D_i, D)$

---

The company takes this recommendation extremely seriously,” it said in a statement.

KABUL, Afghanistan — An Afghan official says a powerful suicide bombing has targeted a U.S. military convoy near the main American Bagram Air Base north of the capital Kabul.

WASHINGTON — The U.S. carried out military strikes in Iraq and Syria targeting a militia blamed for an attack that killed an American contractor, a Defense Department spokesman said Sunday. Mr. Causey, who reported his concern to authorities, was not charged in the indictment, which a grand jury returned last month, and did not immediately comment.

His trial has not yet started.

MEXICO CITY — A fiery freeway accident involving a bus and a tractor-trailer killed 21 people in the Mexican state of Veracruz on Wednesday, according to the authorities and local news outlets.

The indictment accuses Mr. Hayes, a former congressman, of helping to route \$250,000 in bribes to the re-election campaign of Mike Causey, the insurance commissioner.

No Kenyans died in the attack, Kenya’s military spokesman Paul Njuguna said Monday.

Mr. Manafort, 70, will most likely be arraigned on the new charges in State Supreme Court in Manhattan later this month and held at Rikers, though his lawyers could seek to have him held at a federal jail in New York, the people with knowledge said.

Officials said attackers fired as many as 30 rockets in Friday’s assault.

KABUL, Afghanistan — Gunmen attacked a remembrance ceremony for a minority Shiite leader in Afghanistan’s capital on Friday, wounding at least 18 people, officials said.

BEIRUT — A senior Turkish official says Turkey has captured the older sister of the slain leader of the Islamic State group in northwestern Syria, calling the arrest an intelligence “gold mine. ”

Paul J. Manafort, President Trump’s former campaign chairman who is serving a federal prison sentence, is expected to be transferred as early as this week to the Rikers Island jail complex in New York City, where he will most likely be held in solitary confinement while facing state fraud charges, people with knowledge of the matter said.

The watchdog, the Securities and Exchange Surveillance Commission, said Tuesday it made the recommendation to the government’s Financial Services Agency on the disclosure documents from 2014 through 2017.

There are no immediate reports of casualties.

It said the U.S. hit three of the militia’s sites in Iraq and two in Syria, including weapon caches and the militia’s command and control bases.

The rebel group did not immediately comment.

Kep provincial authorities later announced a total of five dead and 18 injured.

QUETTA, Pakistan — Attackers used a remotely-controlled bomb and assault rifles to ambush a convoy of Pakistani troops assigned to protect an oil and gas facility in the country’s restive southwest, killing six soldiers and wounding four, officials said Tuesday.

WASHINGTON — Senator Bernie Sanders of Vermont raised \$18.2 million over the first six weeks of his presidential bid, his campaign announced Tuesday, a display of financial strength that cements his status as one of the top fund-raisers in the sprawling Democratic field.

---

Table 11: Sample of the most likely fact-update sentences, as judged by our top-performing model. Top predictions reflect a combination of statistics, recent or upcoming events, and waiting for quotes.



---

Lowest Predictions for Content Evolution Prediction,  $p(l = \text{Fact Update} | D_i, D)$

---

Sir Anthony Seldon, vice-chancellor of the University of Buckingham, said: "Cheating should be tackled and the problem should not be allowed to fester any longer. "

He added: "This shows the extent to which a party which had such a proud record of fighting racism has been poisoned under Jeremy Corbyn. "

But he said his dream of making it in the game had turned into a nightmare. "

Adam Price, Plaid Cymru leader, said: "There is now no doubt that Wales should be able to hold an independence referendum. "

Others told how excited they had been when they were scouted by Higgins. "

The former Conservative deputy prime minister said it was "complete nonsense" to suggest Brexit could be done by Christmas. "

He said the QAA identified 17,000 academic offences in 2016 - but it was impossible to know how many cases had gone undetected. "

Nationalism leads a "false trail" in "exactly the opposite direction", he argued, "one that pits working people against each other, based on the accident of geography".

He also suggested that universities should adopt "honour codes", in which students formally commit to not cheating, and also recognise the consequences facing students who are subsequently caught.

He added: "But my experience is, if you make that threat, you don't actually need to follow through with the dreaded milkshake tax. "

He said: "There's an anger inside of me, a feeling of disgust that turns my stomach. "

Damian Hinds says it is "unethical for these companies to profit from this dishonest business".

She added: "His plan to hold another two referendums next year – and all the chaos that will bring – will mean that his government will not have time to focus on the people's priorities. "

We would be happy to talk to the Department of Education about their concerns. " "

I am determined to beat the cheats who threaten the integrity of our system and am calling on online giants, such as PayPal, to block payments or end the advertisement of these services - it is their moral duty to do so," said Mr Hinds.

The chief executive of Action on Smoking and Health, Deborah Arnott, also warned it would be a "grave error" to move away from taxing cigarettes. "

Rather than just taxing people more, we should look at how effective the so-called 'sin taxes' really are, and if they actually change behaviour. "

He added: "How many more red lines will be laid down by sensible Labour MPs, only for the leadership to trample right over them?

This shows that the complaints process is a complete sham," she tweeted. "

Mr Hinds added that such firms are "exploiting young people and it is time to stamp them out". "

One said he was abused by Higgins in a gym.

---

Table 12: Sample of the least likely fact-update sentences, as judged by our best-performing model. Predictions represent a combination of opinion quotes or anecdotes, projects and longer-term plans.