

Multi-Sentence Argument Linking

Seth Ebner* Patrick Xia* Ryan Culkin Kyle Rawlins Benjamin Van Durme

Johns Hopkins University
{seth, paxia}@cs.jhu.edu
{rculkin, kgr, vandurme}@jhu.edu

Abstract

We present a novel document-level model for finding argument spans that fill an event’s roles, connecting related ideas in sentence-level semantic role labeling and coreference resolution. Because existing datasets for cross-sentence linking are small, development of our neural model is supported through the creation of a new resource, **Roles Across Multiple Sentences (RAMS)**, which contains 9,124 annotated events across 139 types. We demonstrate strong performance of our model on RAMS and other event-related datasets.¹

1 Introduction

Textual event descriptions may span multiple sentences, yet large-scale datasets predominately annotate for events and their arguments at the sentence level. This has driven researchers to focus on sentence-level tasks such as semantic role labeling (SRL), even though perfect performance at such tasks would still enable a less than complete understanding of an event at the document level.

In this work, we approach event understanding as a form of *linking*, more akin to coreference resolution than sentence-level SRL. An event trigger *evokes* a set of roles regarded as latent arguments, with these implicit arguments then potentially linked to explicit mentions in the text.

Consider the example in Figure 1: the `AirstrikeMissileStrike` event (triggered by “bombarding”) gives rise to a frame or set of type-level roles (`attacker`, `target`, `instrument`, `place`) with the referents (“Russians”, “rebel outpost”, “aircraft”, “Syria”).² Intuitively we recognize the possible existence of fillers for these roles, for example, the `place` of the particular Air-

When Russian aircraft bombed a remote garrison in southeastern Syria last month, alarm bells sounded at the Pentagon and the Ministry of Defense in London.

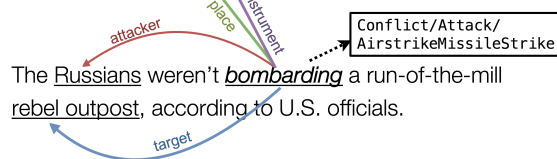


Figure 1: A passage annotated for an event’s type, *trigger*, and *arguments*. Each arc points from the trigger to the argument that fills the labeled role.

`strikeMissileStrike` event. These implicit arguments are linked to explicit arguments in the document (i.e., text spans). We refer to the task of finding explicit argument(s) to fill each role for an event as *argument linking*.

Prior annotation of cross-sentence argument links has produced small datasets, with a focus either on a small number of predicate types (Gerber and Chai, 2010, 2012; Feizabadi and Padó, 2014) or on a small number of documents (Ruppenhofer et al., 2010). To enable the development of a neural model for argument linking, we produce **Roles Across Multiple Sentences (RAMS)**, a dataset of 9,124 annotated events from news based on an ontology of 139 event types and 65 roles. In a 5-sentence window around each event trigger, we annotate the closest argument span for each role.

Our model builds on recent ideas in span selection models (Lee et al., 2018; He et al., 2018; Ouchi et al., 2018), used in this work for the multi-sentence argument linking task for RAMS and for several other event-based datasets (Gerber and Chai, 2012; Pradhan et al., 2013; Pavlick et al., 2016, AIDA Phase 1). On RAMS our best model achieves 68.3 F_1 , and it achieves 73.3 F_1 when event types are also known, outperforming strong baselines. We also demonstrate effective use of RAMS as pre-training for a related dataset.

*Equal Contribution

¹Data and code at <http://nlp.jhu.edu/rams/>.

² ϵ would indicate there is no explicit referent in the text.

Our main contributions are a novel model for argument linking and a new large-scale dataset for the task. Our dataset is annotated for arguments across multiple sentences and has broader coverage of event types and more examples than similar work. Our experiments highlight our model’s adaptability to multiple datasets. Together, these contributions further the automatic understanding of events at the document level.

2 Non-local Arguments

We are not the first to consider non-local event arguments; here we review prior work and refer to O’Gorman (2019) for further reading. Whereas local (sentence-level) event arguments are well-studied as semantic role labeling—utilizing large datasets such as OntoNotes 5.0 (Weischedel et al., 2013; Pradhan et al., 2013)—existing datasets annotated for non-local arguments are too small for training neural models.

Much of the effort on non-local arguments, sometimes called *implicit* SRL, has focused on two datasets: SemEval-2010 Task 10 (Ruppenhofer et al., 2010) and Beyond NomBank (henceforth **BNB**) (Gerber and Chai, 2010, 2012). These datasets are substantially smaller than RAMS: the SemEval Task 10 training set contains 1,370 frame instantiations over 438 sentences, while BNB contains 1,247 examples covering just 10 nominal predicate types. Multi-sentence AMR (MS-AMR) (O’Gorman et al., 2018; Knight et al., 2020) contains 293 documents annotated with a document-level adaptation of the Abstract Meaning Representation (AMR) formalism. O’Gorman (2019) notes that the relatively small size of the MS-AMR and SemEval datasets hinders supervised training. In contrast to these datasets, RAMS contains 9,124 annotated examples covering a wide range of nominal and verbal triggers.

Under the DARPA AIDA program, the Linguistic Data Consortium (LDC) has annotated document-level event arguments under a three-level hierarchical event ontology (see Figure 2) influenced by prior LDC-supported ontologies such as ERE and ACE. These have been packaged as the AIDA Phase 1 Practice³ and Eval⁴ releases (henceforth **AIDA-1**), currently made available to performers in the AIDA program and participants

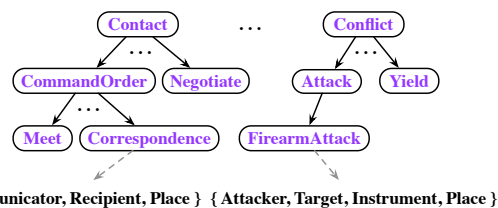


Figure 2: Subset of the AIDA-1 ontology illustrating the three-level Type/Subtype/Sub-subtype event hierarchy. Dashed gray edges point to roles for two event nodes, which have one role in common (Place).

in related NIST evaluations.⁵ AIDA-1 documents focus on recent geopolitical events relating to interactions between Russia and Ukraine. Unless otherwise noted, statistics about AIDA-1 pertain only to the Practice portion of the dataset.

For each document in LDC’s collection, only AIDA-salient events are annotated. This protocol does not guarantee coverage over the event ontology: 1,559 event triggers are annotated in the text portion of the collection, accounting for only 88 of the 139 distinct event sub-subtypes in the ontology. Our dataset, RAMS, employs the same annotation ontology but is substantially larger and covers all 139 types in the ontology. Figure 3 (§3) compares the two datasets.

Across multiple datasets, a substantial number of event arguments are observed to be non-local. For example, Gerber and Chai (2012) found that their annotation of non-local arguments added 71% (relative) role coverage to NomBank annotations. Additionally, 38.1% of the annotated events in AIDA-1 have an argument outside the sentence containing the trigger. This phenomenon is not surprising in light of the analysis of zero anaphora and definite null complements by Fillmore (1986) and the distinction between “core” and “non-core” frame elements or roles in FrameNet (Baker et al., 1998) and PropBank (Palmer et al., 2005).

As previous datasets have been small, various approaches have been taken to handle scarcity. To obtain more training data, Silberer and Frank (2012) created artificial instances from data annotated jointly for coreference and semantic roles. Roth and Frank (2013) automatically induced implicit arguments from pairs of comparable texts, but recovered a proportionally small set of additional arguments. Feizabadi and Padó (2015)

³LDC2019E04 (data); LDC2019E07 (annotations)

⁴LDC2019E42 (data); LDC2019E77 (annotations)

⁵While rarely freely released, historically such collections are eventually made available under a license to anyone, under some timeline established within a program.

combined existing corpora to increase and diversify sources of model supervision. Cheng and Erk (2018, 2019) approached the data scarcity problem by recasting implicit SRL as a cloze task and as a reading comprehension task, for which data can be generated automatically.

The TAC KBP event argument extraction task also seeks arguments from document contexts. However, in our work we are concerned with reified events (explicit mentions) and links between event mentions and argument mentions rather than entity-level arguments (coreference clusters).

3 RAMS

Motivated by the scarcity of data for training neural models to predict non-local arguments, we constructed **Roles Across Multiple Sentences** (RAMS), a crowd-sourced dataset with annotations for 9,124 events following the AIDA ontology. We employed the AIDA ontology in RAMS so-as to be most similar to an existing corpus already being investigated by various members of the community. Each example consists of a typed *trigger* span and 0 or more *argument* spans in an English document. A trigger span is a word or phrase that evokes a certain event type in context, while argument spans denote role-typed participants in the event (e.g., the *Recipient*). Trigger and argument spans are token-level $[start, end]$ offsets into a tokenized document.

Typically, event and relation datasets annotate only the argument spans that are in the same sentence as the trigger, but we present annotators with a *multi-sentence* context window surrounding the trigger. Annotators may select argument spans in any sentence in the context window.

3.1 Dataset Description

Data Source We used Reddit, a popular internet forum, to filter a collection of news articles to be topically similar to AIDA-1. After applying a set of criteria based on keywords, time period, and popularity (listed in Appendix A.1) we identified approximately 12,000 news articles with an average length of approximately 40 sentences.

Annotation We manually constructed a mapping from each event ((sub-)sub)type to a list of lexical units (LUs) likely to evoke that type.⁶ This mapping was designed to give high precision and

⁶For example, *Conflict/Attack/SetFire* is evoked by *inferno*, *blaze*, and *arson* (and word forms).

	Train	Dev	Test	Total
Docs	3,194	399	400	3,993
Examples	7,329	924	871	9,124
Event Types	139	131	–	139
Roles	65	62	–	65
Arguments	17,026	2,188	2,023	21,237

Table 1: Sizes and coverage of RAMS splits. RAMS covers all of the 139 event types and 65 roles types in the AIDA Phase 1 ontology.

low recall, in that for a given (Type, LUs) pair, the items in LUs are all likely to evoke the Type, although LUs can omit items that also evoke the Type. On average, $|LUs| = 3.9$.

We performed a soft match⁷ between every LU and every word in our text collection to select candidate sentences for each event type. This matching procedure produced approximately 94,000 candidates, which we balanced by sampling the same number of sentences for each LU.

Candidate sentences were then vetted by crowd-sourcing to ensure that they evoked their associated event type and had positive factuality. We collected judgments on approximately 17,500 candidate sentences, of which 52% were determined to satisfy these constraints, yielding 9,124 sentences containing a LU trigger. Using these sentences we then collected multi-sentence annotations, presenting annotators with a 5-sentence window containing two sentences of context before the sentence with the trigger and two sentences after.⁸ Annotators then selected in the context window a span to fill each of the event’s roles.

A window size of five sentences was chosen based on internal pilots and supported by our finding that 90% of event arguments in AIDA-1 are recoverable in this window size. Similarly, Gerber and Chai (2010) found that in their data almost 90% of implicit arguments can be resolved in the two sentences preceding the trigger.⁹ Arguments fall close to the trigger in RAMS as well: 82% of arguments occur in the same sentence as the trigger. On average, we collected 66 full annotations (trigger and arguments) per event type. Table 1 shows dataset size and coverage. All aspects of the protocol, including the annotation interface and instructions, are included in Appendix A.

⁷We stem all words and ignore case.

⁸If fewer than two sentences appeared before/after the trigger, annotators were shown as many sentences as were available.

⁹Arguments following the trigger were not annotated.

Inter-Annotator Agreement We randomly selected 93 tasks for redundant annotation in order to measure inter-annotator agreement, collecting five responses per task from distinct users. 68.5% of the time, all annotators mark the role as either absent or present. Less frequently (21.7%), four of the five annotators agree, and rarely (9.8%) is there strong disagreement.

We compute pairwise agreement for span boundaries. For each annotated (event, role) combination, we compare pairs of spans for which both annotators believe the role is present. 55.3% of the pairs agree exactly. Allowing for a fuzzier match, such as to account for whether one includes a determiner, spans whose boundaries differ by one token have a much higher agreement of 69.9%. Fewer spans agree on the start boundary (59.8%) than on the end (73.5%), while 78.0% match at least one of the two boundaries. We demonstrate data quality in §5.2 by showing its positive impact on a downstream task.

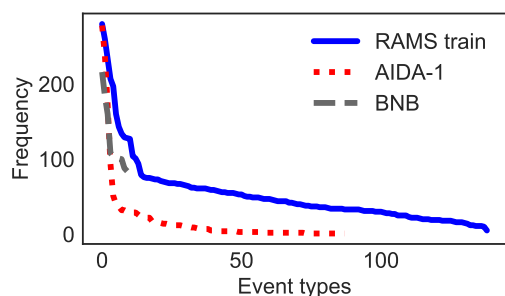


Figure 3: Comparison of frequency of event types in various datasets sorted by decreasing frequency in that dataset. RAMS has a heavier tail than AIDA-1 and BNB and broader coverage of events.

Comparisons to Related Datasets Comparisons of event type coverage among RAMS, AIDA-1, and BNB (Gerber and Chai, 2010, 2012) are given in Figure 3. RAMS provides larger and broader coverage of event types than do AIDA-1 and BNB. By design, BNB focuses on only a few predicate types, but we include its statistics for reference. More figures regarding type and role coverage are included in Appendix A.4.

Related Protocols Feizabadi and Padó (2014) also considered the case of crowdsourcing annotations for cross-sentence arguments. Like us, they provided annotators with a context window rather than the whole document, annotating two frames each with four roles over 384 predicates. Annota-

tors in that work were shown the sentence containing the predicate and the three previous sentences, unlike ours which shows two preceding and two following sentences.

Rather than instructing annotators to highlight spans in the text (“marking”), Feizabadi and Padó (2014) directed annotators to fill in blanks in templatic sentences (“gap filling”). We in contrast require annotators to highlight mention spans directly in the text.

Our protocol of event type verification followed by argument finding is similar to the protocol supported by interfaces such as SALTO (Burchardt et al., 2006) and that of Fillmore et al. (2002).

4 Model

We formulate *argument linking* as follows, similar to the formulation in Das et al. (2010). Assume a document \mathcal{D} contains a set of described events \mathcal{E} , each designated by a trigger—a text span in \mathcal{D} . The type of an event e determines the set of roles the event’s arguments may take, denoted \mathcal{R}_e . For each $e \in \mathcal{E}$, the task is to link the event’s roles with arguments—text spans in \mathcal{D} —if they are attested. Specifically, one must find for each e all (r, a) pairs such that $r \in \mathcal{R}_e$ and $a \in \mathcal{D}$. This formulation does not restrict each role to be filled by only one argument, nor does it restrict each explicit argument to take at most one role.

4.1 Architecture

Our model architecture is related to recent models for SRL (He et al., 2018; Ouchi et al., 2018). Contextualized text embeddings are used to form candidate argument span representations, \mathcal{A} . These are then pruned and scored alongside the trigger span and learned role embeddings to determine the best argument span (possibly none) for each event and role, i.e., $\operatorname{argmax}_{a \in \mathcal{A}} P(a \mid e, r)$ for each event $e \in \mathcal{E}$ and role $r \in \mathcal{R}_e$.

Representations To represent text spans, we adopt the convention from Lee et al. (2017) that has been used for a broad suite of core NLP tasks (Swayamdipta et al., 2018; He et al., 2018; Tenney et al., 2019b). A bidirectional LSTM encodes each sentence’s contextualized embeddings (Peters et al., 2018; Devlin et al., 2018). The hidden states at the start and end of the span are concatenated along with a feature vector for the size of the span and a soft head word vector produced by a learned attention mask over the word vectors

(GloVe embeddings (Pennington et al., 2014) and character-level convolutions) within the span.

We use this method to form representations of trigger spans, \mathbf{e} , and of candidate argument spans, \mathbf{a} . We learn a separate embedding, \mathbf{r} , for each role in the ontology, $r \in \mathcal{R}$. Since our objective is to link candidate arguments to event-role pairs, we construct an event-role representation¹⁰ by applying a feed-forward neural network ($F_{\tilde{a}}$) to the event trigger span and role embedding:

$$\tilde{\mathbf{a}}_{e,r} = F_{\tilde{a}}([\mathbf{e}; \mathbf{r}]) \quad (1)$$

This method is similar to one for forming edge representations for cross-sentence relation extraction (Song et al., 2018), but contrasts with prior work which limits the interaction between r and e (He et al., 2018; Tenney et al., 2019b).

Pruning Given a document with n tokens, there are $O(n^2)$ candidate argument text spans, which leads to intractability for large documents. Following Lee et al. (2017) and He et al. (2018), we consider within-sentence spans up to a certain width (giving $O(n)$ spans) and score each span, a , using a learned unary function of its representation: $s_A(a) = \mathbf{w}_A^\top F_A(\mathbf{a})$. We keep the top $\lambda_A n$ spans (λ_A is a hyperparameter) and refer to this set of high-scoring candidate argument spans as \mathcal{A} .

In an unpruned model, we need to create at least $\sum_e |\mathcal{R}_e|$ event-role representations and evaluate $\Omega(n \sum_e |\mathcal{R}_e|)$ combinations of events, roles, and arguments, which can become prohibitively large when there are numerous events and roles. Assuming the number of events is linear in document length, the number of combinations would be quadratic in document length (rather than quadratic in sentence length as in He et al. (2018)).

Lee et al. (2018) addressed this issue in coreference resolution, a different document-level task, by implementing a coarse pruner to limit the number of candidate spans that are subsequently scored. For our model, any role can potentially be filled (if the event type is not known). Thus, we do not wish to prematurely prune (e, r) pairs, so we must further prune \mathcal{A} . Rather than scoring $a \in \mathcal{A}$ with every event-role pair (e, r) , we assign a score between a and every event e . This relaxation reflects a loose notion of how likely an

¹⁰As a role for an event evokes an *implicit discourse referent*, this can be regarded as an implicit discourse referent representation.

argument span is to participate in an event, which can be determined irrespective of a role:

$$s_c(e, a) = \mathbf{e}^\top \mathbf{W}_c \mathbf{a} + s_A(a) + s_E(e) + \phi_c(e, a)$$

where \mathbf{W}_c is learned and $\phi_c(e, a)$ are task-specific features. We use $\mathcal{A}_e \subseteq \mathcal{A}$ to refer to the top- k -scoring candidate argument spans in relation to e .

Scoring We introduce a link scoring function, $l(a, \tilde{\mathbf{a}}_{e,r})$, between candidate spans $a \in \mathcal{A}_e$ and event-role pairs $\tilde{\mathbf{a}}_{e,r} = (e, r) \in \mathcal{E} \times \mathcal{R}$.¹¹ The scoring function decomposes as:

$$\begin{aligned} l(a, \tilde{\mathbf{a}}_{e,r}) &= s_{E,R}(e, r) + s_{A,R}(a, r) \\ &\quad + s_l(a, \tilde{\mathbf{a}}_{e,r}) + s_c(e, a), \quad a \neq \epsilon \quad (2) \\ s_E(e) &= \mathbf{w}_E^\top F_E(\mathbf{e}) \\ s_{E,R}(e, r) &= \mathbf{w}_{E,R}^\top F_{E,R}([\mathbf{e}; \mathbf{r}]) \\ s_{A,R}(a, r) &= \mathbf{w}_{A,R}^\top F_{A,R}([\mathbf{a}; \mathbf{r}]) \\ s_l(a, \tilde{\mathbf{a}}_{e,r}) &= \mathbf{w}_l^\top F_l([\mathbf{a}; \tilde{\mathbf{a}}_{e,r}; \mathbf{a} \circ \tilde{\mathbf{a}}_{e,r}; \\ &\quad \phi_l(a, \tilde{\mathbf{a}}_{e,r})]) \quad (3) \end{aligned}$$

where $\phi_l(a, \tilde{\mathbf{a}}_{e,r})$ is a feature vector containing information such as the (bucketed) token distance between e and a .¹² F_x are feed-forward neural networks, and \mathbf{w}_x are learned weights. The decomposition is inspired by Lee et al. (2017) and He et al. (2018), while the direct scoring of candidate arguments against event-role pairs, $s_l(a, \tilde{\mathbf{a}}_{e,r})$, bears similarities to the approach taken by Schenk and Chiarcos (2016), which finds the candidate argument whose representation is most similar to the prototypical filler of a frame element (role).

Learning We denote “no explicit argument” by ϵ and assign it link score $l(\epsilon, \tilde{\mathbf{a}}_{e,r}) \triangleq 0$, which acts as a threshold for the link function. For every event-role-argument triple (e, r, a) , we maximize

$$P(a \mid e, r) = \frac{\exp\{l(a, \tilde{\mathbf{a}}_{e,r})\}}{\sum_{a' \in \mathcal{A}_e \cup \{\epsilon\}} \exp\{l(a', \tilde{\mathbf{a}}_{e,r})\}}$$

Decoding We experiment with three decoding strategies: *argmax*, *greedy*, and *type-constrained*. If we assume each role is satisfied by exactly one argument (potentially ϵ), we can perform *argmax* decoding independently for each role:

$$\hat{a} = \operatorname{argmax}_{a \in \mathcal{A}_e \cup \{\epsilon\}} P(a \mid e, r)$$

¹¹If the type of e is known, then we could restrict $r \in \mathcal{R}_e$.

¹²Distance = $\max(e_{start} - a_{end}, a_{start} - e_{end})$.

To instead predict multiple non-overlapping arguments per role, we could use $P(\epsilon | e, r)$ as a threshold in *greedy* decoding (Ouchi et al., 2018).

We may know the gold event types and the mapping between events e and their permitted roles, \mathcal{R}_e . While this information can be used during training, we take a simpler approach of using it for *type-constrained* decoding (TCD). If an event type allows m_r arguments for role r , we keep only the top-scoring m_r arguments based on link scores.

4.2 Related Models

Our model is inspired by several recent span selection models (He et al., 2018; Lee et al., 2018; Ouchi et al., 2018), as well as the long line of neural event extraction models (Chen et al., 2015; Nguyen et al., 2016, *inter alia*). O’Gorman (2019) speculates a joint coreference and SRL model in which implicit discourse referents are generated for each event predicate and subsequently clustered with the discovered referent spans using a model for coreference, which is similar to the approach of Silberer and Frank (2012). O’Gorman (2019) further claims that span selection models would be difficult to scale to the document level, which is the regime we are most interested in. We focus on the implicit discourse referents (i.e., the event-role representations) for an event and link them to argument mentions, rather than cluster them using a coreference resolution system or aggregate event structures across multiple events and documents (Wolfe et al., 2015). Our approach is also similar to the one used by Das et al. (2010) for FrameNet parsing.

CoNLL 2012 SRL As our model bears similarities to the SRL models proposed by He et al. (2018) and Ouchi et al. (2018), we evaluate our model on the sentence-level CoNLL 2012 dataset as a sanity check. Based on a small hyperparameter sweep, our model achieves 81.4 F_1 when given gold predicate spans and 81.2 F_1 when not given gold predicates.¹³ Our model’s recall is harmed because our span pruning occurs at the document level rather than at the sentence level, which leads to overpruning in some sentences. Although our model is designed to accommodate cross-sentence links, it maintains competitive performance on sentence-level SRL.

¹³We use ELMo (Peters et al., 2018) in these experiments. He et al. (2018) achieve 85.5 F_1 with gold predicates and 82.9 F_1 without gold predicates, and Ouchi et al. (2018) achieve 86.2 F_1 with gold predicates.

Model	Dev. F_1	P	R	F_1
Our model	69.9	62.8	74.9	68.3
Our model ^{TCD}	75.1	78.1	69.2	73.3
Most common	17.3	15.7	15.7	15.7
Fixed trigger ^{TCD}	60.2	83.7	41.9	55.8
Context as trigger ^{TCD}	62.1	80.5	45.8	58.4
Distractor args	24.3	60.5	15.1	24.2
Distractor args ^{TCD}	24.2	68.8	14.3	23.7
No given args	8.7	20.2	3.5	6.0
No given args ^{TCD}	8.4	26.6	3.1	5.5

Table 2: P(recision), R(ecall), and F_1 on RAMS development and test data. TCD designates the use of ontology-aware type-constrained decoding.

5 RAMS Experiments and Results

In the following experiments, for each event the model is given the (gold) trigger span and the (gold) spans of the arguments. The model finds for each role the best argument(s) to fill it. Predictions are returned as trigger-role-argument triples.

We use feature-based BERT-base (Devlin et al., 2018)—mixing layers 9 through 12—by splitting the documents into segments of size 512 subtokens and encoding each segment separately.¹⁴

We perform preliminary sweeps across hyperparameter values, which are then fixed while we perform a more exhaustive sweep across scoring features. We also compare argmax decoding with greedy decoding during training. The best model is selected based on F_1 on the development set, and ablations are reported in Table 3. Our final model uses greedy decoding, $s_{A,R}$, and s_l and omits $s_{E,R}$ and s_c (see Equation 2). More details can be found in Appendix B.

The results using our model with greedy decoding and TCD are reported in Table 2. We also report performance of the following baselines: 1) choosing for each link the most common role (place), 2) using the same fixed trigger representation across examples, and 3) using the full context window as the trigger. Additionally, we experiment with two other data conditions: 1) linking the correct argument(s) from among a set of distractor candidate arguments provided by a constituency parser (Kitaev and Klein, 2018),¹⁵ and 2) finding the correct argument(s) from among all possible spans up to a fixed length.

¹⁴0.2% of the training documents span multiple segments.

¹⁵We take as the distractor arguments all (potentially overlapping) NPs predicted by the parser. On average, this yields 44 distractors per training document.

Model	Greedy	TCD
Our model	69.9	75.1
- distance score	69.0	74.3
- $s_l(a, \tilde{a}_{e,r})$	54.9	58.4
- $s_{A,R}(a, r)$	68.6	73.8
+ $s_{E,R}(e, r)$	69.5	74.4
+ $s_c(e, a)$	65.9	70.6
w/ argmax decoding	69.9	75.1
BERT 6–9	69.6	75.3
ELMo	68.5	75.2

Table 3: F₁ on RAMS dev data when link score components are separately included/excluded (Equation 2) or other contextualized encoders are used in the best performing model. TCD = type-constrained decoding.

For the distractor experiment, we use the same hyperparameters as for the main experiment. When not given gold argument spans, we consider all spans up to 5 tokens long and change only the hyperparameters that would prune less aggressively. We hypothesize that the low performance in this setting is due to the sparsity of annotated spans compared to the set of all enumerated spans. In contrast, datasets such as CoNLL 2012 are more densely annotated, so the training signal is not as affected when the model must determine argument spans in addition to linking them.

Finally, we examine the effect of TCD to see whether the model effectively uses gold event types if they are given. TCD filters out illegal predictions, boosting precision. Recall is still affected by this decoding strategy because the model may be more confident in the wrong argument for a given role, thus filtering out the less confident, correct one. Nevertheless, using gold types at test time generally leads to gains in performance.

5.1 Analysis

Ablations Ablation studies on development data for components of the link score as well as the contextualized encoder and decoding strategy are shown in Table 3. Type-constrained decoding based on knowledge of gold event types improves F₁ in all cases because it removes predictions that are invalid with respect to the ontology.

The most important link score component is the score between a combined event-role and a candidate argument. This result follows intuitions that s_l is the primary component of the link score since it directly captures the compatibility of the explicit argument and the implicit argument represented by the event-role pair.

Dist.	# Gold	# Predict	P	R	F ₁
-2	79 (26)	69 (21)	81.2	70.9	75.7
-1	164 (33)	151 (27)	76.8	70.7	73.7
0	1,811 (61)	1,688 (51)	77.7	72.4	75.0
1	87 (24)	83 (22)	78.3	74.7	76.5
2	47 (18)	39 (14)	87.2	72.3	79.1
Total	2,189 (62)	2,030 (52)	78.0	72.3	75.1

Table 4: Performance breakdown by distance (number of sentences) between argument and event trigger for our model using TCD over the development data. Negative distances indicate that the argument occurs before the trigger. # Gold and # Predict list the number of arguments (and unique roles) at that distance.

We also experiment with both ELMo (Peters et al., 2018) and BERT layers 6–9, which were found to have the highest mixture weights for SRL by Tenney et al. (2019a). We found that BERT generally improves over ELMo and layers 9–12 often perform better than layers 6–9.

Argument–Trigger Distance One of the differentiating components of RAMS compared to SRL datasets is its non-local annotation of arguments. At the same time, RAMS uses naturally occurring text so arguments are still heavily distributed within the same sentence as the trigger (Figure 5). This setting allows us to ask whether our model accurately finds arguments outside of the sentence containing the trigger despite the non-uniform distribution. In Table 4, we report F₁ based on distance on the development set and find that performance on distant arguments is comparable to performance on local arguments, demonstrating the model’s ability to handle non-local arguments.

Role Embeddings and Confusion We present in Figure 4 the cosine similarities between the learned 50-dimensional role embeddings in our model and also the errors made by the model under argmax decoding on the dev set.¹⁶ Some roles are highly correlated. For example, *origin* and *destination* have the most similar embeddings, possibly because they co-occur frequently and have the same entity type. Conversely, negatively correlated roles have different entity types or occur in different events, such as *communicator* compared to *destination* and *artifact*. We also observe that incorrect predictions are made more often between highly correlated roles and err

¹⁶Analysis of the confusion matrix with type-constrained decoding is less meaningful because the constraints, which rely on gold event types, filter out major classes of errors.

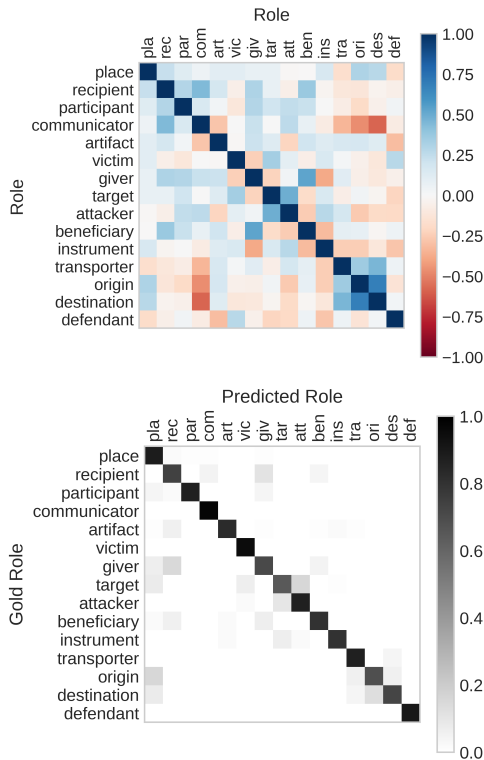


Figure 4: Embedding similarity (top) and row-normalized confusion (bottom) between roles for the 15 most frequent roles with our model. The full figures are included in Appendix C. Best viewed in color.

on the side of the more frequent role, as most errors occur below the diagonal.

Examples We present predictions from the development set which demonstrate some phenomena of interest. These are made without TCD, illustrating the model’s predictions without knowledge of gold event types.

In Table 5, the first example demonstrates the model’s ability to link a non-local argument which occurs in the sentence before the trigger. Greedy decoding helps the model find multiple arguments satisfying the same participant role, which also appear on either side of the trigger. In the second example, the model correctly predicts the driverpassenger, one of the rarer roles in RAMS (17 instances in the training set), consistent with the gold AccidentCrash event type.

In Table 6, the model fills roles corresponding to both the Death and the gold JudicialConsequences event types, thereby mixing roles from different event types. The predictions are plausible when interpreted in context and would be more accurate under TCD.

The EU’s leaders in Brussels are expected to play hard-
PARTICIPANT
 ball in negotiating Britain’s exit, to send a message to other states that might be contemplating a similar move. “Informal meeting of EU 27 next week without PM in the room to decide common negotiating position vs UK
PARTICIPANT
 on exit negotiations” —Faisal Islam.

SPEAKER: I’m Mary Ann Mendoza, the mother of Sergeant Brandon Mendoza, who was killed in a violent
DRIVERPASSENGER
 head-on collision in Mesa.
PLACE

Table 5: Two examples of correct predictions on the development set.

“Many people are saying that the Iranians killed the sci-
KILLER
 entist who helped the US because of Hillary Clinton’s hacked emails.” —8 August, Twitter. Shahran Amiri, the
VICTIM, DEFENDANT
 nuclear scientist executed in Iran last week, ...
PLACE

“Many people are saying that the Iranians killed the
JUDGECOURT
 scientist who helped the US because of Hillary Clinton’s hacked emails.” —8 August, Twitter. Shahran Amiri, the
CRIME
 nuclear scientist executed in Iran last week, ...
DEFENDANT
PLACE

Table 6: A partially correct prediction (top) and its corresponding gold annotations (bottom).

5.2 AIDA Phase 1

We also investigate how well RAMS serves as pre-training data for AIDA-1. A model using the hyperparameters of our best-performing RAMS model and trained on just English AIDA-1 Practice data achieves 19.1 F_1 on the English AIDA-1 Eval data under greedy decoding and 18.2 F_1 with TCD. When our best-performing RAMS model is fine-tuned to the AIDA task by further training on the AIDA-1 data, performance is improved to 24.4 F_1 under greedy decoding and 24.8 F_1 with TCD. The crowdsourced annotations in RAMS are therefore of sufficient quality to serve as augmentation to LDC’s AIDA-1. Experimental details are available in Appendix D.

6 Other Datasets

6.1 Beyond NomBank

The Beyond NomBank (BNB) dataset collected by Gerber and Chai (2010) and refined by Gerber and Chai (2012) contains nominal predicates (event triggers) and multi-sentence arguments, both of which are properties shared with RAMS.

To accommodate our formulation of the argu-

Field	Baseline*	Our Model
Victim Name	9.3 (54.1)	62.2 (69.6)
Shooter Name	4.7 (24.1)	53.1 (57.8)
Location	12.2 (18.9)	34.9 (63.3)
Time	68.1 (69.3)	62.9 (69.4)
Weapon	1.1 (17.9)	32.5 (49.6)

Table 7: Strict (and approximate) match F_1 on GVDB. Due to the different data splits and evaluation conditions, we are not directly comparable to the baseline (Pavlick et al., 2016), provided only for reference.

ment linking task, we modify the BNB data in two ways: 1) we merge “split” arguments, which in all but one case are already contiguous spans; and 2) we reduce each cluster of acceptable argument fillers to a set containing only the argument closest to the trigger. We also make modifications to the data splits for purposes of evaluation. Gerber and Chai (2012) suggest evaluation be done using cross-validation on shuffled data, but this may cause document information to leak between the train and evaluation folds. To prevent such leakage and to have a development set for hyperparameter tuning, we separate the data into train, dev, and test splits with no document overlap. Additional data processing details and hyperparameters are given in Appendix E. When given gold triggers and argument spans, our model achieves 75.4 F_1 on dev data and 76.6 F_1 on test data.

6.2 Gun Violence Database

The Gun Violence Database (GVDB) (Pavlick et al., 2016) is a collection of news articles from the early 2000s to 2016 with annotations specifically related to a *gun violence* event. We split the corpus chronologically into a training set of 5,056 articles, a development set of 400, and a test set of 500. We use this dataset to perform a MUC-style information extraction task (Sundheim, 1992). While GVDB’s schema permits any number of shooters or victims, we simply predict the first mention of each type. Pavlick et al. (2016) perform evaluation in two settings: a *strict* match is awarded if the predicted string matches the gold string exactly, while an *approximate* match is awarded if either string contains the other.

Assuming each document contains a single gun violence event triggered by the full document, our goal is to predict the value (argument) for each slot (role) for the event. As each slot is filled by exactly one value, we use argmax decoding.

While the baseline experiments of Pavlick et al. (2016) made sentence-level predictions focusing on five attributes, we make document-level predictions and consider the larger set of attributes. Table 7 shows our model’s performance on the shared subset of attributes, but the numerical values are not directly comparable because the prior work makes predictions on the full dataset and also combines some roles. Our results show that our model is suitable for information extraction tasks like slot filling. Appendix F contains information on hyperparameters and performance on the full set of roles. To our knowledge, our results are a substantial improvement over prior attempts to predict attributes of gun violence event reports, and we make our models available in the hopes of assisting social scientists in their corpus studies.

7 Conclusion

We introduced a novel model for document-level argument linking. Because of the small amount of existing data for the task, to support training our neural framework we constructed the RAMS dataset consisting of 9,124 events covering 139 event types. Our model outperforms strong baselines on RAMS, and we also illustrated its applicability to a variety of related datasets. We hope that RAMS will stimulate further work on multi-sentence argument linking.

Acknowledgments

We thank Craig Harman for his help in developing the annotation interface. We also thank Tongfei Chen, Yunmo Chen, members of JHU CLSP, and the anonymous reviewers for their helpful discussions and feedback. This work was supported in part by DARPA AIDA (FA8750-18-2-0015) and IARPA BETTER (#2019-19051600005). The views and conclusions contained in this work are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, or endorsements of DARPA, ODNI, IARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright annotation therein.

References

Collin F. Baker, Charles J. Fillmore, and John B. Lowe. 1998. *The Berkeley FrameNet project*. In *COLING*

- 1998 Volume 1: *The 17th International Conference on Computational Linguistics*.
- Aljoscha Burchardt, Katrin Erk, Anette Frank, Andrea Kowalski, and Sebastian Pado. 2006. **SALTO - a versatile multi-level annotation tool**. In *Proceedings of the Fifth International Conference on Language Resources and Evaluation (LREC'06)*, Genoa, Italy. European Language Resources Association (ELRA).
- Yubo Chen, Liheng Xu, Kang Liu, Daojian Zeng, and Jun Zhao. 2015. **Event extraction via dynamic multi-pooling convolutional neural networks**. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 167–176, Beijing, China. Association for Computational Linguistics.
- Pengxiang Cheng and Katrin Erk. 2018. **Implicit argument prediction with event knowledge**. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 831–840, New Orleans, Louisiana. Association for Computational Linguistics.
- Pengxiang Cheng and Katrin Erk. 2019. Implicit argument prediction as reading comprehension. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6284–6291.
- Dipanjan Das, Nathan Schneider, Desai Chen, and Noah A. Smith. 2010. **Probabilistic frame-semantic parsing**. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 948–956, Los Angeles, California. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. **BERT: pre-training of deep bidirectional transformers for language understanding**. *CoRR*, abs/1810.04805.
- Parvin Sadat Feizabadi and Sebastian Padó. 2014. **Crowdsourcing annotation of non-local semantic roles**. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics, volume 2: Short Papers*, pages 226–230, Gothenburg, Sweden. Association for Computational Linguistics.
- Parvin Sadat Feizabadi and Sebastian Padó. 2015. **Combining seemingly incompatible corpora for implicit semantic role labeling**. In *Proceedings of the Fourth Joint Conference on Lexical and Computational Semantics*, pages 40–50, Denver, Colorado. Association for Computational Linguistics.
- Charles J Fillmore. 1986. Pragmatically controlled zero anaphora. In *Annual Meeting of the Berkeley Linguistics Society*, volume 12, pages 95–107.
- Charles J. Fillmore, Collin F. Baker, and Hiroaki Sato. 2002. **The FrameNet database and software tools**. In *Proceedings of the Third International Conference on Language Resources and Evaluation (LREC'02)*, Las Palmas, Canary Islands - Spain. European Language Resources Association (ELRA).
- Matthew Gerber and Joyce Chai. 2010. **Beyond NomBank: A study of implicit arguments for nominal predicates**. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 1583–1592, Uppsala, Sweden. Association for Computational Linguistics.
- Matthew Gerber and Joyce Y. Chai. 2012. **Semantic role labeling of implicit arguments for nominal predicates**. *Computational Linguistics*, 38(4):755–798.
- Luheng He, Kenton Lee, Omer Levy, and Luke Zettlemoyer. 2018. **Jointly predicting predicates and arguments in neural semantic role labeling**. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 364–369, Melbourne, Australia. Association for Computational Linguistics.
- Diederik P. Kingma and Jimmy Ba. 2015. **Adam: A method for stochastic optimization**. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Nikita Kitaev and Dan Klein. 2018. **Constituency parsing with a self-attentive encoder**. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2676–2686, Melbourne, Australia. Association for Computational Linguistics.
- Kevin Knight, Bianca Badarau, Laura Baranescu, Claire Bonial, Madalina Bardocz, Kira Griffitt, Ulf Hermjakob, Daniel Marcu, Martha Palmer, Tim O’Gorman, and Nathan Schneider. 2020. Abstract Meaning Representation (AMR) annotation release 3.0 LDC2020T02. *Linguistic Data Consortium, Philadelphia, PA*.
- Kenton Lee, Luheng He, Mike Lewis, and Luke Zettlemoyer. 2017. **End-to-end neural coreference resolution**. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 188–197, Copenhagen, Denmark. Association for Computational Linguistics.
- Kenton Lee, Luheng He, and Luke Zettlemoyer. 2018. **Higher-order coreference resolution with coarse-to-fine inference**. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 687–692, New Orleans, Louisiana. Association for Computational Linguistics.
- Thien Huu Nguyen, Kyunghyun Cho, and Ralph Grishman. 2016. **Joint event extraction via recurrent**

- neural networks. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 300–309, San Diego, California. Association for Computational Linguistics.
- Tim O’Gorman, Michael Regan, Kira Griffitt, Ulf Her-mjakob, Kevin Knight, and Martha Palmer. 2018. [AMR beyond the sentence: the multi-sentence AMR corpus](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 3693–3702, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Timothy J O’Gorman. 2019. *Bringing Together Computational and Linguistic Models of Implicit Role Interpretation*. PhD dissertation, University of Colorado at Boulder.
- Hiroki Ouchi, Hiroyuki Shindo, and Yuji Matsumoto. 2018. [A span selection model for semantic role labeling](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1630–1642, Brussels, Belgium. Association for Computational Linguistics.
- Martha Palmer, Daniel Gildea, and Paul Kingsbury. 2005. [The proposition bank: An annotated corpus of semantic roles](#). *Computational Linguistics*, 31(1):71–106.
- Ellie Pavlick, Heng Ji, Xiaoman Pan, and Chris Callison-Burch. 2016. [The gun violence database: A new task and data set for NLP](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1018–1024, Austin, Texas. Association for Computational Linguistics.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. [Glove: Global vectors for word representation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. [Deep contextualized word representations](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.
- Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Hwee Tou Ng, Anders Björkelund, Olga Uryupina, Yuchen Zhang, and Zhi Zhong. 2013. [Towards robust linguistic analysis using OntoNotes](#). In *Proceedings of the Seventeenth Conference on Computational Natural Language Learning*, pages 143–152, Sofia, Bulgaria. Association for Computational Linguistics.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Michael Roth and Anette Frank. 2013. [Automatically identifying implicit arguments to improve argument linking and coherence modeling](#). In *Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 1: Proceedings of the Main Conference and the Shared Task: Semantic Textual Similarity*, pages 306–316, Atlanta, Georgia, USA. Association for Computational Linguistics.
- Josef Ruppenhofer, Caroline Sporleder, Roser Morante, Collin Baker, and Martha Palmer. 2010. [SemEval-2010 task 10: Linking events and their participants in discourse](#). In *Proceedings of the 5th International Workshop on Semantic Evaluation*, pages 45–50, Uppsala, Sweden. Association for Computational Linguistics.
- Niko Schenk and Christian Chiarcos. 2016. [Un-supervised learning of prototypical fillers for implicit semantic role labeling](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1473–1479, San Diego, California. Association for Computational Linguistics.
- Carina Silberer and Anette Frank. 2012. [Casting implicit role linking as an anaphora resolution task](#). In **SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 1–10, Montréal, Canada. Association for Computational Linguistics.
- Linfeng Song, Yue Zhang, Zhiguo Wang, and Daniel Gildea. 2018. [N-ary relation extraction using graph-state LSTM](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2226–2235, Brussels, Belgium. Association for Computational Linguistics.
- Beth M. Sundheim. 1992. [Overview of the fourth message understanding evaluation and conference](#). In *FOURTH MESSAGE UNDERSTANDING CONFERENCE (MUC-4), Proceedings of a Conference Held in McLean, Virginia, June 16-18, 1992*.
- Swabha Swayamdipta, Sam Thomson, Kenton Lee, Luke Zettlemoyer, Chris Dyer, and Noah A. Smith. 2018. [Syntactic scaffolds for semantic structures](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3772–3782, Brussels, Belgium. Association for Computational Linguistics.
- Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019a. [BERT rediscovered the classical NLP pipeline](#). In

Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 4593–4601, Florence, Italy. Association for Computational Linguistics.

Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim, Benjamin Van Durme, Sam Bowman, Dipanjan Das, and Ellie Pavlick. 2019b. [What do you learn from context? Probing for sentence structure in contextualized word representations](#). In *International Conference on Learning Representations*.

Ralph Weischedel, Martha Palmer, Mitchell Marcus, Eduard Hovy, Sameer Pradhan, Lance Ramshaw, Nianwen Xue, Ann Taylor, Jeff Kaufman, Michelle Franchini, Mohammed El-Bachouti, Robert Belvin, and Ann Houston. 2013. OntoNotes release 5.0 LDC2013T19. *Linguistic Data Consortium, Philadelphia, PA*.

Travis Wolfe, Mark Dredze, and Benjamin Van Durme. 2015. [Predicate argument alignment using a global coherence model](#). In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 11–20, Denver, Colorado. Association for Computational Linguistics.

A RAMS Data

A.1 Collection

On Reddit, users make submissions containing links to news articles, images, videos, or other kinds of documents, and other users may then vote or comment on the submitted content. We collected news articles matching the following criteria: 1) Posted to the *r/politics* sub-forum between January and October 2016; 2) Resulted in threads with at least 25 comments; and 3) Contained at least one mention of the string “Russia”. The resulting subset of articles tended to describe geopolitical events and relations like the ones in the AIDA ontology. In order to filter out low-quality, fake, or disreputable news articles, we treat the number of comments in the discussion as a signal of information content. Our approach of gathering user-submitted and curated content through Reddit is similar to those used for creating large datasets for language model pre-training (Radford et al., 2019). Documents were split into sentences using NLTK 3.4.3, and sentences were split into tokens using SpaCy 2.1.4.

A.2 Annotation

To assess whether a lexical unit (LU) evoked an event with positive factuality, the vetting task contained an event definition and several candidate sentences, each with a highlighted LU. Annotators were asked to judge how well each highlighted LU, in the context of its sentence, matched the provided event definition. In the same task, they were also asked to assess the factuality of the sentence. Annotation instructions and examples are shown in Figure 9 and Figure 10.

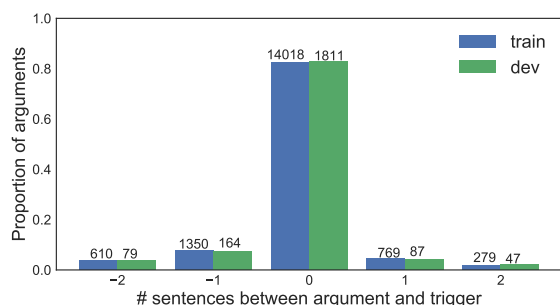


Figure 5: Distances between triggers and arguments in RAMS and proportion of arguments at that distance (counts are shown above each bar). Negative distances indicate that the argument occurs before the trigger.

Each argument selection task contained five

tokenized sentences, a contiguous set of tokens marking the trigger, a definition of the event type, and a list of roles and their associated definitions. For each role, annotators were asked whether a corresponding argument was present in the 5-sentence window, and if so, to highlight the argument span that was closest to the event trigger, as there could be multiple. In cases near the beginning or end of a document, annotators were shown up to two sentences before or after the sentence containing the trigger. Annotators were allowed to highlight any set of (within-sentence) contiguous tokens within the 5-sentence window aside from the trigger tokens. The distribution of distances between triggers and arguments is shown in Figure 5. Annotation instructions and an example are shown in Figure 11 and Figure 12.

A.3 Agreement

We additionally compute the frequency with which annotators agreed a given role was or was not present in the context window. To measure the frequency with which annotators agree whether a given role is present, we treat the majority annotation as the gold standard. Then, we calculated the precision, recall, and F_1 of the annotations. Across the set of redundantly annotated tasks, there were 83 false negatives, 60 false positives, and 892 true positives, giving a precision of 93.7, recall of 91.5, and an F_1 of 92.6.

Threshold	Conjunctive	Disjunctive	Start	End
0	55.3	78.0	59.8	73.5
1	69.9	80.3	74.9	75.3
2	73.9	82.0	78.2	77.8
3	76.4	83.6	80.9	79.1
4	78.8	84.3	82.7	80.4

Table 8: Pairwise span boundary inter-annotator agreement statistics for various span difference thresholds.

We consider a wider range of span difference thresholds, where span difference is calculated by using the absolute difference of the $(start, end)$ token indices from each pair. These are presented in Table 8. In conjunctive agreement, both $|start_1 - start_2|$ and $|end_1 - end_2|$ must be less than the given threshold; therefore, conjunctive agreement at threshold 0 is the percent of pairs that exactly agree (55.3%). Disjunctive agreement is less strict, requiring that either the absolute difference of start offsets or end offsets must be less than the threshold. Start and end agreement is deter-

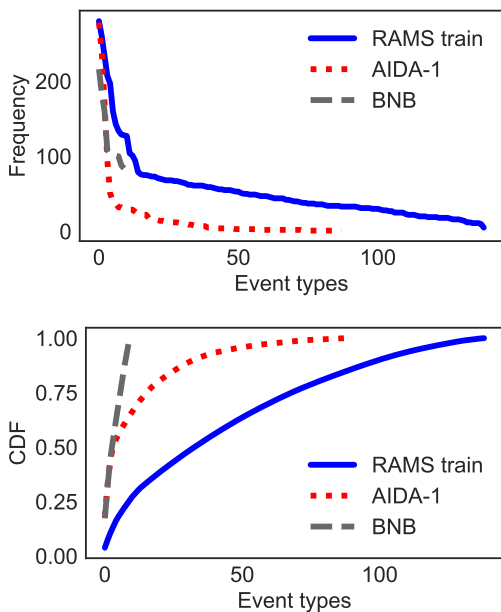


Figure 6: Comparison of frequency (top) and amount of dataset covered (bottom) of event types sorted by decreasing frequency. RAMS has more annotations for a more diverse set of event types than do AIDA Phase 1 and Beyond NomBank.

mined by considering whether the absolute difference of the pair’s start or end offsets (respectively) is within the given threshold.

A.4 Event and Role Type Coverage

Event type and role type coverage are shown in Figure 6 and Figure 7. Figure 6 illustrates that RAMS contains more annotations for a larger set of event types than does AIDA-1. In addition, the distribution of annotations in RAMS is less skewed (more entropic) than in AIDA-1, in that in order to cover a given percentage of the dataset, more event types must be considered in RAMS than in AIDA-1. Figure 7 shows a similar pattern for role type coverage.

Figure 8 shows role coverage per event type, a measure of how much of each event type’s role set is annotated on average. Role coverage per event type is calculated as the average number of filled roles per instance of the event type divided by the number of roles specified for that event type by the ontology. For the RAMS training set, the 25th percentile is 55.6%, the 50th percentile is 61.9%, and the 75th percentile is 68.6% coverage.

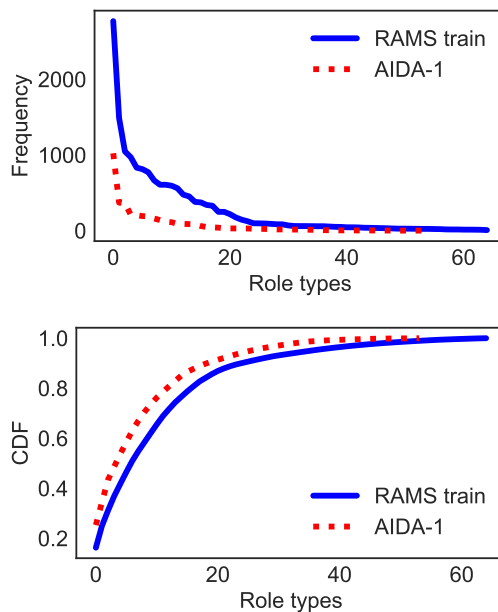


Figure 7: Comparison of frequency (top) and amount of dataset covered (bottom) of roles sorted by decreasing frequency. RAMS has more annotations for a more diverse set of role types than the AIDA Phase 1 data.

B RAMS Hyperparameters

Table 9 lists the numerical hyperparameters shared by all models discussed in this paper. Models may ignore some link score components if they were found to be unhelpful during our sweep of Equation 2 and Equation 3. For our model, we learn a linear combination of the top layers (9, 10, 11, 12) of BERT-base cased, while we use the middle layers (6, 7, 8, 9) for the 6–9 ablation. For ELMo, we use all three layers and encode each sentence separately. We apply a lexical dropout of 0.5 to these embeddings.

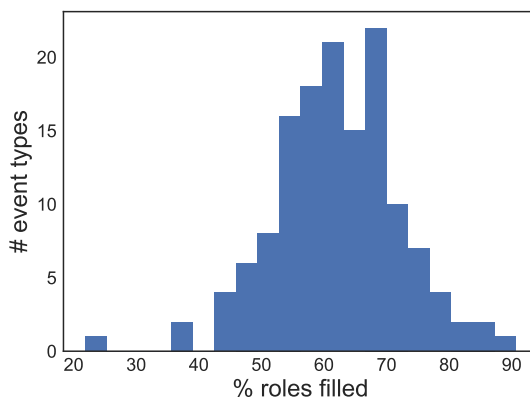


Figure 8: Number of event types for which a given percentage of roles are filled in RAMS train set.

Below, you will be presented with a **description of an event** and several **news article sentences**. Each sentence contains a piece of highlighted text. In short, we are interested in whether the highlighted text "matches" the provided event description.

For each sentence, you will answer up to three questions. First, we establish whether the text contains *any* kind of event:

1. Does the highlighted piece of text refer to an event?

Important:

- Most verbs refer to events, but nouns can refer to events too! For example, "destruction" or "robbery".
- It's fine if just one or two words are highlighted, as long as those words refer to an event.
 - For example, "the **crash** was relatively minor".
- If the sentence contains an event but it is not highlighted, answer "no" to this question.

Then, if the highlighted text *does* refer to an event, we ask:

2. How well does the highlighted event match the provided event description?

Important:

- The event can match even if the event **did not occur!** Imagine we have the event description: "an object is involved in a collision".
 - In "the **crash** was narrowly avoided", the event matches the description even though the crash didn't occur!
 - We ask a separate question to determine whether the event happened or not.
- You may have to distinguish between several senses of a word.
 - In "the stock market **crash** of 1987", the "crash" is not a "collision" event.
- If the sentence contains a matching event, but it is not highlighted, answer "no" to this question.

Finally, we ask:

3. How strongly does the sentence imply that the highlighted event actually happened?

Important:

- Pay attention to words like "definitely", "undoubtedly", "unlikely", "maybe", etc.
- In the absence of any obvious clues, you should lean towards "yes".
- Only use the information in the text! To the extent possible, don't use any background knowledge.
 - For example, in "Ben **crashed** into a unicorn", the text implies the event occurred, even though we know it couldn't have happened.
- If the sentence implies the event **happened** or is presently **happening**, lean towards "yes". If the event will happen in the future, lean towards "no".

Figure 9: Annotation instructions for determining whether a lexical unit (in context) evokes an event type.

Event Name: Collision
Event Description: An object collides violently with an obstacle or another object

Sentence: The car crash was **narrowly** avoided.

Does the highlighted piece of text refer to an event?
 Yes
 No

How well does the highlighted event match the provided event description?
 Very Weakly Very Strongly

How strongly does the sentence imply that the highlighted event actually happened?
 Definitely didn't happen Definitely Happened

Sentence: The car **crash** was narrowly avoided.

Does the highlighted piece of text refer to an event?
 Yes
 No

How well does the highlighted event match the provided event description?
 Very Weakly Very Strongly

How strongly does the sentence imply that the highlighted event actually happened?
 Definitely didn't happen Definitely Happened

Figure 10: Annotation interface for determining whether a lexical unit (in context) evokes an event type.

This task is helping us build systems to better understand documents, thank you for your help.

You will be presented with up to five sentences that have been extracted from a larger document. In one of those sentences an **event trigger** has been highlighted: some word or phrase that refers to an event. We need your help to highlight the **arguments** of that event, somewhere in the provided text. "Arguments" are just the entities that play a certain role in an event; for example, in "John **ate** the sandwich" we could say that "John" is the "Consumer" and "the sandwich" is the "ConsumedObject".

The arguments are defined in two ways:

1. Each argument has a **role name**, which sometimes can be informative such as *Transporter*, and other times may not be very helpful, such as *GPE*
2. We provide a **gloss** for the event, where we have written a sentence based on the event definition, with the arguments initially filled in with words like *someone* or *someplace*.

Additionally, you will be provided with the *name* and *definition* of the event. If you find it hard to interpret a particular argument, we suggest re-reading the event name and definition and seeing how that argument fits in with the event.

As you highlight argument spans, the gloss will update based on your selection. The gloss is meant to be a helpful tool: it can happen that as you select arguments it will not always result in a perfectly grammatical gloss. That situation is OK, just do your best.

We would like you to select arguments even if the event did not happen! For example, in "John might **crash** his car tomorrow", no "Crash" event occurred yet, but we'd still like you to highlight arguments.

Finally, since there may be several text spans that refer to the same entity, we ask that you select the one that is **closest** to the highlighted trigger word. For example, if we are trying to find the "Consumer" in the sentence: "John and Mary like candy. They **ate** some today", we ask that you select "they", rather than "John and Mary", since "they" is closer to the trigger.

Sometimes not all of the arguments will be present in the text; in that case, we ask you explicitly check the box that a given argument can not be found. The following is an example of a completed task where the "Place" argument is not mentioned:

Figure 11: Annotation instructions for selecting arguments for an event.

Gloss

Cohn died non-violently at someplace

Context

" Because he saw these mob guys as pathways to money , and Donald is all about money . "

From a \$ 400 million tax giveaway on his first big project , to getting a casino license , to collecting fees for putting his name on everything from bottled water and buildings to neckties and steaks , Trump 's life has been dedicated to the next big score .

Through Cohn , Trump made choices that - gratuitously , it appears - resulted in his first known business dealings with mob - controlled companies and unions , a pattern that continued long after Cohn died .

What Trump has to say about the reasons for his long , close and wide - ranging dealings with organized crime figures , with the role of mobsters in cheating Trump Tower workers , his dealings with Felix Sater and Trump 's seeming leniency for Weichselbaum , are questions that voters deserve full answers about before casting their ballots .

Event type: Life.Die.NonviolentDeath

Event definition: The non-violent end of the life of a person entity

Arguments

Victim	per	Cohn	<input type="checkbox"/> Argument not present
Place	loc, gpe, fac		<input checked="" type="checkbox"/> Argument not present

Figure 12: Annotation interface for selecting arguments for an event.

Hyperparameter		Value
Embeddings	role size	50
	feature (ϕ_l) size	20
	size	200
LSTM	layers	3
	dropout	0.4
argument (F_A)	size	150
	layers	2
event-role ($F_{E,R}$)	size	150
	layers	2
$F_{\tilde{a}}$ (Eqn. 1)	layers	2
	size	150
arg-role ($F_{A,R}$)	layers	2
	size	150
F_l	layers	2
	size	150
distance FFNN	layers	2
	# buckets	10
Pruning	k	10
Memory Limits	training doc size	1000
	batch size	1
	learning rate	0.001
Training	decay	$\frac{0.999}{100 \text{ steps}}$
	patience	10

Table 9: Hyperparameters of the model trained on RAMS. Sizes of learned weights that are omitted from the table can be determined from these hyperparameters. As the argument spans are given to the model in our experiments, we skip the first pass of pruning. We do not clip gradients.

In our best model, we use learned bucketed distance embeddings (Lee et al., 2017). These embeddings are scored as part of ϕ_c in computing $s_c(e, a)$ in Equation 2 and are also scored as a part of ϕ_l in s_l (Equation 3). Since span boundaries are given in our primary experiments, we do not include a score s_A or s_E in s_c . Our best model uses both $s_{A,R}$ and $s_l(a, \tilde{a}_{e,r})$ in Equation 2. These features were chosen as the result of a sweep over possible features, with other ablations reported in Table 3.

We adopt the span embedding approach by Lee et al. (2017), which uses character convolutions (50 8-dimensional filters of sizes 3, 4, and 5) and 300-dimensional GloVe embeddings. The default dropout applied to all connections is 0.2. We optimize using Adam (Kingma and Ba, 2015) with patience-based early stopping, resulting in the best checkpoint after 19 epochs (9 hours on an NVIDIA 1080Ti), using F_1 as the evaluation metric.

Hyperparameters for the condition with distractor candidate arguments are the same as those in Table 9. For the condition with no given argument spans, we consider all intrasentential spans

up to 5 tokens in length. We include the score of each candidate argument span when pruning to encourage the model to keep correct spans. We modify hyperparameters in Table 9 to prune less aggressively, setting $k = 100$ and $\lambda_A = 1.0$ (defined in §4.1).

C Full Role Confusion and Similarity Matrices

Figure 13 shows the similarity between all 65 role embeddings, while Figure 14 visualizes all the errors made by the model on the development set. These are expansions of the per-role results from §5.1.

Since argument linking is not a one-to-one labeling problem, we need to perform a modified procedure for visualizing a confusion matrix. For example, an argument span may take on multiple roles for the same event. To compute the errors, we first align the correct prediction(s) and subsequently compute the errors for the remaining gold and predicted label(s). For example, if the correct set of roles is {destination, origin} and the model predicts {origin, place}, then we only mark place as an error for destination.

D AIDA Phase 1

D.1 Data Processing

We filter and process the AIDA-1 Practice and Eval data in the following way. Because annotations are available for only a subset of the documents in AIDA-1, we consider only the documents that have textual event triggers. We then take from this set only the English documents, which, due to noisy language ID in the original annotations, were selected by manual inspection of the first 5 sentences of each document by one of the authors of this work.

In addition, the argument spans in each example are only those that participate in events. In other words, arguments of relations (that are not also arguments of events) are not included. Additionally, a document may contain multiple events, unlike in RAMS.

The training and development set come from AIDA-1 Practice, and the test set comes from AIDA-1 Eval. As the AIDA-1 Eval documents are about different topics than the Practice documents are, we emulate the mismatch in topic distribution by using a development set that is about a different

Strategy	Dev. F_1	P	R	F_1
No pre-training	25.0	36.6	12.9	19.1
No pre-training ^{TCD}	27.1	53.5	11.0	18.2
RAMS pre-training	34.1	43.9	16.9	24.4
RAMS pre-training ^{TCD}	34.2	62.5	15.4	24.8

Table 10: P(recision), R(ecall), and F_1 on AIDA-1 English development and test data. TCD designates the use of ontology-aware type-constrained decoding.

topic than the training set is. We use Practice topics R103 and R107 for training and R105 for development because R105 is the smallest of the three practice topics both by number of documents and by number of annotations. The test set consists of all 3 topics (E101, E102, E103) from the (unsequenced) Eval set. After the filtering process described above, we obtain a training set of 46 documents, a development set of 17 documents, and a test set of 69 documents. There are 389 events in the training set, and the training documents have an average length of 50 sentences.

D.2 Hyperparameters

We use the same hyperparameters as the best model for RAMS, shown in Table 9.

D.3 Pre-training on RAMS

Both the models with and without pre-training on RAMS were trained on AIDA-1 for 100 epochs with an early-stopping patience of 50 epochs using the same hyperparameters as the best RAMS model. All parameters were updated during fine-tuning (none were frozen). The vocabulary of the pre-trained model was not expanded when trained on AIDA-1.

The models’ lower performance on AIDA-1 than on RAMS may be in part explained by the presence of distractors in AIDA-1. Moving from RAMS (one trigger per example) to AIDA-1 (many triggers per example) introduces distractor “negative” links: an argument for one event might not participate in a different event in the same document. When given gold argument spans, a model learns from RAMS that every argument gets linked to the trigger, but there are many negative links in the AIDA-1 data, which the model must learn to not predict.

Full results are given in Table 10. Type-constrained decoding does not improve performance on AIDA-1 as much as it did in Table 3, possibly because the AIDA-1 data often does not

adhere to the multiplicity constraints of the ontology. For example, many attack events have more than one annotated attacker or target. Under TCD, correct predictions made in excess of what the ontology allows are deleted, hurting recall.

Interestingly, type-constrained decoding *hurts* performance on AIDA-1 Eval when there is no pre-training. As discussed in §5, type-constrained decoding tends to improve precision and lower recall. Despite the same behavior here, F_1 is nonetheless decreased.

We see similar behavior in this experiment to the RAMS experiment involving distractor candidate arguments: low performance which is reduced further when using TCD.

E BNB Data Processing and Hyperparameters

E.1 Data Processing

We use the data from Gerber and Chai (2012).¹⁷ We processed the data in the following way. The annotations were first aligned to text in the Penn Treebank. Because our model assumes that arguments are contiguous spans, we then manually merged all “split” arguments, which with one exception were already contiguous spans of text. For the one split argument that was not a contiguous span, we replaced it with its maximal span.¹⁸ We then removed special parsing tokens such as “trace” terminals from the text and realigned the spans. While BNB gives full credit as long as one argument in each argument “cluster” is found, our training objective assumes one argument per role. We therefore automatically reduced each argument cluster to a singleton set containing the argument closest to the trigger. This reformulation of the problem limits our ability to compare to prior work.

Once all the data had been processed, we created training, development, and test splits. To avoid leaking information across splits, we bucketed examples by document and randomly assigned documents to the splits so that the splits contained instances in the proportions 80% (train), 10% (dev), and 10% (test).

¹⁷http://lair.cse.msu.edu/projects/implicit_argument_annotations.zip. Information about the data and its fields is available at http://lair.cse.msu.edu/projects/implicit_annotations.html.

¹⁸The instance is a quote broken by speaker attribution, where the split argument consists of the two halves of the quote. This example appears in our training set.

Hyperparameter		Value
Embeddings	role size	50
	feature (ϕ_l) size	20
	size	200
LSTM	layers	3
	dropout	0.4
argument (F_A)	size	150
	layers	2
event-role ($F_{E,R}$)	size	150
	layers	2
$F_{\bar{a}}$ (Eqn. 1)	layers	2
	size	150
F_l	layers	2
	size	150
positional FFNN	layers	2
	# buckets	10
	λ_A	0.8
Pruning	k	45
	training doc size	600
Memory Limits	span width	15
	batch size	1
	learning rate	0.0005
Training	decay	$\frac{0.999}{200 \text{ steps}}$
	patience	20
	gradient clipping	10.0

Table 11: Hyperparameters of the model trained on GVDB.

E.2 Hyperparameters

We use the same hyperparameters as the best model for RAMS, shown in Table 9.

F GVDB Hyperparameters and Additional Results

The entire GVDB corpus consists of 7,366 articles. We exclude articles that do not have a reliable publication date or lack annotated spans for the roles we are interested in. Additionally, a buffer of 100 articles spanning roughly one week between the dev and test set is discarded, limiting the possibility of events occurring in both the development and test sets. We also filter out spans whose start and end boundaries are in different sentences, as these are unlikely to be well-formed argument spans. For evaluation, a slot’s value is marked as correct under the *strict* setting if any of the predictions for that slot match the string of the correct answer exactly, while an *approximate* match is awarded if either a prediction contains the correct answer or if the correct answer contains the predicted string. The approximate setting is necessary due to inconsistent annotations (e.g., omitting first or last names).

We experiment with the feature-based version of BERT-base and with ELMo as our contextualized encoder. Table 11 lists the numerical hyper-

parameters for this model. Since there is only one event per document and no explicit trigger, e is represented by a span embedding of the full document. We use the top four layers (9–12) of BERT-base cased (all three layers for ELMo) with a lexical dropout of 0.5. Everywhere else, we apply a dropout of 0.4. We train with the Adam optimizer (Kingma and Ba, 2015) and use patience-based early stopping. Our best checkpoint was after 8 epochs (roughly 9 hours on a single NVIDIA 1080Ti). Even though the official evaluation is string based, we used a span-based micro F_1 metric for early stopping.

For this model, ϕ_l corresponds to a learned (bucketed) positional embedding of the argument span (i.e., distance from the start of the document). In computing the coarse score, we omit ϕ_c . When computing Equation 2, we omit $s_{A,R}$ but keep all other terms in Equation 2. We adopt the character convolution of 50 8-dimensional filters of window sizes 3, 4, and 5 (Lee et al., 2017).

With the same hyperparameters and feature choices, we perform an identical evaluation using ELMo instead of BERT. As the original documents are not tokenized, we use SpaCy 2.1.4 for finding sentence boundaries and tokenization. The complete list of annotated fields are VICTIM (name, age, race), SHOOTER (name, age, race), LOCATION (specific location¹⁹ or city), TIME (time of day or clock time) and WEAPON (weapon type, number of shots fired). While Pavlick et al. (2016) only make predictions for VICTIM.NAME, SHOOTER.NAME, LOCATION.(CITY|LOCATION), TIME.(TIME|CLOCK), and WEAPON.WEAPON, we perform predictions over all annotated span-based fields. The full results for both BERT and ELMo are reported in Table 12 and Table 13, respectively. BERT generally improves over ELMo across the board, but not by a sizeable margin. Despite the inability to directly compare, we nonetheless present a stronger and more comprehensive baseline for future work with GVDB.

¹⁹For example, a park or a laundromat.

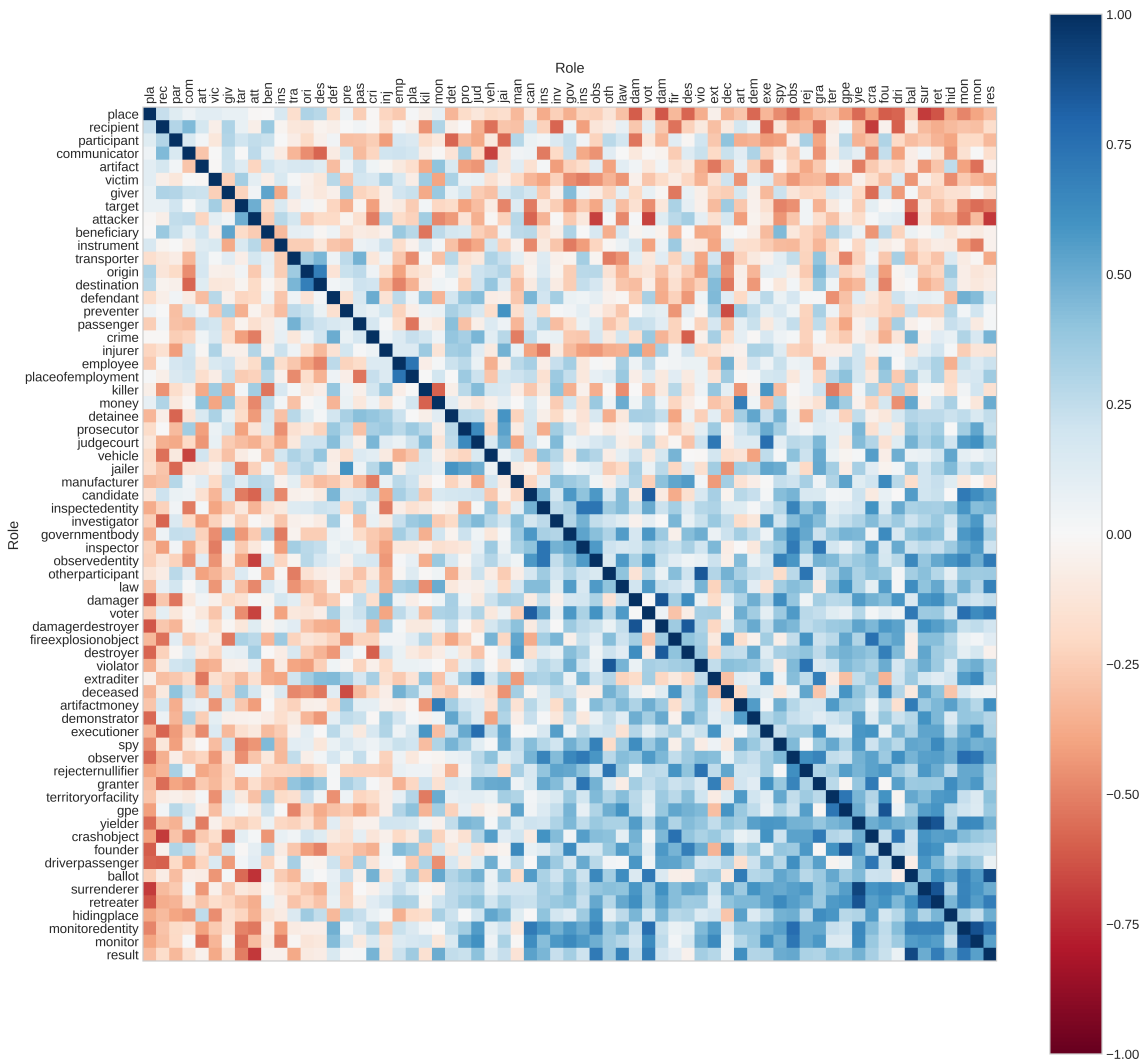


Figure 13: Full version of Figure 4, showing cosine similarity between role embeddings. Best viewed in color.

Field		Strict						Partial					
		Baseline			Us			Baseline			Us		
		P	R	F1	P	R	F1	P	R	F1	P	R	F1
VICTIM	Name	10.2	8.5	9.3	61.2	63.3	62.2	59.5	49.6	54.1	68.4	70.9	69.6
	Age	-	-	-	19.4	24.2	21.5	-	-	-	67.3	84.1	74.8
	Race	-	-	-	75.5	74.1	74.8	-	-	-	75.5	74.1	74.8
SHOOTER	Name	5.8	3.9	4.7	55.3	51.1	53.1	30.2	20.1	24.1	60.2	55.6	57.8
	Age	-	-	-	34.1	32.6	33.3	-	-	-	69.0	65.9	67.4
	Race	-	-	-	72.7	55.2	62.7	-	-	-	81.8	62.1	70.6
LOCATION	City				67.4	66.2	66.8				72.2	70.9	71.5
	Location	19.9	8.8	12.2	36.1	33.8	34.9	30.8	13.6	18.9	65.4	61.2	63.3
TIME	Time				57.2	69.7	62.9				63.2	76.9	69.4
	Clock	69.3	66.9	68.1	44.0	47.6	45.7	70.5	68.1	69.3	84.0	90.8	87.2
WEAPON	Weapon	2.1	0.7	1.1	33.3	31.7	32.5	36.8	11.8	17.9	50.9	48.3	49.6
	Num Shots	-	-	-	40.6	11.2	17.6	-	-	-	62.5	17.2	27.0

Table 12: P(recision), R(ecall), and F_1 on event-based slot filling (GVDB) using BERT as the document encoder. Due to the different data splits and evaluation conditions, the results are not directly comparable to the baseline (Pavlick et al., 2016), which is provided only for reference. Fields that were aggregated in the baseline are predicted separately in our model. ‘-’ indicates result is not reported in the baseline.

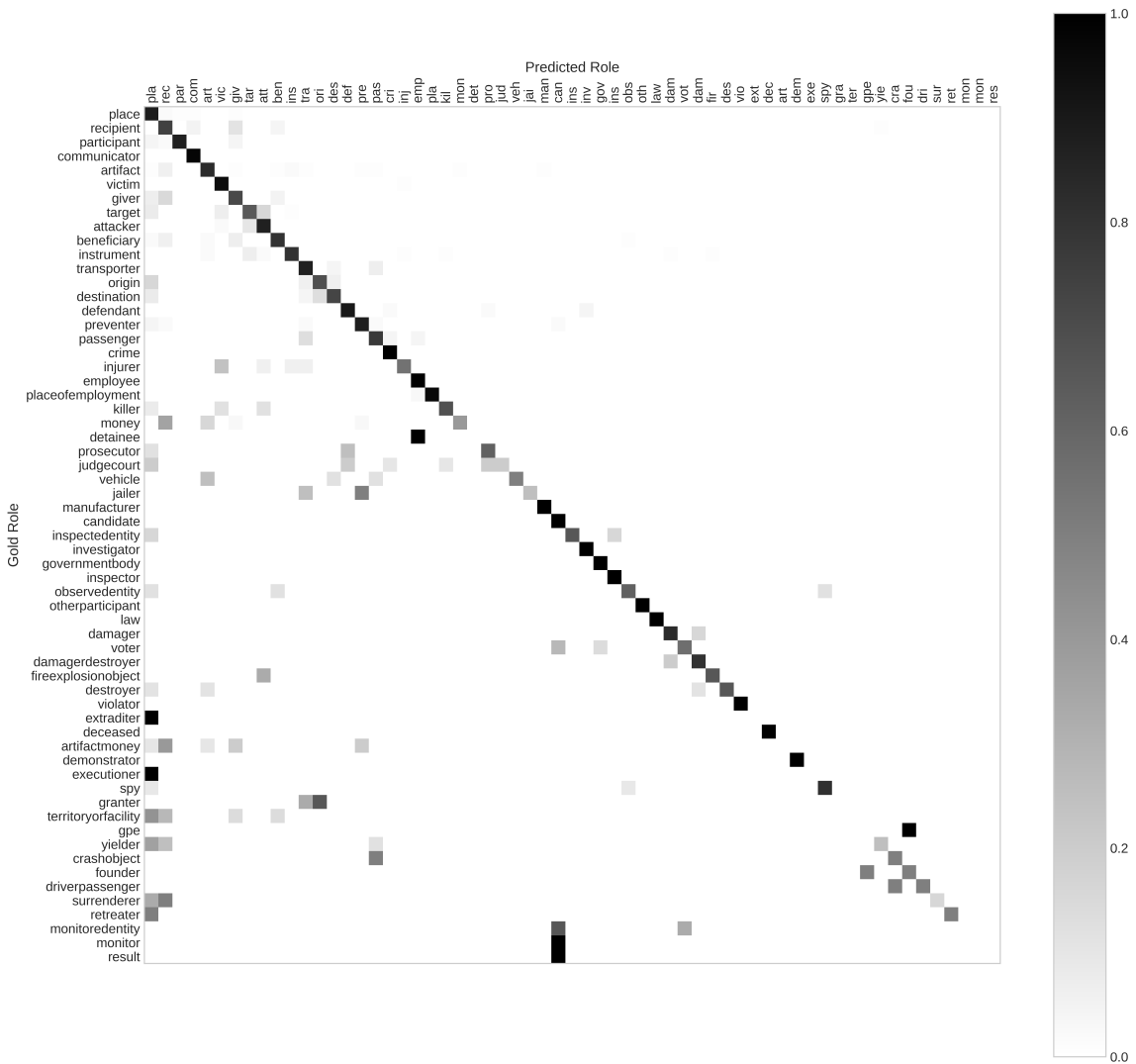


Figure 14: Full version of Figure 4, showing row-normalized confusion between roles. Note that roles not predicted at all would result in empty rows and so are omitted from the table.

Field		Strict						Partial					
		Baseline			Us			Baseline			Us		
		P	R	F1	P	R	F1	P	R	F1	P	R	F1
VICTIM	Name	10.2	8.5	9.3	56.2	56.8	56.5	59.5	49.6	54.1	62.7	63.3	63.0
	Age	-	-	-	29.5	33.9	31.6	-	-	-	64.4	74.0	68.9
	Race	-	-	-	73.2	75.9	74.5	-	-	-	75.0	77.8	76.4
SHOOTER	Name	5.8	3.9	4.7	53.7	60.2	56.7	30.2	20.1	24.1	56.4	63.2	59.6
	Age	-	-	-	27.3	31.8	29.4	-	-	-	53.2	62.1	57.3
	Race	-	-	-	55.9	65.5	60.3	-	-	-	58.8	69.0	63.5
LOCATION	City				59.1	61.1	60.1				64.1	66.2	65.1
	Location	19.9	8.8	12.2	36.6	34.7	35.6	30.8	13.6	18.9	59.1	56.0	57.5
TIME	Time				57.7	64.7	61.0				64.5	72.4	68.2
	Clock	69.3	66.9	68.1	44.6	45.8	45.2	70.5	68.1	69.3	83.5	85.6	84.5
WEAPON	Weapon	2.1	0.7	1.1	32.7	26.7	29.4	36.8	11.8	17.9	44.9	36.7	40.4
	Num Shots	-	-	-	23.3	18.1	20.4	-	-	-	42.2	32.8	36.9

Table 13: P(recision), R(ecall), and F_1 on event-based slot filling (GVDB) using ELMo at the sentence level. On average, the performance is outperformed by BERT.