

Decomposing and Recomposing Event Structure

William Gantt

University of Rochester, USA
wgantt@cs.rochester.edu

Lelia Glass

Georgia Institute of Technology, USA
lelia.glass@modlangs.gatech.edu

Aaron Steven White

University of Rochester, USA
aaron.white@rochester.edu

Abstract

We present an event structure classification empirically derived from inferential properties annotated on sentence- and document-level Universal Decompositional Semantics (UDS) graphs. We induce this classification jointly with semantic role, entity, and event-event relation classifications using a document-level generative model structured by these graphs. To support this induction, we augment existing annotations found in the UDS1.0 dataset, which covers the entirety of the English Web Treebank, with an array of inferential properties capturing fine-grained aspects of the temporal and aspectual structure of events. The resulting dataset (available at decomp.io) is the largest annotation of event structure and (partial) event coreference to date.

1 Introduction

Natural language provides myriad ways of communicating about complex events. For instance, one and the same event can be described at a coarse grain, using a single clause (1), or at a finer grain, using an entire document (2).

- (1) The contractors built the house.
- (2) They started by laying the house’s foundation. They then framed the house before installing the plumbing. After that [...]

Further, descriptions of the same event at different granularities can be interleaved within the same document—for example, (2) might well directly follow (1) as an elaboration on the house-building process.

Consequently, extracting knowledge about complex events from text involves determining the structure of the events being referred to: what their parts are, how those parts are laid out in time,

who participates in them and how, and so forth. Determining this structure requires an event classification whose elements are associated with event structure representations. A number of such classifications and annotated corpora exist: FrameNet (Baker et al., 1998), VerbNet (Kipper Schuler, 2005), PropBank (Palmer et al., 2005), Abstract Meaning Representation (Banarescu et al., 2013), and Universal Conceptual Cognitive Annotation (Abend and Rappoport, 2013), among others.

Similar in spirit to this prior work, but different in method, our work aims to develop an *empirically derived* event structure classification. Where prior work takes a top-down approach—hand-engineering an event classification before deploying it for annotation—we take a bottom-up approach—*decomposing* event structure into a wide variety of theoretically informed, cross-cutting semantic properties, annotating for those properties, then *recomposing* an event classification from them by induction. The properties on which our categories rest target (i) the substructure of an event (e.g., that the building described in (1) consists of a sequence of subevents resulting in the creation of some artifact); (ii) the superstructure in which an event takes part (e.g., that laying a house’s foundation is part of building a house, alongside framing the house, installing the plumbing, etc.); (iii) the relationship between an event and its participants (e.g., that the contractors in (1) build the house collectively through their joint efforts); and (iv) properties of the event’s participants (e.g., that the contractors in (1) are animate while the house is not).

To derive our event structure classification, we extend the Universal Decompositional Semantics dataset (UDS; White et al., 2016, 2020). UDS annotates for a subset of key event structure properties,

but a range of key properties remain to be captured. After motivating the need for these additional properties (§2), we develop annotation protocols for them (§3). We validate our protocols (§4) and use them to collect annotations for the entire Universal Dependencies (Nivre et al., 2016) English Web Treebank (§5; Bies et al., 2012), resulting in the UDS-EventStructure dataset (UDS-E). To derive an event structure classification from UDS-E and existing UDS annotations, we develop a document-level generative model that jointly induces event, entity, semantic role, and event-event relation types (§6). Finally, we compare these types to those found in existing event structure classifications (§7). We make UDS-E and our code available at decomp.io.

2 Background

Contemporary theoretical treatments of event structure tend to take as their starting point Vendler’s (1957) seminal four-way classification. We briefly discuss this classification and elaborations thereon before turning to other event structure classifications developed for annotating corpora.¹ We then contrast these with the fully decompositional approach we take in this paper.

Theoretical Approaches Vendler categorizes event descriptions into four classes: *statives* (3), *activities* (4), *achievements* (5), and *accomplishments* (6). As theoretical constructs, these classes are used to explain both the distributional characteristics of event descriptions as well as inferences about how an event progresses over time.

- (3) Jo was in the park.

$$\textit{stative} = [+DUR, -DYN, -TEL]$$

- (4) Jo ran around in the park.

$$\textit{activity} = [+DUR, +DYN, -TEL]$$

- (5) Jo arrived at the park.

$$\textit{achievement} = [-DUR, +DYN, +TEL]$$

- (6) Jo ran to the park.

$$\textit{accomplishment} = [+DUR, +DYN, +TEL]$$

Work building on Vendler’s discovered that these classes can be decomposed into the now well-accepted component properties in (7)–(9) (Kenny, 1963; Lakoff, 1965; Verkuyl, 1972; Bennett and Partee, 1978; Mourelatos, 1978; Dowty, 1979).

¹The theoretical literature on event structure is truly vast. See Truswell (2019) for a collection of overview articles.

- (7) $DUR(ATIVITY)$: whether the event happens at an instant or extends over time
 (8) $DYN(AMICITY)$: whether the event involves change, broadly construed
 (9) $TEL(ICITY)$: whether the event culminates in a participant changing state or location, being created or destroyed, and so forth.

Later work further expanded these properties and, therefore, the possible classes. Expanding on DYN , Taylor (1977) suggests a distinction between dynamic predicates that refer to events with dynamic subevents (e.g., the individual strides in a running) and ones that do not (e.g., the gliding in (10)) (see also Bach, 1986; Smith, 2003).

- (10) The pelican glided through the air.

Dynamic events with dynamic subevents can be further distinguished based on whether the subevents are similar (e.g., the strides in a running) or dissimilar (e.g., the subevents in a house-building) (Piñón, 1995). In the case where the subevents are similar and a participant itself has subparts (e.g., when the participant is a group), there may be a bijection from participant subparts to subevents. In (11), there is a smiling for each child that makes up the composite smiling—*smile* is *distributive*. In (12), the meeting presumably has some structure, but there is no bijection from members to subevents—*meet* is *collective* (see Champollion, 2010, for a review).

- (11) {The children, Jo and Bo} smiled.

- (12) {The committee, Jo and Bo} met.

Expanding on TEL , Dowty (1991) argues for a distinction among telics in which the culmination comes about incrementally (13) or abruptly (14) (see also Tenny, 1987; Krifka, 1989, 1992, 1998; Levin and Hovav, 1991; Rappaport Hovav and Levin, 1998, 2001; Croft, 2012).

- (13) The gardener mowed the lawn.

- (14) The climber summited at 5pm.

This notion of incrementality is intimately tied up with the notion of $DUR(ATIVITY)$. For instance, Moens and Steedman (1988) point out that certain event structures can be systematically transformed into others—for example, whereas (14) describes the summiting as something that happens at an instant (and is thus abrupt), (15) describes it as a process that culminates in having reached the top of the mountain (see also Pustejovsky, 1995).

(15) The climber was summiting.

Such cases of *aspectual coercion* highlight the importance of grammatical factors in determining the structure of an event. More general contextual factors are also at play when determining event structure: *I ran* can describe a telic event (e.g., when it is known that I run the same distance or to the same place every day) or an atelic event (e.g., when the destination and/or distance is irrelevant in context) (Dowty, 1979; Olsen, 1997). This context-sensitivity strongly suggests that annotating event structure is not simply a matter of building a type-level lexical resource and projecting its labels onto text: Actual text must be annotated.

Resources Early, broad-coverage lexical resources, such as the Lexical Conceptual Structure lexicon (LCS; Dorr, 1993), attempt to directly encode an elaboration of the core Vendler classes in terms of a hand-engineered graph representation proposed by Jackendoff (1990). VerbNet (Kipper Schuler, 2005) further elaborates on LCS by building on the fine-grained syntax-based classification of Levin (1993) and links her classes to LCS-like representations. More recent versions of VerbNet (v3.3+; Brown et al., 2018) update these representations to ones based on the Dynamic Event Model (Pustejovsky, 2013).

COLLIE-V, which expands the TRIPS lexicon and ontology (Ferguson and Allen, 1998, et seq), takes a similar tack of producing hand-engineered event structures, combining this hand-engineering with a procedure for bootstrapping event structures (Allen et al., 2020). FrameNet also contains hand-engineered event structures, though they are significantly more fine-grained than those found in LCS or VerbNet (Baker et al., 1998).

VerbNet, COLLIE-V, and FrameNet are not directly annotated on text, though annotations for at least VerbNet and FrameNet can be obtained by using SemLink to project FrameNet and VerbNet annotations onto PropBank annotations (Palmer et al., 2005). PropBank frames have been enriched in a variety of other ways. One such enrichment can be found in Abstract Meaning Representation (AMR; Banarescu et al., 2013; Donatelli et al., 2018). Another can be found in Richer Event Descriptions (RED; O’Gorman et al., 2016), which annotates events and entities for factuality (whether an event actually happened)

and genericity (whether an event/entity is a particular or generic) as well as annotating for causal, temporal, sub-event, and co-reference relations between events (see also Chklovski and Pantel, 2004; Hovy et al., 2013; Cybulska and Vossen, 2014).

Additional less fine-grained event classifications exist in TimeBank (Pustejovsky et al., 2006), Universal Conceptual Cognitive Annotation (UCCA; Abend and Rappoport, 2013), and the Situation Entities dataset (SitEnt; Friedrich and Palmer, 2014b; Friedrich et al., 2016). Of these, the closest to capturing the standard Vendler classification and decompositions thereof is SitEnt. The original version of SitEnt annotates only for a state-event distinction (alongside related, non-event structural distinctions), but later elaborations further annotate for telicity (Friedrich and Gateva, 2017). Because of this close alignment to the standard Vendler classes, we use SitEnt annotations as part of validating our own annotation protocol in §3.

Universal Decompositional Semantics In contrast to the hand-engineered event structure classifications discussed above, our aim is to derive event structure representations directly from semantic annotations. To do this, we extend the existing annotations in the Universal Decompositional Semantics dataset (UDS; White et al., 2016, 2020) with key annotations for the event structural distinctions discussed above. Our aim is not necessarily to reconstruct any previous classification, though we do find in §6 that our event type classification approximates Vendler’s to some extent.

UDS is a semantic annotation framework and dataset based on the principles that (i) the semantics of words or phrases can be *decomposed* into sets of simpler semantic properties and (ii) these properties can be annotated by asking straightforward questions intelligible to non-experts. UDS comprises two layers of annotations on top of the Universal Dependencies (UD) syntactic graphs in the English Web Treebank (EWT): (i) predicate-argument graphs with mappings into the syntactic graphs, derived using the PredPatt tool (White et al., 2016; Zhang et al., 2017); and (ii) crowd-sourced annotations for properties of events (on the *predicate nodes* of the predicate-argument graph), entities (on the

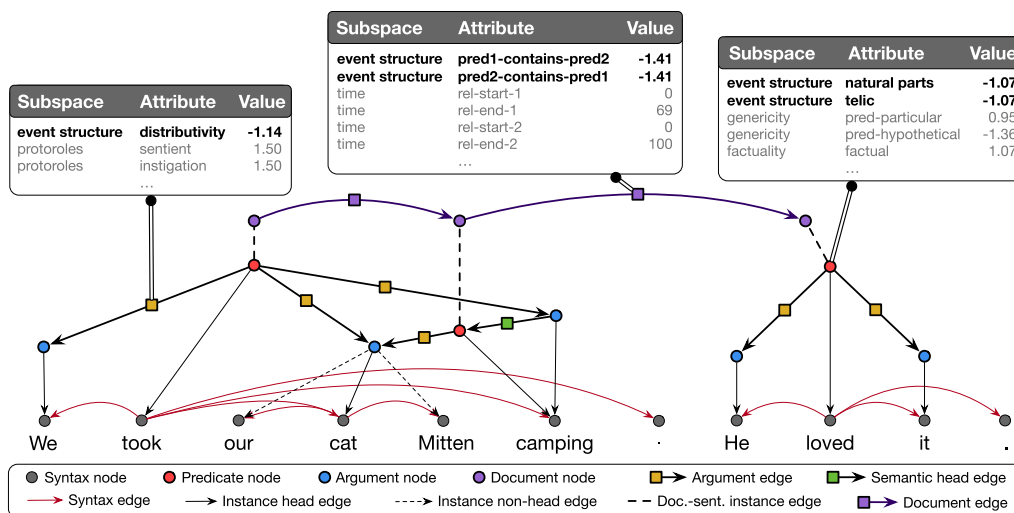


Figure 1: Example UDS semantics and syntax graphs with select properties (see Footnote 2 on the property values). Bolded properties are ones we collect in this paper, and our new *document-level* graph is also shown in purple.

argument nodes), and their relationship (on the *predicate-argument edges*).

The UDS properties are organized into three *predicate subspaces* with five properties in total:

- **FACTUALITY** (Rudinger et al., 2018)
factual: did the event happen?
- **GENERICITY** (Govindarajan et al., 2019)
kind: is the event generic?
hypothetical: is the event hypothetical?
dynamic: is the event dynamic or stative?
- **TIME** (Vashishtha et al., 2019)
duration: how long did/will the event last?

Two *argument subspaces* with four properties:

- **GENERICITY** (Govindarajan et al., 2019)
particular: is the entity a particular?
kind: is the entity a kind?
abstract: is the entity abstract or concrete?
- **WORDSENSE** (White et al., 2016)
Which coarse entity types (WordNet super-sense) does the entity have?

And one *predicate-argument subspace* with 16 properties (see White et al., 2016, for full list):

- **PROTOROLES** (Reisinger et al., 2015)
instigation: did participant cause event?
change of state: did participant change state during or as a consequence of event?
change of location: did participant change location during event?

existed {before, during, after} did participant exist {before, during, after} the event?

Figure 1 shows an example UDS1.0 graph (White et al., 2020) augmented with (i) a subset of the properties we add in bold (see §3); and (ii) *document-level* edges in purple (see §6).²

The UDS annotations and associated toolkit have supported research in a variety of areas, including syntactic and semantic parsing (Stengel-Eskin et al., 2020, 2021), semantic role labeling (Teichert et al., 2017) and induction (White et al., 2017), event factuality prediction (Rudinger et al., 2018), temporal relation extraction (Vashishtha et al., 2019), among others. For our purposes, the annotations do cover some event structural distinctions—for example, *dynamicity*, specific cases of *telicity* (in the form of *change of state*, *change of location*, and *existed {before, during, after}*), and *durativity*. In this sense, UDS provides an alternative, decompositional event representation that distinguishes it from more traditional categorical ones like SitEnt. However, the existing annotations fail to capture a number of the core distinctions above—a lacuna this work aims to fill.

²Following White et al. (2020), the property values in Figure 1 are derived from raw annotations using mixed effects models (MEMs; Gelman and Hill, 2014), which enable one to adjust for differences in how annotators approach a particular annotation task (see also Gantt et al., 2020). In §6, we similarly use MEMs in our event structure induction model, allowing us to work directly with the raw annotations.

3 Annotation Protocol

We annotate for the core event structural distinction not currently covered by UDS, breaking our annotation into three subprotocols. For all questions, annotators report confidence in their response to each question on a scale from 1 (*not at all confident*) to 5 (*totally confident*).³

Event-subevent Annotators are presented with a sentence containing a single highlighted predicate followed by four questions about the internal structure of the event it describes. Q1 asks whether the event described by the highlighted predicate has natural subparts. Q2 asks whether the event has a natural endpoint.

The final questions depend on the response to Q1. If an annotator responds that the highlighted predicate refers to an event that *has* natural parts, they are asked (i) whether the parts are similar to one another and (ii) how long each part lasts on average. If an annotator instead responds that the event referred to does *not* have natural parts, they are asked (i) whether the event is dynamic, and (ii) how long the event lasts.

All questions are binary except those concerning duration, for which answers are supplied as one of twelve ordinal values (see Vashishtha et al., 2019): *effectively no time at all*, *fractions of a second*, *seconds*, *minutes*, *hours*, *days*, *weeks*, *months*, *years*, *decades*, *centuries*, or *effectively forever*. Together, these questions target the three Vendler-inspired features (DYN, DUR, TEL), plus a fourth dimension for subtypes of dynamic predicates. In the context of UDS, these properties form a predicate node subspace, alongside FACTUALITY, GENERICITY, and TIME.

Event-event Annotators are presented with either a single sentence or a pair of adjacent sentences, with the two predicates of interest highlighted in distinct colors. For a predicate pair (p_1, p_2) describing an event pair (e_1, e_2) , annotators are asked whether e_1 is a mereological part of e_2 , and vice versa. Both questions are binary: A positive response to both indicates that e_1 and e_2 are the same event; and a positive response to exactly one of the questions indicates proper parthood. Prior versions of UDS do not contain any predicate-predicate edge subspaces, so we add

³The annotation interfaces for all three subprotocols, including instructions, are available at decomp.io.

document-level graphs to UDS (§6) to capture the relation between adjacently described events.

This subprotocol targets generalized event coreference, identifying *constituency* in addition to strict identity. It also augments the information collected in the event-subevent protocol: insofar as a proper subevent relation holds between e_1 and e_2 , we obtain additional fine-grained information about the subevents of the containing event—for example, an explicit description of at least one subevent.

Event-entity The final subprotocol focuses on the relation between the event described by a predicate and its plural or conjoined arguments, asking whether the predicate is distributive or collective with respect to that argument. This property accordingly forms a predicate-argument subspace in UDS, similar to PROTOROLES.

4 Validation Experiments

We validate our annotation protocol (i) by assessing interannotator agreement (IAA) among both experts and crowd-sourced annotators for each subprotocol on a small sample of items drawn from existing annotated corpora (§4.1-4.2); and (ii) by comparing annotations generated using our protocol against existing annotations that cover (a subset of) the phenomena that ours does and are generated by highly trained annotators (§4.3).

4.1 Item Selection

For each of the three subprotocols, one of the authors selected 100 sentences for inclusion in the pilot for that subprotocol. This author did not consult with the other authors on their selection, so that annotation could be blind.

For the event-subevent subprotocol, the 100 sentences come from the portion of the MASC corpus (Ide et al., 2008) that Friedrich et al. (2016) annotate for eventivity (EVENT v. STATE) and that Friedrich and Gateva (2017) annotate for telicity (TELIC v. ATELIC). For the event-event subprotocol, the 100 sentences come from the portions of the Richer Event Descriptions corpus (RED; O’Gorman et al., 2016) that are annotated for event subpart relations. To our knowledge, no existing annotations cover distributivity, and so for our event-entity protocol, we select 100 sentences (distinct from those used

for the event-subevent subprotocol) and compute IAA, but do not compare against existing annotations.

4.2 Interannotator Agreement

We compute two forms of IAA: (i) IAA among expert annotators (the three authors); and (ii) IAA between experts and crowd-sourced annotators. In both cases, we use Krippendorff’s α as our measure of (dis)agreement (Krippendorff, 2004). For the binary responses, we use the nominal form of α ; for the ordinal responses, we use the ordinal.

Expert Annotators For each subprotocol, the three authors independently annotated the 100 sentences selected for that subprotocol.

Prior to analysis, we rigit score the confidence ratings by annotator to normalize them for differences in annotator scale use (see Govindarajan et al., 2019 for discussion of rigit scoring confidence ratings in a similar annotation protocol). This method maps ordinal labels to (0, 1) on the basis of the empirical CDF of each annotator’s responses—with values closer to 0 implying lower confidence and those nearer 1 implying higher confidence. For questions that are dynamically revealed on the basis of the answer to the *natural parts* question (i.e., *part similarity*, *average part duration*, *dynamicity*, and *situation duration*) we use the average of the rigit scored confidence for *natural parts* and that question.

Figure 2 shows α when including only items that the expert annotators rated with a particular rigit scored confidence or higher. The agreement for the event-event protocol (mereology) is given in two forms: given that e_1 temporally contains e_2 , (i) *directed*: the agreement on whether e_2 is a subevent of e_1 ; and (ii) *undirected*: the agreement on whether e_2 is a subevent of e_1 and whether e_1 is a subevent of e_2 .

The error bars are computed by a nonparametric bootstrap over items. A threshold of 0.0 corresponds to computing α for all annotations, regardless of confidence; a threshold of $t > 0.0$ corresponds to computing α only for annotations associated with a rigit scored confidence of greater than t . When this thresholding results in less than $\frac{1}{3}$ of items having an annotation for at least two annotators, α is not plotted. This situation occurs only for questions that are revealed based on the answer to a previous question.

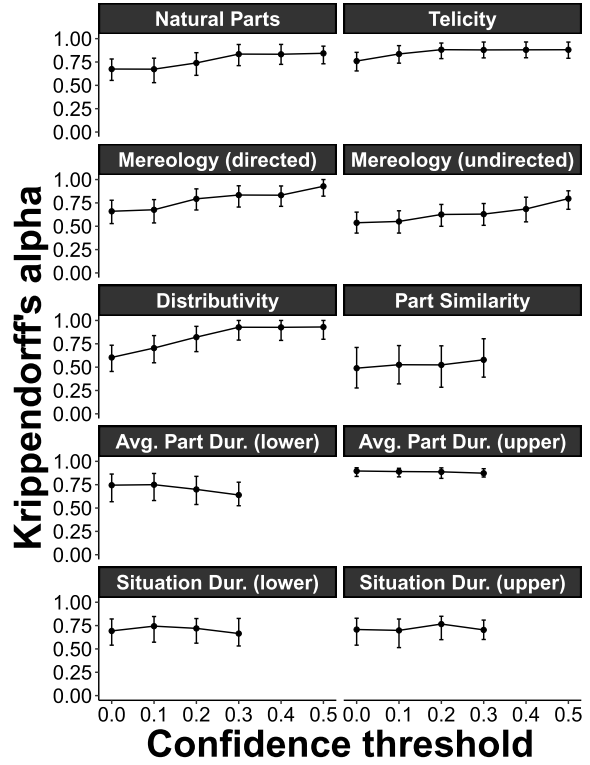


Figure 2: IAA among experts for each property, filtering annotations with rigit-scored confidence ratings below different thresholds. Confidence threshold 0.0 implies no filtering. Error bars show 95% confidence intervals computed by a nonparametric bootstrap.

For *natural parts*, *telicity*, *mereology*, and *distributivity*, agreement is high, even without filtering any responses on the basis of confidence, and that agreement improves with confidence. For *part similarity*, *average part duration*, and *situation duration*, we see more middling, but still reasonable, agreement, though this agreement does not reliably increase with confidence. The fact that it does not increase may have to do with interactions between confidence on the *natural parts* question and its dependent questions that we do not capture by taking the mean of these two confidences.

Crowd-Sourced Annotators We recruit crowd-sourced annotators in two stages. First, we select a small set of items from the 100 we annotate in the expert annotation that have high agreement among experts to create a qualification task. Second, based on performance in this qualification task, we construct a pool of trusted annotators who are allowed to participate in pilot annotations for each of the three subprotocols.⁴

⁴During all validation stages as well as bulk annotation (§5), we targeted an average hourly pay equivalent to that

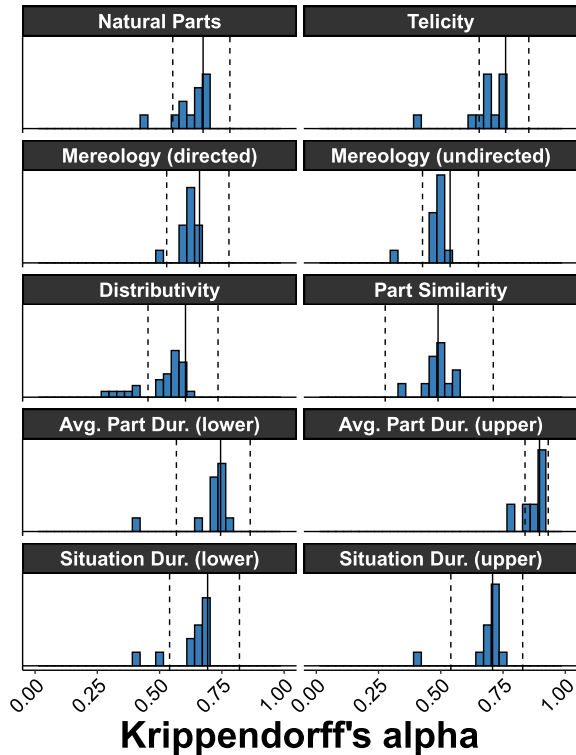


Figure 3: Per-property histograms of alphas for IAA between each crowd-sourced annotator and all experts. Black lines show the experts-only alpha, with dashed lines for the 95% CI. (See Figure 2).

Qualification For the qualification task, we selected eight of the sentences collected from MASC for the event-subevent subprotocol on which expert agreement was very high and which were diverse in the types of events described. We then obtained event-subevent annotations for these sentences from 400 workers on Amazon Mechanical Turk (AMT), and selected the top 200 among them on the basis of their agreement with expert responses on the same items. These workers were then permitted to participate in the pilot tasks.

Pilot We conducted one pilot for each subprotocol, using the items described in §4.1. Sentences were presented in lists of 10 per Human Intelligence Task (HIT) on AMT for the event-event and event-entity subprotocols and in lists of 5 per HIT for event-subevent. We collected annotations from 10 distinct workers for each sentence, and workers were permitted to annotate up to the full

for undergraduate research assistants performing corpus annotation at the first and final author’s institution: \$12.50 per hour. A third-party estimate from TurkerView shows our actual average hourly compensation when data collection was completed to be \$14.71 per hour.

100 sentences in each pilot. Thus, all pilots were guaranteed to include a minimum of 10 distinct workers (all workers do all HITs), up to a maximum of 100 for the subprotocols with 10 sentences per HIT or 200 for the subprotocol with 5 per HIT (each worker does one HIT). All top-200 workers from the qualification were allowed to participate.

Figure 3 shows IAA between all pilot annotators and experts for individual questions across the three pilots. More specifically, it shows the distribution of α scores by question for each annotator when IAA is computed among the three experts and that annotator only. Solid vertical lines show the expert-only IAA and dashed vertical lines show the 95% confidence interval.

4.3 Protocol Comparison

To further validate the event-event and event-subevent subprotocols, we evaluate how well our pilot data predicts the corresponding CONTAINS v. CONTAINS-SUBEVENT annotations from RED in the former case, as well as the EVENT v. STATE and TELIC v. ATELIC annotations from SitEnt in the latter. In both cases, we used the (ridit-scored) confidence-weighted average response across annotators for a particular item as features in a simple SVM classifier with linear kernel. In a leave-one-out cross-validation on the binary classification task for RED, we achieve a micro-averaged F1 score of 0.79—exceeding the reported human F1 agreement for both the CONTAINS (0.640) and CONTAINS-SUBEVENT (0.258) annotations reported by O’Gorman et al. (2016).

For SitEnt, we evaluate on a three-way classification task for STATIVE, EVENTIVE-TELIC, and EVENTIVE-ATELIC, achieving a micro-averaged F1 of 0.68 using the same leave-one-out cross-validation. As Friedrich and Palmer (2014a) do not report interannotator agreement for this class breakdown, we further compute Krippendorff’s alpha from their raw annotations and again find that agreement between our predicted annotations and the gold ones (0.48) slightly exceeds the interannotator agreement among humans (0.47).

These results suggest that our subprotocols capture relevant event structural phenomena as well as linguistically trained annotators can and that they may serve as effective alternatives to existing protocols while not requiring any linguistic expertise.

| | Annotation | Count (%) | Example |
|----------------|-------------------------------------|--------------------------|---|
| Event-subevent | Has natural parts | 6,903 (23%) | The eighteen steps of the dance are <u>done</u> rhythmically |
| | Parts similar | 4,498 (15%) | Israel resumed its policy of <u>targeting</u> militant leaders |
| | Parts dissimilar (Part duration) | 2,158 (7%) (-) | Fish are probably the easiest to <u>take</u> care of (ordinal; not shown) |
| | No natural parts | 23,069 (77%) | It <u>had</u> better nutritional value |
| | Dynamic | 13,903 (48%) | I would like to <u>informally</u> get together with you |
| | Not dynamic (Full duration) | 8,839 (29%) (-) | I assume this <u>is 12:30</u> Central Time? (ordinal; not shown) |
| | Natural endpoint | 6,031 (20%) | I will <u>deliver</u> it to you |
| | No natural endpoint | 23,941 (80%) | If you <u>know</u> or work there could you enlighten me? |
| total | 29,984 | (all event descriptions) | |
| Event-event | P1, P2 identical | 2,435 (6%) | All horses [. . .] are <u>happy</u> ₁ & <u>healthy</u> ₂ when they arrive |
| | P1, P2 disjoint | 30,247 (80%) | I am often <u>stopped</u> ₁ on the street and asked, ‘Who does your hair . . . I <u>LOVE</u> ₂ it’ |
| | P1 \subset P2 | 1,832 (5%) | The office is shared with a foot doctor and it’s <u>very sterile</u> ₁ and <u>medical feeling</u> ₂ , which I liked |
| | P2 \subset P1 | 3,029 (8%) | It is a <u>very cruel death</u> ₁ with bodies <u>dismembered</u> ₂ |
| | total | 37,719 | (pairs of event descriptions w/ temporal overlap) |
| Event-entity | Distributive | 4,812 (50%) | the <u>pics</u> turned out <u>ok</u> |
| | Collective | 4,876 (50%) | <u>we</u> <u>draw</u> on our many faith traditions to arrive at a common conviction |
| | total | 9,710 | (event descriptions with plural arguments) |

Table 1: Descriptive statistics and examples from Train and Dev data. Each item was annotated by a single annotator in Train; and by three annotators in Dev, of which this table reports the majority opinion.

5 Corpus Annotation

We collect crowd-sourced annotations for the entirety of UD-EWT. Predicate and argument spans are obtained from the PredPatt predicate-argument graphs for UD-EWT available in UDS1.0. The total number of items annotated for each subprotocol is presented in Table 1.

Event-subevent These annotations cover all predicates headed by verbs (as identified by UD POS tag), as well as copular constructions with nominal and adjectival complements. In the former case, only the verb token is highlighted in the task; in the latter, the highlighting spans from the copula to the complement head.

Event-event Pairs for the event-event subprotocol were drawn from the UDS-Time dataset,

which features pairs of verbal predicates, either within the same sentence or in adjacent sentences, each annotated with its start- and endpoint relative to the other. We additionally included predicate-argument pairs in cases where the argument is annotated in UDS as having a WordNet supersense of EVENT, STATE, or PROCESS. To our knowledge, this represents the largest publicly available (partial) event coreference dataset to date.

Event-entity For the event-entity subprotocol, we identify predicate-argument pairs in which the argument is plural or conjoined. Plural arguments are identified by the UD NUMBER attribute, and conjoined ones by a conj dependency between an argument head and another noun. We consider only predicate-argument pairs with a UD dependency of nsubj, nsubjpass, dobj, or iobj.

6 Event Structure Induction

Our goal in inducing event structural categories is to learn representations of those categories on the basis of annotated UDS graphs, augmented with the new UDS-E annotations. We aim to learn four sets of interdependent classifications grounded in UDS properties: event types, entity types, semantic role types, and event-event relation types. These classifications are interdependent in that we assume a generative model that incorporates both sentence- and document-level structure.⁵

Document-level UDS Semantics edges in UDS1.0 represent only sentence-internal semantic relations. This constraint implies that annotations for cross-sentential semantic relations—a significant subset of our event-event annotations—cannot be represented in the graph structure. To remedy this, we extend UDS1.0 by adding *document edges* that connect semantics nodes either within a sentence or in two distinct sentences, and we associate our event-event annotations with their corresponding document edge (see Figure 1). Because UDS1.0 does not have a notion of document edge, it does not contain Vashishtha et al.’s (2019) fine-grained temporal relation annotations, which are highly relevant to event-event relations. We additionally add those attributes to their corresponding document edges.

Generative Model Algorithm 1 gives the generative story for our event structure induction model. We assume some number of types of events $\mathcal{T}_{\text{event}}$, roles $\mathcal{R}_{\text{role}}$, entities \mathcal{T}_{ent} , and relations \mathcal{R}_{rel} . Figure 4 shows the resulting factor graph for the semantic graphs shown in Figure 1.

Annotation Likelihoods The distribution f_p^a on the annotations themselves is implemented as a mixed model (Gelman and Hill, 2014) dependent on property p being annotated with annotator random intercepts \mathbf{R} , where the random intercepts for annotator a are $\rho_a \sim \mathcal{N}(\mathbf{0}, \Sigma_{\text{ann}})$ with unknown Σ_{ann} . When p receives binary annotations, a simple logistic mixed model is assumed, where $f_p^a = \text{Bern}(\text{logit}^{-1}(\mu_{i_p} + \rho_{ai_p}))$ and i_p is the index corresponding to property p in the expected annotation μ . When p receives nominal annotations, $f_p^a = \text{Cat}(\text{softmax}(\mu_{i_p} + \rho_{ai_p}))$ and i_p is a set with cardinality of the number of nominal

⁵See Ferraro and Van Durme, 2016 for a related model that uses FrameNet’s ontology, rather than inducing its own.

```

Initialize queue  $I$ ;
for sentence  $s \in S$  do
  Initialize queue  $J$ ;
  Enqueue  $J \rightarrow I$ ;
  if length( $I$ ) >  $W$  then
    Dequeue  $I$ 
  for predicate node  $v \in \text{predicates}(s)$  do
    Sample event type  $t_{sv} \sim \text{Cat}(\theta^{(\text{event})})$ ;
    for property  $p \in \mathcal{P}_{\text{event}}$  do
      for annotator  $i \in \mathcal{A}_{svp}^{(\text{event})}$  do
        Sample  $x_{svpi}^{(\text{event})} \sim f_p^i(\mu_{t_{sv}}^{(\text{event})})$ 
    Enqueue  $\langle s, v \rangle \rightarrow J$ ;
  for argument node  $v' \in \text{arguments}(s, v)$  do
    Sample ent. type  $t_{sv'} \sim \text{Cat}(\theta^{(\text{entity})})$ ;
    for property  $p \in \mathcal{P}_{\text{ent}}$  do
      for annotator  $i \in \mathcal{A}_{sv'p}^{(\text{ent})}$  do
        Sample  $x_{sv'pi}^{(\text{ent})} \sim f_p^i(\mu_{t_{sv'}}^{(\text{ent})})$ 
    if  $v'$  is eventive then
      Enqueue  $\langle s, v' \rangle \rightarrow J$ ;
    Sample role type  $r_{svv'} \sim \text{Cat}(\theta_{t_{sv}t_{sv'}}^{(\text{role})})$ ;
    for property  $p \in \mathcal{P}_{\text{role}}$  do
      for annotator  $i \in \mathcal{A}_{svv'p}^{(\text{role})}$  do
        Sample  $x_{svv'pi}^{(\text{role})} \sim f_p^i(\mu_{r_{svv'}}^{(\text{role})})$ 
    for index pair  $\langle s', v' \rangle \in \text{flatten}(I)$  do
      Sample rel. type  $q \sim \text{Cat}(\theta_{t_{sv}t_{s'v'}}^{(\text{rel})})$ ;
      for property  $p \in \mathcal{P}_{\text{rel}}$  do
        for annotator  $i \in \mathcal{A}_{svs'v'p}^{(\text{rel})}$  do
          Sample  $x_{svs'v'pi}^{(\text{rel})} \sim f_p^i(\mu_q^{(\text{rel})})$ 

```

Algorithm 1: Generative story of event structure induction model for a single document with sentence window W .

categories. And when p receives ordinal annotations, we follow White et al. (2020) in using an ordinal (linked logit) mixed effects model where ρ_a defines the cutpoints between response values in the cumulative density function for annotator a :

$$\mathbb{P}(x_{ai_p} \leq j) = \text{logit}^{-1}(\mu_{i_p} - \rho_{ai_p})$$

$$f_p^a(x_{ai_p} = j) = \mathbb{P}(x_{ai_p} \leq j) - \mathbb{P}(x_{ai_p} \leq j - 1)$$

Conditional Properties For both our dataset and UDS-Protoroles, certain annotations are conditioned on others, owing to the fact that whether some questions are asked at all depends upon annotator responses to previous ones. Following White et al. (2017), we model the likelihoods for these properties using hurdle models (Agresti, 2014): For a given property, a Bernoulli distribution determines whether the property applies; if it does, the property value is determined using a second distribution of the appropriate type.

Temporal Relations Temporal relations annotations from UDS-Time consist of 4-tuples

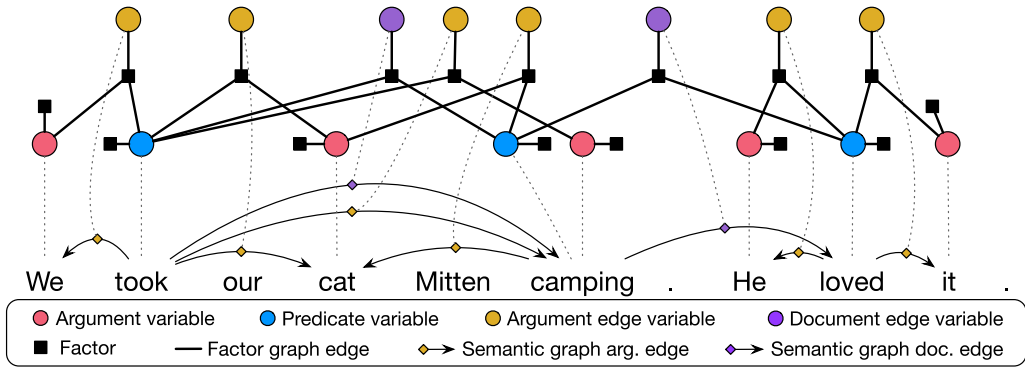


Figure 4: The factor graph for the pair of sentences shown in Figure 1 based on the generative story given in Algorithm 1. Each node or edge annotated in the semantics graphs becomes a variable node in the factor graph, as indicated by the dotted lines. Only factors for the prior distributions over types are shown; the annotation likelihood factors associated with each variable node are omitted for space.

$(\overleftarrow{e}_1, \overleftarrow{e}_2, \overrightarrow{e}_1, \overrightarrow{e}_2)$ of real values on the unit interval, representing start- and endpoints of two event-referring predicates or arguments, e_1 and e_2 . Each tuple is normalized such that the earlier of $(\overleftarrow{e}_1, \overleftarrow{e}_2)$ is always locked to the left end of the scale (0) and the later of $(\overrightarrow{e}_1, \overrightarrow{e}_2)$ to the right end (1). The likelihood for these annotations must consider the different possible orderings of the two events. To do so, we first determine whether \overleftarrow{e}_1 is locked, \overleftarrow{e}_2 is, or both are, according to $\text{Cat}(\text{softmax}(\boldsymbol{\mu}_{\text{lock}^{\leftarrow}} + \boldsymbol{\rho}_{a_{i_{\text{lock}^{\leftarrow}}}}))$. We do likewise for \overrightarrow{e}_1 and \overrightarrow{e}_2 , using a separate distribution $\text{Cat}(\text{softmax}(\boldsymbol{\mu}_{\text{lock}^{\rightarrow}} + \boldsymbol{\rho}_{a_{i_{\text{lock}^{\rightarrow}}}}))$. Finally, if the start point from one event and the endpoint from the other are free (i.e., not locked), we determine their relative ordering using a third distribution $\text{Cat}(\text{softmax}(\boldsymbol{\mu}_{\text{lock}^{\leftrightarrow}} + \boldsymbol{\rho}_{a_{i_{\text{lock}^{\leftrightarrow}}}}))$.

Implementation We fit our model to the training data using expectation-maximization. We use loopy belief propagation to obtain the posteriors over event, entity, role, and relation types in the expectation step and the Adam optimizer to estimate the parameters of the distributions associated with each type in the maximization step.⁶ As a stopping criterion, we compute the evidence that the model assigns to the development data, stopping when this quantity begins to decrease.

To make use of the (ridit-scored) confidence response $c_{a_{i_p}} \in (0, 1)$ associated with each annotation $x_{a_{i_p}}$, we weight the log-likelihood of $x_{a_{i_p}}$ by $c_{a_{i_p}}$ when computing the evidence of the annotations. This weighting encourages the model to

⁶The variable and factor nodes for the relation types can introduce cycles into the factor graph for a document, which is what necessitates the loopy variant of belief propagation.

explain annotations that an annotator was highly confident in, penalizing the model less if it assigns low likelihood to a low confidence annotation.

To select $|\mathcal{T}_{\text{event}}|$, $|\mathcal{T}_{\text{ent}}|$, $|\mathcal{R}_{\text{role}}|$, and $|\mathcal{R}_{\text{rel}}|$ for Algorithm 1, we fit separate mixture models for each classification—i.e., removing all factor nodes—using the same likelihood functions f_p^i as in Algorithm 1. We then compute the evidence that the simplified model assigns to the development data given some number of types, choosing the smallest number such that there is no reliable increase in the evidence for any larger number. To determine reliability, we compute 95% confidence intervals using nonparametric bootstraps. Importantly, this simplified model is only used to select $|\mathcal{T}_{\text{event}}|$, $|\mathcal{T}_{\text{ent}}|$, $|\mathcal{R}_{\text{role}}|$, and $|\mathcal{R}_{\text{rel}}|$: all analyses below are conducted on full model fits.

Types The selection procedure described above yields $|\mathcal{T}_{\text{event}}| = 4$, $|\mathcal{T}_{\text{ent}}| = 8$, $|\mathcal{R}_{\text{role}}| = 2$, and $|\mathcal{R}_{\text{rel}}| = 5$. To interpret these classes, we inspect the property means $\boldsymbol{\mu}_t$ associated with each type t and give examples from UD-EWT for which the posterior probability of that type is high.

Event Types While our goal was not necessarily to reconstruct any particular classification from the theoretical literature, the four event types align fairly well with those proposed by Vendler (1957): statives (16), activities (17), achievements (18), and accomplishments (19). We label our clusters based on these interpretations (Figure 5).

(16) I have finally found a mechanic I **trust**!!

(17) his agency is still **reviewing** the decision.

(18) A suit against [. . .] Kristof was **dismissed**.

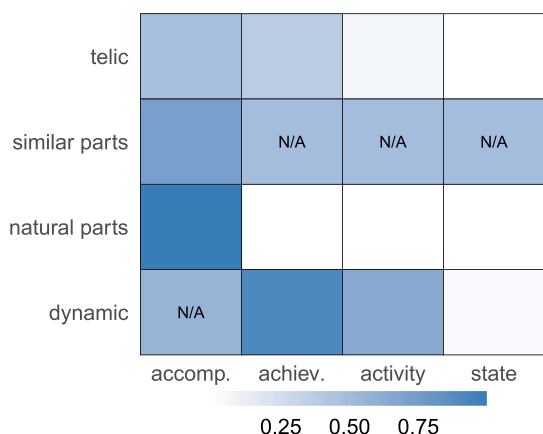


Figure 5: Probability of binary properties from the event-subevent protocol by event type. Cells marked with ‘‘N/A’’ indicate that the property generally does not apply for the corresponding type because of the conditional dependence on natural parts.

(19) a consortium [. . .] **established** in 1997

One difference between Vendler’s classes and our own is that our ‘‘activities’’ correspond primarily to those without dynamic subevents, while our ‘‘accomplishments’’ encompass both his accomplishments and activities with dynamic subevents (see discussion of Taylor, 1977 in §2).

Even if approximate, this alignment is surprising given that Vendler’s classification was not developed with actual language use in mind and thus abstracts away from complexities that arise when dealing with, for example, non-factual or generic events. Nonetheless, there do arise cases where a particular predicate has a wider distribution across types than we might expect based on prior work. For instance, *know* is prototypically stative; and while it does get classed that way by our model, it also gets classed as an accomplishment or achievement (though rarely an activity)—for example, when it is used to talk about coming to know something, as in (20).

(20) Please let me **know** how[. . .]to proceed.

Entity Types Our entity types are: person/group (21), concrete artifact (22), contentful artifact (23), particular state/event (24), generic state/event (25), time (26), kind of concrete objects (27), and particular concrete objects (28).

(21) Have a real **mechanic** check[...]

(22) I have a [. . .] cockatiel, and there are 2 **eggs** in the bottom of the cage[. . .]

(23) Please find attached a credit **worksheet**[. . .]

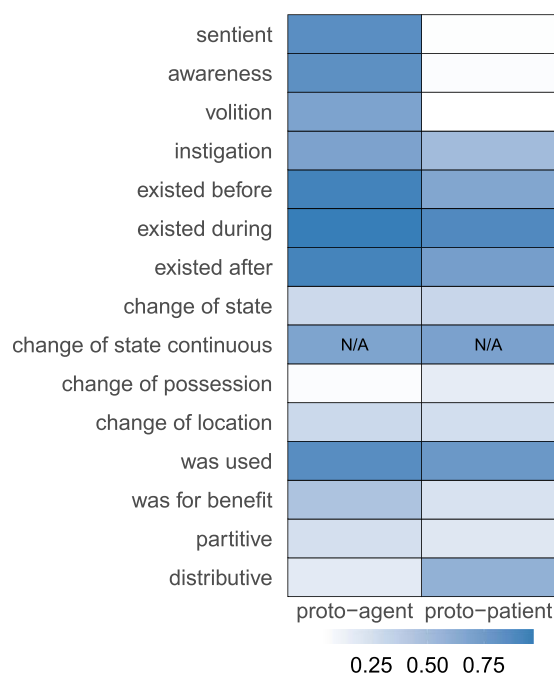


Figure 6: Probability of role type properties. These include existing UDS protoroles properties, along with the *distributive* property from the event-entity sub-protocol. We have labeled the role types with our proto-agent/proto-patient interpretation given below.

(24) He didn’t take a **dislike** to the kids[...]

(25) They require a **lot of attention** [. . .]

(26) Every move Google makes brings this particular **future** closer.

(27) And what is their big / main **meal** of the day.

(28) Find him before he finds the dog **food**.

Role Types The optimality of two role types is consistent with Dowty’s (1991) proposal that there are only two abstract role prototypes—*proto-agent* and *proto-patient*—into which individual thematic roles (i.e., those specific to particular predicates) cluster. Further, the means for the two role types we find very closely track those predicted by Dowty (see Figure 6), with clear proto-agents (29) and proto-patients (30) (see also White et al., 2017).

(29) **they** don’t **press** their sandwiches.

(30) you don’t ever feel like you **ate** too **much**.

Relation Types The relation types we obtain track closely with approaches that use sets of underspecified temporal relations (Cassidy et al., 2014; O’Gorman et al., 2016; Zhou et al., 2019, 2020; Wang et al., 2020): e_1 starts before e_2 (31),

e_2 starts before e_1 (32), e_2 ends after e_1 (33), e_1 contains e_2 (34), and $e_1 = e_2$ (35).

- (31) [. . .]the Spanish, Thai and other contingents are already **committed to leaving** [. . .]
- (32) And I have to **wonder**: Did he **forget** that he already has a memoir[. . .]
- (33) no, i am not **kidding** and no i don't want it b/c of the taco bell dog. i want it b/c it is really **small** and cute.
- (34) they **offer** cheap air tickets to their country [. . .] you may get excellent discount airfare, which may even **surprise** you.
- (35) the food is good, however the tables are so **close together** that it feels very **cramped**.

Type Consistency To assess the impacts of sentence/document-level structure (Algorithm 1) and confidence weighting on the types we induce, we investigate how the distributions over role and event types shift when comparing models fit with and without structure and confidence weighting. Figure 7 visualizes these shifts as row-normalized confusion matrices computed from the posteriors across items derived from each model. It compares (top) the simplified model used for model selection (rows) against the full model without confidence weighting (columns), and (bottom) the full model without confidence weighting (rows) against the one with (columns).⁷

First, we find that the interaction between types afforded by incorporating the full graph structure (top plot) produces small shifts in both the event and role type distributions, suggesting that the added structure may help chiefly in resolving boundary cases, which is exactly what we might hope additional model structure would do. Second, weighting likelihoods by annotator confidence (bottom) yields somewhat larger shifts as well as more entropic posteriors (0.22 average normalized entropy for events; 0.30 for roles) than without weighting (0.02 for events; 0.22 for roles).

Higher entropy is expected (and to some extent, desirable) here: Introducing a notion of confidence should make the model less confident about items that annotators were less confident about. Further, among event types, the distribution of posterior entropy across items is driven by a minority of high uncertainty items, as evidenced by a very low

⁷The distributional shifts for entity and relation types were extremely small, and so we do not discuss them here.

| | | [+struct, -conf] | | | | | |
|--------------------|------------|--------------------|---------|----------|-------|-----------|------------|
| | | accomp. | achiev. | activity | state | proto-ag. | proto-pat. |
| [-struct, -conf] | accomp. | 0.996 | 0.002 | 0.002 | 0.001 | | |
| | achiev. | 0.001 | 0.885 | 0.1 | 0.013 | | |
| | activity | 0.003 | 0.078 | 0.917 | 0.002 | | |
| | state | 0.001 | 0.064 | 0.004 | 0.93 | | |
| | proto-ag. | | | | | 0.932 | 0.068 |
| | proto-pat. | | | | | 0.008 | 0.992 |

| | | [+struct, +conf] | | | | | |
|--------------------|------------|--------------------|---------|----------|-------|-----------|------------|
| | | accomp. | achiev. | activity | state | proto-ag. | proto-pat. |
| [+struct, -conf] | accomp. | 0.986 | 0.008 | 0.004 | 0.002 | | |
| | achiev. | 0.001 | 0.75 | 0.074 | 0.176 | | |
| | activity | 0.002 | 0.125 | 0.839 | 0.034 | | |
| | state | 0.001 | 0.071 | 0.05 | 0.879 | | |
| | proto-ag. | | | | | 0.877 | 0.123 |
| | proto-pat. | | | | | 0.35 | 0.65 |

Figure 7: Confusion matrices for event and role types.

median normalized entropy for event types (0.02). The opposite appears to be true among the role types, for which the median is high (0.60). This latter pattern is perhaps not surprising in light of theoretical accounts of semantic roles, such as Dowty's: The entire point of such accounts is that it is very difficult to determine sharp role categories, suggesting the need for a more continuous notion.

7 Comparison to Existing Ontologies

To explore the relationship between our induced classification and existing event and role ontologies, we ask how well our event, role, and entity types map onto those found in PropBank and VerbNet. Importantly, the goal here is not perfect alignment between our types and PropBank and VerbNet types, but rather to compare other classifications that reflect top-down assumptions to the one we derive bottom-up.

Implementation To carry out these comparisons, we use the parameters of the posterior distributions over event types $\theta_p^{(ev)}$ for each predicate p , over role types $\theta_{pa}^{(role)}$ for each argument a of each predicate p , and over entity types $\theta_{pa}^{(ent)}$ for each argument a of each predicate p as features in an SVM with RBF kernel predicting the event and role types found in PropBank and VerbNet. We take this route, over direct comparison of types, to account for the possibility that information encoded in role or event types within VerbNet or PropBank is distributed differently in our more abstract classification. We tune L2 regularization ($\lambda \in \{1, 0.5, 0.2, 0.1, 0.01, 0.001\}$) and bandwidth ($\gamma \in \{1e-2, 1e-3, 1e-4, 1e-5\}$) using grid search, selecting the best model based on performance on the standard UD-EWT development set. All metrics reflect UD-EWT test set performance.

| | Role | P | R | F | Micro F |
|---------|---------|------|------|------|---------|
| argnum | A0 | 0.58 | 0.63 | 0.60 | 0.67 |
| | A1 | 0.72 | 0.78 | 0.75 | |
| functag | pag | 0.57 | 0.59 | 0.58 | 0.62 |
| | ppt | 0.65 | 0.77 | 0.71 | |
| verbnet | agent | 0.64 | 0.54 | 0.59 | NA |
| | patient | 0.20 | 0.14 | 0.16 | |
| | theme | 0.55 | 0.58 | 0.57 | |

Table 2: Test set results for all role types that are labeled on at least 5% of the development data.

Role Type Comparison We first obtain a mapping from UDS predicates and arguments to the PropBank predicates and arguments annotated in EWT. Each such argument in PropBank is annotated with an argument number (A0-A4) as well as a function tag (PAG = *agent*, PPT = *patient*, etc.). We then compose this mapping with the mapping given in the PropBank frame files from PropBank rolesets to sets of VerbNet classes and from PropBank roles to sets of VerbNet roles (AGENT, PATIENT, THEME, etc.) to obtain a mapping from UDS arguments to sets of VerbNet roles. Because a particular argument maps to a set of VerbNet roles, we treat predicting VerbNet roles as a multi-label problem, fitting one SVM per role. For each argument a of predicate p , we use as predictors $[\theta_p^{(ev)}; \theta_{pa}^{(role)}; \theta_{pa}^{(ent)}; \theta_{p-a}^{(role)}; \theta_{p-a}^{(ent)}]$, with $\theta_{p-a}^{(role/ent)} = [\max_{a' \neq a} \theta_{pa'j}^{(role/ent)}, \text{mean}_{a' \neq a} \theta_{pa'j}^{(role/ent)}]$.

Table 2 gives the test set results for all role types labeled on at least 5% of the development data. For comparison, a majority guessing baseline obtains micro F1s of 0.58 (argnum) and 0.53 (functag).⁸ Our roles tend to align well with agentive roles (PAG, AGENT, and A0) and some non-agentive roles (PPT, THEME, and A1), but they align less well with other non-agentive roles (PATIENT). This result suggests that our two-role classification aligns fairly closely with the agentivity distinctions in PropBank and VerbNet, as we would expect if our roles in fact captured something like Dowty’s coarse distinction among prototypical agents and patients.

Event Type Comparison The PropBank roleset and VerbNet class ontologies are extremely

⁸A majority baseline for VerbNet roles always yields an F1 of 0 in our multi-label setup, since no role is assigned to more than half of arguments.

| Predicate | P | R | F |
|----------------|------|------|------|
| cause | 0.51 | 0.95 | 0.66 |
| do | 0.30 | 0.25 | 0.27 |
| has_possession | 0.23 | 0.18 | 0.20 |
| has_location | 0.11 | 0.14 | 0.12 |
| motion | 0.09 | 0.10 | 0.09 |

Table 3: Test set results for all VerbNet predicates that are labeled on five most frequent predicates.

fine-grained, with PropBank capturing specific predicate senses and VerbNet capturing very fine-grained syntactic behavior of a generally small set of predicates. Since our event types are intended to be more general than either, we do not compare it directly to PropBank rolesets or VerbNet classes.

Instead, we compare to the generative lexicon-inspired variant of VerbNet’s semantics layer (Brown et al., 2018). An example of this layer for the predicate `give-13.1` is `has_possession(e1, Ag, Th) & transfer(e2, Ag, Th, Rec) & cause(e2, e3) & has_possession(e3, Rec, Th)`. We predict only the abstract predicates in this decomposition (e.g., `transfer` or `cause`), treating the problem as multi-label and fitting one SVM per predicate. For each predicate p , we use as predictors $[\theta_p^{(ev)}; \theta_p^{(role)}; \theta_p^{(ent)}]$, with $\theta_{p:j}^{(role/ent)} = [\max_a \theta_{pa,j}^{(role/ent)}, \text{mean}_a \theta_{pa,j}^{(role/ent)}]$.

Table 3 gives the test set results for the five most frequent predicates in the corpus. For comparison, a majority guessing baseline would yield the same F (0.66) as our model for CAUSE, but since none of the other classes are assigned to more than half of events, majority guessing for those would yield an F of 0. This result suggests that, while there may be some agreement between our classification and VerbNet’s semantics layer, the two representations are relatively distinct.

8 Conclusion

We have presented an event structure classification derived from inferential properties annotated on sentence- and document-level semantic graphs. We induced this classification jointly with semantic role, entity, and event-event relation types using a document-level generative model. Our model identifies types that approximate theoretical predictions—notably, four event types like

Vendler’s, as well as proto-agent and proto-patient role types like Dowty’s. We hope this work encourages greater interest in computational approaches to event structural understanding while also supporting work on adjacent problems in NLU, such as temporal information extraction and (partial) event coreference, for which we provide the largest publicly available dataset to date.

Acknowledgments

We would like to thank Emily Bender, Dan Gildea, and three anonymous reviewers for detailed comments on this paper. We would also like to thank members of the Formal and Computational Semantics lab at the University of Rochester for feedback on the annotation protocols. This work was supported in part by the National Science Foundation (BCS-2040820/2040831, *Collaborative Research: Computational Modeling of the Internal Structure of Events*) as well as by DARPA AIDA and DARPA KAIROS. 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 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

- Omri Abend and Ari Rappoport. 2013. Universal Conceptual Cognitive Annotation (UCCA). In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 228–238, Sofia, Bulgaria. Association for Computational Linguistics.
- Alan Agresti. 2014. *Categorical Data Analysis*, John Wiley & Sons.
- James Allen, Hannah An, Ritwik Bose, Will de Beaumont, and Choh Man Teng. 2020. A broad-coverage deep semantic lexicon for verbs. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 3243–3251, Marseille, France. European Language Resources Association.
- Emmon Bach. 1986. The algebra of events. *Linguistics and Philosophy*, 9(1):5–16.
- Collin F. Baker, Charles J. Fillmore, and John B. Lowe. 1998. The Berkeley Framenet Project. In *Proceedings of the 17th International Conference on Computational Linguistics*, volume 1, pages 86–90. Association for Computational Linguistics.
- Laura Banarescu, Claire Bonial, Shu Cai, Madalina Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin Knight, Philipp Koehn, Martha Palmer, and Nathan Schneider. 2013. Abstract Meaning Representation for sembanking. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pages 178–186, Sofia, Bulgaria. Association for Computational Linguistics.
- Michael Ruisdael Bennett and Barbara Hall Partee. 1978. *Towards the Logic of Tense and Aspect in English*. Indiana University Linguistics Club, Bloomington, IN.
- Ann Bies, Justin Mott, Colin Warner, and Seth Kulick. 2012. English web treebank. *Linguistic Data Consortium, Philadelphia, PA*.
- Susan Windisch Brown, James Pustejovsky, Annie Zaenen, and Martha Palmer. 2018. Integrating Generative Lexicon event structures into VerbNet. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Taylor Cassidy, Bill McDowell, Nathanael Chambers, and Steven Bethard. 2014. An annotation framework for dense event ordering. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 501–506, Baltimore, Maryland. Association for Computational Linguistics. <https://doi.org/10.3115/v1/P14-2082>
- Lucas Champollion. 2010. *Parts of a Whole: Distributivity as a Bridge between Aspect and Measurement*. Ph.D. thesis, University of Pennsylvania, Philadelphia.
- Timothy Chklovski and Patrick Pantel. 2004. Verbocean: Mining the web for fine-grained semantic verb relations. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 33–40.
- William Croft. 2012. *Verbs: Aspect and causal structure*. Oxford University Press.

- Agata Cybulska and Piek Vossen. 2014. Using a sledgehammer to crack a nut? Lexical diversity and event coreference resolution. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 4545–4552, Reykjavik, Iceland. European Language Resources Association (ELRA).
- Lucia Donatelli, Michael Regan, William Croft, and Nathan Schneider. 2018. Annotation of tense and aspect semantics for sentential AMR. In *Proceedings of the Joint Workshop on Linguistic Annotation, Multiword Expressions and Constructions (LAW-MWE-CxG-2018)*, pages 96–108.
- Bonnie Jean Dorr. 1993. *Machine Translation: A View from the Lexicon*. MIT Press.
- David Dowty. 1979. *Word Meaning and Montague Grammar: The Semantics of Verbs and Times in Generative Semantics and in Montague's PTQ*, volume 7, Springer Science & Business Media.
- David Dowty. 1991. Thematic proto-roles and argument selection. *Language*, 67(3):547–619. <https://doi.org/10.2307/415037>, <https://doi.org/10.1353/lan.1991.0021>
- George Ferguson and James F. Allen. 1998. TRIPS: An integrated intelligent problem-solving assistant. In *Proceedings of the Fifteenth National/Tenth Conference on Artificial Intelligence/Innovative Applications of Artificial Intelligence*, AAAI '98/IAAI '98, pages 567–572, USA. American Association for Artificial Intelligence.
- Francis Ferraro and Benjamin Van Durme. 2016. A unified Bayesian model of scripts, frames and language. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30.
- Annemarie Friedrich and Damyana Gateva. 2017. Classification of telicity using cross-linguistic annotation projection. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2559–2565, Copenhagen, Denmark. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D17-1271>
- Annemarie Friedrich and Alexis Palmer. 2014a. Automatic prediction of aspectual class of verbs in context. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 517–523, Baltimore, Maryland. Association for Computational Linguistics. <https://doi.org/10.3115/v1/P14-2085>
- Annemarie Friedrich and Alexis Palmer. 2014b. Situation entity annotation. In *Proceedings of LAW VIII - The 8th Linguistic Annotation Workshop*, pages 149–158, Dublin, Ireland. Association for Computational Linguistics and Dublin City University. <https://doi.org/10.3115/v1/W14-4921>
- Annemarie Friedrich, Alexis Palmer, and Manfred Pinkal. 2016. Situation entity types: Automatic classification of clause-level aspect. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1757–1768, Berlin, Germany. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P16-1166>
- William Gantt, Benjamin Kane, and Aaron Steven White. 2020. Natural language inference with mixed effects. In *Proceedings of the Ninth Joint Conference on Lexical and Computational Semantics*, pages 81–87, Barcelona, Spain (Online). Association for Computational Linguistics.
- Andrew Gelman and Jennifer Hill. 2014. *Data Analysis Using Regression and Multi-level/Hierarchical Models*. Cambridge University Press, New York City.
- Venkata Govindarajan, Benjamin Van Durme, and Aaron Steven White. 2019. Decomposing generalization: Models of generic, habitual, and episodic statements. *Transactions of the Association for Computational Linguistics*, 7:501–517. https://doi.org/10.1162/tacl_a_00285
- Eduard Hovy, Teruko Mitamura, Felisa Verdejo, Jun Araki, and Andrew Philpot. 2013. Events are not simple: Identity, non-identity, and quasi-identity. In *Workshop on Events: Definition, Detection, Coreference, and Representation*, pages 21–28.
- Nancy Ide, Collin Baker, Christiane Fellbaum, Charles Fillmore, and Rebecca Jane Passonneau. 2008. MASC: The manually annotated sub-corpus

- of American English. In *6th International Conference on Language Resources and Evaluation, LREC 2008*, pages 2455–2460. European Language Resources Association (ELRA).
- Ray Jackendoff. 1990. *Semantic Structures*, volume 18. MIT Press.
- Anthony Kenny. 1963. *Action, Emotion and Will*, Humanities Press, London.
- Karin Kipper Schuler. 2005. *VerbNet: A Broad-coverage, Comprehensive Verb Lexicon*. Ph.D. thesis, University of Pennsylvania.
- Manfred Krifka. 1989. Nominal reference, temporal constitution and quantification in event semantics. In Renate Bartsch, Johan van Benthem, and Peter van Emde Boas, editors, *Semantics and Contextual Expressions*, pages 75–115. Foris, Dordrecht. <https://doi.org/10.1515/9783110877335-005>
- Manfred Krifka. 1992. Thematic relations as links between nominal reference and temporal constitution. *Lexical Matters*, 2953.
- Manfred Krifka. 1998. The origins of telicity. In Susan Rothstein, editor, *Events and Grammar*, Studies in Linguistics and Philosophy, pages 197–235. Springer Netherlands, Dordrecht. https://doi.org/10.1007/978-94-011-3969-4_9
- K. Krippendorff. 2004. *Content Analysis: An Introduction to Its Methodology*. Sage.
- George Lakoff. 1965. *On the Nature of Syntactic Irregularity*. Ph.D. thesis, Massachusetts Institute of Technology.
- Beth Levin. 1993. *English Verb Classes and Alternations: A Preliminary Investigation*. University of Chicago Press, Chicago.
- Beth Levin and Malka Rappaport Hovav. 1991. Wiping the slate clean: A lexical semantic exploration. *Cognition*, 41(1-3):123–151. [https://doi.org/10.1016/0010-0277\(91\)90034-2](https://doi.org/10.1016/0010-0277(91)90034-2)
- Marc Moens and Mark Steedman. 1988. Temporal ontology and temporal reference. *Computational Linguistics*, 14(2):15–28.
- Alexander P. D. Mourelatos. 1978. Events, processes, and states. *Linguistics and Philosophy*, 2(3):415–434. <https://doi.org/10.1007/BF00149015>
- Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Yoav Goldberg, Jan Hajič, Christopher D. Manning, Ryan McDonald, Slav Petrov, Sampo Pyysalo, Natalia Silveira, Reut Tsarfaty, and Daniel Zeman. 2016. Universal Dependencies v1: A multilingual treebank collection. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 1659–1666, Portorož, Slovenia. European Language Resources Association (ELRA).
- Tim O’Gorman, Kristin Wright-Bettner, and Martha Palmer. 2016. Richer event description: Integrating event coreference with temporal, causal and bridging annotation. In *Proceedings of the 2nd Workshop on Computing News Storylines (CNS 2016)*, pages 47–56, Austin, Texas. Association for Computational Linguistics. <https://doi.org/10.18653/v1/W16-5706>
- Mari Broman Olsen. 1997. *A Semantic and Pragmatic Model of Lexical and Grammatical Aspect*. Outstanding Dissertations in Linguistics. Garland. <https://doi.org/10.1162/0891201053630264>
- Martha Palmer, Daniel Gildea, and Paul Kingsbury. 2005. The proposition bank: An annotated corpus of semantic roles. *Computational Linguistics*, 31(1):71–106.
- Christopher Piñón. 1995. *An Ontology for Event Semantics*. Ph.D. thesis, Stanford University, Palo Alto.
- James Pustejovsky. 1995. *The Generative Lexicon*. MIT Press, Cambridge, MA.
- James Pustejovsky. 2013. Dynamic event structure and habitat theory. In *Proceedings of the 6th International Conference on Generative Approaches to the Lexicon (GL2013)*, pages 1–10, Pisa, Italy. Association for Computational Linguistics.
- James Pustejovsky, Marc Verhagen, Roser Saurí, Jessica Littman, Robert Gaizauskas, Graham Katz, Inderjeet Mani, Robert Knippen, and Andrea Setzer. 2006. TimeBank 1.2. *Linguistic Data Consortium*, 40.
- Malka Rappaport Hovav and Beth Levin. 1998. Building verb meanings. *The Projection of Arguments: Lexical and Compositional Factors*, pages 97–134.

- Malka Rappaport Hovav and Beth Levin. 2001. An event structure account of english resultatives. *Language*, 77(4):766–797.
- Drew Reisinger, Rachel Rudinger, Francis Ferraro, Craig Harman, Kyle Rawlins, and Benjamin Van Durme. 2015. Semantic proto-roles. *Transactions of the Association for Computational Linguistics*, 3:475–488. https://doi.org/10.1162/tacl_a_00152
- Rachel Rudinger, Adam Teichert, Ryan Culkin, Sheng Zhang, and Benjamin Van Durme. 2018. Neural-Davidsonian semantic proto-role labeling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 944–955, Brussels, Belgium. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D18-1114>
- Carlota S. Smith. 2003. *Modes of Discourse: The Local Structure of Texts*, volume 103. Cambridge University Press. <https://doi.org/10.1017/CBO9780511615108>
- Elias Stengel-Eskin, Kenton Murray, Sheng Zhang, Aaron Steven White, and Benjamin Van Durme. 2021. Joint universal syntactic and semantic parsing. *arXiv preprint arXiv:2104.05696*
- Elias Stengel-Eskin, Aaron Steven White, Sheng Zhang, and Benjamin Van Durme. 2020. Universal decompositional semantic parsing. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8427–8439, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.acl-main.746>
- Barry Taylor. 1977. Tense and continuity. *Linguistics and Philosophy*, 1(2):199–220
- Adam Teichert, Adam Poliak, Benjamin Van Durme, and Matthew Gormley. 2017. Semantic proto-role labeling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31.
- Carol Lee Tenny. 1987. *Grammaticalizing Aspect and Affectedness*. Ph.D. thesis, Massachusetts Institute of Technology.
- Robert Truswell, editor. 2019. *The Oxford Handbook of Event Structure*, Oxford University Press, Oxford. Publication Title: The Oxford Handbook of Event Structure. <https://doi.org/10.1093/oxfordhob/9780199685318.001.0001>
- Siddharth Vashishtha, Benjamin Van Durme, and Aaron Steven White. 2019. Fine-grained temporal relation extraction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2906–2919, Florence, Italy. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P19-1280>
- Zeno Vendler. 1957. Verbs and times. *Philosophical Review*, 66(2):143–160. <https://doi.org/10.2307/2182371>
- Henk J. Verkuyl. 1972. *On The Compositional Nature Of The Aspects*, volume 15 of *Foundations of Language*. D. Reidel Publishing Company, Dordrecht. <https://doi.org/10.1007/978-94-017-2478-4>
- Haoyu Wang, Muhao Chen, Hongming Zhang, and Dan Roth. 2020. Joint constrained learning for event-event relation extraction. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 696–706, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.emnlp-main.51>
- Aaron Steven White, Kyle Rawlins, and Benjamin Van Durme. 2017. The semantic proto-role linking model. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 92–98, Valencia, Spain. Association for Computational Linguistics.
- Aaron Steven White, Drew Reisinger, Keisuke Sakaguchi, Tim Vieira, Sheng Zhang, Rachel Rudinger, Kyle Rawlins, and Benjamin Van Durme. 2016. Universal decompositional semantics on Universal Dependencies. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1713–1723, Austin, Texas. Association for Computational Linguistics.
- Aaron Steven White, Elias Stengel-Eskin, Siddharth Vashishtha, Venkata Subrahmanyam Govindarajan, Dee Ann Reisinger, Tim Vieira, Keisuke Sakaguchi, Sheng Zhang, Francis Ferraro, Rachel Rudinger, Kyle Rawlins, and Benjamin Van Durme. 2020. The universal decompositional semantics dataset and decomp toolkit. In

Proceedings of the 12th Language Resources and Evaluation Conference, pages 5698–5707, Marseille, France. European Language Resources Association.

Sheng Zhang, Rachel Rudinger, and Benjamin Van Durme. 2017. An Evaluation of PredPatt and Open IE via Stage 1 Semantic Role Labeling. In *IWCS 2017—12th International Conference on Computational Semantics—Short papers*.

Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth. 2019. “going on a vacation” takes longer than “going for a walk”: A study of temporal commonsense understanding. In *Proceedings of the 2019 Conference on*

Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3363–3369, Hong Kong, China. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D19-1332>

Ben Zhou, Qiang Ning, Daniel Khashabi, and Dan Roth. 2020. Temporal common sense acquisition with minimal supervision. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7579–7589, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.acl-main.678>